

Spokoiny, Vladimir

Working Paper

Parametric estimation: Finite sample theory

SFB 649 Discussion Paper, No. 2011-081

Provided in Cooperation with:

Collaborative Research Center 649: Economic Risk, Humboldt University Berlin

Suggested Citation: Spokoiny, Vladimir (2011) : Parametric estimation: Finite sample theory, SFB 649 Discussion Paper, No. 2011-081, Humboldt University of Berlin, Collaborative Research Center 649 - Economic Risk, Berlin

This Version is available at:

<https://hdl.handle.net/10419/56754>

Standard-Nutzungsbedingungen:

Die Dokumente auf EconStor dürfen zu eigenen wissenschaftlichen Zwecken und zum Privatgebrauch gespeichert und kopiert werden.

Sie dürfen die Dokumente nicht für öffentliche oder kommerzielle Zwecke vervielfältigen, öffentlich ausstellen, öffentlich zugänglich machen, vertreiben oder anderweitig nutzen.

Sofern die Verfasser die Dokumente unter Open-Content-Lizenzen (insbesondere CC-Lizenzen) zur Verfügung gestellt haben sollten, gelten abweichend von diesen Nutzungsbedingungen die in der dort genannten Lizenz gewährten Nutzungsrechte.

Terms of use:

Documents in EconStor may be saved and copied for your personal and scholarly purposes.

You are not to copy documents for public or commercial purposes, to exhibit the documents publicly, to make them publicly available on the internet, or to distribute or otherwise use the documents in public.

If the documents have been made available under an Open Content Licence (especially Creative Commons Licences), you may exercise further usage rights as specified in the indicated licence.

Parametric estimation. Finite sample theory

Vladimir Spokoiny*

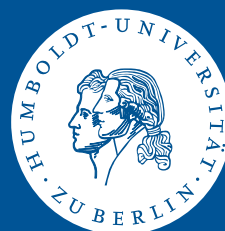


* Weierstrass Institute (WIAS) Berlin, Germany

This research was supported by the Deutsche Forschungsgemeinschaft through the SFB 649 "Economic Risk".

<http://sfb649.wiwi.hu-berlin.de>
ISSN 1860-5664

SFB 649, Humboldt-Universität zu Berlin
Spandauer Straße 1, D-10178 Berlin



Parametric estimation. Finite sample theory

Vladimir Spokoiny *

Weierstrass-Institute,
Mohrenstr. 39, 10117 Berlin, Germany
`spokoiny@wias-berlin.de`

Abstract

The paper aims at reconsidering the famous Le Cam LAN theory. The main features of the approach which make it different from the classical one are: (1) the study is non-asymptotic, that is, the sample size is fixed and does not tend to infinity; (2) the parametric assumption is possibly misspecified and the underlying data distribution can lie beyond the given parametric family.

The main results include a large deviation bounds for the (quasi) maximum likelihood and the local quadratic majorization of the log-likelihood process. The latter yields a number of important corollaries for statistical inference: concentration, confidence and risk bounds, expansion of the maximum likelihood estimate, etc. All these corollaries are stated in a non-classical way admitting a model misspecification and finite samples. However, the classical asymptotic results including the efficiency bounds can be easily derived as corollaries of the obtained non-asymptotic statements. The general results are illustrated for the i.i.d. set-up as well as for generalized linear and median estimation. The results apply for any dimension of the parameter space and provide a quantitative lower bound on the sample size yielding the root-n accuracy.

AMS 2000 Subject Classification: Primary 62F10. Secondary 62J12, 62F25, 62H12

Keywords: maximum likelihood, local quadratic approximation, concentration, coverage, deficiency

JEL-Classification: C13, C14

*Financial support by the German Research Foundation (DFG) through the Collaborative Research Center 649 “Economic Risk” is gratefully acknowledged

1 Introduction

One of the most popular approaches in statistics is based on the parametric assumption (PA) that the distribution \mathbb{P} of the observed data \mathbf{Y} belongs to a given parametric family $(\mathbb{P}_{\boldsymbol{\theta}}, \boldsymbol{\theta} \in \Theta \subseteq \mathbb{R}^p)$, where p states for the number of parameters. This assumption allows to reduce the problem of statistical inference about \mathbb{P} to recovering the parameter $\boldsymbol{\theta}$. The theory of parameter estimation and inference is nicely developed in a quite general set-up. There is a vast literature on this issue. We only mention the book by Ibragimov and Khas'minskij (1981), which provides a comprehensive study of asymptotic properties of maximum likelihood and Bayesian estimators. The theory is essentially based on two major assumptions: (1) the underlying data distribution follows the PA; (2) the sample size or the amount of available information is large relative to the number of parameters.

In many practical applications, both assumptions can be very restrictive and limiting the scope of applicability for the whole approach. Indeed, the PA is usually only an approximation of real data distribution and in the most of statistical problems it is too restrictive to assume that the PA is exactly fulfilled. Many modern statistical problems deal with very complex high dimensional data where a huge number of parameters are involved. In such situations, the applicability of large sample asymptotics is questionable. These two issues partially explain why the parametric and nonparametric theory are almost isolated from each other. Relaxing these restrictive assumptions can be viewed as an important challenge of the modern statistical theory. The present paper attempts at developing a unified approach which does not require the restrictive parametric assumptions but still enjoys the main benefits of the parametric theory. The main feature of the presentation is the *non-asymptotic* framework. The notions like asymptotic normality, convergence or tightness are meaningless in the non-asymptotic setup, the arguments based on compactness of the parameter space are not really helpful. Instead some exact exponential bounds and concentration results are systematically used. The main steps of the approach are similar to the classical local asymptotic normality (LAN) theory; see e.g. Chapters 1–3 in the monograph Ibragimov and Khas'minskij (1981): first we establish a kind of large deviation bound allowing to localize the problem into a neighborhood of the target parameter. Then we use a local quadratic expansion of the log-likelihood to solve the corresponding estimation problem.

Let \mathbf{Y} stand for the available data. Everywhere below we assume that the observed data \mathbf{Y} follow the distribution \mathbb{P} on a metric space \mathcal{Y} . We do not specify any particular structure of \mathbf{Y} . In particular, no assumption like independence or weak dependence of individual observations is imposed. The basic parametric assumption is that \mathbb{P} can be

approximated by a parametric distribution \mathbb{P}_θ from a given parametric family $(\mathbb{P}_\theta, \theta \in \Theta \subseteq \mathbb{R}^p)$. Our approach allows that the PA can be misspecified, that is, in general, $\mathbb{P} \notin (\mathbb{P}_\theta)$.

Let $L(\mathbf{Y}, \theta)$ be the log-likelihood for the considered parametric model: $L(\mathbf{Y}, \theta) = \log \frac{d\mathbb{P}_\theta}{d\mu_0}(\mathbf{Y})$, where μ_0 is any dominating measure for the family (\mathbb{P}_θ) . The classical likelihood principle suggests to estimate θ by maximizing the corresponding log-likelihood function $L(\mathbf{Y}, \theta)$:

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

Our ultimate goal is to study the properties of the quasi MLE $\tilde{\theta}$. It turns out that such properties can be naturally described in terms of the maximum of the process $L(\theta)$ rather than the point of maximum $\tilde{\theta}$. To avoid technical burdens it is assumed that the maximum is attained leading to the identity $\max_{\theta} L(\theta) = L(\tilde{\theta})$. However, the point of maximum needs not to be unique. If there are many such points we take $\tilde{\theta}$ as any of them. Basically, the notation $\tilde{\theta}$ is used for the identity $L(\tilde{\theta}) = \sup_{\theta \in \Theta} L(\theta)$.

If $\mathbb{P} \notin (\mathbb{P}_\theta)$, then the (quasi) MLE estimate $\tilde{\theta}$ from (1.1) is still meaningful and it appears to be an estimate of the value θ^* defined by maximizing the expected value of $L(\mathbf{Y}, \theta)$:

$$\theta^* \stackrel{\text{def}}{=} \operatorname{argmax}_{\theta \in \Theta} \mathbb{E} L(\mathbf{Y}, \theta) \quad (1.2)$$

which is the true value in the parametric situation and can be viewed as the parameter of the best parametric fit in the general case.

We focus on the properties of the process $L(\mathbf{Y}, \theta)$ as a function of the parameter θ . Therefore, we suppress the argument \mathbf{Y} there and write $L(\theta)$ instead of $L(\mathbf{Y}, \theta)$. One has to keep in mind that $L(\theta)$ is random and depends on the observed data \mathbf{Y} . We also define the excess or maximum log-likelihood $L(\theta, \theta^*) = L(\theta) - L(\theta^*)$. The results below show that the main properties of the quasi MLE $\tilde{\theta}$ like concentration or coverage probability can be described in terms of the quasi maximum likelihood $L(\tilde{\theta}) - L(\theta^*) = \max_{\theta \in \Theta} L(\theta) - L(\theta^*)$, which is the difference between the maximum of the process $L(\theta)$ and its value at the “true” point θ^* .

The established results can be split into two big groups. A large deviation bound states some concentration properties of the estimate $\tilde{\theta}$. For specific local sets $\Theta_0(\mathbf{r})$ with elliptic shape, the deviation probability $\mathbb{P}(\tilde{\theta} \notin \Theta_0(\mathbf{r}))$ is exponentially small in \mathbf{r} . This concentration bound allows for restricting the parameter space to a properly selected vicinity $\Theta_0(\mathbf{r})$. Our main results describe local properties of the process $L(\theta)$ within $\Theta_0(\mathbf{r})$. They can be viewed as a non-asymptotic version of the Le Cam LAN

theory.

The paper is organized as follows. Section 2 presents the list of conditions which are systematically used in the text. The conditions only concern the properties of the quasi likelihood process $L(\boldsymbol{\theta})$.

Section 3 appears to be central in the whole approach and it focuses on local properties of the process $L(\boldsymbol{\theta})$ within $\Theta_0(\mathbf{r})$. The idea is to sandwich the underlying (quasi) log-likelihood process $L(\boldsymbol{\theta})$ for $\boldsymbol{\theta} \in \Theta_0(\mathbf{r})$ between two quadratic (in parameter) expressions. Then the maximum of $L(\boldsymbol{\theta})$ over $\Theta_0(\mathbf{r})$ will be sandwiched as well by the maxima of the lower and upper processes. The quadratic structure of these processes help to compute these maxima explicitly yielding the bounds for the value of the original problem. This approximation result is used to derive a number of corollaries including the concentration and coverage probability, expansion of the estimate $\tilde{\boldsymbol{\theta}}$, polynomial risk bounds, etc. In the contrary to the classical theory, all the results are non-asymptotic and do not involve any small values of the form $o(1)$, all the terms are specified explicitly. Also the results are stated under possible model misspecification.

Section 4 accomplishes the local results with the concentration property which bounds the probability that $\tilde{\boldsymbol{\theta}}$ deviates from the local set $\Theta_0(\mathbf{r})$. In the modern statistical literature there is a number of studies considering maximum likelihood or more generally minimum contrast estimators in a general i.i.d. situation, when the parameter set Θ is a subset of some functional space. We mention the papers Van de Geer (1993), Birgé and Massart (1993), Birgé and Massart (1998), Birgé (2006) and references therein. The established results are based on deep probabilistic facts from the empirical process theory; see e.g. Talagrand (1996, 2001, 2005), van der Vaart and Wellner (1996), Boucheron et al. (2003). The general result presented in Section 7 follows the generic chaining idea due to Talagrand (2005); cf. Bednorz (2006). However, we do not assume any specific structure of the model. In particular, we do not assume independent observations and thus, cannot apply the most developed concentration bounds from the empirical process theory.

Section 5 illustrates the applicability of the general results to the classical case of an i.i.d. sample. The previously established general results apply under rather mild conditions. Basically we assume some smoothness of the log-likelihood process and some minimal number of observations pro parameter: the sample size should be at least of order of the dimensionality p of the parameter space. We also consider the examples of generalized linear modeling and of median regression.

It is important to mention that the non-asymptotic character of our study yields an almost complete change of the mathematical tools: the notions of convergence and tightness become meaningless, the arguments based on compactness of the parameter space do not apply, etc. Instead we utilize the tools of the empirical process theory based

on the ideas of concentration of measures and nonasymptotic entropy bounds. Section 6 presents an exponential bound for a general quadratic form which is very essential for getting the sharp risk bounds for the quasi MLE. This bound is an important step in the concentration results for the quasi MLE. Section 7 explains how generic chaining and majorizing measure device by Talagrand (2005) refined in Bednorz (2006) can be used for obtaining a general exponential bound for the log-likelihood process.

2 Conditions

Below we collect the list of conditions which are systematically used in the text. It seems to be an advantage of the whole approach that all the results are stated in a unified way under the same conditions which are quite general and not very much restrictive. We do not try to formulate the conditions and the results in the most general form. In some cases we sacrifice generality in favor of readability and ease of presentation. In some cases we indicate possible extensions of the results under more general conditions. It is important to stress that all the conditions only concern the properties of the quasi likelihood process $L(\boldsymbol{\theta})$.

The imposed conditions on the process can be classified into the following groups by their meaning:

- smoothness conditions on $L(\boldsymbol{\theta})$ allowing the second order Taylor expansion;
- exponential moment conditions;
- identifiability and regularity conditions;

We also distinguish between local and global conditions. The global conditions concern the global behavior of the process $L(\boldsymbol{\theta})$ while the local conditions focus on its behavior in the vicinity of the central point $\boldsymbol{\theta}^*$.

2.1 Global conditions

The first global condition (E) assumes some exponential moments for the quasi log-likelihood $L(\boldsymbol{\theta})$ for each $\boldsymbol{\theta} \in \Theta$. The second condition (ED) assumes some smoothness of $L(\boldsymbol{\theta})$ and requires exponential moments of its gradient. The important identifiability condition is stated later; see Section 4.2.

The formulation involves a subset \mathbb{M} of \mathbb{R}_+ describing all possible exponents in the moment conditions.

(E) For each $\boldsymbol{\theta} \in \Theta$, there exists a positive value $\mu \in \mathbb{M}$ such that

$$E \exp\{\mu L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)\} < \infty.$$

Note that this condition is automatically fulfilled if $\mathbb{P} = \mathbb{P}_{\boldsymbol{\theta}^*}$ and all the $\mathbb{P}_{\boldsymbol{\theta}}$'s are absolutely continuous w.r.t. $\mathbb{P}_{\boldsymbol{\theta}^*}$ with $\mu \leq 1$. Indeed, $L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) = \log(d\mathbb{P}_{\boldsymbol{\theta}}/d\mathbb{P}_{\boldsymbol{\theta}^*})$ and $\log \mathbb{E}_{\boldsymbol{\theta}^*}(d\mathbb{P}_{\boldsymbol{\theta}}/d\mathbb{P}_{\boldsymbol{\theta}^*}) = 0$. For $\mu < 1$, it holds by the Jensen inequality that $-\log \mathbb{E}_{\boldsymbol{\theta}^*}(d\mathbb{P}_{\boldsymbol{\theta}}/d\mathbb{P}_{\boldsymbol{\theta}^*})^\mu \geq 0$.

Condition (E) enables us to define the function

$$\mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*) \stackrel{\text{def}}{=} -\log \mathbb{E} \exp\{\mu L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)\} < \infty. \quad (2.1)$$

This definition can be extended to all $\mu \in \mathbb{M}$ by letting $\mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*) = -\infty$ when the exponential moment of $\mu L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ does not exist. The main observation behind condition (E) is that

$$\mathbb{E} \exp\{\mu L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) + \mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*)\} = 1$$

provided that $\mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*)$ is finite. Note that $\mathfrak{M}(0, \boldsymbol{\theta}, \boldsymbol{\theta}^*) = 0$. The concentration results established in Section 4 require some *identification* properties. A pointwise *identifiability* means that one can separate the target measure $\mathbb{P}_{\boldsymbol{\theta}^*}$ and the measure $\mathbb{P}_{\boldsymbol{\theta}}$ corresponding to another point $\boldsymbol{\theta}$ in the parameter space. This condition can be expressed in terms of the function $\mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*)$ from (2.1). Namely, this value has to be significantly positive for some $\mu \in \mathbb{M}$:

$$\mathfrak{M}^*(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \stackrel{\text{def}}{=} \sup_{\mu \in \mathbb{M}} \mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*) > 0.$$

This condition, however, only ensures a pointwise separation between $\boldsymbol{\theta}^*$ and $\boldsymbol{\theta}$. A *strict* identification requires a quantitative lower bound on the value $\mathfrak{M}^*(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$: it has to grow at least logarithmically with the norm $\|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|$. A precise formulation involves some more notation and it is given in Corollary 4.3 below; see condition (4.9).

To bound local fluctuations of the process $L(\boldsymbol{\theta})$, we introduce an exponential moment condition on the stochastic component $\zeta(\boldsymbol{\theta})$:

$$\zeta(\boldsymbol{\theta}) \stackrel{\text{def}}{=} L(\boldsymbol{\theta}) - \mathbb{E} L(\boldsymbol{\theta}).$$

Suppose that the random function $\zeta(\boldsymbol{\theta})$ is differentiable in $\boldsymbol{\theta}$ and its gradient $\nabla \zeta(\boldsymbol{\theta}) = \partial \zeta(\boldsymbol{\theta}) / \partial \boldsymbol{\theta} \in \mathbb{R}^p$ fulfills the following condition:

(ED) *There exist some constant ν_0 , a positive symmetric matrix V^2 , and constant $g > 0$ such that for all $|\lambda| \leq g$*

$$\sup_{\boldsymbol{\gamma} \in \mathbb{R}^p} \sup_{\boldsymbol{\theta} \in \Theta} \log \mathbb{E} \exp\left\{\lambda \frac{\boldsymbol{\gamma}^\top \nabla \zeta(\boldsymbol{\theta})}{\|V \boldsymbol{\gamma}\|}\right\} \leq \nu_0^2 \lambda^2 / 2.$$

This condition effectively means that the gradient $\nabla\zeta(\boldsymbol{\theta})$ normalized by the matrix V has bounded exponential moments. It can be relaxed by allowing the matrix V^2 and/or the value \mathbf{g} to be dependent of $\boldsymbol{\theta}$ in a uniformly continuous way.

2.2 Local conditions

Local conditions describe the properties of $L(\boldsymbol{\theta})$ in a vicinity of the central point $\boldsymbol{\theta}^*$ from (1.2). First we refine condition (ED). It is fulfilled for all $\boldsymbol{\theta}$ including $\boldsymbol{\theta} = \boldsymbol{\theta}^*$. However, the matrix V can be improved if the only point $\boldsymbol{\theta}^*$ is concerned.

(ED₀) *There exist a positive symmetric matrix V_0^2 , and constants $\mathbf{g} > 0$, $\nu_0 \geq 1$ such that $\text{Var}\{\nabla\zeta(\boldsymbol{\theta}^*)\} \leq V_0^2$ and for all $|\lambda| \leq \mathbf{g}$*

$$\sup_{\boldsymbol{\gamma} \in \mathbb{R}^p} \log \mathbb{E} \exp \left\{ \lambda \frac{\boldsymbol{\gamma}^\top \nabla\zeta(\boldsymbol{\theta}^*)}{\|V_0 \boldsymbol{\gamma}\|} \right\} \leq \nu_0^2 \lambda^2 / 2.$$

In typical situation, the matrix V_0^2 can be defined as the covariance matrix of the gradient vector $\nabla\zeta(\boldsymbol{\theta}^*)$: $V_0^2 = \text{Var}(\nabla\zeta(\boldsymbol{\theta}^*)) = \text{Var}(\nabla L(\boldsymbol{\theta}^*))$. If $L(\boldsymbol{\theta})$ is the log-likelihood for a correctly specified model, then $\boldsymbol{\theta}^*$ is the true parameter value and V_0^2 coincides with the corresponding Fisher information matrix.

The matrix V_0 shown in this condition determines the local geometry in the vicinity of $\boldsymbol{\theta}^*$. In particular, define the local elliptic neighborhoods of $\boldsymbol{\theta}^*$ as

$$\Theta_0(\mathbf{r}) \stackrel{\text{def}}{=} \{\boldsymbol{\theta} \in \Theta : \|V_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\| \leq \mathbf{r}\}. \quad (2.2)$$

The further conditions are restricted to such defined neighborhoods $\Theta_0(\mathbf{r})$. In fact, they quantify local smoothness properties of the log-likelihood function $L(\boldsymbol{\theta})$.

(ED₁) *For some R and each $\mathbf{r} \leq R$, there exist a constant $\omega(\mathbf{r}) \leq 1/2$ such that it holds for all $\boldsymbol{\theta} \in \Theta_0(\mathbf{r})$*

$$\sup_{\boldsymbol{\gamma} \in \mathbb{S}^p} \log \mathbb{E} \exp \left\{ \lambda \frac{\boldsymbol{\gamma}^\top \{\nabla\zeta(\boldsymbol{\theta}) - \nabla\zeta(\boldsymbol{\theta}^*)\}}{\omega(\mathbf{r}) \|V_0 \boldsymbol{\gamma}\|} \right\} \leq \nu_0^2 \lambda^2 / 2, \quad |\lambda| \leq \mathbf{g}.$$

The main majorization result also requires second order smoothness of the expected log-likelihood $\mathbb{E}L(\boldsymbol{\theta})$. By definition, $L(\boldsymbol{\theta}^*, \boldsymbol{\theta}^*) \equiv 0$ and $\nabla \mathbb{E}L(\boldsymbol{\theta}^*) = 0$ because $\boldsymbol{\theta}^*$ is the extreme point of $\mathbb{E}L(\boldsymbol{\theta})$. Therefore, $-\mathbb{E}L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ can be approximated by a quadratic function of $\boldsymbol{\theta} - \boldsymbol{\theta}^*$ in the neighborhood of $\boldsymbol{\theta}^*$. The *local identifiability* condition qualifies this quadratic approximation from above and from below on the set $\Theta_0(\mathbf{r})$ from (2.2).

(\mathcal{L}_0) There are a positive matrix D_0 and for each $\mathbf{r} \leq R$ and a constant $\delta(\mathbf{r}) \leq 1/2$, such that it holds on the set $\Theta_0(\mathbf{r}) = \{\boldsymbol{\theta} : \|V_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\| \leq \mathbf{r}\}$

$$\left| \frac{-2\mathbb{E}L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)}{\|D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2} - 1 \right| \leq \delta(\mathbf{r}).$$

Note that if $L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ is the log-likelihood ratio and $\mathbb{P} = \mathbb{P}_{\boldsymbol{\theta}^*}$ then $-\mathbb{E}L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) = \mathbb{E}_{\boldsymbol{\theta}^*} \log(d\mathbb{P}_{\boldsymbol{\theta}^*}/d\mathbb{P}_{\boldsymbol{\theta}}) = \mathcal{K}(\mathbb{P}_{\boldsymbol{\theta}^*}, \mathbb{P}_{\boldsymbol{\theta}})$, the Kullback-Leibler divergence between $\mathbb{P}_{\boldsymbol{\theta}^*}$ and $\mathbb{P}_{\boldsymbol{\theta}}$. Then condition (\mathcal{L}_0) with $D_0 = V_0$ follows from the usual regularity conditions on the family $(\mathbb{P}_{\boldsymbol{\theta}})$; cf. Ibragimov and Khas'minskij (1981).

If the log-likelihood process $L(\boldsymbol{\theta})$ is sufficiently smooth in $\boldsymbol{\theta}$, e.g. three times stochastically differentiable, then the quantities $\omega(\mathbf{r})$ and $\delta(\mathbf{r})$ are proportional to the radius $\varrho(\mathbf{r})$ of the set $\Theta_0(\mathbf{r})$ defined as

$$\varrho(\mathbf{r}) \stackrel{\text{def}}{=} \max_{\boldsymbol{\theta} \in \Theta_0(\mathbf{r})} \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|.$$

In the important special case of an i.i.d. model one can take $\omega(\mathbf{r}) = \omega^* \mathbf{r}/n^{1/2}$ and $\delta(\mathbf{r}) = \delta^* \mathbf{r}/n^{1/2}$ for some constants ω^*, δ^* ; see Section 5.1.

3 Local non-asymptotic quadraticity

The *Local Asymptotic Normality* (LAN) condition since introduced by L. Le Cam in Le Cam (1960) became one of the central notions in the statistical theory. It postulates a kind of local approximation of the log-likelihood of the original model by the log-likelihood of a Gaussian shift experiment. The LAN property being once checked yields a number of important corollaries for statistical inference. In words, if you can solve a statistical problem for the Gaussian shift model, the result can be translated under the LAN condition to the original setup. We refer to Ibragimov and Khas'minskij (1981) for a nice presentation of the LAN theory including asymptotic efficiency of MLE and Bayes estimators. The LAN properties was extended to *mixed LAN* or *Local Asymptotic Quadraticity* (LAQ); see e.g. Le Cam and Yang (2000). All these notions are very much asymptotic and very much local. The LAN theory also requires that $L(\boldsymbol{\theta})$ is the correctly specified log-likelihood. The strict localization does not allow for considering a growing or infinite parameter dimension and limits applications of the LAN theory to nonparametric estimation.

Our approach tries to avoid asymptotic constructions and attempts to include a possible model misspecification and a large dimension of the parameter space. The presentation below shows that such an extension of the LAN theory can be made essentially by no price: all the major asymptotic results like Fisher and Cramér-Rao information

bounds, as well as the Wilks phenomenon can be derived as corollaries of the obtained non-asymptotic statements simply by letting the sample size to infinity. At the same time, it applies to a high dimensional parameter space.

The LAN property states that the considered process $L(\boldsymbol{\theta})$ can be approximated by a quadratic in $\boldsymbol{\theta}$ expression in a vicinity of the central point $\boldsymbol{\theta}^*$. This property is usually checked using the second order Taylor expansion. The main problem arising here is that the error of the approximation grows too fast with the local size of the neighborhood. Section 3.1 presents the non-asymptotic version of the LAN property in which the local quadratic approximation of $L(\boldsymbol{\theta})$ is replaced by a local quadratic majorization of this process from above and from below by two different quadratic in $\boldsymbol{\theta}$ processes. More precisely, we apply the *sandwiching* idea: the difference $L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) = L(\boldsymbol{\theta}) - L(\boldsymbol{\theta}^*)$ is put between two quadratic processes $\mathbb{L}_{\underline{\epsilon}}(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ and $\mathbb{L}_{\epsilon}(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$:

$$\mathbb{L}_{\underline{\epsilon}}(\boldsymbol{\theta}, \boldsymbol{\theta}^*) - \diamond_{\underline{\epsilon}} \leq L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \leq \mathbb{L}_{\epsilon}(\boldsymbol{\theta}, \boldsymbol{\theta}^*) + \diamond_{\epsilon}, \quad \boldsymbol{\theta} \in \Theta_0(\mathbf{r}) \quad (3.1)$$

where ϵ is a numerical parameter, $\underline{\epsilon} = -\epsilon$, and $\diamond_{\underline{\epsilon}}$ and \diamond_{ϵ} are stochastic errors which only depends on the selected vicinity $\Theta_0(\mathbf{r})$. The upper process $\mathbb{L}_{\epsilon}(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ and the lower process $\mathbb{L}_{\underline{\epsilon}}(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ can deviate substantially from each other, however, the errors $\diamond_{\epsilon}, \diamond_{\underline{\epsilon}}$ remain small even if the value \mathbf{r} describing the size of the local neighborhood $\Theta_0(\mathbf{r})$ is large.

The sandwiching result (3.1) naturally leads to two important notions: value of the problem and deficiency. It turns out that the most of statements like confidence and concentration probability rely upon the maximum of $L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ over $\boldsymbol{\theta}$ which we call *the value of the problem*. Due to (3.1) this value can be bounded from above and from below using the similar quantities $\max_{\boldsymbol{\theta}} \mathbb{L}_{\underline{\epsilon}}(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ and $\max_{\boldsymbol{\theta}} \mathbb{L}_{\epsilon}(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ which can be called the *values of the lower and upper problems*. Note that $\max_{\boldsymbol{\theta}} \{\mathbb{L}_{\epsilon}(\boldsymbol{\theta}, \boldsymbol{\theta}^*) - \mathbb{L}_{\underline{\epsilon}}(\boldsymbol{\theta}, \boldsymbol{\theta}^*)\} = \infty$. However, this is not crucial. What really matters is the difference between the upper and the lower values. The *deficiency* Δ_{ϵ} can be defined as the width of the interval bounding the value of the problem due to (3.1), that is, as the sum of the approximation errors and of this difference:

$$\Delta_{\epsilon} \stackrel{\text{def}}{=} \diamond_{\epsilon} + \diamond_{\underline{\epsilon}} + \left\{ \max_{\boldsymbol{\theta}} \mathbb{L}_{\epsilon}(\boldsymbol{\theta}, \boldsymbol{\theta}^*) - \max_{\boldsymbol{\theta}} \mathbb{L}_{\underline{\epsilon}}(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \right\}.$$

The range of applicability of this approach can be described by the following mnemonic rule: “The value of the upper problem is larger in order than the deficiency.” The further sections explain in details the meaning and content of this rule. Section 3.1 presents the key bound (3.1) and derives it from the general results on empirical processes. Section 3.2 presents some straightforward corollaries of the bound (3.1) including the coverage and concentration probabilities, expansion of the MLE and the risk bounds. It also indicates

how the classical results on asymptotic efficiency of the MLE follow from the obtained non-asymptotic bounds.

3.1 Local quadratic majorization

This section presents the key result about local quadratic approximation of the quasi log-likelihood process given by Theorem 3.1 below.

Let the radius \mathbf{r} of the local neighborhood $\Theta_0(\mathbf{r})$ be fixed in a way that the deviation probability $\mathbb{P}(\tilde{\boldsymbol{\theta}} \notin \Theta_0(\mathbf{r}))$ is sufficiently small. Precise results about the choice of \mathbf{r} which ensures this property are postponed until Section 4.2. In this neighborhood $\Theta_0(\mathbf{r})$ we aim to build a quadratic majorization of the process $L(\boldsymbol{\theta})$. The first step is the usual decomposition of this process into deterministic and stochastic components:

$$L(\boldsymbol{\theta}) = \mathbb{E}L(\boldsymbol{\theta}) + \zeta(\boldsymbol{\theta}),$$

where $\zeta(\boldsymbol{\theta}) = L(\boldsymbol{\theta}) - \mathbb{E}L(\boldsymbol{\theta})$. Condition (\mathcal{L}_0) allows for approximating the smooth deterministic function $\mathbb{E}L(\boldsymbol{\theta}) - \mathbb{E}L(\boldsymbol{\theta}^*)$ around the point of maximum $\boldsymbol{\theta}^*$ by the quadratic form $-\|D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2/2$. The smoothness properties of the stochastic component $\zeta(\boldsymbol{\theta})$ given by conditions (ED_0) and (ED_1) leads to linear approximation $\zeta(\boldsymbol{\theta}) - \zeta(\boldsymbol{\theta}^*) \approx (\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top \nabla \zeta(\boldsymbol{\theta}^*)$. Putting these two approximations together yields the following approximation of the process $L(\boldsymbol{\theta})$ on $\Theta_0(\mathbf{r})$:

$$L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \approx \mathbb{L}(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \stackrel{\text{def}}{=} (\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top \nabla \zeta(\boldsymbol{\theta}^*) - \|D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2/2.$$

This expansion is used in the most of asymptotic statistical calculus. However, it does not suit our purposes because the error of approximation grows quadratically with the radius \mathbf{r} and starts to dominate at some critical value of \mathbf{r} . We slightly modify the construction by introducing two different approximating processes. They only differ in the deterministic quadratic terms which is either shrunk or stretched relative to the term $\|D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2/2$ in $\mathbb{L}(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$.

Introduce for a vector $\boldsymbol{\epsilon} = (\delta, \varrho)$ the following notation:

$$\begin{aligned} \mathbb{L}_{\boldsymbol{\epsilon}}(\boldsymbol{\theta}, \boldsymbol{\theta}^*) &\stackrel{\text{def}}{=} (\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top \nabla L(\boldsymbol{\theta}^*) - \|D_{\boldsymbol{\epsilon}}(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2/2 \\ &= \boldsymbol{\xi}_{\boldsymbol{\epsilon}}^\top D_{\boldsymbol{\epsilon}}(\boldsymbol{\theta} - \boldsymbol{\theta}^*) - \|D_{\boldsymbol{\epsilon}}(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2/2, \end{aligned} \tag{3.2}$$

where

$$D_{\boldsymbol{\epsilon}}^2 = D_0^2(1 - \delta) - \varrho V_0^2, \quad \boldsymbol{\xi}_{\boldsymbol{\epsilon}} \stackrel{\text{def}}{=} D_{\boldsymbol{\epsilon}}^{-1} \nabla L(\boldsymbol{\theta}^*).$$

Here we implicitly assume that with the proposed choice of the constants δ and ϱ , the matrix $D_{\boldsymbol{\epsilon}}^2$ is non-negative: $D_{\boldsymbol{\epsilon}}^2 \geq 0$. The representation (3.2) indicates that the process

$\mathbb{L}_\epsilon(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ has the geometric structure of log-likelihood of a linear Gaussian model. We do not require that the vector $\boldsymbol{\xi}_\epsilon$ is Gaussian and hence, it is not the Gaussian log-likelihood. However, the geometric structure of this process appears to be more important than its distributional properties.

One can see that if δ, ϱ are positive, the quadratic drift component of the process $\mathbb{L}_\epsilon(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ is shrunk relative to $\mathbb{L}(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ and stretched if δ, ϱ are negative. Now, given \mathbf{r} , define $\delta = \delta(\mathbf{r})$, $\varrho = \omega(\mathbf{r})/(3\nu_0)$ with the value $\delta(\mathbf{r})$ from condition (\mathcal{L}_0) and $\omega(\mathbf{r})$ from condition (ED_1) . Finally set $\underline{\epsilon} = -\epsilon$, so that

$$D_{\underline{\epsilon}}^2 = D_0^2(1 + \delta) + \varrho V_0^2.$$

Theorem 3.1. *Assume (ED_1) and (\mathcal{L}_0) . Let for some \mathbf{r} , the values $\varrho \geq 3\nu_0 \omega(\mathbf{r})$ and $\delta \geq \delta(\mathbf{r})$ be such that $D_0^2(1 - \delta) - \varrho V_0^2 \geq 0$. Set $\epsilon = (\delta, \varrho)$, $\underline{\epsilon} = -\epsilon = (-\delta, -\varrho)$. Then*

$$\mathbb{L}_{\underline{\epsilon}}(\boldsymbol{\theta}, \boldsymbol{\theta}^*) - \diamond_{\underline{\epsilon}}(\mathbf{r}) \leq L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \leq \mathbb{L}_{\epsilon}(\boldsymbol{\theta}, \boldsymbol{\theta}^*) + \diamond_{\epsilon}(\mathbf{r}), \quad \boldsymbol{\theta} \in \Theta_0(\mathbf{r}), \quad (3.3)$$

with $\mathbb{L}_{\epsilon}(\boldsymbol{\theta}, \boldsymbol{\theta}^*), \mathbb{L}_{\underline{\epsilon}}(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ defined by (3.2). Moreover, the random variable $\diamond_{\epsilon}(\mathbf{r})$ fulfills

$$\mathbb{P}\{\varrho^{-1} \diamond_{\epsilon}(\mathbf{r}) \geq \mathfrak{z}_0(\mathbf{x}, p)\} \leq \exp(-\mathbf{x}) \quad (3.4)$$

with $\mathfrak{z}_0(\mathbf{x}, p)$ given for $\mathfrak{g}_0 = \mathfrak{g}\nu_0 \geq 3$ by

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

where $\mathfrak{c}_1 = 2$ for $p \geq 2$ and $\mathfrak{c}_1 = 2.4$ for $p = 1$. Similarly for $\diamond_{\underline{\epsilon}}(\mathbf{r})$.

Proof. Consider for fixed \mathbf{r} and $\epsilon = (\delta, \varrho)$ a quantity

$$\diamond_{\epsilon}(\mathbf{r}) \stackrel{\text{def}}{=} \sup_{\boldsymbol{\theta} \in \Theta_0(\mathbf{r})} \left\{ L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) - \mathbb{E}L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) - (\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top \nabla L(\boldsymbol{\theta}^*) - \frac{\varrho}{2} \|V_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2 \right\}.$$

In view of $\mathbb{E}\nabla L(\boldsymbol{\theta}^*) = 0$, this definition can be rewritten in the form

$$\diamond_{\epsilon}(\mathbf{r}) \stackrel{\text{def}}{=} \sup_{\boldsymbol{\theta} \in \Theta_0(\mathbf{r})} \left\{ \zeta(\boldsymbol{\theta}, \boldsymbol{\theta}^*) - (\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top \nabla \zeta(\boldsymbol{\theta}^*) - \frac{\varrho}{2} \|V_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2 \right\}.$$

Similarly define

$$\begin{aligned} \diamond_{\underline{\epsilon}}(\mathbf{r}) &\stackrel{\text{def}}{=} \sup_{\boldsymbol{\theta} \in \Theta_0(\mathbf{r})} \left\{ L(\boldsymbol{\theta}^*, \boldsymbol{\theta}) - \mathbb{E}L(\boldsymbol{\theta}^*, \boldsymbol{\theta}) - (\boldsymbol{\theta}^* - \boldsymbol{\theta})^\top \nabla L(\boldsymbol{\theta}^*) - \frac{\varrho}{2} \|V_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2 \right\} \\ &= \sup_{\boldsymbol{\theta} \in \Theta_0(\mathbf{r})} \left\{ \zeta(\boldsymbol{\theta}^*, \boldsymbol{\theta}) - (\boldsymbol{\theta}^* - \boldsymbol{\theta})^\top \nabla \zeta(\boldsymbol{\theta}^*) - \frac{\varrho}{2} \|V_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2 \right\}. \end{aligned}$$

Now the claim of the theorem can be easily reduced to an exponential bound for the quantities $\diamond_{\epsilon}(\mathbf{r}), \diamond_{\underline{\epsilon}}(\mathbf{r})$. We apply Theorem 7.9 to the process

$$\mathcal{U}(\boldsymbol{\theta}, \boldsymbol{\theta}^*) = \frac{1}{\omega(\mathbf{r})} \{ \zeta(\boldsymbol{\theta}, \boldsymbol{\theta}^*) - (\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top \nabla \zeta(\boldsymbol{\theta}^*) \}, \quad \boldsymbol{\theta} \in \Theta_0(\mathbf{r}),$$

and $H_0 = V_0$. Condition $(\mathcal{E}D)$ follows from (ED_1) with the same ν_0 and \mathbf{g} in view of $\nabla \mathcal{U}(\boldsymbol{\theta}, \boldsymbol{\theta}^*) = \{ \nabla \zeta(\boldsymbol{\theta}) - \nabla \zeta(\boldsymbol{\theta}^*) \} / \omega(\mathbf{r})$. So, the conditions of Theorem 7.9 are fulfilled yielding (3.4) in view of $\varrho \geq 3\nu_0 \omega(\mathbf{r})$. \square

Remark 3.1. If \mathbf{x} is not too big, then the value $\mathfrak{z}_0(\mathbf{x}, p)$ is close to $\mathbf{x} + \mathbf{c}_1 p$; cf. (7.7). The bound (3.4) tells us that the errors $\diamond_{\epsilon}(\mathbf{r})$ and $\diamond_{\underline{\epsilon}}(\mathbf{r})$ are of order $\omega(\mathbf{r})p$.

3.2 Local inference. Deficiency

This section presents a list of corollaries from the basic approximation bounds of Theorem 3.1. The idea is to replace the original problem by a similar one for the approximating upper and lower models. It is important to stress once again that all the corollaries only rely on the *majorization* (3.1) and the *geometric structure* of the processes \mathbb{L}_{ϵ} and $\mathbb{L}_{\underline{\epsilon}}$.

The random quantity $\sup_{\boldsymbol{\theta}} \mathbb{L}_{\epsilon}(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ can be called the *value of the upper problem*, while $\sup_{\boldsymbol{\theta}} \mathbb{L}_{\underline{\epsilon}}(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ is the *value of the lower problem*. The quadratic (in $\boldsymbol{\theta}$) structure of the functions $\mathbb{L}_{\epsilon}(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ and $\mathbb{L}_{\underline{\epsilon}}(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ enables us to explicitly solve the problem of maximizing the corresponding function w.r.t. $\boldsymbol{\theta}$. The value of the original problem which is the maximum of the original process $L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ over $\boldsymbol{\theta}$ is sandwiched between the similar expressions for the two approximating processes in view of the approximation bound (3.3). This suggests to measure the quality of approximation by the difference between values of the upper and lower approximating problems. The approximating quality is sufficiently good if this difference is smaller in order than the value itself.

3.2.1 Upper and lower values

First consider the maximization problem for the upper approximating processes $\mathbb{L}_{\epsilon}(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$. This is a quadratic optimization with the closed form solution $\|\boldsymbol{\xi}_{\epsilon}\|^2/2$. Moreover, the Euclidean norm of the random vector $\boldsymbol{\xi}_{\epsilon}$ behaves nearly as a chi-squared random variable. The results below make these statements more concise. Define for

$$\mathbb{B}_{\epsilon} \stackrel{\text{def}}{=} \frac{D_{\epsilon}^{-1} V_0^2 D_{\epsilon}^{-1}}{\lambda_{\max}(D_{\epsilon}^{-1} V_0^2 D_{\epsilon}^{-1})}, \quad \mathbf{p}_{\epsilon} \stackrel{\text{def}}{=} \text{tr}(\mathbb{B}_{\epsilon}), \quad \mathbf{v}_{\epsilon}^2 \stackrel{\text{def}}{=} 2 \text{tr}(\mathbb{B}_{\epsilon}^2). \quad (3.5)$$

Moreover, with the constant \mathbf{g} from (ED_0) , define also $\mu_c = \mathbf{g}^2/(\mathbf{p}_\epsilon + \mathbf{g}^2)$, and

$$\begin{aligned} \mathbf{y}_c^2 &\stackrel{\text{def}}{=} \mathbf{g}^2/\mu_c^2 - \mathbf{p}_\epsilon/\mu_c, \\ \mathbf{g}_c &\stackrel{\text{def}}{=} \mu_c \mathbf{y}_c = \sqrt{\mathbf{g}^2 - \mu_c \mathbf{p}_\epsilon}, \\ 2\mathbf{x}_c &\stackrel{\text{def}}{=} \mathbf{g}_c \mathbf{y}_c + \log \det(\mathbf{I}_p - \mu_c \mathbb{B}_\epsilon^2). \end{aligned}$$

To gain some feeling of these quantities consider a special case with $\mathbf{g}^2 = \mathbf{p}_\epsilon$. Then $\mu_c = 1/2$, and the inequality $x + \log(1 - x) \geq x/3$ for $0 \leq x \leq 1/2$ implies

$$\begin{aligned} \mathbf{y}_c &= 4\mathbf{g}^2 - 2\mathbf{p}_\epsilon = 2\mathbf{p}_\epsilon, \\ \mathbf{g}_c^2 &= (\mathbf{y}_c/2)^2 = \mathbf{p}_\epsilon/2, \\ 2\mathbf{x}_c &= \sqrt{\mathbf{p}_\epsilon/2 \cdot 2\mathbf{p}_\epsilon} + \log \det(\mathbf{I}_p - \mathbb{B}_\epsilon/2) \geq \mathbf{p}_\epsilon/3. \end{aligned}$$

Theorem 3.2. *Assume (ED_0) with $\nu_0 = 1$ and $\mathbf{g}^2 \geq 2\mathbf{p}_\epsilon$. It holds*

$$\sup_{\boldsymbol{\theta} \in \mathbb{R}^p} \mathbb{L}_\epsilon(\boldsymbol{\theta}, \boldsymbol{\theta}^*) = \|\boldsymbol{\xi}_\epsilon\|^2/2. \quad (3.6)$$

Moreover, $\mathbb{E}\|\boldsymbol{\xi}_\epsilon\|^2 \leq \mathbf{p}_\epsilon$, and for each $\mathbf{x} > 0$

$$\mathbb{P}(\|\boldsymbol{\xi}_\epsilon\|^2 \geq \mathfrak{z}(\mathbf{x}, \mathbb{B}_\epsilon)) \leq 2e^{-\mathbf{x}} + 8.4e^{-\mathbf{x}_c}, \quad (3.7)$$

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

$$\mathfrak{z}(\mathbf{x}, \mathbb{B}_\epsilon) \stackrel{\text{def}}{=} \begin{cases} \mathbf{p}_\epsilon + \sqrt{2\mathbf{x}\mathbf{v}_\epsilon}, & \mathbf{x} \leq \mathbf{v}_\epsilon/18, \\ \mathbf{p}_\epsilon + 6\mathbf{x} & \mathbf{v}_\epsilon/18 < \mathbf{x} \leq \mathbf{x}_c, \\ |y_c + 2(\mathbf{x} - \mathbf{x}_c)/\mathbf{g}_c|^2, & \mathbf{x} > \mathbf{x}_c. \end{cases}$$

Proof. The unconstrained maximum of the quadratic form $\mathbb{L}_\epsilon(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ w.r.t. $\boldsymbol{\theta}$ is attained at $\tilde{\boldsymbol{\theta}}$ yielding the expression (3.6). The second moment of $\boldsymbol{\xi}_\epsilon$ can be bounded from condition (ED_0) . Indeed,

$$\begin{aligned} \mathbb{E}\|\boldsymbol{\xi}_\epsilon\|^2 &= \mathbb{E} \operatorname{tr} \boldsymbol{\xi}_\epsilon \boldsymbol{\xi}_\epsilon^\top \\ &= \operatorname{tr} D_\epsilon^{-1} [\mathbb{E} \nabla L(\boldsymbol{\theta}^*) \{\nabla L(\boldsymbol{\theta}^*)\}^\top] D_\epsilon^{-1} = \operatorname{tr} [D_\epsilon^{-2} \operatorname{Var}\{\nabla L(\boldsymbol{\theta}^*)\}] \end{aligned}$$

and (ED_0) implies $\boldsymbol{\gamma}^\top \operatorname{Var}\{\nabla L(\boldsymbol{\theta}^*)\} \boldsymbol{\gamma} \leq \boldsymbol{\gamma}^\top V_0^2 \boldsymbol{\gamma}$ and thus, $\mathbb{E}\|\boldsymbol{\xi}_\epsilon\|^2 \leq \mathbf{p}_\epsilon$. The deviation bound (3.7) is proved in Corollary 6.12. \square

Next result describes the *lower value* $\sup_{\boldsymbol{\theta}} \mathbb{L}_\epsilon(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ and the *difference* between upper and lower values. The proof is straightforward.

Theorem 3.3. *On the random set $\{\|\xi_\epsilon\| \leq \mathbf{r}\}$, it holds*

$$\sup_{\boldsymbol{\theta} \in \Theta_0(\mathbf{r})} \mathbb{L}_\epsilon(\boldsymbol{\theta}, \boldsymbol{\theta}) = \|\xi_\epsilon\|^2/2.$$

Moreover, $\|\xi_\epsilon\| \leq \|\xi_\epsilon\|$ and

$$\|\xi_\epsilon\|^2 - \|\xi_\epsilon\|^2 = \xi_\epsilon^\top (I_p - D_\epsilon D_\epsilon^{-2} D_\epsilon) \xi_\epsilon \leq \alpha_\epsilon \|\xi_\epsilon\|^2,$$

with

$$\alpha_\epsilon \stackrel{\text{def}}{=} \|I_p - D_\epsilon D_\epsilon^{-2} D_\epsilon\|_\infty = \lambda_{\max}(I_p - D_\epsilon D_\epsilon^{-2} D_\epsilon). \quad (3.8)$$

If the value α_ϵ is small then the difference $\|\xi_\epsilon\|^2 - \|\xi_\epsilon\|^2$ is automatically smaller than the upper value $\|\xi_\epsilon\|^2$.

3.2.2 Deficiency

In view of the results of Theorems 3.2 and 3.3, the sandwiching approach (3.3) bounds the value $L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) = \sup_{\boldsymbol{\theta} \in \Theta} L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ in the interval $[\|\xi_\epsilon\|^2 - \diamond_\epsilon(\mathbf{r}), \|\xi_\epsilon\|^2 + \diamond_\epsilon(\mathbf{r})]$. The width of this interval describes the accuracy of the approach. By analogy to the general Le Cam theory of statistical experiments, this value will be called the *deficiency*. Define the value $\Delta_\epsilon(\mathbf{r})$ by

$$\Delta_\epsilon(\mathbf{r}) \stackrel{\text{def}}{=} \diamond_\epsilon(\mathbf{r}) + \diamond_\epsilon(\mathbf{r}) + (\|\xi_\epsilon\|^2 - \|\xi_\epsilon\|^2)/2. \quad (3.9)$$

This quantity is random but it can be easily evaluated under the considered conditions. Indeed, the approximation errors $\diamond_\epsilon(\mathbf{r}), \diamond_\epsilon(\mathbf{r})$ can be bounded by $\varrho \mathfrak{z}_0(\mathbf{x}, p)$ with the probability at least $1 - 2e^{-\mathbf{x}}$; see (3.4). Also $\|\xi_\epsilon\|^2 \leq \mathfrak{z}(\mathbf{x}, \mathcal{B}_\epsilon)$ with a probability of order $1 - 2e^{-\mathbf{x}}$; see (3.7). This yields for the deficiency $\Delta_\epsilon(\mathbf{r})$ with a probability about $1 - 4e^{-\mathbf{x}}$

$$\Delta_\epsilon(\mathbf{r}) \leq 2\varrho \mathfrak{z}_0(\mathbf{x}, p) + \alpha_\epsilon \mathfrak{z}(\mathbf{x}, \mathcal{B}_\epsilon).$$

3.2.3 The regular case

The bound (3.9) can be further specified in the so called *regular case* under the condition

$$V_0 \leq \mathfrak{a} D_0. \quad (3.10)$$

If the parametric assumption is correct, that is, $\mathcal{P} = \mathcal{P}_{\boldsymbol{\theta}^*}$ for a regular parametric family, then the both matrices coincide with the total Fisher information matrix, and the regularity condition is fulfilled automatically with $\mathfrak{a} = 1$. Otherwise, the regularity

means that the local variability of the process $L(\boldsymbol{\theta})$ measured by the matrix V_0 is not significantly larger than the local information measured by the matrix D_0 .

Theorem 3.4. *Suppose (3.10). Then*

$$D_\epsilon^2 \geq (1 - \delta - \varrho \mathfrak{a}^2) D_0^2, \quad \alpha_\epsilon = \|I_p - D_\epsilon D_\epsilon^{-2} D_\epsilon\|_\infty \leq \frac{2(\delta + \varrho \mathfrak{a}^2)}{1 - \delta - \varrho \mathfrak{a}^2}.$$

Moreover, $B_\epsilon^2 = D_\epsilon^{-1} V_0^2 D_\epsilon^{-1}$ satisfies with $\mathfrak{a}_\epsilon \stackrel{\text{def}}{=} \mathfrak{a}(1 - \delta - \varrho \mathfrak{a}^2)$:

$$B_\epsilon^2 \leq \mathfrak{a}_\epsilon^{-2} I_p, \quad \mathfrak{p}_\epsilon \leq \mathfrak{a}_\epsilon^{-2} p, \quad \mathfrak{v}_\epsilon^2 \leq 2\mathfrak{a}_\epsilon^{-4} p, \quad \lambda_\epsilon \leq \mathfrak{a}_\epsilon^{-2}. \quad (3.11)$$

Proof. The results follow directly from the definition of D_ϵ and D_ϵ and (3.5) by making use of (3.10). \square

In particular, the matrices D_ϵ and D_ϵ are close to each other if δ and $\varrho \mathfrak{a}^2$ are small. So, all we need in the regular case, is a large deviation bound for the probability $\mathbb{P}(\tilde{\boldsymbol{\theta}} \notin \Theta_0(\mathbf{r}))$ and that the quantities $\omega(\mathbf{r})$ and $\delta(\mathbf{r})$ are small.

3.2.4 Local coverage probability

Now we state some immediate corollaries of the exponential bound from Theorem 3.1. First we study the probability of covering $\boldsymbol{\theta}^*$ by the random set $\mathcal{E}(\mathfrak{z}) = \{\boldsymbol{\theta} : 2L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}) \leq \mathfrak{z}\}$.

Theorem 3.5. *Suppose (ED_0) , (ED_1) , and (\mathcal{L}_0) on $\Theta_0(\mathbf{r})$, and let $\varrho \geq 3\nu_0 \omega(\mathbf{r})$, $\delta \geq \delta(\mathbf{r})$, and $D_0^2(1 - \delta) - \varrho V_0^2 \geq 0$. Then for any $\mathfrak{z} > 0$, it holds*

$$\begin{aligned} \mathbb{P}\{\mathcal{E}(\mathfrak{z}) \not\ni \boldsymbol{\theta}^*, \tilde{\boldsymbol{\theta}} \in \Theta_0(\mathbf{r})\} &= \mathbb{P}\{2L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) > \mathfrak{z}, \tilde{\boldsymbol{\theta}} \in \Theta_0(\mathbf{r})\} \\ &\leq \mathbb{P}\{\|\boldsymbol{\xi}_\epsilon\|^2 \geq \mathfrak{z} - \diamond_\epsilon(\mathbf{r})\}. \end{aligned} \quad (3.12)$$

Proof. The bound (3.12) follows from the upper bound of Theorem 3.1 and the statement (3.6) of Lemma 3.2. \square

The exponential bound (3.7) helps to answer a very important question about a proper choice of the critical value \mathfrak{z} providing the prescribed covering probability. Namely, this probability starts to decrease gradually when \mathfrak{z} grows over $\mathfrak{z}(\mathbf{x}, B_\epsilon)$.

3.2.5 Local concentration

Now we describe local concentration properties of $\tilde{\boldsymbol{\theta}}$ assuming that $\tilde{\boldsymbol{\theta}}$ is restricted to $\Theta_0(\mathbf{r})$. More precisely, we bound the probability that $\tilde{\boldsymbol{\theta}}$ does not belong to the set

$\mathcal{A}_\epsilon(z)$ of the form

$$\mathcal{A}_\epsilon(z) = \{\boldsymbol{\theta} : \|D_\epsilon(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\| \leq z\}.$$

It is obvious that

$$\{\tilde{\boldsymbol{\theta}} \notin \mathcal{A}_\epsilon(z)\} = \{\|D_\epsilon(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\| > z\}.$$

Theorem 3.6. Assume (ED_0) , (ED_1) , and (\mathcal{L}_0) on $\Theta_0(\mathbf{r})$, and let $\varrho \geq 3\nu_0\omega(\mathbf{r})$, $\delta \geq \delta(\mathbf{r})$, and $D_0^2(1 - \delta) - \varrho V_0^2 \geq 0$. Then for any $z > 0$, it holds with α_ϵ from (3.8)

$$\mathbb{P}\{\|D_\epsilon(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\| > z, \tilde{\boldsymbol{\theta}} \in \Theta_0(\mathbf{r})\} \leq \mathbb{P}\{\|\boldsymbol{\xi}_\epsilon\| > z - \sqrt{2\Delta_\epsilon(\mathbf{r})}\} \quad (3.13)$$

$$\leq \mathbb{P}\{(1 - \sqrt{\alpha_\epsilon})\|\boldsymbol{\xi}_\epsilon\| > z - \sqrt{2\Diamond_\epsilon(\mathbf{r}) + 2\Diamond_\epsilon(\mathbf{r})}\}. \quad (3.14)$$

Proof. It obviously holds on the set $\{\tilde{\boldsymbol{\theta}} \in \Theta_0(\mathbf{r})\}$

$$\begin{aligned} \{\tilde{\boldsymbol{\theta}} \notin \mathcal{A}_\epsilon(z)\} &= \left\{ \sup_{\boldsymbol{\theta} \notin \mathcal{A}_\epsilon(z)} L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) = \sup_{\boldsymbol{\theta}} L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \right\} \\ &\subseteq \left\{ \sup_{\boldsymbol{\theta} \notin \mathcal{A}_\epsilon(z)} \mathbb{L}_\epsilon(\boldsymbol{\theta}, \boldsymbol{\theta}^*) + \Diamond_\epsilon(\mathbf{r}) \geq \sup_{\boldsymbol{\theta}} \mathbb{L}_\epsilon(\boldsymbol{\theta}, \boldsymbol{\theta}^*) - \Diamond_\epsilon(\mathbf{r}) \right\}. \end{aligned}$$

As $\mathbb{L}_\epsilon(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ is a quadratic function of $D_\epsilon(\boldsymbol{\theta} - \boldsymbol{\theta}^*)$; cf. (3.2), its maximum on the complement $\mathcal{A}_\epsilon^c(z)$ of the set $\mathcal{A}_\epsilon(z)$ is attained at the point $\boldsymbol{\theta}$ satisfying $D_\epsilon(\boldsymbol{\theta} - \boldsymbol{\theta}^*) = \gamma\boldsymbol{\xi}_\epsilon$ with $\gamma = z/\|\boldsymbol{\xi}_\epsilon\|$. This implies for all $\boldsymbol{\theta} \notin \mathcal{A}_\epsilon(z)$

$$\mathbb{L}_\epsilon(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \leq \gamma\|\boldsymbol{\xi}_\epsilon\|^2 - \gamma^2\|\boldsymbol{\xi}_\epsilon\|^2/2 = z\|\boldsymbol{\xi}_\epsilon\| - z^2/2.$$

By Lemma 3.2 $\sup_{\boldsymbol{\theta}} \mathbb{L}_\epsilon(\boldsymbol{\theta}, \boldsymbol{\theta}^*) = \|\boldsymbol{\xi}_\epsilon\|^2/2$. Therefore,

$$\begin{aligned} \{\tilde{\boldsymbol{\theta}} \notin \mathcal{A}_\epsilon(z)\} &\subseteq \{z\|\boldsymbol{\xi}_\epsilon\| - z^2/2 \geq \|\boldsymbol{\xi}_\epsilon\|^2/2 - \Diamond_\epsilon(\mathbf{r}) - \Diamond_\epsilon(\mathbf{r})\} \\ &= \{z^2/2 - z\|\boldsymbol{\xi}_\epsilon\| + \|\boldsymbol{\xi}_\epsilon\|^2/2 \leq \Delta_\epsilon(\mathbf{r})\} \end{aligned}$$

and (3.13) follows. Further, the bound $\|\boldsymbol{\xi}_\epsilon\|^2 - \|\boldsymbol{\xi}_\epsilon\|^2 \leq \alpha_\epsilon\|\boldsymbol{\xi}_\epsilon\|^2$ implies

$$\sqrt{2\Delta_\epsilon(\mathbf{r})} \leq \sqrt{2\Diamond_\epsilon(\mathbf{r}) + 2\Diamond_\epsilon(\mathbf{r})} + \sqrt{\alpha_\epsilon}\|\boldsymbol{\xi}_\epsilon\|$$

which yields (3.14). \square

An interesting and important question is for which z the probability of the event $\{\|D_\epsilon(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\| > z\}$ becomes small. We use that $\max\{\Diamond_\epsilon(\mathbf{r}), \Diamond_\epsilon(\mathbf{r})\} \leq \varrho\mathfrak{z}_0(\mathbf{x}, p)$ and on a set of probability at least $1 - 2e^{-\mathbf{x}}$. This and the bound (3.14) imply

$$\sqrt{2\Delta_\epsilon(\mathbf{r})} \leq 2\sqrt{\varrho\mathfrak{z}_0(\mathbf{x}, p)} + \sqrt{\alpha_\epsilon}\|\boldsymbol{\xi}_\epsilon\|.$$

Now it follows from the local concentration result of Theorem 3.6:

$$\begin{aligned} & \mathbb{P}\{\|D_\epsilon(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\| > z, \tilde{\boldsymbol{\theta}} \in \Theta_0(\mathbf{r})\} \\ & \leq \mathbb{P}\left\{(1 - \sqrt{\alpha_\epsilon})\|\boldsymbol{\xi}_\epsilon\| > z - 2\sqrt{\varrho \mathfrak{z}_0(\mathbf{x}, p)}\right\} + 2e^{-x}. \end{aligned}$$

The probability $\mathbb{P}(\|\boldsymbol{\xi}_\epsilon\| > z)$ starts to vanish when z^2 significantly exceeds $\mathbf{p}_\epsilon = \mathbb{E}\|\boldsymbol{\xi}_\epsilon\|^2$. Under the regularity conditions, the value \mathbf{p}_ϵ is of order p . Moreover, $\mathfrak{z}_0(\mathbf{x}, p)$ is also of order p for moderate \mathbf{x} . If ϱ and α_ϵ are small then the latter deviation probability can be bounded by $\mathbb{P}(\|\boldsymbol{\xi}_\epsilon\| > z')$ with $z/z' \approx 1$ which can be evaluated by (3.7).

3.2.6 Local expansions

Now we show how the bound (3.3) can be used for obtaining a local expansion of the quasi MLE $\tilde{\boldsymbol{\theta}}$. The basic idea is to plug $\tilde{\boldsymbol{\theta}}$ in place of $\boldsymbol{\theta}$ in the definition of $\diamond_\epsilon(\mathbf{r})$.

Theorem 3.7. *Assume (ED_0) , (ED_1) , and (\mathcal{L}_0) on $\Theta_0(\mathbf{r})$ and let $\varrho \geq 3\nu_0\omega(\mathbf{r})$, $\delta \geq \delta(\mathbf{r})$, and $D_0^2(1 - \delta) - \varrho V_0^2 \geq 0$. Then the following approximation holds on the random set $\mathcal{C}_\epsilon(\mathbf{r}) = \{\tilde{\boldsymbol{\theta}} \in \Theta_0(\mathbf{r}), \|\boldsymbol{\xi}_\epsilon\| \leq \mathbf{r}\}$:*

$$\|\boldsymbol{\xi}_\epsilon\|^2/2 - \diamond_\epsilon(\mathbf{r}) \leq L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) \leq \|\boldsymbol{\xi}_\epsilon\|^2/2 + \diamond_\epsilon(\mathbf{r}). \quad (3.15)$$

Moreover, it holds on the same random set $\mathcal{C}_\epsilon(\mathbf{r})$

$$\|D_\epsilon(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) - \boldsymbol{\xi}_\epsilon\|^2 \leq 2\Delta_\epsilon(\mathbf{r}). \quad (3.16)$$

Proof. The bound (3.3) together with Lemma 3.2 yield on $\mathcal{C}_\epsilon(\mathbf{r})$

$$\begin{aligned} L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) &= \sup_{\boldsymbol{\theta} \in \Theta_0(\mathbf{r})} L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \\ &\geq \sup_{\boldsymbol{\theta} \in \Theta_0(\mathbf{r})} \mathbb{L}_\epsilon(\boldsymbol{\theta}, \boldsymbol{\theta}^*) - \diamond_\epsilon(\mathbf{r}) = \|\boldsymbol{\xi}_\epsilon\|^2/2 - \diamond_\epsilon(\mathbf{r}). \end{aligned} \quad (3.17)$$

Similarly

$$L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) \leq \sup_{\boldsymbol{\theta} \in \Theta_0(\mathbf{r})} \mathbb{L}_\epsilon(\boldsymbol{\theta}, \boldsymbol{\theta}^*) + \diamond_\epsilon(\mathbf{r}) \leq \|\boldsymbol{\xi}_\epsilon\|^2/2 + \diamond_\epsilon(\mathbf{r}) \quad (3.18)$$

yielding (3.15). For getting (3.16), we again apply the inequality $L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \leq \mathbb{L}_\epsilon(\boldsymbol{\theta}, \boldsymbol{\theta}^*) + \diamond_\epsilon(\mathbf{r})$ from Theorem 3.1 for $\boldsymbol{\theta}$ equal to $\tilde{\boldsymbol{\theta}}$. With $\boldsymbol{\xi}_\epsilon = D_\epsilon^{-1}\nabla L(\boldsymbol{\theta}^*)$ and $\mathbf{u}_\epsilon \stackrel{\text{def}}{=} D_\epsilon(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)$, this gives

$$L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) - \boldsymbol{\xi}_\epsilon^\top \mathbf{u}_\epsilon + \|\mathbf{u}_\epsilon\|^2/2 \leq \diamond_\epsilon(\mathbf{r}).$$

Therefore, by (3.17)

$$\|\xi_\epsilon\|^2/2 - \diamond_\epsilon(\mathbf{r}) - \xi_\epsilon^\top \mathbf{u}_\epsilon + \|\mathbf{u}_\epsilon\|^2/2 \leq \diamond_\epsilon(\mathbf{r})$$

or, equivalently

$$\|\xi_\epsilon\|^2/2 - \xi_\epsilon^\top \mathbf{u}_\epsilon + \|\mathbf{u}_\epsilon\|^2/2 \leq \diamond_\epsilon(\mathbf{r}) + \diamond_\epsilon(\mathbf{r}) + (\|\xi_\epsilon\|^2 - \|\xi_\epsilon\|^2)/2$$

and the definition of $\Delta_\epsilon(\mathbf{r})$ implies $\|\mathbf{u}_\epsilon - \xi_\epsilon\|^2 \leq 2\Delta_\epsilon(\mathbf{r})$. \square

3.2.7 A local risk bound

Below we also bound the moments of the excess $L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*)$ when $\tilde{\boldsymbol{\theta}}$ is restricted to the local vicinity $\Theta_0(\mathbf{r})$ of $\boldsymbol{\theta}^*$.

Theorem 3.8. *Assume (ED_0) , (ED_1) , and (\mathcal{L}_0) on $\Theta_0(\mathbf{r})$ and let $\varrho \geq 3\nu_0\omega(\mathbf{r})$, $\delta \geq \delta(\mathbf{r})$, and $D_0^2(1 - \delta) - \varrho V_0^2 \geq 0$. Then for $u > 0$*

$$\mathbb{E} L^u(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) \mathbb{I}(\tilde{\boldsymbol{\theta}} \in \Theta_0(\mathbf{r})) \leq \mathbb{E} [\|\xi_\epsilon\|^2/2 + \diamond_\epsilon(\mathbf{r})]^u. \quad (3.19)$$

Moreover, for $\mathcal{C}_\epsilon(\mathbf{r}) = \{\tilde{\boldsymbol{\theta}} \in \Theta_0(\mathbf{r}), \|\xi_\epsilon\| \leq \mathbf{r}\}$, it holds

$$\mathbb{E} \|D_\epsilon(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\|^u \mathbb{I}(\mathcal{C}_\epsilon(\mathbf{r})) \leq \mathbb{E} \{\|\xi_\epsilon\| + \sqrt{2\Delta_\epsilon(\mathbf{r})}\}^u. \quad (3.20)$$

Proof. The bound (3.19) follows from (3.18). Next, the expansion (3.16) yields on $\mathcal{C}_\epsilon(\mathbf{r})$

$$\|D_\epsilon(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\| \leq \|\xi_\epsilon\| + \sqrt{2\Delta_\epsilon(\mathbf{r})}$$

and (3.20) follows. \square

3.2.8 Range of applicability

The whole proposed approach relies implicitly on the two groups of assumptions: global and local. These assumptions are linked to each other by the value \mathbf{r} . From one side, the *global* assumptions listed in Theorem 4.1 and its corollaries should ensure a sensitive bound for the deviation probability $\mathbb{P}\{\tilde{\boldsymbol{\theta}} \notin \Theta_0(\mathbf{r})\}$. This particularly requires that \mathbf{r} is sufficiently large. In the contrary, the *local* conditions are based on the assumption that the local set $\Theta_0(\mathbf{r})$ is sufficiently small to guarantee that errors of approximation $\diamond_\epsilon(\mathbf{r})$ and $\diamond_\epsilon(\mathbf{r})$ and the deficiency $\Delta_\epsilon(\mathbf{r})$ from (3.9) are small as well in a probabilistic sense.

More precisely, the obtained results are sharp and meaningful if the deficiency $\Delta_\epsilon(\mathbf{r})$ is smaller in order than the value of the upper problem $\|\xi_\epsilon\|^2$. The general results of Section 6 state that $\|\xi_\epsilon\|^2$ concentrates around its expected value $\mathbf{p}_\epsilon \stackrel{\text{def}}{=} \mathbb{E}\|\xi_\epsilon\|^2$. Therefore, the latter condition about the deficiency can be decomposed into two others:

- the values $\diamond_{\epsilon}(\mathbf{r}), \diamond_{\underline{\epsilon}}(\mathbf{r})$ are smaller in order than \mathbf{p}_{ϵ} ;
- the difference $\|\xi_{\epsilon}\|^2 - \|\xi_{\underline{\epsilon}}\|^2$ is smaller in order than \mathbf{p}_{ϵ} ;

Due to the result of Theorem 3.1, the random quantity $\diamond_{\epsilon}(\mathbf{r})$ can be bounded with a probability larger than $1 - e^{-x}$ by $\varrho \mathfrak{z}_0(\mathbf{x}, p)$, where $\mathfrak{z}_0(\mathbf{x}, p) \approx p + \mathbf{x}$ if \mathbf{x} is not too large. So, the first conditions requires that ϱp is smaller in order than \mathbf{p}_{ϵ} . If \mathbf{p}_{ϵ} is of order p (see the regular case in the next section) then the condition “ $\diamond_{\epsilon}(\mathbf{r})$ is small” only requires that ϱ , or equivalently, $\omega(\mathbf{r})$ is small. The same holds for $\diamond_{\underline{\epsilon}}(\mathbf{r})$.

Summarizing the above discussion yields that the local results apply if, for a fixed \mathbf{r} :

- $\omega(\mathbf{r})p/\mathbf{p}_{\epsilon}$ is small;
- α_{ϵ} is small.

In the regular case studied in Section 3.2.3, these two conditions simplify to “ $\omega(\mathbf{r}), \delta(\mathbf{r})$ are small”.

3.2.9 Non-asymptotic efficiency

This section discusses the efficiency issues. The famous Cramér-Rao result describes the lower bound for the estimation risk of an unbiased estimate. For linear models this result implies that the true MLE is efficient while a quasi MLE for a misspecified noise covariance is not. The Le Cam LAN theory transfers this result on the general statistical model under the LAN condition; see e.g. Chapter 3 in Ibragimov and Khas'minskij (1981). The results obtained in the previous sections provide a non-asymptotic version of the LAN approach. As already mentioned in Section 2.2, if $L(\theta)$ is the true log-likelihood function of a regular parametric family and $\mathbb{P} = \mathbb{P}_{\theta^*}$, then both matrices D_0^2 and V_0^2 are equal to the Fisher information matrix of this family. This implies the regularity condition with $\mathfrak{a} = 1$; see (3.10). Suppose in addition that the values $\omega(\mathbf{r}), \delta(\mathbf{r})$ are small, so that

$$\mathfrak{a}_{\epsilon} = 1 - \delta - \varrho = 1 - \delta(\mathbf{r}) - 3\nu_0\omega(\mathbf{r})$$

is close to one. By (3.11) this implies that $D_{\epsilon} \approx D_0$, the matrix \mathbb{B}_{ϵ} is close to the identity matrix I_p and $\mathbf{p}_{\epsilon} \approx p$, $\mathbf{v}_{\epsilon}^2 \approx 2p$, $\lambda_{\epsilon} \approx 1$. This in turn implies that the twice upper value $\|\xi_{\epsilon}\|^2$ behaves as a χ_p^2 random variable and the deficiency $\Delta_{\epsilon}(\mathbf{r})$ is small in probability. In particular, all the corollaries about confidence and coverage probability for $\tilde{\theta}$ reproduce the similar statements for the correct linear Gaussian model. Moreover, the decomposition (3.16) can be rewritten as

$$D_0(\tilde{\theta} - \theta^*) \approx D_0^{-1} \nabla L(\theta^*)$$

and this is the famous expansion of the MLE in the LAN situation yielding the asymptotic normality and all other asymptotic properties of $\tilde{\boldsymbol{\theta}}$ including its asymptotic efficiency. We present a precise statement in Section 5.1 when studying the i.i.d. case.

4 Deviation bounds and concentration of the qMLE

A very important step in the analysis of the qMLE $\tilde{\boldsymbol{\theta}}$ is *localization*. This property means that $\tilde{\boldsymbol{\theta}}$ concentrates in a small vicinity of the central point $\boldsymbol{\theta}^*$. This section states such a concentration bound under the global conditions of Section 2.

Given a local vicinity Θ_0 of Θ , the concentration result describes the deviation probability $\mathbb{P}(\tilde{\boldsymbol{\theta}} \notin \Theta_0)$. The key step in this large deviation bound is made in terms of a *multiscale upper function* for the process $L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \stackrel{\text{def}}{=} L(\boldsymbol{\theta}) - L(\boldsymbol{\theta}^*)$. Namely, we build a deterministic function $\mathcal{C}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*)$ such that the probability

$$\mathbb{P}\left(\sup_{\boldsymbol{\theta} \notin \Theta_0} \sup_{\mu \in \mathbb{M}} [\mu L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) + \mathcal{C}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*)] \geq \mathbf{x}\right) \quad (4.1)$$

is exponentially small in \mathbf{x} . Here μ is a positive *scale* parameter and \mathbb{M} is the discrete set of considered scale values. Concentration sets for $\tilde{\boldsymbol{\theta}}$ can be naturally defined via level sets of the function $\mathcal{C}^*(\boldsymbol{\theta}, \boldsymbol{\theta}^*) = \max_{\mu \in \mathbb{M}} \mathcal{C}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*)$. The bound (4.1) is established by some rather crude methods of the theory of empirical processes; see Section 7. Once established, the concentration properties of $\tilde{\boldsymbol{\theta}}$ can be refined in the local vicinity Θ_0 using local majorization technique; see Section 3.2.

The other important result describes an *upper function* $\mathfrak{b}(\boldsymbol{\theta}, \mathbf{x})$ for the non-scaled process $L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ providing that the probability

$$\mathbb{P}\left(\sup_{\boldsymbol{\theta} \notin \Theta_0} \{L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) + \mathfrak{b}(\boldsymbol{\theta}, \mathbf{x})\} > 0\right)$$

is exponentially small in \mathbf{x} . Such bounds are usually called for in the analysis of the posterior measure in the Bayes approach.

4.1 A multiscale upper function for the log-likelihood

This section presents a construction of the multiscale upper function $\mathcal{C}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*)$. Then it will be used for controlling the deviation probability of the estimate $\tilde{\boldsymbol{\theta}}$. Below we suppose that the scaling factor μ runs over some discrete set \mathbb{M} and consider the process $\mu L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ for $\boldsymbol{\theta} \in \Theta$ and $\mu \in \mathbb{M}$. Assume that for each $\boldsymbol{\theta}$, there is some $\mu \in \mathbb{M}$ such that the exponential moment of $\mu L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ is finite. This enables us to define for each $\boldsymbol{\theta}$ the function $\mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*)$ ensuring the identity

$$\mathbb{E} \exp\{\mu L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) + \mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*)\} = 1.$$

This means that the process $\mu L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) + \mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*)$ is pointwise stochastically bounded in a rather strict sense. For each $\mu > 0$ and $\mathfrak{z} \geq 0$ by the Markov inequality

$$\begin{aligned} \mathbb{P}(L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \geq \mathfrak{z}) &\leq \exp\{-\mu\mathfrak{z}\} \mathbb{E} \exp\{\mu L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)\} \\ &= \exp\{-\mu\mathfrak{z} - \mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*)\}. \end{aligned} \quad (4.2)$$

In particular, with $\mathfrak{z} = 0$,

$$\mathbb{P}(L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \geq 0) \leq \exp\left\{-\max_{\mu \in \mathbb{M}} \mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*)\right\}.$$

So, a reasonable choice of μ can be made via maximization of $\mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*)$ w.r.t. μ . The famous Chernoff result describes the asymptotic separation rate between these two measures $\mathbb{P}_{\boldsymbol{\theta}^*}$ and $\mathbb{P}_{\boldsymbol{\theta}}$ in terms of the value $\mathfrak{M}^*(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ with

$$\mathfrak{M}^*(\boldsymbol{\theta}, \boldsymbol{\theta}^*) = \sup_{\mu \in \mathbb{M}} \mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*). \quad (4.3)$$

The larger $\mathfrak{M}^*(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ is, the stronger is the pointwise identification.

If Θ° is a subset of Θ not containing $\boldsymbol{\theta}^*$, then the event $\tilde{\boldsymbol{\theta}} \in \Theta^\circ$ is only possible if $\sup_{\boldsymbol{\theta} \in \Theta^\circ} L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \geq 0$, because $L(\boldsymbol{\theta}^*, \boldsymbol{\theta}^*) \equiv 0$. It is intuitively clear that the uniform identification over Θ° requires that the rate function $\mathfrak{M}^*(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ is bounded away from zero on Θ° . The result of Theorem 4.1 below quantifies this condition.

We present a result for the smooth case which assumes the conditions (E) , (ED) , and (ED_0) to be fulfilled with the corresponding matrices V and V_0 . The set Θ° is taken as the complement of the local set $\Theta_0(\mathbf{r})$ with a sufficiently large \mathbf{r} . The result also involves a value \mathbf{r}_1 entering into the pilling device; see Section 7.3 for details. The choice of \mathbf{r}_1 is done by the equality $\det(\mathbf{r}_1^{-1}V) = \det(V_0)$ relating two matrices V_0 and V to each others.

For stating the results, some further notations have to be introduced. Define $\zeta(\boldsymbol{\theta}, \boldsymbol{\theta}^*) = \zeta(\boldsymbol{\theta}) - \zeta(\boldsymbol{\theta}^*)$ and

$$\mathfrak{N}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*) \stackrel{\text{def}}{=} \log \mathbb{E} \exp\{\mu \zeta(\boldsymbol{\theta}, \boldsymbol{\theta}^*)\}.$$

Then it holds for the function $\mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*)$:

$$\mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*) = -\log \mathbb{E} \exp\{\mu L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)\} = -\mu \mathbb{E} L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) - \mathfrak{N}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*).$$

Further, consider for each $\boldsymbol{\theta}^\circ \in \Theta$ a local ball $\mathcal{B}_\mu(\boldsymbol{\theta}^\circ) \stackrel{\text{def}}{=} \{\boldsymbol{\theta} \in \Theta : \|V(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\| \leq \mathbf{r}_1/\mu\}$. Next, for the Lebesgue measure π on \mathbb{R}^p , define the smoothing operator \mathbb{S}_μ by

$$\mathbb{S}_\mu f(\boldsymbol{\theta}^\circ) \stackrel{\text{def}}{=} \frac{1}{\pi_\mu(\boldsymbol{\theta}^\circ)} \int_{\mathcal{B}_\mu(\boldsymbol{\theta}^\circ)} f(\boldsymbol{\theta}) \pi(d\boldsymbol{\theta}).$$

Let ν_0 and \mathbf{g} be the constants from (ED). Define \mathbf{c}_1 by $\mathbf{c}_1 = 2$ for $p \geq 2$ and $\mathbf{c}_1 = 2.4$ for $p = 1$. Let a constant s be selected under the condition $3\nu_0\mathbf{r}_1/s \leq \mathbf{g} \wedge \sqrt{2\mathbf{c}_1 p}$. The value $s = 1$ is a proper candidate in typical situations.

For shortening the notation assume that the discrete set \mathbb{M} is fixed under the condition $M \equiv \sum_{\mu \in \mathbb{M}} \mu^{p+2} \leq e/2$.

Theorem 4.1. *Suppose that (E), (ED), and (ED₀) hold with some \mathbf{g}, ν_0 and with matrices V and V_0 . Let \mathbf{r}_1 be such that $\det(\mathbf{r}_1^{-1}V) \leq \det(V_0)$, and s be such that $3\nu_0\mathbf{r}_1/s \leq \mathbf{g} \wedge \sqrt{2\mathbf{c}_1 p}$. For $\mathbf{r} \geq \sqrt{p/2}$ and $\mathbf{x} > 0$, it holds*

$$\mathbb{P}\left\{\sup_{\boldsymbol{\theta} \notin \Theta_0(\mathbf{r})} \sup_{\mu \in \mathbb{M}} [\mu L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) + \mathcal{C}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*)] > \mathfrak{z}_1(\mathbf{x})\right\} \leq e^{-\mathbf{x}+1}, \quad (4.4)$$

with $\mathfrak{z}_1(\mathbf{x}) \stackrel{\text{def}}{=} 2s\mathbf{c}_1 p + (1+s)\mathbf{x}$ and

$$\begin{aligned} \mathcal{C}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*) &\stackrel{\text{def}}{=} -\mu \mathbb{E}L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) - \mathbb{S}_\mu \mathfrak{N}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*) - (1+s)\mathfrak{t}(\boldsymbol{\theta}) \\ \mathfrak{t}(\boldsymbol{\theta}) &\stackrel{\text{def}}{=} (p+2) \log(\|V_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|). \end{aligned} \quad (4.5)$$

Proof. First fix $\mu \in \mathbb{M}$ and apply the general results of Theorem 7.10 to the process $\mathcal{U}(\boldsymbol{\theta}) = \zeta(\boldsymbol{\theta}, \boldsymbol{\theta}^*) = L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) - \mathbb{E}L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$, $M(\boldsymbol{\theta}) = -\mathbb{E}L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$, and $H_0 = V_0$ on $\mathcal{R}^\circ = \Theta_0^c(\mathbf{r})$. The condition (E) is fulfilled, condition (ED) implies (ED) with $H(\boldsymbol{\theta}) = V$, and Theorem 7.10 yields (4.4) in view of $\|\mathbf{r}_1^{-1}V(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\| \leq \|V_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|$. \square

Remark 4.1. The construction of the multiscale upper function $\mathcal{C}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*)$ in (4.5) deserves some discussion. More precisely, it is interesting to compare this construction with the pointwise upper function $\mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*) = -\mu \mathbb{E}L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) - \mathfrak{N}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*)$. The smoothing operator in the term $\mathbb{S}_\mu \mathfrak{N}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*)$ from (4.5) is only used for technical reasons and it can be handled by simple rescaling arguments. The other term $(1+s)\mathfrak{t}(\boldsymbol{\theta})$ in (4.5) is really important and it can be viewed as the price for uniform concentration. Simple white noise examples show that this penalty term is nearly sharp if s is taken small and p large. However, our aim is only to obtain a rough exponential bound because really sharp results are stated by mean of local quadratic approximation; see Section 3.2.5.

4.2 Concentration sets and deviation probability

This section describes so called *concentration sets* for the estimate $\tilde{\boldsymbol{\theta}}$. Any such set is deterministic ensuring that $\tilde{\boldsymbol{\theta}}$ deviates from this set with a small probability. Such concentration sets are usually used in theoretical studies and their construction typically depends on the unknown quantities the moment generation function $\mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*)$. Let

$\mathcal{C}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*)$ be the upper multiscale function from (4.5). Define

$$\mu(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \stackrel{\text{def}}{=} \operatorname{argmax}_{\mu \in \mathbb{M}} \mathcal{C}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*), \quad (4.6)$$

$$\mathcal{C}^*(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \stackrel{\text{def}}{=} \max_{\mu \in \mathbb{M}} \mathcal{C}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*) = \mathcal{C}(\mu(\boldsymbol{\theta}, \boldsymbol{\theta}^*), \boldsymbol{\theta}, \boldsymbol{\theta}^*). \quad (4.7)$$

For simplicity we assume that the maximum in (4.7) is attained and $\mu(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ is well defined. In the considered case of a discrete set \mathbb{M} this is always fulfilled. The value $\mathcal{C}^*(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ shown in (4.7) replaces $\mathfrak{M}^*(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ from (4.3) if one is concerned with uniform non-asymptotic bounds. For any fixed subset $\mathcal{A} \subset \Theta$ define the value $\mathfrak{g}(\mathcal{A})$ by

$$\mathfrak{g}(\mathcal{A}) \stackrel{\text{def}}{=} \inf_{\boldsymbol{\theta} \notin \mathcal{A}} \mathcal{C}^*(\boldsymbol{\theta}, \boldsymbol{\theta}^*). \quad (4.8)$$

A particular choice of the set \mathcal{A} is given by the level set of $\mathcal{C}^*(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$: for every $\mathbf{x} > 0$, define $\mathcal{A}(\mathbf{x}, \boldsymbol{\theta}^*)$ with

$$\mathcal{A}(\mathbf{x}, \boldsymbol{\theta}^*) = \{\boldsymbol{\theta} : \mathcal{C}^*(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \leq \mathfrak{z}_1(\mathbf{x})\} = \left\{ \boldsymbol{\theta} : \sup_{\mu \in \mathbb{M}} \mathcal{C}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*) \leq \mathfrak{z}_1(\mathbf{x}) \right\},$$

where $\mathfrak{z}_1(\mathbf{x}) = 2s\mathfrak{c}_1p + (1+s)\mathbf{x}$. Obviously $\mathfrak{g}(\mathcal{A}(\mathbf{x}, \boldsymbol{\theta}^*)) \geq \mathfrak{z}_1(\mathbf{x})$.

Below we show that the bound (4.4) yields some concentration properties of the estimate $\tilde{\boldsymbol{\theta}}$ in terms of the function $\mathfrak{g}(\cdot)$ from (4.8).

Corollary 4.2. *Under (4.4), it holds for any set \mathcal{A} with $\mathfrak{g}(\mathcal{A}) \geq \mathfrak{z}_1(\mathbf{x})$*

$$\mathbb{P}(\tilde{\boldsymbol{\theta}} \notin \mathcal{A}) \leq e^{-\mathbf{x}+1}.$$

In particular,

$$\log \mathbb{P}(\tilde{\boldsymbol{\theta}} \notin \mathcal{A}(\mathbf{x}, \boldsymbol{\theta}^*)) \leq -\mathbf{x} + 1.$$

Proof. Denote $\tilde{\mu} = \mu(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*)$ with $\mu(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ from (4.6). Then in view of $L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) \geq 0$ and $\mathcal{C}(\tilde{\mu}, \boldsymbol{\theta}, \boldsymbol{\theta}^*) \geq \mathfrak{g}(\mathcal{A})$ for $\tilde{\boldsymbol{\theta}} \notin \mathcal{A}$

$$\begin{aligned} \{\tilde{\boldsymbol{\theta}} \notin \mathcal{A}\} &\subseteq \{\mathcal{C}(\tilde{\mu}, \tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) \geq \mathfrak{g}(\mathcal{A})\} \\ &\subseteq \{\tilde{\mu}L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) + \mathcal{C}(\tilde{\mu}, \tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) > \mathfrak{g}(\mathcal{A})\}, \end{aligned}$$

and the result follows from (4.4). \square

Corollary 4.2 presents an upper bound for the probability that $\tilde{\boldsymbol{\theta}}$ deviates from a set \mathcal{A} defined via the function $\mathcal{C}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*)$. The local approach of Section 3.1 requires also to bound the deviation probability for a local set $\Theta_0(\mathbf{r}) \stackrel{\text{def}}{=} \{\boldsymbol{\theta} : \|V_0(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\| > \mathbf{r}\}$ for

\mathbf{r} sufficiently large. It suffices to evaluate the quantity $\mathfrak{g}(\mathcal{A})$ for $\mathcal{A} = \Theta_0(\mathbf{r})$. This in turns requires to bound from below the value $\mathfrak{C}^*(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ for $\boldsymbol{\theta} \notin \mathcal{A}$.

The next result presents sufficient conditions ensuring a sensible large deviation probability bound in terms of the rate function $\mathfrak{M}^*(\boldsymbol{\theta}, \boldsymbol{\theta}^*) = \max_{\mu} \{-\log \mathbb{E} \exp \mu L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)\}$. Define

$$\delta_{\mu}(\boldsymbol{\theta}) \stackrel{\text{def}}{=} |\mathfrak{N}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*) - \mathbb{S}_{\mu} \mathfrak{N}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*)|.$$

Corollary 4.3. *Suppose the conditions of Theorem 4.1. Let, given $\mathbf{x} > 0$, there be a value $\mathbf{r} = \mathbf{r}(\mathbf{x})$, such that it holds*

$$\mathfrak{M}^*(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \geq (1 + s)\mathfrak{t}(\boldsymbol{\theta}) + \delta_{\mu}(\boldsymbol{\theta}) + \mathfrak{z}_1(\mathbf{x}), \quad \boldsymbol{\theta} \notin \Theta_0(\mathbf{r}), \quad (4.9)$$

with $\mathfrak{t}(\boldsymbol{\theta}) = (p + 2) \log(\|V_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|)$ and $\mathfrak{z}_1(\mathbf{x}) \stackrel{\text{def}}{=} 2s\mathfrak{c}_1 p + (1 + s)\mathbf{x}$. Then

$$\mathbb{P}(\|V_0(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\| > \mathbf{r}(\mathbf{x})) \leq e^{-\mathbf{x}+1}.$$

Proof. The result follows from Corollary 4.2 by making use of the inequality

$$-\mu \mathbb{E} L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) - \mathbb{S}_{\mu} \mathfrak{N}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*) \geq \mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*) - \delta_{\mu}(\boldsymbol{\theta})$$

for any $\mu \in \mathbb{M}$. □

Remark 4.2. Due to this result, a logarithmic growth of $\mathfrak{M}^*(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ as function of the distance $\|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|$ ensures a sensitive large deviation bound for the process $L(\boldsymbol{\theta})$.

4.3 Probability bounds for the quasi log-likelihood

This section presents a uniform upper bound for process $L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ without scaling. Our starting point is again a pointwise bound on $L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ for a fixed $\boldsymbol{\theta} \neq \boldsymbol{\theta}^*$. Namely, given \mathbf{x} , we first try to find $b(\boldsymbol{\theta}, \mathbf{x})$ providing

$$\mathbb{P}(L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) + b(\boldsymbol{\theta}, \mathbf{x}) \geq 0) \leq e^{-\mathbf{x}}. \quad (4.10)$$

Define

$$b(\boldsymbol{\theta}, \mathbf{x}) \stackrel{\text{def}}{=} \max_{\mu \in \mathbb{M}} \mu^{-1} \{-\mathbf{x} + \mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*)\}.$$

Then applying (4.2) with the corresponding value μ yields (4.10). Now we aim to establish an extension of this pointwise bound to a uniform bound on a subset Θ° of Θ , that is, to build a function $\mathfrak{b}(\boldsymbol{\theta}, \mathbf{x})$ such that

$$\mathbb{P}\left(\sup_{\boldsymbol{\theta} \in \Theta^{\circ}} \{L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) + \mathfrak{b}(\boldsymbol{\theta}, \mathbf{x})\} > 0\right) \leq e^{-\mathbf{x}+1}. \quad (4.11)$$

Such bounds are naturally called for in the analysis of the posterior measure in the Bayes approach. The function $\mathfrak{b}(\boldsymbol{\theta}, \mathbf{x})$ can be described via the multiscale upper function $\mathcal{C}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*)$. Namely, for each $\boldsymbol{\theta}$ and $\mathbf{x} > 0$ define

$$\mathfrak{b}(\boldsymbol{\theta}, \mathbf{x}) \stackrel{\text{def}}{=} \max_{\mu \in \mathbb{M}} \{-\mu^{-1} \mathfrak{z}_1(\mathbf{x}) + \mu^{-1} \mathcal{C}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*)\}. \quad (4.12)$$

Corollary 4.4. *Assume the bound (4.4) for a function $\mathcal{C}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*)$. Then (4.11) holds with $\mathfrak{b}(\boldsymbol{\theta}, \mathbf{x})$ from (4.12).*

Proof. The bound (4.4) implies

$$\mathbb{P}\left(\sup_{\boldsymbol{\theta} \in \Theta^\circ} \min_{\mu \in \mathbb{M}} \{\mu L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) + \mathcal{C}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*)\} > \mathfrak{z}_1(\mathbf{x})\right) \leq e^{-\mathbf{x}+1}.$$

This yields that there exists a random set $\Omega(\mathbf{x})$ with $\mathbb{P}(\Omega(\mathbf{x})) \geq 1 - e^{-\mathbf{x}+1}$ such that for any $\mu \in \mathbb{M}$ and any $\boldsymbol{\theta} \in \Theta^\circ$, it holds $\mu L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) + \mathcal{C}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*) \leq \mathfrak{z}_1(\mathbf{x})$ on $\Omega(\mathbf{x})$. Let $\mu(\boldsymbol{\theta}, \mathbf{x})$ fulfill

$$\mu(\boldsymbol{\theta}, \mathbf{x}) \stackrel{\text{def}}{=} \operatorname{argmax}_{\mu \in \mathbb{M}} \{-\mu^{-1} \mathfrak{z}_1(\mathbf{x}) + \mu^{-1} \mathcal{C}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*)\}.$$

Then particular choice $\mu = \mu(\boldsymbol{\theta}, \mathbf{x})$ yields on $\Omega(\mathbf{x})$

$$L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) + \frac{1}{\mu(\boldsymbol{\theta}, \mathbf{x})} \{\mathcal{C}(\mu(\boldsymbol{\theta}, \mathbf{x}), \boldsymbol{\theta}, \boldsymbol{\theta}^*) - \mathfrak{z}_1(\mathbf{x})\} \leq 0$$

or equivalently $L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) + \mathfrak{b}(\boldsymbol{\theta}, \mathbf{x}) \leq 0$. □

5 Examples

The model with independent identically distributed (i.i.d.) observations is one of the most popular setups in statistical literature and in statistical applications. The essential and the most developed part of the statistical theory is designed for the i.i.d. modeling. Especially, the classical asymptotic parametric theory is almost complete including asymptotic root-n normality and efficiency of the MLE and Bayes estimators under rather mild assumptions; see e.g. Chapter 2 and 3 in Ibragimov and Khas'minskij (1981). So, the i.i.d. model can naturally serve as a benchmark for any extension of the statistical theory: being applied to the i.i.d. setup, the new approach should lead to essentially the same conclusions as in the classical theory. Similar reasons apply to the regression model and its extensions. Below we try demonstrate that the proposed non-asymptotic viewpoint is able to reproduce the existing brilliant and well established results of the classical parametric theory. With some surprise, the majority of classical efficiency results can be easily derived from the obtained general non-asymptotic bounds.

The next question is whether there is any added value or benefits of the new approach being restricted to the i.i.d. situation relative to the classical one. Two important issues have been already mentioned: the new approach applies to the situation with finite samples and survives under model misspecification. One more important question is whether the obtained results remain applicable and informative if the dimension of the parameter space is high – this is one of the main challenge in the modern statistics. We show that the dimensionality p naturally appears in the risk bounds and the results apply as long as the sample size exceeds in order this value p . All these questions are addressed in Section 5.1 for the i.i.d. setup, Section 5.2 focuses on generalized linear modeling, while Section 5.3 discusses linear median regression.

5.1 Quasi MLE in an i.i.d. model

The basic i.i.d. parametric model means that the observations $\mathbf{Y} = (Y_1, \dots, Y_n)$ are independent identically distributed from a distribution P from a given parametric family $(P_\theta, \theta \in \Theta)$ on the observation space \mathcal{Y}_1 . Each $\theta \in \Theta$ clearly yields the product data distribution $\mathbb{P}_\theta = P_\theta^{\otimes n}$ on the product space $\mathcal{Y} = \mathcal{Y}_1^n$. This section illustrates how the obtained general results can be applied to this type of modeling under possible model misspecification. Different types of misspecification can be considered. Each of the assumptions, namely, data independence, identical distribution, parametric form of the marginal distribution can be violated. To be specific, we assume the observations Y_i independent and identically distributed. However, we admit that the distribution of each Y_i does not necessarily belong to the parametric family (P_θ) . The case of non-identically distributed observations can be done similarly at cost of more complicated notation.

In what follows the parametric family (P_θ) is supposed to be dominated by a measure μ_0 , and each density $p(y, \theta) = dP_\theta/d\mu_0(y)$ is two times continuously differentiable in θ for all y . Denote $\ell(y, \theta) = \log p(y, \theta)$. The parametric assumption $Y_i \sim P_{\theta^*} \in (P_\theta)$ leads to the log-likelihood

$$L(\theta) = \sum \ell(Y_i, \theta),$$

where the summation is taken over $i = 1, \dots, n$. The quasi MLE $\tilde{\theta}$ maximizes this sum over $\theta \in \Theta$:

$$\tilde{\theta} \stackrel{\text{def}}{=} \operatorname{argmax}_{\theta \in \Theta} L(\theta) = \operatorname{argmax}_{\theta \in \Theta} \sum \ell(Y_i, \theta).$$

The target of estimation θ^* maximizes the expectation of $L(\theta)$:

$$\theta^* \stackrel{\text{def}}{=} \operatorname{argmax}_{\theta \in \Theta} \mathbb{E} L(\theta) = \operatorname{argmax}_{\theta \in \Theta} \sum \mathbb{E} \ell(Y_i, \theta).$$

Let $\zeta_i(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \ell(Y_i, \boldsymbol{\theta}) - \mathbb{E}\ell(Y_i, \boldsymbol{\theta})$. Then $\zeta(\boldsymbol{\theta}) = \sum \zeta_i(\boldsymbol{\theta})$. The equation $\mathbb{E}\nabla L(\boldsymbol{\theta}^*) = 0$ implies

$$\nabla \zeta(\boldsymbol{\theta}^*) = \sum \nabla \zeta_i(\boldsymbol{\theta}^*) = \sum \nabla \ell_i(\boldsymbol{\theta}^*). \quad (5.1)$$

I.i.d. structure of the Y_i 's allows for rewriting the conditions (E) , (ED) , (ED_0) , (ED_1) , and (\mathcal{L}_0) in terms of the marginal distribution. In the following conditions the index i runs from 1 to n .

(e) For each $\boldsymbol{\theta} \in \Theta$, there exists a positive value $\mu \in \mathbb{M}$ such that

$$\mathfrak{m}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*) \stackrel{\text{def}}{=} -\log \mathbb{E} \exp\{\mu[\ell(Y_i, \boldsymbol{\theta}) - \ell(Y_i, \boldsymbol{\theta}^*)]\}$$

is finite.

(ed) There exist some constants ν_0 , and $\mathfrak{g}_1 > 0$, and a positive symmetric $p \times p$ matrix \mathbf{v} , such that for all $|\lambda| \leq \mathfrak{g}_1$

$$\sup_{\boldsymbol{\gamma} \in \mathbb{S}^p} \sup_{\boldsymbol{\theta} \in \Theta} \log \mathbb{E} \exp\left\{\lambda \frac{\boldsymbol{\gamma}^\top \nabla \zeta_i(\boldsymbol{\theta})}{\|\mathbf{v}\boldsymbol{\gamma}\|}\right\} \leq \nu_0^2 \lambda^2 / 2.$$

(ed₀) There exists a positive symmetric matrix \mathbf{v}_0 , such that for all $|\lambda| \leq \mathfrak{g}_1$

$$\sup_{\boldsymbol{\gamma} \in \mathbb{S}^p} \log \mathbb{E} \exp\left\{\lambda \frac{\boldsymbol{\gamma}^\top \nabla \zeta_i(\boldsymbol{\theta}^*)}{\|\mathbf{v}_0\boldsymbol{\gamma}\|}\right\} \leq \nu_0^2 \lambda^2 / 2.$$

A natural candidate on \mathbf{v}_0^2 is given by the variance of the gradient $\nabla \ell(Y_1, \boldsymbol{\theta}^*)$, that is, $\mathbf{v}_0^2 = \text{Var} \nabla \ell(Y_1, \boldsymbol{\theta}) = \text{Var} \nabla \zeta_1(\boldsymbol{\theta})$.

Next consider the local sets

$$\Theta_0(\mathbf{r}) = \{\boldsymbol{\theta} : \|\mathbf{v}_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\| \leq \mathbf{r}/n^{1/2}\}.$$

The local smoothness conditions (ED_1) and (\mathcal{L}_0) require to specify the functions $\delta(\mathbf{r})$ and $\varrho(\mathbf{r})$. If the log-likelihood function $\ell(y, \boldsymbol{\theta})$ is sufficiently smooth in $\boldsymbol{\theta}$, these functions can be selected proportional to \mathbf{r} .

(ed₁) For each $\mathbf{r} \leq R$, there exists a constant ω^* such that for all $i = 1, \dots, n$ and $|\lambda| \leq \mathfrak{g}_1$

$$\sup_{\boldsymbol{\gamma} \in \mathbb{S}^p} \sup_{\boldsymbol{\theta} \in \Theta_0(\mathbf{r})} \log \mathbb{E} \exp\left\{\lambda \frac{\boldsymbol{\gamma}^\top [\nabla \zeta_i(\boldsymbol{\theta}) - \nabla \zeta_i(\boldsymbol{\theta}^*)]}{\omega^* \mathbf{r} \|\mathbf{v}_0\boldsymbol{\gamma}\|}\right\} \leq \nu_0^2 \lambda^2 / 2.$$

Further we restate the local identifiability condition (\mathcal{L}_0) in terms of the expected value $\bar{\ell}(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \mathbb{E}\ell(Y_i, \boldsymbol{\theta})$ of each $\ell(Y_i, \boldsymbol{\theta})$. We suppose that $\bar{\ell}(\boldsymbol{\theta})$ is two times differentiable w.r.t. $\boldsymbol{\theta}$ and define the matrix $\mathbf{F}_0 = -\nabla^2 \bar{\ell}(\boldsymbol{\theta}^*)$.

(ℓ_0) For each $\mathbf{r} \leq R$, there is a constant δ^* , such that it holds on $\Theta_0(\mathbf{r})$

$$\left| \frac{\bar{\ell}(\boldsymbol{\theta}) - \bar{\ell}(\boldsymbol{\theta}^*) - (\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top \nabla \bar{\ell}(\boldsymbol{\theta}^*)}{(\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top \mathbb{F}_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)/2} - 1 \right| \leq \delta^* \mathbf{r}.$$

In the regular parametric case with $\mathbb{P} \in (P_{\boldsymbol{\theta}})$, the matrices \mathbf{v}_0^2 and \mathbb{F}_0 coincide with the Fisher information matrix $\mathbb{F}(\boldsymbol{\theta}^*)$ of the family $(P_{\boldsymbol{\theta}})$ at the point $\boldsymbol{\theta}^*$.

Lemma 5.1. Let Y_1, \dots, Y_n be i.i.d. Then (e) , (ed) , (ed_0) , $(ed)_1$, and (ℓ_0) imply (E) , (ED) , (ED_0) , $(ED)_1$, and (\mathcal{L}_0) with $\mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*) = n\mathfrak{m}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*)$, $V^2 = n\mathbf{v}^2$, $V_0^2 = n\mathbf{v}_0^2$, $D_0^2 = n\mathbb{F}_0$, $\omega(\mathbf{r}) = \omega^* \mathbf{r}$, $\delta(\mathbf{r}) = \delta^* \mathbf{r}$, the same constant ν_0 , and

$$\mathbf{g} \stackrel{\text{def}}{=} \mathbf{g}_1 \sqrt{n}.$$

Proof. The identities $\mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*) = n\mathfrak{m}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*)$, $V^2 = n\mathbf{v}^2$, $V_0^2 = n\mathbf{v}_0^2$, $D_0^2 = n\mathbb{F}_0$ follow from the i.i.d. structure of the observations Y_i . We briefly comment on condition (ED) . The use once again the i.i.d. structure yields by (5.1) in view of $V^2 = n\mathbf{v}^2$

$$\log \mathbb{E} \exp \left\{ \lambda \frac{\boldsymbol{\gamma}^\top \nabla \zeta(\boldsymbol{\theta})}{\|V\boldsymbol{\gamma}\|} \right\} = n \mathbb{E} \exp \left\{ \frac{\lambda}{n^{1/2}} \frac{\boldsymbol{\gamma}^\top \nabla \zeta_1(\boldsymbol{\theta})}{\|\mathbf{v}\boldsymbol{\gamma}\|} \right\} \leq \nu_0^2 \lambda^2 / 2$$

as long as $\lambda \leq n^{1/2} \mathbf{g}_1 \leq \mathbf{g}$. Similarly one can check (ED_0) and $(ED)_1$. \square

Below we specify the general results of Sections 3 and 4 to the i.i.d. setup.

5.1.1 A large deviation bound

First we describe the large deviation probability for the event $\{\tilde{\boldsymbol{\theta}} \notin \Theta_0(\mathbf{r})\}$ for a fixed \mathbf{r} . Corollary 4.3 provides a sufficient condition for such a bound: the rate function $\mathfrak{M}^*(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ should grow at least logarithmic with the distance $\|V(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|$. Define

$$\mathfrak{m}^*(\boldsymbol{\theta}, \boldsymbol{\theta}^*) = \max_{\mu} \mathfrak{m}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*).$$

(ld_1) There exist constants $b > 0$, such that it holds on the set $\Theta_0^c(\mathbf{r})$

$$\mathfrak{m}^*(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \geq b \log(1 + \|\mathbf{v}_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2).$$

This condition implies by $\mathfrak{M}^*(\boldsymbol{\theta}, \boldsymbol{\theta}^*) = n\mathfrak{m}^*(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$

$$\mathfrak{M}^*(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \geq nb \log(1 + \|\mathbf{v}_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2), \quad \boldsymbol{\theta} \notin \Theta_0(\mathbf{r}).$$

For the next result we assume that the value $\mathbf{r}_1 \geq 1$ is fixed by $\det(\mathbf{r}_1^{-1} \mathbf{v}) \leq \det(\mathbf{v}_0)$. Further, the constant s has to be fixed such that $3\nu_0 \mathbf{r}_1 / s \leq \mathbf{g} \wedge \sqrt{2\mathbf{c}_1 p}$.

Theorem 5.2. *Suppose (e) , (ed) , (ld_1) . For each \mathbf{x} and $\mathbf{r} = \mathbf{r}(\mathbf{x})$ such that*

$$(2/3)b\mathbf{r}^2 \geq (1+s)(p/2+1)\log(1+\mathbf{r}^2) + \mathfrak{z}_1(\mathbf{x}), \quad (5.2)$$

it holds for $n \geq \mathbf{r}^2$

$$\mathbb{P}(\sqrt{n}\|\mathbf{v}_0(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\| \geq \mathbf{r}) \leq e^{-\mathbf{x}+1}.$$

Proof. Given \mathbf{r} and \mathbf{x} , (5.2) implies for all $u \geq \mathbf{r}$

$$nb \log(1 + u^2/n) \geq (1+s)(p/2+1)\log(1 + u^2) + \mathfrak{z}_1(\mathbf{x}).$$

because $n \log(1 + u^2/n) \geq 2\mathbf{r}^2/3$ for $\mathbf{r}^2 \leq n$. Now it follows from (ld_1) for any $\boldsymbol{\theta} \notin \Theta_0(\mathbf{r})$ that

$$\begin{aligned} \mathfrak{M}^*(\boldsymbol{\theta}, \boldsymbol{\theta}^*) &\geq nb \log(1 + \|\mathbf{v}_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2) \\ &\geq (1+s)(p/2+1)\log(1 + \|V_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2) + \mathfrak{z}_1(\mathbf{x}). \end{aligned}$$

Now the result follows from Corollary 4.3.

(to be done) treatment of δ_μ . □

Remark 5.1. The presented result helps to quantify two important values \mathbf{r} and n providing a sensitive deviation probability bound: the radius \mathbf{r} of the local neighborhood should be large enough to ensure (5.2), while the sample size n should be larger than \mathbf{r}^2 . It is straightforward to see that (5.2) starts to hold for $\mathbf{r}^2 \geq \text{Const. } p \log p$ for some fixed constant Const. Therefore, for any \mathbf{r} exceeding this value, a deviation probability $\mathbb{P}(\tilde{\boldsymbol{\theta}} \notin \Theta_0(\mathbf{r}))$ is negligible when $n \geq \mathbf{r}^2 \geq \text{Const. } p \log p$.

5.1.2 Local inference

Now we restate the general local bounds of Section 3 for the i.i.d. case. First we describe the approximating linear models. The matrices \mathbf{v}_0^2 and \mathbb{F}_0 from conditions (ed_0) , (ed_1) , and (ℓ_0) determine their drift and variance components. Define

$$\mathbb{F}_\epsilon \stackrel{\text{def}}{=} \mathbb{F}_0(1 - \delta) - \varrho \mathbf{v}_0^2.$$

Then $D_\epsilon^2 = n\mathbb{F}_\epsilon$ and

$$\boldsymbol{\xi}_\epsilon \stackrel{\text{def}}{=} D_\epsilon^{-1} \nabla \zeta(\boldsymbol{\theta}^*) = (n\mathbb{F}_\epsilon)^{-1/2} \sum \nabla \ell(Y_i, \boldsymbol{\theta}^*). \quad (5.3)$$

The upper approximating process reads as

$$\mathbb{L}_\epsilon(\boldsymbol{\theta}, \boldsymbol{\theta}^*) = (\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top D_\epsilon \boldsymbol{\xi}_\epsilon - \|D_\epsilon(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2/2.$$

This expression appears as log-likelihood for the linear model $\boldsymbol{\xi}_\epsilon = D_\epsilon \boldsymbol{\theta} + \boldsymbol{\varepsilon}$ for a standard normal error $\boldsymbol{\varepsilon}$. The (quasi) MLE $\tilde{\boldsymbol{\theta}}_\epsilon$ for this model is of the form $\tilde{\boldsymbol{\theta}}_\epsilon = D_\epsilon^{-1} \boldsymbol{\xi}_\epsilon$.

Theorem 5.3. *Suppose (ed_0) . Given \mathbf{r} , assume (ed_1) , and (ℓ_0) on $\Theta_0(\mathbf{r})$, and let $\varrho = 3\nu_0 \omega^* \mathbf{r}/n^{1/2}$, $\delta = \delta(\mathbf{r}) = \delta^* \mathbf{r}/n^{1/2}$, and $\mathbb{F}_\epsilon \stackrel{\text{def}}{=} \mathbb{F}_0(1 - \delta) - \varrho \mathbf{v}_0^2 \geq 0$. Then the results of Theorem 3.1 through 3.8 apply to the case of i.i.d. modeling. In particular, for any $z > 0$, it holds*

$$\begin{aligned} & \mathbb{P}\{\|\sqrt{n\mathbb{F}_\epsilon}(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\| > z, \tilde{\boldsymbol{\theta}} \in \Theta_0(\mathbf{r})\} \\ & \leq \mathbb{P}\left\{(1 - \sqrt{\alpha_\epsilon})\|\boldsymbol{\xi}_\epsilon\| > z - \sqrt{2\Diamond_\epsilon(\mathbf{r}) + 2\Diamond_{\underline{\epsilon}}(\mathbf{r})}\right\}, \end{aligned}$$

where $\Diamond_\epsilon(\mathbf{r}), \Diamond_{\underline{\epsilon}}(\mathbf{r})$ follow the bound (3.4) and α_ϵ is defined by

$$\alpha_\epsilon \stackrel{\text{def}}{=} \|I_p - \mathbb{F}_\epsilon^{1/2} \mathbb{F}_{\underline{\epsilon}}^{-1} \mathbb{F}_\epsilon^{1/2}\|_\infty = \lambda_{\max}(I_p - \mathbb{F}_\epsilon^{1/2} \mathbb{F}_{\underline{\epsilon}}^{-1} \mathbb{F}_\epsilon^{1/2}).$$

Moreover, on the random set $\mathcal{C}_\epsilon(\mathbf{r}) = \{\tilde{\boldsymbol{\theta}} \in \Theta_0(\mathbf{r}), \|\boldsymbol{\xi}_{\underline{\epsilon}}\| \leq \mathbf{r}\}$, it holds

$$\|\sqrt{n\mathbb{F}_\epsilon}(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) - \boldsymbol{\xi}_\epsilon\|^2 \leq 2\Delta_\epsilon(\mathbf{r}).$$

The presented results are stated via the probability bound for the squared norm of the vector $\boldsymbol{\xi}_\epsilon$ from (5.3). One can apply the general results of Section 6.4. For ease of notation we consider the vector $\boldsymbol{\xi}$ instead of $\boldsymbol{\xi}_\epsilon$:

$$\boldsymbol{\xi} \stackrel{\text{def}}{=} (n\mathbb{F}_0)^{-1/2} \nabla \zeta(\boldsymbol{\theta}^*) = (n\mathbb{F}_0)^{-1/2} \sum \nabla \ell(Y_i, \boldsymbol{\theta}^*). \quad (5.4)$$

The necessary condition (6.18) coincides with (ED_0) . Now Corollary 6.13 implies that the probability $\mathbb{P}(\|\boldsymbol{\xi}\|^2 > \mathfrak{z}(\mathbf{x}, \mathbb{B}))$ is of order $2e^{-\mathbf{x}}$ for all moderate \mathbf{x} and $\mathbb{B} = \mathbb{F}_0^{-1/2} \mathbf{v}_0^2 \mathbb{F}_0^{-1/2}$. In particular, this probability starts to degenerate when \mathfrak{z} significantly exceeds the value $\mathbf{p} = \text{tr}(\mathbb{B})$.

5.1.3 The regular parametric case and asymptotic efficiency

The conditions and the results become even more transparent in the regular situation when $\mathbb{F}_0 \geq \mathfrak{a}^2 \mathbf{v}_0^2$. An important special case corresponds to the correct parametric specification with $\mathfrak{a} = 1$ and $\mathbb{F}_0 = \mathbf{v}_0^2$. Under regularity one can use $\mathbb{F}_\epsilon = (1 - \delta - \varrho \mathfrak{a}^2) \mathbb{F}_0$. Here we briefly discuss the corollaries of Theorem 5.3 for the classical asymptotic setup when n tends to infinity.

Under the imposed conditions, the quantities δ and ϱ can be taken of order $\mathbf{r}/n^{1/2}$, where $\mathbf{r}^2 \gg p \log(p)$. If n grows then δ and ϱ decreases to zero, and the matrix \mathbb{F}_ϵ is close to \mathbb{F}_0 . The value $\Delta_\epsilon(\mathbf{r})$ is close to zero in probability. Moreover, the random vector $\boldsymbol{\xi}$ from (5.4) fulfills $\text{Var}(\boldsymbol{\xi}) \leq \mathbb{F}_0^{-1/2} \mathbf{v}_0^2 \mathbb{F}_0^{-1/2} \stackrel{\text{def}}{=} \mathbb{B}^2$ and by the central limit theorem $\boldsymbol{\xi}$

is asymptotically normal $\mathcal{N}(0, \mathcal{B}^2)$. This yields by Theorem 5.3 that $\sqrt{n\mathcal{F}_0}(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)$ is asymptotically normal $\mathcal{N}(0, \mathcal{B}^2)$ as well. The correct model specification implies $\mathcal{B} \equiv \mathcal{I}_p$ and hence $\tilde{\boldsymbol{\theta}}$ is asymptotically efficient; see Ibragimov and Khas'minskij (1981). Also $2L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) \approx \|\boldsymbol{\xi}\|^2$ which is nearly χ^2 r.v. with p degrees of freedom. This result is known as asymptotic Wilks theorem.

5.2 Generalized linear modeling

Now we consider a generalized linear modeling (GLM) which is often used for describing some categorical data. Let $\mathcal{P} = (P_w, w \in \mathcal{Y})$ be an exponential family with a canonical parametrization; see e.g. McCullagh and Nelder (1989). The corresponding log-density can be represented as $\ell(y, w) = yw - d(w)$ for a convex function $d(w)$. The popular examples are given by the binomial (binary response, logistic) model with $d(w) = \log(e^w + 1)$, the Poisson model with $d(w) = e^w$, the exponential model with $d(w) = -\log(w)$. Note that linear Gaussian regression is a special case with $d(w) = w^2/2$.

A GLM specification means that every observation Y_i has a distribution from the family \mathcal{P} with the parameter w_i which linearly depends on the regressor $\Psi_i \in \mathbb{R}^p$:

$$Y_i \sim P_{\Psi_i^\top \boldsymbol{\theta}^*}. \quad (5.5)$$

The corresponding log-density of a GLM reads as

$$L(\boldsymbol{\theta}) = \sum \{Y_i \Psi_i^\top \boldsymbol{\theta} - d(\Psi_i^\top \boldsymbol{\theta})\}.$$

First we specify the data distribution allowing that the parametric model (5.5) is misspecified. Misspecification of the first kind means that the vector $\mathbf{f} \stackrel{\text{def}}{=} \mathbb{E}\mathbf{Y}$ cannot be represented in the form $\boldsymbol{\Psi}^\top \boldsymbol{\theta}$ whatever $\boldsymbol{\theta}$ is. In this situation, the target of estimation $\boldsymbol{\theta}^*$ is defined by

$$\boldsymbol{\theta}^* \stackrel{\text{def}}{=} \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \mathbb{E}L(\boldsymbol{\theta}).$$

The other sort of misspecification concerns the data distribution. The model (5.5) assumes that the Y_i 's are independent and the marginal distribution belongs to the given parametric family \mathcal{P} . In what follows, we only assume independent data having certain exponential moments. The quasi MLE $\tilde{\boldsymbol{\theta}}$ is defined by maximization of $L(\boldsymbol{\theta})$:

$$\tilde{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} L(\boldsymbol{\theta}) = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \sum \{Y_i \Psi_i^\top \boldsymbol{\theta} - d(\Psi_i^\top \boldsymbol{\theta})\}.$$

Convexity of $d(\cdot)$ implies that $L(\boldsymbol{\theta})$ is a concave function of $\boldsymbol{\theta}$, so that the optimization problem has a unique solution and can be effectively solved. However, a closed form

solution is only available for the constant regression or for the linear Gaussian regression. The corresponding target $\boldsymbol{\theta}^*$ is the maximizer of the expected log-likelihood:

$$\boldsymbol{\theta}^* = \operatorname{argmax}_{\boldsymbol{\theta}} \mathbb{E}L(\boldsymbol{\theta}) = \operatorname{argmax}_{\boldsymbol{\theta}} \sum \{f_i \Psi_i^\top \boldsymbol{\theta} - d(\Psi_i^\top \boldsymbol{\theta})\}$$

with $f_i = \mathbb{E}Y_i$. The function $\mathbb{E}L(\boldsymbol{\theta})$ is concave as well and the vector $\boldsymbol{\theta}^*$ is also well defined.

Define the individual errors (residuals) $\varepsilon_i = Y_i - \mathbb{E}Y_i$. Below we assume that these errors fulfill some exponential moment conditions.

(e₁) *There exist some constants ν_0 and $\mathbf{g}_1 > 0$, and for every i a constant \mathbf{n}_i such that $\mathbb{E}(\varepsilon_i/\mathbf{n}_i)^2 \leq 1$ and for all $|\lambda| \leq \mathbf{g}_1$*

$$\log \mathbb{E} \exp(\lambda \varepsilon_i / \mathbf{n}_i) \leq \nu_0^2 \lambda^2 / 2, \quad |\lambda| \leq \mathbf{g}_1. \quad (5.6)$$

A natural candidate for \mathbf{n}_i is σ_i where $\sigma_i^2 = \mathbb{E}\varepsilon_i^2$ is the variance of ε_i ; see Lemma 7.12. Under (5.6), introduce a $p \times p$ matrix V_0 defined by

$$V_0^2 \stackrel{\text{def}}{=} \sum \mathbf{n}_i^2 \Psi_i \Psi_i^\top. \quad (5.7)$$

Condition (e₁) effectively means that each error term $\varepsilon_i = Y_i - \mathbb{E}Y_i$ has some bounded exponential moments: for $\lambda = \mathbf{g}_1$, it holds $f(\lambda) \stackrel{\text{def}}{=} \log \mathbb{E} \exp(\lambda \varepsilon_i / \mathbf{n}_i) < \infty$. This implies the quadratic upper bound for the function $f(\lambda)$ for $|\lambda| \leq \mathbf{g}_1$; see Lemma 7.12. In words, condition (e₁) requires light (exponentially decreasing) tail for the marginal distribution of each ε_i .

Define also

$$N^{-1/2} \stackrel{\text{def}}{=} \max_i \sup_{\boldsymbol{\gamma} \in \mathbb{R}^p} \frac{\mathbf{n}_i |\Psi_i^\top \boldsymbol{\gamma}|}{\|V_0 \boldsymbol{\gamma}\|}. \quad (5.8)$$

Lemma 5.4. *Assume (e₁) and let V_0 be defined by (5.7) and N by (5.8). Then conditions (ED₀) and (ED) follow with $V = V_0$, $\mathbf{g} = \mathbf{g}_1 N^{1/2}$, and with the constant ν_0 from (e₁). Moreover, the stochastic component $\zeta(\boldsymbol{\theta})$ is linear in $\boldsymbol{\theta}$ and the condition (ED₁) is fulfilled with $\omega(\mathbf{r}) \equiv 0$.*

Proof. The gradient of the stochastic component $\zeta(\boldsymbol{\theta})$ of $L(\boldsymbol{\theta})$ does not depend on $\boldsymbol{\theta}$:

$$\nabla \zeta(\boldsymbol{\theta}) = \sum \Psi_i \varepsilon_i$$

with $\varepsilon_i = Y_i - \mathbb{E}Y_i$. Now, for any unit vector $\boldsymbol{\gamma} \in \mathbb{R}^p$ and $\lambda \leq \mathbf{g}$, independence of the ε_i 's implies that

$$\log \mathbb{E} \exp \left\{ \frac{\lambda}{\|V_0 \boldsymbol{\gamma}\|} \boldsymbol{\gamma}^\top \sum \Psi_i \varepsilon_i \right\} = \sum \log \mathbb{E} \exp \left\{ \frac{\lambda \mathbf{n}_i \Psi_i^\top \boldsymbol{\gamma}}{\|V_0 \boldsymbol{\gamma}\|} \varepsilon_i / \mathbf{n}_i \right\}. \quad (5.9)$$

By definition $\mathfrak{n}_i |\Psi_i^\top \gamma| / \|V_0 \gamma\| \leq N^{-1/2}$ and therefore, $\lambda \mathfrak{n}_i |\Psi_i^\top \gamma| / \|V_0 \gamma\| \leq \mathfrak{g}_1$. Hence, (5.6) implies

$$\log \mathbb{E} \exp \left\{ \frac{\lambda}{\|V_0 \gamma\|} \gamma^\top \sum \Psi_i \varepsilon_i \right\} \leq \frac{\nu_0^2 \lambda^2}{2 \|V_0 \gamma\|^2} \sum \mathfrak{n}_i^2 |\Psi_i^\top \gamma|^2 = \frac{\nu_0^2 \lambda^2}{2}, \quad (5.10)$$

and (ED_0) follows. \square

It remains only to bound the quality of quadratic approximation for the mean of the process $L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ in a vicinity of $\boldsymbol{\theta}^*$. An interesting feature of the GLM is that the effect of model misspecification disappears in the expectation of $L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$.

Lemma 5.5. *It holds*

$$\begin{aligned} -\mathbb{E} L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) &= \sum \{d(\Psi_i^\top \boldsymbol{\theta}) - d(\Psi_i^\top \boldsymbol{\theta}^*) - d'(\Psi_i^\top \boldsymbol{\theta}^*) \Psi_i^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^*)\} \\ &= \mathcal{K}(\mathbb{P}_{\boldsymbol{\theta}^*}, \mathbb{P}_{\boldsymbol{\theta}}), \end{aligned} \quad (5.11)$$

where $\mathcal{K}(\mathbb{P}_{\boldsymbol{\theta}^*}, \mathbb{P}_{\boldsymbol{\theta}})$ is the Kullback-Leibler divergence between measures $\mathbb{P}_{\boldsymbol{\theta}^*}$ and $\mathbb{P}_{\boldsymbol{\theta}}$. Moreover,

$$-\mathbb{E} L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) = \frac{1}{2} \|D_{\boldsymbol{\theta}^\circ}(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2, \quad (5.12)$$

where $\boldsymbol{\theta}^\circ \in [\boldsymbol{\theta}^*, \boldsymbol{\theta}]$ and

$$D^2(\boldsymbol{\theta}^\circ) = \sum d''(\Psi_i^\top \boldsymbol{\theta}^\circ) \Psi_i \Psi_i^\top.$$

Proof. The definition implies

$$\mathbb{E} L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) = \sum \{f_i \Psi_i^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^*) - d(\Psi_i^\top \boldsymbol{\theta}) + d(\Psi_i^\top \boldsymbol{\theta}^*)\}.$$

As $\boldsymbol{\theta}^*$ is the extreme point of $\mathbb{E} L(\boldsymbol{\theta})$, it holds $\nabla \mathbb{E} L(\boldsymbol{\theta}^*) = \sum [f_i - d'(\Psi_i^\top \boldsymbol{\theta}^*)] \Psi_i = 0$ and (5.11) follows. The Taylor expansion of the second order around $\boldsymbol{\theta}^*$ yields the expansion (5.12). \square

Define now the matrix D_0 by

$$D_0^2 \stackrel{\text{def}}{=} D^2(\boldsymbol{\theta}^*) = \sum d''(\Psi_i^\top \boldsymbol{\theta}^*) \Psi_i \Psi_i^\top.$$

Let also V_0 be defined by (5.7). Note that the matrices D_0 and V_0 coincide if the model $Y_i \sim P_{\Psi_i^\top \boldsymbol{\theta}^*}$ is correctly specified and $\mathfrak{n}_i^2 = d''(\Psi_i^\top \boldsymbol{\theta}^*)$. The matrix V_0 describes a local elliptic neighborhood of the central point $\boldsymbol{\theta}^*$ in the form $\Theta_0(\mathbf{r}) = \{\boldsymbol{\theta} : \|V_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\| \leq \mathbf{r}\}$. If the matrix function $D^2(\boldsymbol{\theta})$ is continuous in this vicinity $\Theta_0(\mathbf{r})$ then the value $\delta(\mathbf{r})$ measuring the approximation quality of $-\mathbb{E} L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ by the quadratic function $\|D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2/2$ is small and the identifiability condition (\mathcal{L}_0) is fulfilled on $\Theta_0(\mathbf{r})$.

Lemma 5.6. *Suppose that*

$$\|I_p - D_0^{-1}D^2(\boldsymbol{\theta})D_0^{-1}\|_\infty \leq \delta(\mathbf{r}), \quad \boldsymbol{\theta} \in \Theta_0(\mathbf{r}). \quad (5.13)$$

Then (\mathcal{L}_0) holds with this $\delta(\mathbf{r})$. Moreover, as the quantities $\omega(\mathbf{r}), \diamond_\epsilon(\mathbf{r}), \diamond_{\underline{\epsilon}}(\mathbf{r})$ vanish, one can take $\varrho = 0$ leading to the following representation for D_ϵ and $\boldsymbol{\xi}_\epsilon$:

$$\begin{aligned} D_\epsilon^2 &= (1 - \delta)D_0^2, & \boldsymbol{\xi}_\epsilon &= (1 + \delta)^{1/2}\boldsymbol{\xi} \\ D_{\underline{\epsilon}}^2 &= (1 + \delta)D_0^2, & \boldsymbol{\xi}_{\underline{\epsilon}} &= (1 - \delta)^{1/2}\boldsymbol{\xi} \end{aligned}$$

with

$$\boldsymbol{\xi} \stackrel{\text{def}}{=} D_0^{-1}\nabla\zeta = D_0^{-1}\sum \Psi_i(Y_i - \mathbb{E}Y_i).$$

Now we are prepared to state the local results for the GLM estimation.

Theorem 5.7. *Let (e_1) hold. Then for $\epsilon = (\delta, 0)$ with $\delta \geq \delta(\mathbf{r})$ and any $\boldsymbol{\theta} \in \Theta_0(\mathbf{r})$*

$$\mathbb{L}_{\underline{\epsilon}}(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \leq L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \leq \mathbb{L}_\epsilon(\boldsymbol{\theta}, \boldsymbol{\theta}^*). \quad (5.14)$$

Moreover, for any $z > 0$ and $\mathfrak{z} > 0$, it holds

$$\begin{aligned} \mathbb{P}(\|D_0(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\| > z, \|V_0(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\| \leq \mathbf{r}) &\leq \mathbb{P}\{\|\boldsymbol{\xi}\|^2 > z^2[1 - \delta(\mathbf{r})]\} \\ \mathbb{P}(L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) > \mathfrak{z}, \|V_0(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\| \leq \mathbf{r}) &\leq \mathbb{P}\{\|\boldsymbol{\xi}\|^2/2 > \mathfrak{z}[1 - \delta(\mathbf{r})]\}. \end{aligned}$$

Linearity of the stochastic component $\zeta(\boldsymbol{\theta})$ in the considered GLM implies important fact that the quantities $\diamond_\epsilon(\mathbf{r}), \diamond_{\underline{\epsilon}}(\mathbf{r})$ in the majorization bound (5.14) vanish for any \mathbf{r} . However, the deterministic component is not quadratic in $\boldsymbol{\theta}$ unless the function $d(w)$ is quadratic. Therefore, the presented bounds are local and have to be accomplished with the large deviation bounds.

An interesting question, similarly to the i.i.d. case, is the minimal radius \mathbf{r} of the local vicinity $\Theta_0(\mathbf{r})$ ensuring the desirable concentration property. We apply the sufficient condition (4.9) of Corollary 4.3 ensuring the concentration property. By Lemma 5.5, the function $\mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*)$ can be decomposed as

$$\mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*) = \mu\mathcal{K}(\mathbb{P}_{\boldsymbol{\theta}^*}, \mathbb{P}_{\boldsymbol{\theta}}) - \mathfrak{N}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*).$$

Let $\mu = \mu(\boldsymbol{\theta})$ be selected such that $\mu|\Psi_i^\top(\boldsymbol{\theta} - \boldsymbol{\theta}^*)|n_i \leq g_1$ for all i . Then the arguments from (5.9) and (5.10) yield

$$\begin{aligned} \mathfrak{N}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*) &\leq \sum \log \mathbb{E} \exp\{\mu n_i \Psi_i^\top(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\varepsilon_i/n_i\} \\ &\leq \frac{\nu_0^2 \mu^2}{2} \sum |n_i \Psi_i^\top(\boldsymbol{\theta} - \boldsymbol{\theta}^*)|^2 = \frac{\nu_0^2 \mu^2}{2} \|V_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2. \end{aligned}$$

Therefore,

$$\mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*) \geq \mu \mathcal{K}(\mathbb{P}_{\boldsymbol{\theta}^*}, \mathbb{P}_{\boldsymbol{\theta}}) - \frac{\nu_0^2 \mu^2}{2} \|V_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2$$

and a reasonable choice of $\mu = \mu(\boldsymbol{\theta})$ is given by $\mu(\boldsymbol{\theta}) = \nu_0^{-2} \mathcal{K}(\mathbb{P}_{\boldsymbol{\theta}^*}, \mathbb{P}_{\boldsymbol{\theta}}) / \|V_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2$ leading to

$$\mathfrak{M}^*(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \geq \frac{\nu_0^{-2} \mathcal{K}^2(\mathbb{P}_{\boldsymbol{\theta}^*}, \mathbb{P}_{\boldsymbol{\theta}})}{2 \|V_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2}.$$

So, the Kullback-Leibler divergence $\mathcal{K}(\mathbb{P}_{\boldsymbol{\theta}^*}, \mathbb{P}_{\boldsymbol{\theta}})$ should of order at least

$$\|V_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\| \sqrt{\|p \log(V_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|)}.$$

One can check that this condition is fulfilled if the convex function $d(\cdot)$ satisfies $d'(t) \geq \text{Const.}/t$ for some $\text{Const} > 0$ and if the effective sample size N from (5.8) is sufficiently large.

5.3 Linear median estimation

This section illustrates how the proposed approach applies to robust estimation in linear models. The target of analysis is the linear dependence of the observed data $\mathbf{Y} = (Y_1, \dots, Y_n)$ on the set of features $\boldsymbol{\Psi}_i \in \mathbb{R}^p$:

$$Y_i = \boldsymbol{\Psi}_i^\top \boldsymbol{\theta} + \varepsilon_i \tag{5.15}$$

where ε_i denotes the i th individual error.

The study of the qMLE in the GLM (5.5) heavily relies on the assumption (e_1) . If this assumption is not verified, then the proposed approach would not apply. An explicit structure of the qMLE, especially in the linear regression case, allows for direct study of the properties of $\tilde{\boldsymbol{\theta}}$ under weaker moment assumptions than (e_1) . However, we aim at establishing some exponential bounds and therefore, the condition of bounded exponential moments for each observation is really necessary within the least squares or generalized linear approach. In the case of heavily tailed data with only polynomial moments, one can obtain some convergence results for the LSE $\tilde{\boldsymbol{\theta}}$, however an exponential bound is not available. In such cases, it is natural to use a robustified version of the contrast, e.g. the least absolute deviation (LAD) method. We consider the linear model (5.15) and suppose for a moment that the errors ε_i are i.i.d. and follow the double exponential (Laplace) distribution with the density $(1/2)e^{-|y|}$. Then the model (5.15) yields the log-likelihood

$$L(\boldsymbol{\theta}) = -\frac{1}{2} \sum |Y_i - \boldsymbol{\Psi}_i^\top \boldsymbol{\theta}|$$

and $\tilde{\boldsymbol{\theta}} \stackrel{\text{def}}{=} \operatorname{argmax}_{\boldsymbol{\theta}} L(\boldsymbol{\theta})$ is called the *least absolute deviation* (LAD) estimate. In the context of linear regression, it is also called the *linear median* estimate. The target of estimation $\boldsymbol{\theta}^*$ is defined as usually by the equation $\boldsymbol{\theta}^* = \operatorname{argmax}_{\boldsymbol{\theta}} \mathbb{E}L(\boldsymbol{\theta})$.

It is useful to define the residuals $\tilde{\varepsilon}_i = Y_i - \Psi_i^\top \boldsymbol{\theta}^*$ and their distributions

$$P_i(A) = \mathbb{P}(\tilde{\varepsilon}_i \in A) = \mathbb{P}(Y_i - \Psi_i^\top \boldsymbol{\theta}^* \in A)$$

for any Borel set A on the real line. If $Y_i = \Psi_i^\top \boldsymbol{\theta}^* + \varepsilon_i$ is the true model then P_i coincides with the distribution of each ε_i . Below we suppose that each $P_i = \mathcal{L}(Y_i - \Psi_i^\top \boldsymbol{\theta}^*)$ has a positive density $p_i(y)$.

Note that the difference $L(\boldsymbol{\theta}) - L(\boldsymbol{\theta}^*)$ is bounded by $\frac{1}{2} \sum |\Psi_i^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^*)|$ and condition (E) is fulfilled automatically. Next we check conditions (ED₀) and (ED₁). Denote $\xi_i(\boldsymbol{\theta}) = \mathbb{I}(Y_i - \Psi_i^\top \boldsymbol{\theta} \leq 0) - q_i(\boldsymbol{\theta})$ for $q_i(\boldsymbol{\theta}) = \mathbb{P}(Y_i - \Psi_i^\top \boldsymbol{\theta} \leq 0)$. This is a centered Bernoulli random variable, and it is easy to check that

$$\nabla \zeta(\boldsymbol{\theta}) = \sum \xi_i(\boldsymbol{\theta}) \Psi_i. \quad (5.16)$$

This expression differs from the similar ones from the linear and generalized linear regression because the error terms ξ_i now depends on $\boldsymbol{\theta}$. First we check the global condition (ED). Fix any $\mathbf{g}_1 < 1$. Then it holds for a Bernoulli r.v. Z with $\mathbb{P}(Z = 1) = q$, $\xi = Z - q$, and $|\lambda| \leq \mathbf{g}_1$

$$\begin{aligned} \log \mathbb{E} \exp(\lambda \xi) &= \log [q \exp\{\lambda(1 - q)\} + (1 - q) \exp(-\lambda q)] \\ &\leq \nu_0^2 q(1 - q) \lambda^2 / 2, \end{aligned} \quad (5.17)$$

where $\nu_0 \geq 1$ depends on \mathbf{g}_1 only. Let now a vector $\boldsymbol{\gamma} \in \mathbb{R}^p$ and $\rho > 0$ be such that $\rho |\Psi_i^\top \boldsymbol{\gamma}| \leq \mathbf{g}_1$ for all $i = 1, \dots, n$. Then

$$\begin{aligned} \log \mathbb{E} \exp\{\rho \boldsymbol{\gamma}^\top \nabla \zeta(\boldsymbol{\theta})\} &\leq \frac{\nu_0^2 \rho^2}{2} \sum q_i(\boldsymbol{\theta}) \{1 - q_i(\boldsymbol{\theta})\} |\Psi_i^\top \boldsymbol{\gamma}|^2 \\ &\leq \nu_0^2 \rho^2 \|V(\boldsymbol{\theta}) \boldsymbol{\gamma}\|^2 / 2, \end{aligned} \quad (5.18)$$

where

$$V^2(\boldsymbol{\theta}) = \sum q_i(\boldsymbol{\theta}) \{1 - q_i(\boldsymbol{\theta})\} \Psi_i \Psi_i^\top. \quad (5.19)$$

Denote also

$$V^2 = \frac{1}{4} \sum \Psi_i \Psi_i^\top.$$

Clearly $V(\boldsymbol{\theta}) \leq V$ for all $\boldsymbol{\theta}$ and condition (ED) is fulfilled globally with the matrix V and $\mathbf{g} = \mathbf{g}_1 N^{1/2}$ for N defined by

$$N^{-1/2} \stackrel{\text{def}}{=} \max_i \sup_{\boldsymbol{\gamma} \in \mathbb{R}^p} \frac{\boldsymbol{\Psi}_i^\top \boldsymbol{\gamma}}{2\|V\boldsymbol{\gamma}\|}; \quad (5.20)$$

cf. (5.9).

5.3.1 A local central bound

Now we restrict ourselves to the elliptic vicinity $\Theta_0(\mathbf{r}) = \{\boldsymbol{\theta} : \|V_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\| \leq \mathbf{r}\}$ of the central point $\boldsymbol{\theta}^*$ for $V_0 = V(\boldsymbol{\theta}^*)$ and some $\mathbf{r} > 0$. Define the matrix $V_0 = V(\boldsymbol{\theta}^*)$. Then condition (ED_0) with the matrix V_0 and $\mathbf{g} = N^{1/2}\mathbf{g}_1$ is fulfilled on $\Theta_0(\mathbf{r})$ due to (5.18). Next, for checking (ED_1) suppose the following regularity condition:

$$V \leq \nu_1 V_0 \quad (5.21)$$

for some $\nu_1 \geq 1$. This condition implies the inequality $|\boldsymbol{\Psi}_i^\top \boldsymbol{\gamma}| \leq \nu_1 N^{-1/2} \|V_0 \boldsymbol{\gamma}\|$ for any vector $\boldsymbol{\gamma} \in \mathbb{R}^p$. By (5.16)

$$\nabla \zeta(\boldsymbol{\theta}) - \nabla \zeta(\boldsymbol{\theta}^*) = \sum \boldsymbol{\Psi}_i \{\xi_i(\boldsymbol{\theta}) - \xi_i(\boldsymbol{\theta}^*)\}.$$

If $\boldsymbol{\Psi}_i^\top \boldsymbol{\theta} \geq \boldsymbol{\Psi}_i^\top \boldsymbol{\theta}^*$, then

$$\xi_i(\boldsymbol{\theta}) - \xi_i(\boldsymbol{\theta}^*) = \mathbb{I}(\boldsymbol{\Psi}_i^\top \boldsymbol{\theta}^* \leq Y_i < \boldsymbol{\Psi}_i^\top \boldsymbol{\theta}) - \mathbb{I}(\boldsymbol{\Psi}_i^\top \boldsymbol{\theta}^* \leq Y_i < \boldsymbol{\Psi}_i^\top \boldsymbol{\theta}).$$

Similarly for $\boldsymbol{\Psi}_i^\top \boldsymbol{\theta} < \boldsymbol{\Psi}_i^\top \boldsymbol{\theta}^*$

$$\xi_i(\boldsymbol{\theta}) - \xi_i(\boldsymbol{\theta}^*) = -\mathbb{I}(\boldsymbol{\Psi}_i^\top \boldsymbol{\theta} \leq Y_i < \boldsymbol{\Psi}_i^\top \boldsymbol{\theta}^*) + \mathbb{I}(\boldsymbol{\Psi}_i^\top \boldsymbol{\theta} \leq Y_i < \boldsymbol{\Psi}_i^\top \boldsymbol{\theta}^*).$$

Define $q_i(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \stackrel{\text{def}}{=} |q_i(\boldsymbol{\theta}) - q_i(\boldsymbol{\theta}^*)|$. Now (5.17) yields similarly to (5.18)

$$\begin{aligned} \log \mathbb{E} \exp\{\rho \boldsymbol{\gamma}^\top \{\nabla \zeta(\boldsymbol{\theta}) - \nabla \zeta(\boldsymbol{\theta}^*)\}\} &\leq \frac{\nu_0^2 \rho^2}{2} \sum q_i(\boldsymbol{\theta}, \boldsymbol{\theta}^*) |\boldsymbol{\Psi}_i^\top \boldsymbol{\gamma}|^2 \\ &\leq 2\nu_0^2 \rho^2 \max_{i \leq n} q_i(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \|V\boldsymbol{\gamma}\|^2 \leq \omega(\mathbf{r}) \nu_0^2 \rho^2 \|V_0 \boldsymbol{\gamma}\|^2 / 2, \end{aligned}$$

with

$$\omega(\mathbf{r}) \stackrel{\text{def}}{=} 4\nu_1 \max_{i \leq n} \sup_{\boldsymbol{\theta} \in \Theta_0(\mathbf{r})} q_i(\boldsymbol{\theta}, \boldsymbol{\theta}^*).$$

If each density function p_i is uniformly bounded by a constant C then

$$|q_i(\boldsymbol{\theta}) - q_i(\boldsymbol{\theta}^*)| \leq C |\boldsymbol{\Psi}_i^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^*)| \leq C \nu_1 N^{-1/2} \|V_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\| \leq C \nu_1 N^{-1/2} \mathbf{r}.$$

Next we check the identifiability condition. We use the following technical lemma.

Lemma 5.8. *It holds for any $\boldsymbol{\theta}$*

$$\frac{\partial^2}{\partial^2 \boldsymbol{\theta}} \mathbb{E}L(\boldsymbol{\theta}) = D^2(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \sum p_i(\boldsymbol{\Psi}_i^\top(\boldsymbol{\theta} - \boldsymbol{\theta}^*)) \boldsymbol{\Psi}_i \boldsymbol{\Psi}_i^\top, \quad (5.22)$$

where $p_i(\cdot)$ is the density of $\tilde{\varepsilon}_i = Y_i - \boldsymbol{\Psi}_i^\top \boldsymbol{\theta}^*$. Moreover, there is $\boldsymbol{\theta}^\circ \in [\boldsymbol{\theta}, \boldsymbol{\theta}^*]$ such that

$$\begin{aligned} -\mathbb{E}L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) &= \frac{1}{2} \sum |\boldsymbol{\Psi}_i^\top(\boldsymbol{\theta} - \boldsymbol{\theta}^*)|^2 p_i(\boldsymbol{\Psi}_i^\top(\boldsymbol{\theta}^\circ - \boldsymbol{\theta}^*)) \\ &= (\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top D^2(\boldsymbol{\theta}^\circ)(\boldsymbol{\theta} - \boldsymbol{\theta}^*)/2. \end{aligned} \quad (5.23)$$

Proof. Obviously

$$\frac{\partial \mathbb{E}L(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \sum \{ \mathbb{P}(Y_i \leq \boldsymbol{\Psi}_i^\top \boldsymbol{\theta}) - 1/2 \} \boldsymbol{\Psi}_i.$$

The identity (5.22) is obtained by one more differentiation. By definition, $\boldsymbol{\theta}^*$ is the extreme point of $\mathbb{E}L(\boldsymbol{\theta})$. The equality $\nabla \mathbb{E}L(\boldsymbol{\theta}^*) = 0$ yields

$$\sum \{ \mathbb{P}(Y_i \leq \boldsymbol{\Psi}_i^\top \boldsymbol{\theta}^*) - 1/2 \} \boldsymbol{\Psi}_i = 0.$$

Now (5.23) follows by the Taylor expansion of the second order at $\boldsymbol{\theta}^*$. \square

Define

$$D_0^2 \stackrel{\text{def}}{=} \sum |\boldsymbol{\Psi}_i^\top(\boldsymbol{\theta} - \boldsymbol{\theta}^*)|^2 p_i(0)$$

Due to this lemma, condition (\mathcal{L}_0) is fulfilled in $\Theta_0(\mathbf{r})$ with this choice D_0 for $\delta(\mathbf{r})$ from (5.13); see Lemma 5.6. Now all the local conditions are fulfilled yielding the general majorizing bound of Theorem 3.1 and all its corollaries.

Theorem 5.9. *Let $V_0 = V(\boldsymbol{\theta}^*)$; see (5.19). Assume (5.21) for $V^2 = (1/4) \sum \boldsymbol{\Psi}_i \boldsymbol{\Psi}_i^\top$. Fix any $\delta \geq \delta(\mathbf{r})$ and $\varrho \geq 3\nu_0\omega(\mathbf{r})$. Then Theorem 3.1 and its corollaries holds for the linear median estimation.*

This example is one more confirmation of the applicability of the the general approach: as soon as the local conditions have been checked the main local statements follow for free. It only remains to accomplish them by a large deviation bound, that is, to describe the local vicinity $\Theta_0(\mathbf{r})$ providing the prescribed concentration bound.

5.3.2 A large deviation bound

A sufficient condition for the concentration property is that the rate function $\mathfrak{M}^*(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ grows at least logarithmic in $\|\boldsymbol{\theta}^* - \boldsymbol{\theta}\|$. For $y > 0$, define for $A_y = [y, \infty]$

$$\lambda_i(y) = -(2y)^{-1} \log[P_i(A_y)] = -(2y)^{-1} \log[\mathbb{P}(Y_i - \boldsymbol{\Psi}_i^\top \boldsymbol{\theta}^* > y)].$$

The case with $\lambda_i(y) \geq \lambda_0 > 0$ corresponds to light tails while $\lambda_i(y) \rightarrow 0$ as $|y| \rightarrow \infty$ means heavy tails of the distribution P_i . Below we focus on the most interesting case when $\lambda_i(y)$ is positive and monotonously decreases to zero in $y > 0$. For simplicity of presentation we also assume that $\lambda_i(y)$ is sufficiently regular and its first derivative $\lambda'_i(y)$ is uniformly continuous on \mathbb{R} . Define

$$\begin{aligned}\xi_i(\boldsymbol{\theta}, \boldsymbol{\theta}^*) &\stackrel{\text{def}}{=} \{|Y_i - \boldsymbol{\Psi}_i^\top \boldsymbol{\theta}| - |Y_i - \boldsymbol{\Psi}_i^\top \boldsymbol{\theta}^*|\}/2 \\ f_i(\mu) &\stackrel{\text{def}}{=} \log \mathbb{E} \exp\{\mu \xi_i(\boldsymbol{\theta}, \boldsymbol{\theta}^*)\}.\end{aligned}$$

As $|\xi_i(\boldsymbol{\theta}, \boldsymbol{\theta}^*)| \leq |\boldsymbol{\Psi}_i^\top(\boldsymbol{\theta} - \boldsymbol{\theta}^*)|/2$ and $f_i(\mu)$ is analytic in μ , it holds for any $\mu \geq 0$ with $\mu \max_i |\boldsymbol{\Psi}_i^\top(\boldsymbol{\theta} - \boldsymbol{\theta}^*)|/2 \leq 1$

$$f_i(\mu) \geq -\mu \mathbb{E} \xi_i(\boldsymbol{\theta}, \boldsymbol{\theta}^*) - \mu^2 \text{Var} \xi_i(\boldsymbol{\theta}, \boldsymbol{\theta}^*).$$

This implies by independence of the ξ_i 's

$$\begin{aligned}\mathfrak{N}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}^*) &\stackrel{\text{def}}{=} -\log \mathbb{E} \exp\{\mu L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)\} \\ &= \sum f_i(\mu) \geq -\mu \mathbb{E} L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) - \mu^2 \text{Var} L(\boldsymbol{\theta}, \boldsymbol{\theta}^*).\end{aligned}$$

The choice $\mu = \mu(\boldsymbol{\theta}) = -\mathbb{E} L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) / \{2 \text{Var} L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)\}$ yields

$$\mathfrak{N}^*(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \geq \frac{|\mathbb{E} L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)|^2}{4 \text{Var} L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)}.$$

Note that the inequality $|\xi_i(\boldsymbol{\theta}, \boldsymbol{\theta}^*)| \leq |\boldsymbol{\Psi}_i^\top(\boldsymbol{\theta} - \boldsymbol{\theta}^*)|/2$ implies

$$\text{Var} L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \leq \frac{1}{2} \max_i |\boldsymbol{\Psi}_i^\top(\boldsymbol{\theta} - \boldsymbol{\theta}^*)| \sum \mathbb{E} |\xi_i(\boldsymbol{\theta}, \boldsymbol{\theta}^*)|.$$

This easily yields that the sufficient condition (4.9) of Corollary 4.3 is fulfilled in this situation if N from (5.20) is sufficiently large.

6 Deviation probability for quadratic forms

The approximation results of the previous sections rely on the probability of the form $\mathbb{P}(\|\boldsymbol{\xi}\| > y)$ for a given random vector $\boldsymbol{\xi} \in \mathbb{R}^p$. The only condition imposed on this vector is that

$$\log \mathbb{E} \exp(\boldsymbol{\gamma}^\top \boldsymbol{\xi}) \leq \nu_0^2 \|\boldsymbol{\gamma}\|^2/2, \quad \boldsymbol{\gamma} \in \mathbb{R}^p, \|\boldsymbol{\gamma}\| \leq \mathbf{g}.$$

To simplify the presentation we rewrite this condition as

$$\log \mathbb{E} \exp(\boldsymbol{\gamma}^\top \boldsymbol{\xi}) \leq \|\boldsymbol{\gamma}\|^2/2, \quad \boldsymbol{\gamma} \in \mathbb{R}^p, \|\boldsymbol{\gamma}\| \leq \mathbf{g}. \quad (6.1)$$

The general case can be reduced to $\nu_0 = 1$ by rescaling $\boldsymbol{\xi}$ and \mathbf{g} :

$$\log \mathbb{E} \exp(\boldsymbol{\gamma}^\top \boldsymbol{\xi} / \nu_0) \leq \|\boldsymbol{\gamma}\|^2 / 2, \quad \boldsymbol{\gamma} \in \mathbb{R}^p, \|\boldsymbol{\gamma}\| \leq \nu_0 \mathbf{g}$$

that is, $\nu_0^{-1} \boldsymbol{\xi}$ fulfills (6.1) with a slightly increased \mathbf{g} . In typical situations like in Section 5, the value \mathbf{g} is large (of order root- n) while the value ν_0 is close to one.

6.1 Gaussian case

Our benchmark will be a deviation bound for $\|\boldsymbol{\xi}\|^2$ for a standard Gaussian vector $\boldsymbol{\xi}$. The ultimate goal is to show that under (6.1) the norm of the vector $\boldsymbol{\xi}$ exhibits behavior expected for a Gaussian vector, at least in the region of moderate deviations. For the reason of comparison, we begin by stating the result for a Gaussian vector $\boldsymbol{\xi}$.

Theorem 6.1. *Let $\boldsymbol{\xi}$ be a standard normal vector in \mathbb{R}^p . Then for any $u > 0$, it holds*

$$\mathbb{P}(\|\boldsymbol{\xi}\|^2 > p + u) \leq \exp\{-(p/2)\phi(u/p)\}$$

with

$$\phi(t) \stackrel{\text{def}}{=} t - \log(1 + t).$$

Let $\phi^{-1}(\cdot)$ stand for the inverse of $\phi(\cdot)$. For any \mathbf{x} ,

$$\mathbb{P}(\|\boldsymbol{\xi}\|^2 > p + \phi^{-1}(2\mathbf{x}/p)) \leq \exp(-\mathbf{x}).$$

This particularly yields with $\varkappa = 6.6$

$$\mathbb{P}(\|\boldsymbol{\xi}\|^2 > p + \sqrt{\varkappa xp} \vee (\varkappa x)) \leq \exp(-x).$$

Proof. The proof utilizes the following well known fact: for $\mu < 1$

$$\log \mathbb{E} \exp(\mu \|\boldsymbol{\xi}\|^2 / 2) = -0.5p \log(1 - \mu).$$

It can be obtained by straightforward calculus. Now consider any $u > 0$. By the exponential Chebyshev inequality

$$\begin{aligned} \mathbb{P}(\|\boldsymbol{\xi}\|^2 > p + u) &\leq \exp\{-\mu(p + u)/2\} \mathbb{E} \exp(\mu \|\boldsymbol{\xi}\|^2 / 2) \\ &= \exp\{-\mu(p + u)/2 - (p/2) \log(1 - \mu)\}. \end{aligned} \tag{6.2}$$

It is easy to see that the value $\mu = u/(u + p)$ maximizes $\mu(p + u) + p \log(1 - \mu)$ w.r.t. μ yielding

$$\mu(p + u) - p \log(1 - \mu) = u - p \log(1 + u/p).$$

Further we use that $x - \log(1+x) \geq a_0 x^2$ for $x \leq 1$ and $x - \log(1+x) \geq a_0 x$ for $x > 1$ with $a_0 = 1 - \log(2) \geq 0.3$. This implies with $x = u/p$ for $u = \sqrt{\varkappa x p}$ or $u = \varkappa x$ and $\varkappa = 2/a_0 < 6.6$ that

$$\mathbb{P}(\|\xi\|^2 \geq p + \sqrt{\varkappa x p} \vee (\varkappa x)) \leq \exp(-x)$$

as required. \square

The message of this result is that the squared norm of the Gaussian vector ξ concentrates around the value p and the deviation over the level $p + \sqrt{x p}$ are exponentially small in x .

A similar bound can be obtained for a norm of the vector $B\xi$ where B is some given matrix. For notational simplicity we assume that B is symmetric. Otherwise one should replace it with $(B^\top A)^{1/2}$.

Theorem 6.2. *Let ξ be standard normal in \mathbb{R}^p . Then for every $x > 0$ and any symmetric matrix B , it holds with $p = \text{tr}(B^2)$, $v^2 = 2 \text{tr}(B^4)$, and $a^* = \|B^2\|_\infty$*

$$\mathbb{P}(\|B\xi\|^2 > p + (2vx^{1/2}) \vee (6a^*x)) \leq \exp(-x).$$

Proof. The matrix B^2 can be represented as $U^\top \text{diag}(a_1, \dots, a_p)U$ for an orthogonal matrix U . The vector $\tilde{\xi} = U\xi$ is also standard normal and $\|B\xi\|^2 = \tilde{\xi}^\top U B^2 U^\top \tilde{\xi}$. This means that one can reduce the situation to the case of a diagonal matrix $B^2 = \text{diag}(a_1, \dots, a_p)$. We can also assume without loss of generality that $a_1 \geq a_2 \geq \dots \geq a_p$. The expressions for the quantities p and v^2 simplifies to

$$p = \text{tr}(B^2) = a_1 + \dots + a_p,$$

$$v^2 = 2 \text{tr}(B^4) = 2(a_1^2 + \dots + a_p^2).$$

Moreover, rescaling the matrix B^2 by a_1 reduces the situation to the case with $a_1 = 1$.

Lemma 6.3. *It holds*

$$\mathbb{E}\|B\xi\|^2 = \text{tr}(B^2), \quad \text{Var}(\|B\xi\|^2) = 2 \text{tr}(B^4).$$

Moreover, for $\mu < 1$

$$\mathbb{E} \exp\{\mu\|B\xi\|^2/2\} = \det(1 - \mu B^2)^{-1/2} = \prod_{i=1}^p (1 - \mu a_i)^{-1/2}. \quad (6.3)$$

Proof. If B^2 is diagonal, then $\|B\xi\|^2 = \sum_i a_i \xi_i^2$ and the summands $a_i \xi_i^2$ are independent. It remains to note that $E(a_i \xi_i^2) = a_i$, $\text{Var}(a_i \xi_i^2) = 2a_i^2$, and for $\mu a_i < 1$,

$$E \exp\{\mu a_i \xi_i^2 / 2\} = (1 - \mu a_i)^{-1/2}$$

yielding (6.3). \square

Given u , fix $\mu < 1$. The exponential Markov inequality yields

$$\begin{aligned} P(\|B\xi\|^2 > p + u) &\leq \exp\left\{-\frac{\mu(p + u)}{2}\right\} E \exp\left(\frac{\mu\|B\xi\|^2}{2}\right) \\ &\leq \exp\left\{-\frac{\mu u}{2} - \frac{1}{2} \sum_{i=1}^p [\mu a_i + \log(1 - \mu a_i)]\right\}. \end{aligned}$$

We start with the case when $x^{1/2} \leq v/3$. Then $u = 2x^{1/2}v$ fulfills $u \leq 2v^2/3$. Define $\mu = u/v^2 \leq 2/3$ and use that $t + \log(1 - t) \geq -t^2$ for $t \leq 2/3$. This implies

$$\begin{aligned} P(\|B\xi\|^2 > p + u) &\leq \exp\left\{-\frac{\mu u}{2} + \frac{1}{2} \sum_{i=1}^p \mu^2 a_i^2\right\} = \exp(-u^2/(4v^2)) = e^{-x}. \end{aligned} \quad (6.4)$$

Next, let $x^{1/2} > v/3$. Set $\mu = 2/3$. It holds similarly to the above

$$\sum_{i=1}^p [\mu a_i + \log(1 - \mu a_i)] \geq -\sum_{i=1}^p \mu^2 a_i^2 \geq -2v^2/9 \geq -2x.$$

Now, for $u = 6x$ and $\mu u/2 = 2x$, (6.4) implies

$$P(\|B\xi\|^2 > p + u) \leq \exp\{-(2x - x)\} = \exp(-x)$$

as required. \square

Below we establish similar bounds for a non-Gaussian vector ξ obeying (6.1).

6.2 A bound for the ℓ_2 -norm

This section presents a general exponential bound for the probability $P(\|\xi\| > y)$ under (6.1). Define the value

$$x_c = 0.5[g^2 - p \log(1 + g^2/p)]. \quad (6.5)$$

Theorem 6.4. *Let $\xi \in \mathbb{R}^p$ fulfill (6.1). Then it holds for each $x \leq x_c$*

$$P(\|\xi\|^2 > p + \sqrt{\kappa x p} \vee (\kappa x), \|\xi\|^2 \leq p + g^2) \leq 2 \exp(-x),$$

where $\varkappa = 6.6$. Moreover, for $y^2 \geq y_c^2 \stackrel{\text{def}}{=} p + g^2$, it holds with $g_c = g^2(p + g^2)^{-1/2}$

$$\begin{aligned} \mathbb{P}(\|\xi\|^2 > y^2) &\leq 8.4 \exp\{-g_c y/2 - (p/2) \log(1 - g_c/y)\} \\ &\leq 8.4 \exp\{-x_c - g_c(y - y_c)/2\}. \end{aligned}$$

Proof. The main step of the proof is the following exponential bound.

Lemma 6.5. *Suppose (6.1). For any $\mu < 1$ with $g^2 > p\mu$, it holds*

$$\mathbb{E} \exp\left(\frac{\mu \|\xi\|^2}{2}\right) \mathbb{I}\left(\|\xi\|^2 \leq \frac{g^2}{\mu^2} - \frac{p}{\mu}\right) \leq 2(1 - \mu)^{-p/2}. \quad (6.6)$$

Proof. Let ε be a standard normal vector in \mathbb{R}^p and $\mathbf{u} \in \mathbb{R}^p$. Given $\mathbf{r} > 0$, define $m(\mathbf{u}, r) = \mathbb{P}(\|\varepsilon - \mathbf{u}\| \leq \mathbf{r})$. This value can be rewritten as

$$\begin{aligned} m(\mathbf{u}, r) &= \mathbb{P}(\|\varepsilon - \mathbf{u}\|^2 \leq \mathbf{r}^2) \\ &= \mathbb{P}(\|\varepsilon\|^2 - \mathbb{E}\|\varepsilon\|^2 - 2\mathbf{u}^\top \varepsilon \leq \mathbf{r}^2 - \|\mathbf{u}\|^2 - \mathbb{E}\|\varepsilon\|^2), \end{aligned}$$

and $m(\mathbf{u}, r) \geq 1/2$ for $\mathbf{r}^2 \geq \|\mathbf{u}\|^2 + \mathbb{E}\|\varepsilon\|^2 = \|\mathbf{u}\|^2 + p$. Let us fix some ξ with $\mu \|\xi\|^2 \leq g^2/\mu - p$ and denote by \mathbb{P}_ξ the conditional probability given ξ . It holds with $c_p = (2\pi)^{-p/2}$

$$\begin{aligned} &c_p \int \exp\left(\gamma^\top \xi - \frac{\|\gamma\|^2}{2\mu}\right) \mathbb{I}(\|\gamma\| \leq g) d\gamma \\ &= c_p \exp(\mu \|\xi\|^2/2) \int \exp\left(-\frac{1}{2} \|\mu^{-1/2} \gamma - \mu^{1/2} \xi\|^2\right) \mathbb{I}(\mu^{-1} \|\gamma\|^2 \leq \mu^{-1} g^2) d\gamma \\ &= \mu^{p/2} \exp(\mu \|\xi\|^2/2) \mathbb{P}_\xi(\|\varepsilon + \mu^{1/2} \xi\|^2 \leq \mu^{-1} g^2) \\ &\geq 0.5 \mu^{p/2} \exp(\mu \|\xi\|^2/2), \end{aligned}$$

because $\|\mu^{1/2} \xi\|^2 + p \leq \mu^{-1} g^2$. This implies in view of $p < g^2/\mu$ that

$$\begin{aligned} &\exp(\mu \|\xi\|^2/2) \mathbb{I}(\|\xi\|^2 \leq g^2/\mu^2 - p/\mu) \\ &\leq 2\mu^{-p/2} c_p \int \exp\left(\gamma^\top \xi - \frac{\|\gamma\|^2}{2\mu}\right) \mathbb{I}(\|\gamma\| \leq g) d\gamma. \end{aligned}$$

Further, by (6.1)

$$\begin{aligned} &c_p \mathbb{E} \int \exp\left(\gamma^\top \xi - \frac{1}{2\mu} \|\gamma\|^2\right) \mathbb{I}(\|\gamma\| \leq g) d\gamma \\ &\leq c_p \int \exp\left(-\frac{\mu^{-1} - 1}{2} \|\gamma\|^2\right) \mathbb{I}(\|\gamma\| \leq g) d\gamma \\ &\leq c_p \int \exp\left(-\frac{\mu^{-1} - 1}{2} \|\gamma\|^2\right) d\gamma \\ &\leq (\mu^{-1} - 1)^{-p/2} \end{aligned}$$

and (6.6) follows. \square

Due to this result, the scaled squared norm $\mu\|\xi\|^2/2$ after a proper truncation possesses the same exponential moments as in the Gaussian case. A straightforward implication is the probability bound $\mathbb{P}(\|\xi\|^2 > p + u)$ for moderate values u . Namely, given $u > 0$, define $\mu = u/(u + p)$. This value optimizes the inequality (6.2) in the Gaussian case. Now we can apply a similar bound under the constraints $\|\xi\|^2 \leq g^2/\mu^2 - p/\mu$. Therefore, the bound is only meaningful if $p + u \leq g^2/\mu^2 - p/\mu$ with $\mu = u/(u + p)$. One can check that the largest value u for which this constraint is still valid, is given by $u = g^2$. Hence, (6.6) yields for $u \leq g^2$

$$\begin{aligned} & \mathbb{P}(\|\xi\|^2 > p + u, \|\xi\|^2 \leq p + g^2) \\ & \leq \exp\left\{-\frac{\mu(p + u)}{2}\right\} \mathbb{E} \exp\left(\frac{\mu\|\xi\|^2}{2}\right) \mathbb{I}\left(\|\xi\|^2 \leq \frac{g^2}{\mu^2} - \frac{p}{\mu}\right) \\ & \leq 2 \exp\{-0.5[\mu(p + u) + p \log(1 - \mu)]\} \\ & = 2 \exp\{-0.5[u - p \log(1 + u/p)]\}. \end{aligned}$$

Similarly to the Gaussian case, this implies with $\varkappa = 6.6$ that

$$\mathbb{P}(\|\xi\|^2 \geq p + \sqrt{\varkappa xp} \vee (\varkappa x), \|\xi\|^2 \leq p + g^2) \leq 2 \exp(-x).$$

The Gaussian case yields (6.1) with $g = \infty$ and the result is done. In the non-Gaussian case with a finite g , we have to accompany the moderate deviation bound with a large deviation bound $\mathbb{P}(\|\xi\| > y)$ for $y^2 \geq p + g^2$. This is done by combining the bound (6.6) with the standard slicing arguments.

Lemma 6.6. *Let $\mu_0 \leq g^2/p$. Define $y_0^2 = g^2/\mu_0^2 - p/\mu_0$ and $g_0^2 = g^2 - \mu_0 p$. It holds for $y \geq y_0$*

$$\mathbb{P}(\|\xi\| > y) \leq 8.4(1 - g_0/y)^{-p/2} \exp(-g_0 y/2) \quad (6.7)$$

$$\leq 8.4 \exp\{-x_0 - g_0(y - y_0)/2\}. \quad (6.8)$$

with x_0 defined by

$$2x_0 = \mu_0 y_0^2 + p \log(1 - \mu_0) = g^2/\mu_0 - p + p \log(1 - \mu_0).$$

Proof. Consider the growing sequence y_k with $y_1 = y$ and $g_0 y_{k+1} = g_0 y + k$. Define also $\mu_k = g_0/y_k$. In particular, $\mu_k \leq \mu_1 = g_0/y$. Obviously

$$\mathbb{P}(\|\xi\| > y) = \sum_{k=1}^{\infty} \mathbb{P}(\|\xi\| > y_k, \|\xi\| \leq y_{k+1}).$$

Now we try to evaluate every slicing probability in this expression. We use that

$$\mu_{k+1}y_k^2 = \frac{(\mathbf{g}_0\mathbf{y} + k - 1)^2}{\mathbf{g}_0\mathbf{y} + k} \geq \mathbf{g}_0\mathbf{y} + k - 2,$$

and also $\mathbf{g}^2/\mu_k^2 - p/\mu_k \geq y_k^2$ because $y \geq y_0$ and

$$\mathbf{g}^2/\mu_k^2 - p/\mu_k - y_k^2 = \mu_k^{-2}(\mathbf{g}^2 - \mu_k p - \mathbf{g}_0^2) \geq \mu_k^{-2}(\mathbf{g}^2 - \mathbf{g}_0 p/y - \mathbf{g}_0^2) \geq 0.$$

Hence by (6.6)

$$\begin{aligned} \mathbb{P}(\|\boldsymbol{\xi}\| > y) &\leq \sum_{k=1}^{\infty} \mathbb{P}(\|\boldsymbol{\xi}\| > y_k, \|\boldsymbol{\xi}\| \leq y_{k+1}) \\ &\leq \sum_{k=1}^{\infty} \exp\left(-\frac{\mu_{k+1}y_k^2}{2}\right) \mathbb{E} \exp\left(\frac{\mu_{k+1}\|\boldsymbol{\xi}\|^2}{2}\right) \mathbb{I}(\|\boldsymbol{\xi}\|^2 \leq y_{k+1}^2) \\ &\leq \sum_{k=1}^{\infty} 2(1 - \mu_{k+1})^{-p/2} \exp\left(-\frac{\mu_{k+1}y_k^2}{2}\right) \\ &\leq 2(1 - \mu_1)^{-p/2} \sum_{k=1}^{\infty} \exp\left(-\frac{\mathbf{g}_0\mathbf{y} + k - 2}{2}\right) \\ &= 2e^{1/2}(1 - e^{-1/2})^{-1}(1 - \mu_1)^{-p/2} \exp(-\mathbf{g}_0\mathbf{y}/2) \\ &\leq 8.4(1 - \mu_1)^{-p/2} \exp(-\mathbf{g}_0\mathbf{y}/2) \end{aligned}$$

and the first assertion follows. For $y = y_0$, it holds

$$\mathbf{g}_0 y_0 + p \log(1 - \mu_0) = \mu_0 y_0^2 + p \log(1 - \mu_0) = 2\mathbf{x}_0$$

and (6.7) implies $\mathbb{P}(\|\boldsymbol{\xi}\| > y_0) \leq 8.4 \exp(-\mathbf{x}_0)$. Now observe that the function $f(y) = \mathbf{g}_0 y/2 + (p/2) \log(1 - \mathbf{g}_0/y)$ fulfills $f(y_0) = \mathbf{x}_0$ and $f'(y) \geq \mathbf{g}_0/2$ yielding $f(y) \geq \mathbf{x}_0 + \mathbf{g}_0(y - y_0)/2$. This implies (6.8). \square

The statements of the theorem are obtained by applying the lemmas with $\mu_0 = \mu_c = \mathbf{g}^2/(p + \mathbf{g}^2)$. This also implies $y_c^2 = p + \mathbf{g}^2$, $\mathbf{g}_c y_c = \mathbf{g}^2$, $1 - \mu_c = p/(p + \mathbf{g}^2)$, and \mathbf{x}_c from (6.5). \square

The statements of Theorem 6.8 can be represented in the form:

Corollary 6.7. *Let $\boldsymbol{\xi}$ fulfill (6.1). Then it holds for $\mathbf{x} \leq \mathbf{x}_c = 0.5[\mathbf{g}^2 - p \log(1 + \mathbf{g}^2/p)]$:*

$$\mathbb{P}(\|\boldsymbol{\xi}\|^2 \geq \mathfrak{z}(\mathbf{x}, p)) \leq 2e^{-\mathbf{x}} + 8.4e^{-\mathbf{x}_c}, \quad (6.9)$$

$$\mathfrak{z}(\mathbf{x}, p) \stackrel{\text{def}}{=} \begin{cases} p + \sqrt{\mathbf{x} p}, & \mathbf{x} \leq p/\mathbf{x}_c, \\ p + \mathbf{x} & p/\mathbf{x}_c < \mathbf{x} \leq \mathbf{x}_c. \end{cases} \quad (6.10)$$

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

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

This result implicitly assumes that $p \leq \varkappa \mathbf{x}_c$ which is fulfilled if $u_0 = \mathbf{g}^2/p \geq 1$:

$$\varkappa \mathbf{x}_c = 0.5\varkappa[u_0 - \log(1 + u_0)]p \geq 3.3[1 - \log(2)]p > p.$$

In the zone $\mathbf{x} \leq p/\varkappa$ we obtain sub-Gaussian behavior of the tail of $\|\boldsymbol{\xi}\|^2 - p$, in the zone $p/\varkappa < \mathbf{x} \leq \mathbf{x}_c$ it becomes sub-exponential. Note that the sub-exponential zone is empty if $\mathbf{g}^2 < p$.

For $\mathbf{x} \leq \mathbf{x}_c$, the function $\mathfrak{z}(\mathbf{x}, p)$ mimics the quantile behavior of the chi-squared distribution χ_p^2 with p degrees of freedom. Moreover, increase the dimension p yields growth of the sub-Gaussian zone.

Finally, in the large deviation zone $\mathbf{x} > \mathbf{x}_c$ the deviation probability decays as $e^{-c\mathbf{x}^{1/2}}$ for some fixed c . However, if the constant \mathbf{g} in the condition (6.1) is sufficiently large relative to p , then \mathbf{x}_c is large as well and the large deviation zone $\mathbf{x} > \mathbf{x}_c$ can be ignored at a small price of $8.4e^{-\mathbf{x}_c}$ and one can focus on the deviation bound described by (6.9) and (6.10).

6.3 A bound for a quadratic form

Now we extend the result to more general bound for $\|\mathbb{B}\boldsymbol{\xi}\|^2 = \boldsymbol{\xi}^\top \mathbb{B}^2 \boldsymbol{\xi}$ with a given matrix \mathbb{B} and a vector $\boldsymbol{\xi}$ obeying the condition (6.1). Similarly to the Gaussian case we assume that \mathbb{B} is symmetric. Define important characteristics of \mathbb{B}

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

For simplicity of formulation we suppose that $\lambda^* = 1$, otherwise one has to replace \mathbf{p} and \mathbf{v}^2 with \mathbf{p}/λ^* and \mathbf{v}^2/λ^* .

Let \mathbf{g} be shown in (6.1). Define similarly to the ℓ_2 -case $\mu_c = \mathbf{g}^2/(\mathbf{p} + \mathbf{g}^2)$. Further define the values $\mathbf{y}_c, \mathbf{g}_c$, and \mathbf{x}_c by

$$\begin{aligned} \mathbf{y}_c^2 &\stackrel{\text{def}}{=} \mathbf{g}^2/\mu_c^2 - \mathbf{p}/\mu_c, \\ \mathbf{g}_c &\stackrel{\text{def}}{=} \mu_c \mathbf{y}_c = \sqrt{\mathbf{g}^2 - \mu_c \mathbf{p}}, \\ 2\mathbf{x}_c &\stackrel{\text{def}}{=} \mathbf{g}_c \mathbf{y}_c + \log \det(\mathbf{I}_p - \mu_c \mathbb{B}^2). \end{aligned} \tag{6.11}$$

Theorem 6.8. *Let a random vector $\boldsymbol{\xi}$ in \mathbb{R}^p fulfill (6.1). Then for each $\mathbf{x} < \mathbf{x}_c$*

$$\mathbb{P}(\|\mathbb{B}\boldsymbol{\xi}\|^2 > \mathbf{p} + (2\mathbf{v}\mathbf{x}^{1/2}) \vee (a_c \mathbf{x}), \|\mathbb{B}\boldsymbol{\xi}\|^2 \leq \mathbf{y}_c^2) \leq 2 \exp(-\mathbf{x})$$

with $a_c \stackrel{\text{def}}{=} 6 \vee (4\mu_c^{-1})$. Moreover, for $y \geq y_c$, it holds

$$\mathbb{P}(\|\mathcal{B}\xi\|^2 > y) \leq 8.4 \exp(-x_c - g_c(y - y_c)/2).$$

Proof. The main steps of the proof are similar to the proof of Theorem 6.4.

Lemma 6.9. *Suppose (6.1). For any $\mu < 1$ with $g^2/\mu \geq p$, it holds*

$$\mathbb{E} \exp(\mu \|\mathcal{B}\xi\|^2/2) \mathbb{I}(\|\mathcal{B}^2\xi\|^2 \leq g^2/\mu^2 - p/\mu) \leq 2\det(I_p - \mu\mathcal{B}^2)^{-1/2}. \quad (6.12)$$

Proof. With $c_p(\mathcal{B}) = (2\pi)^{-p/2} \det(\mathcal{B}^{-1})$

$$\begin{aligned} c_p(\mathcal{B}) \int \exp\left(\gamma^\top \xi - \frac{1}{2\mu} \|\mathcal{B}^{-1}\gamma\|^2\right) \mathbb{I}(\|\gamma\| \leq g) d\gamma \\ = c_p(\mathcal{B}) \exp\left(\frac{\mu \|\mathcal{B}\xi\|^2}{2}\right) \int \exp\left(-\frac{1}{2} \|\mu^{1/2}\mathcal{B}\xi - \mu^{-1/2}\mathcal{B}^{-1}\gamma\|^2\right) \mathbb{I}(\|\gamma\| \leq g) d\gamma \\ = \mu^{p/2} \exp\left(\frac{\mu \|\mathcal{B}\xi\|^2}{2}\right) \mathbb{P}_\xi(\|\mathcal{B}\varepsilon - \mu^{1/2}\mathcal{B}^2\xi\| \leq g\mu^{-1/2}), \end{aligned}$$

where ε denotes a standard normal vector in \mathbb{R}^p and \mathbb{P}_ξ means the conditional expectation given ξ . Moreover, for any $u \in \mathbb{R}^p$ and $r \geq p + \|u\|^2$, it holds in view of $\mathbb{E}\|\mathcal{B}\varepsilon\|^2 = p$

$$\begin{aligned} \mathbb{P}(\|\mathcal{B}\varepsilon - u\| \leq r) &= \mathbb{P}(\|\mathcal{B}\varepsilon\|^2 - p - 2u^\top \mathcal{B}\varepsilon \leq r^2 - \|u\|^2 - p) \\ &\geq \mathbb{P}(\|\mathcal{B}\varepsilon\|^2 - \mathbb{E}\|\mathcal{B}\varepsilon\|^2 - 2u^\top \mathcal{B}\varepsilon \leq 0) \geq 1/2. \end{aligned}$$

This implies

$$\begin{aligned} \exp(\mu \|\mathcal{B}\xi\|^2/2) \mathbb{I}(\|\mathcal{B}^2\xi\|^2 \leq g^2/\mu^2 - p/\mu) \\ \leq 2\mu^{-p/2} c_p(\mathcal{B}) \int \exp\left(\gamma^\top \xi - \frac{1}{2\mu} \|\mathcal{B}^{-1}\gamma\|^2\right) \mathbb{I}(\|\gamma\| \leq g) d\gamma. \end{aligned}$$

Further, by (6.1)

$$\begin{aligned} c_p(\mathcal{B}) \mathbb{E} \int \exp\left(\gamma^\top \xi - \frac{1}{2\mu} \|\mathcal{B}^{-1}\gamma\|^2\right) \mathbb{I}(\|\gamma\| \leq g) d\gamma \\ \leq c_p(\mathcal{B}) \int \exp\left(\frac{\|\gamma\|^2}{2} - \frac{1}{2\mu} \|\mathcal{B}^{-1}\gamma\|^2\right) d\gamma \\ \leq \det(\mathcal{B}^{-1}) \det(\mu^{-1}\mathcal{B}^{-2} - I_p)^{-1/2} = \mu^{p/2} \det(I_p - \mu\mathcal{B}^2)^{-1/2} \end{aligned}$$

and (6.12) follows. \square

Now we evaluate the probability $\mathbb{P}(\|\mathcal{B}\xi\| > y)$ for moderate values of y .

Lemma 6.10. *Let $\mu_0 < 1 \wedge (\mathbf{g}^2/\mathbf{p})$. With $y_0^2 = \mathbf{g}^2/\mu_0^2 - \mathbf{p}/\mu_0$, it holds for any $u > 0$*

$$\begin{aligned} \mathbb{P}(\|\mathcal{B}\boldsymbol{\xi}\|^2 > \mathbf{p} + u, \|\mathcal{B}^2\boldsymbol{\xi}\| \leq y_0) \\ \leq 2 \exp\{-0.5\mu_0(\mathbf{p} + u) - 0.5 \log \det(\mathcal{I}_p - \mu_0 \mathcal{B}^2)\}. \end{aligned} \quad (6.13)$$

In particular, if \mathcal{B}^2 is diagonal, that is, $\mathcal{B}^2 = \text{diag}(a_1, \dots, a_p)$, then

$$\begin{aligned} \mathbb{P}(\|\mathcal{B}\boldsymbol{\xi}\|^2 > \mathbf{p} + u, \|\mathcal{B}^2\boldsymbol{\xi}\| \leq y_0) \\ \leq 2 \exp\left\{-\frac{\mu_0 u}{2} - \frac{1}{2} \sum_{i=1}^p [\mu_0 a_i + \log(1 - \mu_0 a_i)]\right\}. \end{aligned} \quad (6.14)$$

Proof. The exponential Chebyshev inequality and (6.12) imply

$$\begin{aligned} \mathbb{P}(\|\mathcal{B}\boldsymbol{\xi}\|^2 > \mathbf{p} + u, \|\mathcal{B}^2\boldsymbol{\xi}\| \leq y_0) \\ \leq \exp\left\{-\frac{\mu_0(\mathbf{p} + u)}{2}\right\} \mathbb{E} \exp\left(\frac{\mu_0 \|\mathcal{B}\boldsymbol{\xi}\|^2}{2}\right) \mathbb{1}\left(\|\mathcal{B}^2\boldsymbol{\xi}\|^2 \leq \frac{\mathbf{g}^2}{\mu_0^2} - \frac{\mathbf{p}}{\mu_0}\right) \\ \leq 2 \exp\{-0.5\mu_0(\mathbf{p} + u) - 0.5 \log \det(\mathcal{I}_p - \mu_0 \mathcal{B}^2)\}. \end{aligned}$$

Moreover, the standard change-of-basis arguments allow us to reduce the problem to the case of a diagonal matrix $\mathcal{B}^2 = \text{diag}(a_1, \dots, a_p)$ where $1 = a_1 \geq a_2 \geq \dots \geq a_p > 0$. Note that $\mathbf{p} = a_1 + \dots + a_p$. Then the claim (6.13) can be written in the form (6.14). \square

Now we evaluate a large deviation probability that $\|\mathcal{B}\boldsymbol{\xi}\| > \mathbf{y}$ for a large \mathbf{y} . Note that the condition $\|\mathcal{B}^2\boldsymbol{\xi}\|_\infty \leq 1$ implies $\|\mathcal{B}^2\boldsymbol{\xi}\| \leq \|\mathcal{B}\boldsymbol{\xi}\|$. So, the bound (6.13) continues to hold when $\|\mathcal{B}^2\boldsymbol{\xi}\| \leq y_0$ is replaced by $\|\mathcal{B}\boldsymbol{\xi}\| \leq y_0$.

Lemma 6.11. *Let $\mu_0 < 1$ and $\mu_0 \mathbf{p} < \mathbf{g}^2$. Define \mathbf{g}_0 by $\mathbf{g}_0^2 = \mathbf{g}^2 - \mu_0 \mathbf{p}$. For any $\mathbf{y} \geq y_0 \stackrel{\text{def}}{=} \mathbf{g}_0/\mu_0$, it holds*

$$\begin{aligned} \mathbb{P}(\|\mathcal{B}\boldsymbol{\xi}\| > \mathbf{y}) &\leq 8.4 \det\{\mathcal{I}_p - (\mathbf{g}_0/\mathbf{y})\mathcal{B}^2\}^{-1/2} \exp(-\mathbf{g}_0 \mathbf{y}/2). \\ &\leq 8.4 \exp(-\mathbf{x}_0 - \mathbf{g}_0(\mathbf{y} - y_0)/2), \end{aligned} \quad (6.15)$$

where \mathbf{x}_0 is defined by

$$2\mathbf{x}_0 = \mathbf{g}_0 y_0 + \log \det\{\mathcal{I}_p - (\mathbf{g}_0/y_0)\mathcal{B}^2\}.$$

Proof. The slicing arguments of Lemma 6.6 apply here in the same manner. One has to replace $\|\boldsymbol{\xi}\|$ by $\|\mathcal{B}\boldsymbol{\xi}\|$ and $(1 - \mu_1)^{-p/2}$ by $\det\{\mathcal{I}_p - (\mathbf{g}_0/\mathbf{y})\mathcal{B}^2\}^{-1/2}$. We omit the details. In particular, with $\mathbf{y} = y_0 = \mathbf{g}_0/\mu_0$, this yields

$$\mathbb{P}(\|\mathcal{B}\boldsymbol{\xi}\| > y_0) \leq 8.4 \exp(-\mathbf{x}_0).$$

Moreover, for the function $f(y) = g_0 y + \log \det\{I_p - (g_0/y)B^2\}$, it holds $f'(y) \geq g_0$ and hence, $f(y) \geq f(y_0) + g_0(y - y_0)$ for $y > y_0$. This implies (6.15). \square

One important feature of the results of Lemma 6.10 and Lemma 6.11 is that the value $\mu_0 < 1 \wedge (g^2/p)$ can be selected arbitrarily. In particular, for $y \geq y_c$, Lemma 6.11 with $\mu_0 = \mu_c$ yields the large deviation probability $\mathbb{P}(\|B\xi\| > y)$. For bounding the probability $\mathbb{P}(\|B\xi\|^2 > p + u, \|B\xi\| \leq y_c)$, we use the inequality $\log(1 - t) \geq -t - t^2$ for $t \leq 2/3$. It implies for $\mu \leq 2/3$ that

$$\begin{aligned} & -\log \mathbb{P}(\|B\xi\|^2 > p + u, \|B\xi\| \leq y_c) \\ & \geq \mu(p + u) + \sum_{i=1}^p \log(1 - \mu a_i) \\ & \geq \mu(p + u) - \sum_{i=1}^p (\mu a_i + \mu^2 a_i^2) \geq \mu u - \mu^2 v^2 / 2. \end{aligned} \quad (6.16)$$

Now we distinguish between $\mu_c \geq 2/3$ and $\mu_c < 2/3$ starting with $\mu_c \geq 2/3$. The bound (6.16) with $\mu_0 = (u/v^2) \wedge (2/3) \leq \mu_c$ and with $u = (2vx^{1/2}) \vee (6x)$ yields

$$\mathbb{P}(\|B\xi\|^2 > p + u, \|B\xi\| \leq y_c) \leq 2 \exp(-x);$$

see the proof of Theorem 6.2 for the Gaussian case.

Now consider $\mu_c < 2/3$. For $x^{1/2} \leq \mu_c v/2$, use $u = 2vx^{1/2}$ and $\mu_0 = u/v^2$. It holds $\mu_0 = u/v^2 \leq \mu_c$ and $u^2/(4v^2) = x$ yielding the desired bound by (6.16). For $x^{1/2} > \mu_c v/2$, we select again $\mu_0 = \mu_c$. It holds with $u = 4\mu_c^{-1}x$ that $\mu_c u/2 - \mu_c^2 v^2/4 \geq 2x - x = x$. This completes the proof. \square

Now we describe the value $\mathfrak{z}(x, B)$ ensuring a small value for the large deviation probability $\mathbb{P}(\|B\xi\|^2 > \mathfrak{z}(x, B))$. For ease of formulation, we suppose that $g^2 \geq 2p$ yielding $\mu_c^{-1} \leq 3/2$. The other case can be easily adjusted.

Corollary 6.12. *Let ξ fulfill (6.1) with $g^2 \geq 2p$. Then it holds for $x \leq x_c$ with x_c from (6.11):*

$$\begin{aligned} \mathbb{P}(\|B\xi\|^2 \geq \mathfrak{z}(x, B)) & \leq 2e^{-x} + 8.4e^{-x_c}, \\ \mathfrak{z}(x, B) & \stackrel{\text{def}}{=} \begin{cases} p + \sqrt{2xv}, & x \leq v/18, \\ p + 6x & v/18 < x \leq x_c. \end{cases} \end{aligned} \quad (6.17)$$

For $x > x_c$

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

6.4 Rescaling and regularity condition

The result of Theorem 6.8 can be extended to a more general situation when the condition (6.1) is fulfilled for a vector ζ rescaled by a matrix V_0 . More precisely, let the random p -vector ζ fulfill for some $p \times p$ matrix V_0 the condition

$$\sup_{\gamma \in \mathbb{R}^p} \log \mathbb{E} \exp \left\{ \lambda \frac{\gamma^\top \zeta}{\|V_0 \gamma\|} \right\} \leq \nu_0^2 \lambda^2 / 2, \quad |\lambda| \leq \mathbf{g}, \quad (6.18)$$

with some constants $\mathbf{g} > 0$, $\nu_0 \geq 1$. Again, a simple change of variables reduces the case of an arbitrary $\nu_0 \geq 1$ to $\nu_0 = 1$. Our aim is to bound the squared norm $\|D_0^{-1} \zeta\|^2$ of a vector $D_0^{-1} \zeta$ for another $p \times p$ positive symmetric matrix D_0^2 . Note that condition (6.18) implies (6.1) for the rescaled vector $\xi = V_0^{-1} \zeta$. This leads to bounding the quadratic form $\|D_0^{-1} V_0 \xi\|^2 = \|\mathbb{B} \xi\|^2$ with $\mathbb{B}^2 = D_0^{-1} V_0^2 D_0^{-1}$. It obviously holds

$$\mathbf{p} = \text{tr}(\mathbb{B}^2) = \text{tr}(D_0^{-2} V_0^2).$$

Now we can apply the result of Corollary 6.12.

Corollary 6.13. *Let ζ fulfill (6.18) with some V_0 and \mathbf{g} . Given D_0 , define $\mathbb{B}^2 = D_0^{-1} V_0^2 D_0^{-1}$, and let $\mathbf{g}^2 \geq 2\mathbf{p}$. Then it holds for $\mathbf{x} \leq \mathbf{x}_c$ with \mathbf{x}_c from (6.11):*

$$\mathbb{P}(\|D_0^{-1} \zeta\|^2 \geq \mathfrak{z}(\mathbf{x}, \mathbb{B})) \leq 2e^{-\mathbf{x}} + 8.4e^{-\mathbf{x}_c},$$

with $\mathfrak{z}(\mathbf{x}, \mathbb{B})$ from (6.17). For $\mathbf{x} > \mathbf{x}_c$

$$\mathbb{P}(\|D_0^{-1} \zeta\|^2 \geq \mathfrak{z}_c(\mathbf{x}, \mathbb{B})) \leq 8.4e^{-\mathbf{x}}, \quad \mathfrak{z}_c(\mathbf{x}, \mathbb{B}) \stackrel{\text{def}}{=} |\mathbf{y}_c + 2(\mathbf{x} - \mathbf{x}_c)/\mathbf{g}_c|^2.$$

Finally we briefly discuss the *regular* case with $D_0 \geq \mathbf{a}V_0$ for some $\mathbf{a} > 0$. This implies $\|\mathbb{B}\|_\infty \leq \mathbf{a}^{-1}$ and

$$\mathbf{v}^2 = 2 \text{tr}(\mathbb{B}^4) \leq 2\mathbf{a}^{-2}\mathbf{p}.$$

This together with $\mathbf{g}^2 \geq 2\mathbf{p}$ yields

$$\mathbf{y}_c^2 \stackrel{\text{def}}{=} \mathbf{g}^2 / \mu_c^2 - \mathbf{p} / \mu_c \geq \mathbf{p} / \mu_c^2,$$

$$\mathbf{g}_c \stackrel{\text{def}}{=} \mu_c \mathbf{y}_c \geq \sqrt{\mathbf{p}},$$

$$2\mathbf{x}_c \stackrel{\text{def}}{=} \mathbf{g}_c \mathbf{y}_c + \log \det(I_p - \mu_c \mathbb{B}^2).$$

7 Some results for empirical processes

This chapter presents some general results of the theory of empirical processes. Under the global conditions from Section 2 one can apply the well developed chaining arguments in Orlic spaces; see e.g. van der Vaart and Wellner (1996), Chapter ???. We follow the more recent approach inspired by the notions of generic chaining and majorizing measures due to M. Talagrand. The chaining arguments are replaced by the *pulling* device; see e.g. Talagrand (1996, 2001, 2005). The results are close to that of Bednorz (2006). We state the results in a slightly different form and present an independent and self-containing proof.

The first result states a bound for local fluctuations of the process $\mathcal{U}(\mathbf{v})$ given on a metric space \mathcal{Y} . Then this result will be used for bounding the maximum of the negatively drifted process $\mathcal{U}(\mathbf{v}) - \mathcal{U}(\mathbf{v}_0) - \rho d^2(\mathbf{v}, \mathbf{v}_0)$ over a vicinity $\mathcal{Y}_\circ(\mathbf{r})$ of the central point \mathbf{v}_0 . The behavior of $\mathcal{U}(\mathbf{v})$ outside of the local central set $\mathcal{Y}_\circ(\mathbf{r})$ is described using the *upper function* method. Namely, we construct a multiscale deterministic function $\mathbf{u}(\mu, \mathbf{v})$ ensuring that with probability at least $1 - e^{-x}$ it holds $\mu \mathcal{U}(\mathbf{v}) + \mathbf{u}(\mu, \mathbf{v}) \leq \mathfrak{z}(\mathbf{x})$ for all $\mathbf{v} \notin \mathcal{Y}_\circ(\mathbf{r})$ and $\mu \in \mathbb{M}$, where $\mathfrak{z}(\mathbf{x})$ grows linearly in \mathbf{x} .

7.1 A bound for local fluctuations

An important step in the whole construction is an exponential bound on the maximum of a random process $\mathcal{U}(\mathbf{v})$ under the exponential moment conditions on its increments. Let $d(\mathbf{v}, \mathbf{v}')$ be a semi-distance on \mathcal{Y} . We suppose the following condition to hold:

(Ed) *There exist $\mathbf{g} > 0$, $\mathbf{r}_1 > 0$, $\nu_0 \geq 1$, such that for any $\lambda \leq \mathbf{g}$ and $\mathbf{v}, \mathbf{v}' \in \mathcal{Y}$ with $d(\mathbf{v}, \mathbf{v}') \leq \mathbf{r}_1$*

$$\log \mathbb{E} \exp \left\{ \lambda \frac{\mathcal{U}(\mathbf{v}) - \mathcal{U}(\mathbf{v}')}{d(\mathbf{v}, \mathbf{v}')} \right\} \leq \nu_0^2 \lambda^2 / 2.$$

Formulation of the result involves a sigma-finite measure π on the space \mathcal{Y} which is often called the *majorizing measure* and used in the *generic chaining* device; see Talagrand (2005). A typical example of choosing π is the Lebesgue measure on \mathbb{R}^p . Let \mathcal{Y}° be a subset of \mathcal{Y} , a sequence \mathbf{r}_k be fixed with $\mathbf{r}_0 = \text{diam}(\mathcal{Y}^\circ)$ and $\mathbf{r}_k = \mathbf{r}_0 2^{-k}$. Let also $\mathcal{B}_k(\mathbf{v}) \stackrel{\text{def}}{=} \{\mathbf{v}' \in \mathcal{Y}^\circ : d(\mathbf{v}, \mathbf{v}') \leq \mathbf{r}_k\}$ be the d -ball centered at \mathbf{v} of radius \mathbf{r}_k and $\pi_k(\mathbf{v})$ denote its π -measure:

$$\pi_k(\mathbf{v}) \stackrel{\text{def}}{=} \int_{\mathcal{B}_k(\mathbf{v})} \pi(d\mathbf{v}') = \int_{\mathcal{Y}^\circ} \mathbb{I}(d(\mathbf{v}, \mathbf{v}') \leq \mathbf{r}_k) \pi(d\mathbf{v}').$$

Denote also

$$M_k \stackrel{\text{def}}{=} \max_{\mathbf{v} \in \mathcal{I}^\circ} \frac{\pi(\mathcal{I}^\circ)}{\pi_k(\mathbf{v})} \quad k \geq 1. \quad (7.1)$$

Finally set $c_1 = 1/3$, $c_k = 2^{-k+2}/3$ for $k \geq 2$, and define the value $\mathbb{Q}(\mathcal{I}^\circ)$ by

$$\mathbb{Q}(\mathcal{I}^\circ) \stackrel{\text{def}}{=} \sum_{k=1}^{\infty} c_k \log(2M_k) = \frac{1}{3} \log(2M_1) + \frac{4}{3} \sum_{k=2}^{\infty} 2^{-k} \log(2M_k).$$

Theorem 7.1. *Suppose (Ed). If \mathcal{I}° is a central set with the center \mathbf{v}° and the radius \mathbf{r}_1 , i.e. $d(\mathbf{v}, \mathbf{v}^\circ) \leq \mathbf{r}_1$ for all $\mathbf{v} \in \mathcal{I}^\circ$, then for $\lambda \leq \mathbf{g}_0 \stackrel{\text{def}}{=} \nu_0 \mathbf{g}$*

$$\log \mathbb{E} \exp \left\{ \frac{\lambda}{3\nu_0 \mathbf{r}_1} \sup_{\mathbf{v} \in \mathcal{I}^\circ} |\mathcal{U}(\mathbf{v}) - \mathcal{U}(\mathbf{v}^\circ)| \right\} \leq \lambda^2/2 + \mathbb{Q}(\mathcal{I}^\circ). \quad (7.2)$$

Proof. A simple change $\mathcal{U}(\cdot)$ with $\nu_0^{-1}\mathcal{U}(\cdot)$ and \mathbf{g} with $\mathbf{g}_0 = \nu_0 \mathbf{g}$ allows for reducing the result to the case with $\nu_0 = 1$ which we assume below. Consider for $k \geq 1$ the smoothing operator \mathbb{S}_k defined as

$$\mathbb{S}_k f(\mathbf{v}^\circ) = \frac{1}{\pi_k(\mathbf{v}^\circ)} \int_{\mathcal{B}_k(\mathbf{v}^\circ)} f(\mathbf{v}) \pi(d\mathbf{v}).$$

Further, define

$$\mathbb{S}_0 \mathcal{U}(\mathbf{v}) \equiv \mathcal{U}(\mathbf{v}^\circ)$$

so that $\mathbb{S}_0 \mathcal{U}$ is a constant function and the same holds for $\mathbb{S}_k \mathbb{S}_{k-1} \dots \mathbb{S}_0 \mathcal{U}$ with any $k \geq 1$. If $f(\cdot) \leq g(\cdot)$ for two non-negative functions f and g , then $\mathbb{S}_k f(\cdot) \leq \mathbb{S}_k g(\cdot)$. Separability of the process \mathcal{U} implies that $\lim_k \mathbb{S}_k \mathcal{U}(\mathbf{v}) = \mathcal{U}(\mathbf{v})$. We conclude that for each $\mathbf{v} \in \mathcal{I}^\circ$

$$\begin{aligned} |\mathcal{U}(\mathbf{v}) - \mathcal{U}(\mathbf{v}^\circ)| &= \lim_{k \rightarrow \infty} |\mathbb{S}_k \mathcal{U}(\mathbf{v}) - \mathbb{S}_k \dots \mathbb{S}_0 \mathcal{U}(\mathbf{v})| \\ &\leq \lim_{k \rightarrow \infty} \sum_{i=1}^k |\mathbb{S}_k \dots \mathbb{S}_i (I - \mathbb{S}_{i-1}) \mathcal{U}(\mathbf{v})| \leq \sum_{i=1}^{\infty} \xi_k^*. \end{aligned}$$

Here $\xi_k^* \stackrel{\text{def}}{=} \sup_{\mathbf{v} \in \mathcal{I}^\circ} \xi_k(\mathbf{v})$ for $k \geq 1$ with

$$\xi_1(\mathbf{v}) \equiv |\mathbb{S}_1 \mathcal{U}(\mathbf{v}) - \mathcal{U}(\mathbf{v}^\circ)|, \quad \xi_k^* \stackrel{\text{def}}{=} |\mathbb{S}_k (I - \mathbb{S}_{k-1}) \mathcal{U}(\mathbf{v})|, \quad k \geq 2$$

For a fixed point \mathbf{v}^\sharp , it holds

$$\xi_k(\mathbf{v}^\sharp) \leq \frac{1}{\pi_k(\mathbf{v}^\sharp)} \int_{\mathcal{B}_k(\mathbf{v}^\sharp)} \frac{1}{\pi_{k-1}(\mathbf{v})} \int_{\mathcal{B}_{k-1}(\mathbf{v})} |\mathcal{U}(\mathbf{v}) - \mathcal{U}(\mathbf{v}')| \pi(d\mathbf{v}') \pi(d\mathbf{v}).$$

For each $\mathbf{v}' \in \mathcal{B}_{k-1}(\mathbf{v})$, it holds $d(\mathbf{v}, \mathbf{v}') \leq \mathbf{r}_{k-1} = 2\mathbf{r}_k$ and

$$|\mathcal{U}(\mathbf{v}) - \mathcal{U}(\mathbf{v}')| \leq \mathbf{r}_{k-1} \frac{|\mathcal{U}(\mathbf{v}) - \mathcal{U}(\mathbf{v}')|}{d(\mathbf{v}, \mathbf{v}')}.$$

This implies for each $\mathbf{v}^\# \in \mathcal{Y}^\circ$ and $k \geq 2$ by the Jensen inequality and (7.1)

$$\begin{aligned} \exp\left\{\frac{\lambda}{\mathbf{r}_{k-1}}\xi_k(\mathbf{v}^\#)\right\} &\leq \int_{\mathcal{B}_k(\mathbf{v}^\#)} \left(\int_{\mathcal{B}_{k-1}(\mathbf{v})} \exp \frac{\lambda|\mathcal{U}(\mathbf{v}) - \mathcal{U}(\mathbf{v}')|}{d(\mathbf{v}, \mathbf{v}')} \frac{\pi(d\mathbf{v}')}{\pi_{k-1}(\mathbf{v})} \right) \frac{\pi(d\mathbf{v})}{\pi_k(\mathbf{v}^\#)} \\ &\leq M_k \int_{\mathcal{Y}^\circ} \left(\int_{\mathcal{B}_{k-1}(\mathbf{v})} \exp \frac{\lambda|\mathcal{U}(\mathbf{v}) - \mathcal{U}(\mathbf{v}')|}{d(\mathbf{v}, \mathbf{v}')} \frac{\pi(d\mathbf{v}')}{\pi_{k-1}(\mathbf{v})} \right) \frac{\pi(d\mathbf{v})}{\pi(\mathcal{Y}^\circ)}. \end{aligned}$$

As the right hand-side does not depend on $\mathbf{v}^\#$, this yields for $\xi_k^* \stackrel{\text{def}}{=} \sup_{\mathbf{v} \in \mathcal{Y}^\circ} \xi_k(\mathbf{v})$ by condition $(\mathcal{E}d)$ in view of $e^{|x|} \leq e^x + e^{-x}$

$$\begin{aligned} \mathbb{E} \exp\left(\frac{\lambda}{\mathbf{r}_{k-1}}\xi_k^*\right) &\leq M_k \int_{\mathcal{Y}^\circ} \left(\int_{\mathcal{B}_{k-1}(\mathbf{v})} \mathbb{E} \exp \frac{\lambda|\mathcal{U}(\mathbf{v}) - \mathcal{U}(\mathbf{v}')|}{d(\mathbf{v}, \mathbf{v}')} \frac{\pi(d\mathbf{v}')}{\pi_{k-1}(\mathbf{v})} \right) \frac{\pi(d\mathbf{v})}{\pi(\mathcal{Y}^\circ)} \\ &\leq 2M_k \exp(\lambda^2/2) \int_{\mathcal{Y}^\circ} \left(\int_{\mathcal{B}_{k-1}(\mathbf{v})} \frac{\pi(d\mathbf{v}')}{\pi_{k-1}(\mathbf{v})} \right) \frac{\pi(d\mathbf{v})}{\pi(\mathcal{Y}^\circ)} \\ &= 2M_k \exp(\lambda^2/2). \end{aligned}$$

Further, the use of $d(\mathbf{v}, \mathbf{v}^\circ) \leq \mathbf{r}_1$ for all $\mathbf{v} \in \mathcal{Y}^\circ$ yields by $(\mathcal{E}d)$

$$\mathbb{E} \exp\left\{\frac{\lambda}{\mathbf{r}_1}|\mathcal{U}(\mathbf{v}) - \mathcal{U}(\mathbf{v}^\circ)|\right\} \leq 2 \exp(\lambda^2/2) \quad (7.3)$$

and thus

$$\begin{aligned} \mathbb{E} \exp\left\{\frac{\lambda}{\mathbf{r}_1}|\mathbb{S}_1\mathcal{U}(\mathbf{v}) - \mathcal{U}(\mathbf{v}^\circ)|\right\} &\leq \frac{1}{\pi_1(\mathbf{v})} \int_{\mathcal{B}_1(\mathbf{v})} \mathbb{E} \exp\left\{\frac{\lambda}{\mathbf{r}_1}|\mathcal{U}(\mathbf{v}') - \mathcal{U}(\mathbf{v}^\circ)|\right\} \pi(d\mathbf{v}') \\ &\leq \frac{M_1}{\pi(\mathcal{Y}^\circ)} \int_{\mathcal{Y}^\circ} \mathbb{E} \exp\left\{\frac{\lambda}{\mathbf{r}_1}|\mathcal{U}(\mathbf{v}') - \mathcal{U}(\mathbf{v}^\circ)|\right\} \pi(d\mathbf{v}'). \end{aligned}$$

This implies by (7.3) for $\xi_1^* \equiv \sup_{\mathbf{v} \in \mathcal{Y}^\circ} |\mathbb{S}_1\mathcal{U}(\mathbf{v}) - \mathcal{U}(\mathbf{v}^\circ)|$

$$\mathbb{E} \exp\left(\frac{\lambda}{\mathbf{r}_1}\xi_1^*\right) \leq 2M_1 \exp(\lambda^2/2).$$

Denote $c_1 = 1/3$ and $c_k = \mathbf{r}_{k-1}/(3\mathbf{r}_1) = 2^{-k+2}/3$ for $k \geq 2$. Then $\sum_{k=1}^\infty c_k = 1$ and it holds by the Hölder inequality; see Lemma 7.11 below:

$$\begin{aligned} \log \mathbb{E} \exp\left(\frac{\lambda}{3\mathbf{r}_1} \sum_{k=1}^\infty \xi_k^*\right) &\leq c_1 \log \mathbb{E} \exp\left(\frac{\lambda}{\mathbf{r}_1}\xi_1^*\right) + \sum_{k=2}^\infty c_k \log \mathbb{E} \exp\left(\frac{\lambda}{\mathbf{r}_{k-1}}\xi_k^*\right) \\ &\leq \lambda^2/2 + c_1 \log(2M_1) + \sum_{k=2}^\infty c_k \log(2M_k) \\ &< \lambda^2/2 + \mathbb{Q}(\mathcal{Y}^\circ). \end{aligned}$$

This implies the result. \square

7.2 A local central bound

Due to the result of Theorem 7.1, the bound for the maximum of $\mathcal{U}(\mathbf{v}, \mathbf{v}_0)$ over $\mathbf{v} \in \mathcal{B}_{\mathbf{r}}(\mathbf{v}_0)$ grows quadratically in \mathbf{r} . So, its applications to situations with $\mathbf{r}^2 \gg \mathbb{Q}(\mathcal{Y}^\circ)$ are limited. The next result shows that introducing a negative quadratic drift helps to state a uniform in \mathbf{r} local probability bound. Namely, the bound for the process $\mathcal{U}(\mathbf{v}, \mathbf{v}_0) - \rho d^2(\mathbf{v}, \mathbf{v}_0)/2$ with some positive ρ over a ball $\mathcal{B}_{\mathbf{r}}(\mathbf{v}_0)$ around the point \mathbf{v}_0 only depends on the drift coefficient ρ but not on \mathbf{r} . Here the generic chaining arguments are accomplished with the *slicing* technique. The idea is for a given $\mathbf{r}^* > 1$ to split the ball $\mathcal{B}_{\mathbf{r}^*}(\mathbf{v}_0)$ into the slices $\mathcal{B}_{\mathbf{r}+1}(\mathbf{v}_0) \setminus \mathcal{B}_{\mathbf{r}}(\mathbf{v}_0)$ and to apply Theorem 7.1 to each slice separately with a proper choice of the parameter λ .

Theorem 7.2. *Let \mathbf{r}^* be such that $(\mathcal{E}d)$ holds on $\mathcal{B}_{\mathbf{r}^*}(\mathbf{v}_0)$. Let also $\mathbb{Q}(\mathcal{Y}^\circ) \leq \mathbb{Q}$ for $\mathcal{Y}^\circ = \mathcal{B}_{\mathbf{r}}(\mathbf{v}_0)$ with $\mathbf{r} \leq \mathbf{r}^*$. If $\rho > 0$ and \mathfrak{z} are fixed to ensure $\sqrt{\rho\mathfrak{z}} \leq \mathfrak{g}_0 = \nu_0 \mathfrak{g}$ and $\rho(\mathfrak{z} - 1) \geq 2$, then it holds*

$$\begin{aligned} \log \mathbb{P} \left(\sup_{\mathbf{v} \in \mathcal{B}_{\mathbf{r}^*}(\mathbf{v}_0)} \left\{ \frac{1}{3\nu_0} \mathcal{U}(\mathbf{v}, \mathbf{v}_0) - \frac{\rho}{2} d^2(\mathbf{v}, \mathbf{v}_0) \right\} > \mathfrak{z} \right) \\ \leq -\rho(\mathfrak{z} - 1) + \log(4\mathfrak{z}) + \mathbb{Q}. \end{aligned} \quad (7.4)$$

Moreover, if $\sqrt{\rho\mathfrak{z}} > \mathfrak{g}_0$, then

$$\begin{aligned} \log \mathbb{P} \left(\sup_{\mathbf{v} \in \mathcal{B}_{\mathbf{r}^*}(\mathbf{v}_0)} \left\{ \frac{1}{3\nu_0} \mathcal{U}(\mathbf{v}, \mathbf{v}_0) - \frac{\rho}{2} d^2(\mathbf{v}, \mathbf{v}_0) \right\} > \mathfrak{z} \right) \\ \leq -\mathfrak{g}_0 \sqrt{\rho(\mathfrak{z} - 1)} + \mathfrak{g}_0^2/2 + \log(4\mathfrak{z}) + \mathbb{Q}. \end{aligned} \quad (7.5)$$

Remark 7.1. Formally the bound applies even with $\mathbf{r}^* = \infty$ provided that $(\mathcal{E}d)$ is fulfilled on the whole set \mathcal{Y}° .

Proof. Denote

$$\mathbf{u}(\mathbf{r}) \stackrel{\text{def}}{=} \frac{1}{3\nu_0 \mathbf{r}} \sup_{\mathbf{v} \in \mathcal{B}_{\mathbf{r}}(\mathbf{v}_0)} \{ \mathcal{U}(\mathbf{v}) - \mathcal{U}(\mathbf{v}_0) \}.$$

Then we have to bound the probability

$$\mathbb{P} \left(\sup_{\mathbf{r} \leq \mathbf{r}^*} \{ \mathbf{r} \mathbf{u}(\mathbf{r}) - \rho \mathbf{r}^2/2 \} > \mathfrak{z} \right).$$

For each $\mathbf{r} \leq \mathbf{r}^*$ and $\lambda \leq \mathfrak{g}_0$, it follows from (7.2) that

$$\log \mathbb{E} \exp \{ \lambda \mathbf{u}(\mathbf{r}) \} \leq \lambda^2/2 + \mathbb{Q}.$$

The choice $\lambda = \sqrt{\rho\mathfrak{z}}$ is admissible in view of $\sqrt{\rho\mathfrak{z}} \leq \mathfrak{g}_0$. This implies by the exponential Chebyshev inequality

$$\begin{aligned} \log \mathbb{P}(\mathbf{r} u(\mathbf{r}) - \rho \mathbf{r}^2/2 \geq \mathfrak{z}) &\leq -\lambda(\mathfrak{z}/\mathbf{r} + \rho \mathbf{r}/2) + \lambda^2/2 + \mathbb{Q} \\ &= -\rho\mathfrak{z}(x + x^{-1} - 1) + \mathbb{Q}, \end{aligned} \quad (7.6)$$

where $u = \sqrt{\rho/(2\mathfrak{z})} \mathbf{r}$. We now apply the slicing arguments w.r.t. $t = \rho \mathbf{r}^2/2 = \mathfrak{z}x^2$. By definition, $\mathbf{r}u(\mathbf{r})$ increases in \mathbf{r} . We use that for any growing function $f(\cdot)$ and any $t \geq 0$, it holds

$$f(t) - t \leq \int_t^{t+1} \{f(s) - s + 1\} ds$$

Therefore, for any $t > 0$, it holds by (7.6) in view of $dt = 2\mathfrak{z} x dx$

$$\begin{aligned} \mathbb{P}\left(\sup_{\mathbf{r} \leq \mathbf{r}^*} \{\mathbf{r} u(\mathbf{r}) - \rho \mathbf{r}^2/2\} > \mathfrak{z}\right) &\leq \int_0^{t^*+1} \mathbb{P}\{\mathbf{r} u(\mathbf{r}) - t \geq \mathfrak{z} - 1\} dt \\ &\leq 2\mathfrak{z} \int_0^{t^*+1} \exp\{-\rho(\mathfrak{z} - 1)(x + x^{-1} - 1) + \mathbb{Q}\} x dx \\ &\leq 2\mathfrak{z} e^{-b+\mathbb{Q}} \int_0^\infty \exp\{-b(x + x^{-1} - 2)\} x dx \end{aligned}$$

with $b = \rho(\mathfrak{z} - 1)$ and $t^* = \rho \mathbf{r}^{*2}/2$. This implies for $b \geq 2$

$$\begin{aligned} \mathbb{P}\left(\sup_{\mathbf{r} \leq \mathbf{r}^*} \{\mathbf{r} u(\mathbf{r}) - \rho \mathbf{r}^2/2\} > \mathfrak{z}\right) &\leq 2\mathfrak{z} e^{-b+\mathbb{Q}} \int_0^\infty \exp\{-2(x + x^{-1} - 2)\} x dx \\ &\leq 4\mathfrak{z} \exp\{-\rho(\mathfrak{z} - 1) + \mathbb{Q}\} \end{aligned}$$

and (7.4) follows.

If $\sqrt{\rho\mathfrak{z}} > \mathfrak{g}_0$, then select $\lambda = \mathfrak{g}_0$. For $\mathbf{r} \leq \mathbf{r}^*$

$$\begin{aligned} \log \mathbb{P}\{\mathbf{r} u(\mathbf{r}) - \rho \mathbf{r}^2/2 \geq \mathfrak{z}\} &= \log \mathbb{P}\{u(\mathbf{r}) > \mathfrak{z}/\mathbf{r} + \rho \mathbf{r}/2\} \\ &\leq -\lambda(\mathfrak{z}/\mathbf{r} + \rho \mathbf{r}/2) + \lambda^2/2 + \mathbb{Q} \\ &\leq -\lambda\sqrt{\rho\mathfrak{z}}(x + x^{-1} - 2)/2 - \lambda\sqrt{\rho\mathfrak{z}} + \lambda^2/2 + \mathbb{Q}, \end{aligned}$$

where $u = \sqrt{\rho/\mathfrak{z}} \mathbf{r}$. This allows to bound in the same way as above

$$\mathbb{P}\left(\sup_{\mathbf{r} \leq \mathbf{r}^*} \{\mathbf{r} u(\mathbf{r}) - \rho \mathbf{r}^2/2\} > \mathfrak{z}\right) \leq 4\mathfrak{z} \exp(-\lambda\sqrt{\rho(\mathfrak{z} - 1)} + \lambda^2/2 + \mathbb{Q})$$

yielding (7.5). □

This result can be used for describing the concentration bound for the maximum of $(3\nu_0)^{-1}\mathcal{U}(\mathbf{v}, \mathbf{v}_0) - \rho d^2(\mathbf{v}, \mathbf{v}_0)/2$. Namely, it suffices to find \mathfrak{z} ensuring the prescribed deviation probability. We state the result for a special case with $\rho = 1$ and $\mathbf{g}_0 \geq 3$ which simplifies the notation.

Corollary 7.3. *Under the conditions of Theorem 7.2, it holds for $\mathbf{x} \geq 0$ with $\mathbf{x} + \mathbb{Q} \geq 4$:*

$$P\left(\sup_{\mathbf{v} \in \mathcal{B}_{\mathbf{r}^*}(\mathbf{v}_0)} \left\{ \frac{1}{3\nu_0} \mathcal{U}(\mathbf{v}, \mathbf{v}_0) - \frac{1}{2} d^2(\mathbf{v}, \mathbf{v}_0) \right\} > \mathfrak{z}_0(\mathbf{x}, \mathbb{Q})\right) \leq \exp(-\mathbf{x}),$$

where with $\mathbf{g}_0 = \nu_0 \mathbf{g} \geq 2$

$$\mathfrak{z}_0(\mathbf{x}, \mathbb{Q}) \stackrel{\text{def}}{=} \begin{cases} (1 + \sqrt{\mathbf{x} + \mathbb{Q}})^2 & \text{if } 1 + \sqrt{\mathbf{x} + \mathbb{Q}} \leq \mathbf{g}_0, \\ 1 + (1 + 2\mathbf{g}_0^{-1})^2 (\mathbf{g}_0^{-1}(\mathbf{x} + \mathbb{Q}) + \mathbf{g}_0/2)^2 & \text{otherwise.} \end{cases} \quad (7.7)$$

Proof. In view of (7.4), it suffices to check that $\mathfrak{z} = (1 + \sqrt{\mathbf{x} + \mathbb{Q}})^2$ ensures

$$\mathfrak{z} - 1 - \log(4\mathfrak{z}) - \mathbb{Q} \geq \mathbf{x}.$$

This follows from the inequality

$$(1 + \mathbf{y})^2 - 1 - 2\log(2 + 2\mathbf{y}) \geq \mathbf{y}^2$$

with $\mathbf{y} = \sqrt{\mathbf{x} + \mathbb{Q}} \geq 2$. Similarly $\mathfrak{z} = 1 + (1 + 2\mathbf{g}_0^{-1})^2 \mathbf{y}^2$ for $\mathbf{y} = \mathbf{g}_0^{-1}(\mathbf{x} + \mathbb{Q}) + \mathbf{g}_0/2$ ensures

$$\mathbf{g}_0 \sqrt{\mathfrak{z} - 1} - \mathbf{g}_0^2/2 - \log(4\mathfrak{z}) - \mathbb{Q} \geq \mathbf{x}$$

in view of $\sqrt{\mathfrak{z} - 1} \geq \mathbf{y}(1 + 2\mathbf{g}_0^{-1})$, $4\mathfrak{z} \leq (1 + 2\mathbf{y} + 4\mathbf{g}_0^{-1}\mathbf{y})^2$, and

$$\mathbf{y} - \log(1 + 2\mathbf{y} + 4\mathbf{g}_0^{-1}\mathbf{y}) \geq 0$$

for $\mathbf{y} \geq 2$ and $\mathbf{g}_0 \geq 3$. □

If $\mathbf{g} \gg \sqrt{\mathbb{Q}}$ and \mathbf{x} is not too big then $\mathfrak{z}_0(\mathbf{x}, \mathbb{Q})$ is of order $\mathbf{x} + \mathbb{Q}$. So, the main message of this result is that with a high probability the maximum of $(3\nu_0)^{-1}\mathcal{U}(\mathbf{v}, \mathbf{v}_0) - d^2(\mathbf{v}, \mathbf{v}_0)/2$ does not significantly exceed the level \mathbb{Q} .

7.3 A global upper function and concentration sets

The result of the previous section can be explained as a local upper function for the process $\mathcal{U}(\cdot)$. Indeed, in a vicinity $\mathcal{B}_{\mathbf{r}^*}(\mathbf{v}_0)$ of the central point \mathbf{v}_0 , it holds $(3\nu_0)^{-1}\mathcal{U}(\mathbf{v}) \leq d^2(\mathbf{v}, \mathbf{v}_0)/2 + \mathfrak{z}$ with a probability exponentially small in \mathfrak{z} . However, an extension of

this result on the whole set \mathcal{T} is only possible under some quite restrictive conditions. This section presents one possible construction of an upper function for the process $\mathcal{U}(\cdot)$ on the complement of the local set $\mathcal{T}_\circ(\mathbf{r}^*)$. For simplifying the notations assume that $\mathcal{U}(\mathbf{v}_0) \equiv 0$. Then $\mathcal{U}(\mathbf{v}, \mathbf{v}_0) = \mathcal{U}(\mathbf{v})$. We say that $\mathbf{u}(\mu, \mathbf{v})$ is a *multiscale upper function* for $\mu\mathcal{U}(\cdot)$ on a subset \mathcal{T}° of \mathcal{T} if

$$\mathbb{P}\left(\sup_{\mu \in \mathbb{M}} \sup_{\mathbf{v} \in \mathcal{T}^\circ} \{\mu\mathcal{U}(\mathbf{v}) - \mathbf{u}(\mu, \mathbf{v})\} \geq \mathfrak{z}(\mathbf{x})\right) \leq e^{-\mathbf{x}},$$

for some fixed function $\mathfrak{z}(\mathbf{x})$. An upper function can be used for describing the concentration sets of the point of maximum $\tilde{\mathbf{v}} = \operatorname{argmax}_{\mathbf{v} \in \mathcal{T}^\circ} \mathcal{U}(\mathbf{v})$; see Section 7.3.1 below.

For constructing such an upper function, the following condition is used which extends condition (Ed):

(E) For any $\mathbf{v} \in \mathcal{T}$ there exists $\mu \in \mathbb{M}$ such that

$$\mathfrak{N}(\mu, \mathbf{v}) \stackrel{\text{def}}{=} \log \mathbb{E} \exp\{\mu\mathcal{U}(\mathbf{v})\} < \infty. \quad (7.8)$$

This condition can be used for building a simple pointwise upper function for $\mu\mathcal{U}(\mathbf{v})$. Indeed, (7.8) implies

$$\mathbb{E} \exp\{\mu\mathcal{U}(\mathbf{v}) - \mathfrak{N}(\mu, \mathbf{v})\} = 1. \quad (7.9)$$

The next step is in extending this pointwise result to a uniform one. The standard approach is based on the notion of a ϵ -net which is a discrete set $\mathcal{T}_\epsilon^\circ$ providing that for any point $\mathbf{v} \in \mathcal{T}$, there exists a point $\mathbf{v}^\circ \in \mathcal{T}_\epsilon^\circ$ with $d(\mathbf{v}, \mathbf{v}^\circ) \leq \epsilon$. The upper function is first constructed on this discrete set $\mathcal{T}_\epsilon^\circ$ using (7.9) by increasing the pointwise bound with $\log N_\epsilon$, where the covering number N_ϵ is the cardinality of $\mathcal{T}_\epsilon^\circ$. Then it is extended on the whole set \mathcal{T}° using stochastic continuity of the process $\mathcal{U}(\cdot)$.

We apply a slightly different construction usually called *pillling*. Let the value \mathbf{r}_1 be fixed. Let also a measure π on \mathcal{T} be fixed. By $\mathcal{B}_\mu(\mathbf{v})$ we denote the ball of radius \mathbf{r}_1/μ at $\mathbf{v} \in \mathcal{T}$, while $\pi_\mu(\mathbf{v})$ denotes its π -measure. In our results the value $1/\pi_\mu(\mathbf{v})$ replaces the covering number N_ϵ with $\epsilon = \mathbf{r}_1/\mu$. Also define the constant ν_1 describing the local variability of $\pi_1(\cdot)$:

$$\nu_1 \stackrel{\text{def}}{=} \sup_{\mu \in \mathbb{M}} \sup_{\mathbf{v}^\circ \in \mathcal{T}^\circ} \sup_{\mathbf{v} \in \mathcal{B}_\mu(\mathbf{v}^\circ)} \frac{\pi_\mu(\mathbf{v})}{\pi_\mu(\mathbf{v}^\circ)}. \quad (7.10)$$

For any fixed point \mathbf{v}° , the local maximum of the process $\mu\mathcal{U}$ over the ball $\mathcal{B}_\mu(\mathbf{v}^\circ)$ can be bounded by combining the pointwise result (7.9) and the result of Theorem 7.1 for local fluctuations of the process $\mu\mathcal{U}$ within $\mathcal{B}_\mu(\mathbf{v}^\circ)$. To get a global bound over \mathcal{T}° , we introduce the so called *penalty* function $\mathbf{t}_\mu(\mathbf{v})$ which accounts for the size of the set

Υ° . In the next result this function is allowed to be μ -dependent. However, in typical situations, one can apply a universal function $\mathbf{t}(\mathbf{v})$; cf. the smooth case in Section 7.4.

Remind the definition of the smoothing operator \mathbb{S}_μ :

$$\mathbb{S}_\mu f(\mathbf{v}^\circ) = \frac{1}{\pi_\mu(\mathbf{v}^\circ)} \int_{\mathcal{B}_\mu(\mathbf{v}^\circ)} f(\mathbf{v}) \pi(d\mathbf{v}).$$

The next result suggests a construction of the upper function $\mathbf{u}(\mu, \mathbf{v})$. The construction involves a constant s which can be selected as the smallest value ensuring the bound $3\nu_0 \mathbf{r}_1 / s \leq \mathbf{g} \wedge \sqrt{2\mathbb{Q}}$.

Theorem 7.4. *Let the process $\mathcal{U}(\cdot)$ fulfill (\mathcal{E}) and $(\mathcal{E}d)$ for all $\mathbf{v}, \mathbf{v}' \in \Upsilon^\circ$ with $d(\mathbf{v}, \mathbf{v}') \leq \mathbf{r}_1$. If s is such that $3\nu_0 \mathbf{r}_1 / s \leq \mathbf{g} \wedge \sqrt{2\mathbb{Q}}$, then it holds for any $\mathbf{x} > 0$*

$$\mathbb{P} \left(\sup_{\mu \in \mathbb{M}} \sup_{\mathbf{v} \in \Upsilon^\circ} \{ \mu \mathcal{U}(\mathbf{v}) - \mathbb{S}_\mu \mathfrak{N}(\mu, \mathbf{v}) - (1+s) \mathbb{S}_\mu \mathbf{t}_\mu(\mathbf{v}) \} \geq \mathfrak{z}_1(\mathbf{x}) \right) \leq 2e^{-\mathbf{x}}, \quad (7.11)$$

where $\mathfrak{z}_1(\mathbf{x})$ is a linear function in \mathbf{x} :

$$\mathfrak{z}_1(\mathbf{x}) \stackrel{\text{def}}{=} (1+s) \{ \mathbf{x} + \log(\nu_1 \mathcal{T}) \} + 2s\mathbb{Q},$$

and

$$\mathcal{T} \stackrel{\text{def}}{=} \sum_{\mu \in \mathbb{M}} \int_{\Upsilon^\circ} \exp\{-\mathbf{t}_\mu(\mathbf{v})\} \frac{\pi(d\mathbf{v})}{\pi_\mu(\mathbf{v})}. \quad (7.12)$$

Proof. We bound $\mu \mathcal{U}(\mathbf{v})$ in two steps: first we evaluate $\mathbb{S}_\mu [\mu \mathcal{U}(\mathbf{v}^\circ) - \mathfrak{N}(\mu, \mathbf{v}^\circ)]$ and then $\mu [\mathcal{U}(\mathbf{v}^\circ) - \mathbb{S}_\mu \mathcal{U}(\mathbf{v}^\circ)]$. Convexity of the exp-function implies by the Jensen inequality

$$\begin{aligned} & \exp\{ \mathbb{S}_\mu [\mu \mathcal{U}(\mathbf{v}^\circ) - \mathfrak{N}(\mu, \mathbf{v}^\circ) - \mathbf{t}_\mu(\mathbf{v}^\circ)] \} \\ & \leq \mathbb{S}_\mu \exp\{ \mu \mathcal{U}(\mathbf{v}^\circ) - \mathfrak{N}(\mu, \mathbf{v}^\circ) - \mathbf{t}_\mu(\mathbf{v}^\circ) \} \\ & \leq \int_{\mathcal{B}_\mu(\mathbf{v}^\circ)} \exp\{ \mu \mathcal{U}(\mathbf{v}) - \mathfrak{N}(\mu, \mathbf{v}) - \mathbf{t}_\mu(\mathbf{v}) \} \frac{\pi(d\mathbf{v})}{\pi_\mu(\mathbf{v}^\circ)} \\ & \leq \nu_1 \int_{\mathcal{B}_\mu(\mathbf{v}^\circ)} \exp\{ \mu \mathcal{U}(\mathbf{v}) - \mathfrak{N}(\mu, \mathbf{v}) - \mathbf{t}_\mu(\mathbf{v}) \} \frac{\pi(d\mathbf{v})}{\pi_\mu(\mathbf{v})} \\ & \leq \nu_1 \int_{\Upsilon^\circ} \exp\{ \mu \mathcal{U}(\mathbf{v}) - \mathfrak{N}(\mu, \mathbf{v}) - \mathbf{t}_\mu(\mathbf{v}) \} \frac{\pi(d\mathbf{v})}{\pi_\mu(\mathbf{v})}. \end{aligned}$$

As the right hand-side does not depend on \mathbf{v}° , the bound applies to the maximum of this expression over \mathbf{v}° . This implies in view of $\mathbb{E} \exp\{ \mu \mathcal{U}(\mathbf{v}) - \mathfrak{N}(\mu, \mathbf{v}) \} = 1$

$$\mathbb{E} \exp \sup_{\mathbf{v} \in \Upsilon^\circ} \{ \mathbb{S}_\mu [\mu \mathcal{U}(\mathbf{v}) - \mathfrak{N}(\mu, \mathbf{v}) - \mathbf{t}_\mu(\mathbf{v})] \} \leq \nu_1 \int_{\Upsilon^\circ} \exp\{-\mathbf{t}_\mu(\mathbf{v})\} \frac{\pi(d\mathbf{v})}{\pi_\mu(\mathbf{v})}.$$

As the sup over $\mu \in \mathbb{M}$ is not larger than the sum of the exponential terms, it holds

$$\mathbb{E} \exp \sup_{\mathbf{v} \in \mathcal{V}^\circ} \sup_{\mu \in \mathbb{M}} \{ \mathbb{S}_\mu [\mu \mathcal{U}(\mathbf{v}) - \mathfrak{N}(\mu, \mathbf{v}) - \mathfrak{t}_\mu(\mathbf{v})] \} \leq \nu_1 \mathcal{T}.$$

This bound implies for each $\mathbf{x} > 0$ with probability at least $1 - e^{-\mathbf{x}}$

$$\mathbb{S}_\mu [\mu \mathcal{U}(\mathbf{v}) - \mathfrak{N}(\mu, \mathbf{v}) - \mathfrak{t}_\mu(\mathbf{v})] \leq \mathbf{x} + \log(\nu_1 \mathcal{T}), \quad \mathbf{v} \in \mathcal{V}^\circ. \quad (7.13)$$

Now define

$$w_\mu(\mathbf{v}^\circ) \stackrel{\text{def}}{=} \sup_{\mathbf{v} \in \mathcal{B}_\mu(\mathbf{v}^\circ)} \mu |\mathcal{U}(\mathbf{v}) - \mathcal{U}(\mathbf{v}^\circ)|.$$

With $\lambda \stackrel{\text{def}}{=} (3\nu_0 \mathbf{r}_1/s) \wedge \sqrt{2\mathbb{Q}}$, it holds by Theorem 7.1 in view of $\lambda^2/2 \leq \mathbb{Q}$

$$\log \mathbb{E} \exp \{ s^{-1} w_\mu(\mathbf{v}) \} \leq 2\mathbb{Q}.$$

Then

$$\begin{aligned} & \exp \{ s^{-1} \mu [\mathcal{U}(\mathbf{v}^\circ) - \mathbb{S}_\mu \mathcal{U}(\mathbf{v}^\circ)] - \mathbb{S}_\mu \mathfrak{t}_\mu(\mathbf{v}) \} \\ &= \exp \int \{ s^{-1} \mu [\mathcal{U}(\mathbf{v}^\circ) - \mathcal{U}(\mathbf{v})] - \mathfrak{t}_\mu(\mathbf{v}) \} \frac{\pi(d\mathbf{v})}{\pi_\mu(\mathbf{v}^\circ)} \\ &\leq \int_{\mathcal{B}_\mu(\mathbf{v}^\circ)} \exp \{ s^{-1} w_\mu(\mathbf{v}) - \mathfrak{t}_\mu(\mathbf{v}) \} \frac{\pi(d\mathbf{v})}{\pi_\mu(\mathbf{v}^\circ)} \\ &\leq \nu_1 \int_{\mathcal{B}_\mu(\mathbf{v}^\circ)} \exp \{ s^{-1} w_\mu(\mathbf{v}) - \mathfrak{t}_\mu(\mathbf{v}) \} \frac{\pi(d\mathbf{v})}{\pi_\mu(\mathbf{v})} \\ &\leq \nu_1 \int_{\mathcal{V}^\circ} \exp \{ s^{-1} w_\mu(\mathbf{v}) - \mathfrak{t}_\mu(\mathbf{v}) \} \frac{\pi(d\mathbf{v})}{\pi_\mu(\mathbf{v})}. \end{aligned}$$

As the right hand-side does not depend on \mathbf{v}° , the bound applies to the maximum of this expression over \mathbf{v}° . This implies

$$\mathbb{E} \exp \sup_{\mathbf{v} \in \mathcal{V}^\circ} \sup_{\mu \in \mathbb{M}} \{ s^{-1} \mu [\mathcal{U}(\mathbf{v}^\circ) - \mathbb{S}_\mu \mathcal{U}(\mathbf{v}^\circ)] - \mathbb{S}_\mu \mathfrak{t}_\mu(\mathbf{v}) \} \leq \nu_1 \mathcal{T} \exp(2\mathbb{Q}).$$

This implies similarly to (7.13) with probability at least $1 - e^{-\mathbf{x}}$:

$$s^{-1} \mu [\mathcal{U}(\mathbf{v}) - \mathbb{S}_\mu \mathcal{U}(\mathbf{v})] - \mathbb{S}_\mu \mathfrak{t}_\mu(\mathbf{v}) \leq \mathbf{x} + 2\mathbb{Q} + \log(\nu_1 \mathcal{T}), \quad \mathbf{v} \in \mathcal{V}^\circ.$$

Combining these two bounds yields (7.11) with probability at least $1 - 2e^{-\mathbf{x}}$. \square

Remark 7.2. It is interesting to compare the uniform bound (7.11) of Theorem 7.4 and the pointwise bound (7.9): at which prise the pointwise result can be extended to a global one. The proposed construction involves two additional terms. One of them is

proportional to the local entropy \mathbb{Q} and it comes from the local bound of Theorem 7.1 as the price for taking the local supremum. The second term is proportional to $\mathfrak{t}_\mu(\mathbf{v}) + \log(\mathcal{T})$ with \mathcal{T} from (7.12) and it is responsible for extending the local maximum into the global one over \mathcal{T}° .

7.3.1 Hitting probability

Let $M(\mathbf{v})$ be a deterministic *boundary* function. We aim at bounding the probability that a process $\mathcal{U}(\mathbf{v})$ on \mathcal{T}° hits this boundary on the set \mathcal{T}° . This precisely means the probability that $\sup_{\mathbf{v} \in \mathcal{T}^\circ} \{\mathcal{U}(\mathbf{v}) - M(\mathbf{v})\} \geq 0$. A particularly interesting problem is to describe for each $\mathbf{x} > 0$ the value $\mathfrak{z}_1(\mathbf{x})$ ensuring that

$$\mathbb{P}\left(\sup_{\mathbf{v} \in \mathcal{T}^\circ} \{\mathcal{U}(\mathbf{v}) - M(\mathbf{v})\} \geq \mathfrak{z}_1(\mathbf{x})\right) \leq e^{-\mathbf{x}}.$$

Let $\mathfrak{u}(\mu, \mathbf{v})$ be the multiscale upper function for $\mathcal{U}(\mathbf{v})$:

$$\mathbb{P}\left\{\sup_{\mu \in \mathbb{M}} \sup_{\mathbf{v} \in \mathcal{T}^\circ} [\mu \mathcal{U}(\mathbf{v}) - \mathfrak{u}(\mu, \mathbf{v})] \geq \mathfrak{z}_1(\mathbf{x})\right\} \leq 2e^{-\mathbf{x}}; \quad (7.14)$$

cf. (7.11). Define

$$\begin{aligned} \mathcal{C}(\mu, \mathbf{v}) &\stackrel{\text{def}}{=} -\mathfrak{u}(\mu, \mathbf{v}) + \mu M(\mathbf{v}), \\ \mu^*(\mathbf{v}) &\stackrel{\text{def}}{=} \operatorname{argmax}_{\mu \in \mathbb{M}} \mathcal{C}(\mu, \mathbf{v}), \\ \mathcal{C}^*(\mathbf{v}) &\stackrel{\text{def}}{=} \max_{\mu \in \mathbb{M}} \mathcal{C}(\mu, \mathbf{v}) = \mathcal{C}(\mu^*(\mathbf{v}), \mathbf{v}). \end{aligned}$$

The studied hitting probability can be described via the value

$$\mathfrak{g}(\mathcal{T}^\circ) \stackrel{\text{def}}{=} \inf_{\mathbf{v} \in \mathcal{T}^\circ} \mathcal{C}^*(\mathbf{v}) = \inf_{\mathbf{v} \in \mathcal{T}^\circ} \max_{\mu \in \mathbb{M}} \{-\mathfrak{u}(\mu, \mathbf{v}) + \mu M(\mathbf{v})\}.$$

The larger this value is, the smaller is the bound for the hitting probability. More precisely, let a fixed \mathbf{x} and the corresponding $\mathfrak{z}_1(\mathbf{x})$ in (7.14) be fixed. If for each $\mathbf{v} \in \mathcal{T}^\circ$, the inequality $\mu M(\mathbf{v}) \geq \mathfrak{u}(\mu, \mathbf{v})$ holds with a properly selected $\mu = \mu(\mathbf{v})$, then the hitting probability is bounded by $2e^{-\mathbf{x}}$.

Theorem 7.5. *Suppose (7.14). If $\mathfrak{g}(\mathcal{T}^\circ) \geq \mathfrak{z}_1(\mathbf{x})$, then*

$$\mathbb{P}\left(\sup_{\mathbf{v} \in \mathcal{T}^\circ} \{\mathcal{U}(\mathbf{v}) - M(\mathbf{v})\} \geq 0\right) \leq 2e^{-\mathbf{x}}.$$

Proof. For each $\mathbf{v} \in \mathcal{T}^\circ$, in view of $\mathcal{C}^*(\mathbf{v}) \geq \mathfrak{g}(\mathcal{T}^\circ)$, it holds

$$\begin{aligned} \{\mathcal{U}(\mathbf{v}) - M(\mathbf{v}) \geq 0\} &= \{\mu^*(\mathbf{v})[\mathcal{U}(\mathbf{v}) - M(\mathbf{v})] \geq 0\} \\ &\subseteq \{\mu^*(\mathbf{v})[\mathcal{U}(\mathbf{v}) - M(\mathbf{v})] + \mathcal{C}^*(\mathbf{v}) \geq \mathfrak{g}(\mathcal{T}^\circ)\} \\ &= \{\mu^*(\mathbf{v})\mathcal{U}(\mathbf{v}) - \mathfrak{u}(\mu^*(\mathbf{v}), \mathbf{v}) \geq \mathfrak{z}_1(\mathbf{x})\} \\ &\subseteq \{\max_{\mu} [\mu \mathcal{U}(\mathbf{v}) - \mathfrak{u}(\mu, \mathbf{v})] \geq \mathfrak{z}_1(\mathbf{x})\}, \end{aligned}$$

and the result follows from (7.14). \square

7.4 Finite-dimensional smooth case

Here we discuss the special case when \mathcal{T} is an open subset in \mathbb{R}^p , the stochastic process $\mathcal{U}(\mathbf{v})$ is absolutely continuous and its gradient $\nabla \mathcal{U}(\mathbf{v}) \stackrel{\text{def}}{=} d\mathcal{U}(\mathbf{v})/d\mathbf{v}$ has bounded exponential moments.

(ED) *There exist $\mathfrak{g} > 0$, $\nu_0 \geq 1$, and for each $\mathbf{v} \in \mathcal{T}$, a symmetric non-negative matrix $H(\mathbf{v})$ such that for any $\lambda \leq \mathfrak{g}$ and any unit vector $\boldsymbol{\gamma} \in \mathbb{R}^p$, it holds*

$$\log \mathbb{E} \exp \left\{ \lambda \frac{\boldsymbol{\gamma}^\top \nabla \mathcal{U}(\mathbf{v})}{\|H(\mathbf{v})\boldsymbol{\gamma}\|} \right\} \leq \nu_0^2 \lambda^2 / 2.$$

A natural candidate for $H^2(\mathbf{v})$ is the covariance matrix $\text{Var}(\nabla \mathcal{U}(\mathbf{v}))$ provided that this matrix is well posed. Then the constant ν_0 can be taken close to one by reducing the value \mathfrak{g} ; see Lemma 7.12 below.

In what follows we fix a subset \mathcal{T}° of \mathcal{T} and establish a bound for the maximum of the process $\mathcal{U}(\mathbf{v}, \mathbf{v}^\circ) = \mathcal{U}(\mathbf{v}) - \mathcal{U}(\mathbf{v}^\circ)$ on \mathcal{T}° for a fixed point \mathbf{v}° . We will assume existence of a dominating matrix $H^* = H^*(\mathcal{T}^\circ)$ such that $H(\mathbf{v}) \preceq H^*$ for all $\mathbf{v} \in \mathcal{T}^\circ$. We also assume that π is the Lebesgue measure on \mathcal{T} . First we show that the differentiability condition (ED) implies (Ed).

Lemma 7.6. *Assume that (ED) holds with some \mathfrak{g} and $H(\mathbf{v}) \preceq H^*$ for $\mathbf{v} \in \mathcal{T}^\circ$. Consider any $\mathbf{v}, \mathbf{v}^\circ \in \mathcal{T}^\circ$. Then it holds for $|\lambda| \leq \mathfrak{g}$*

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

Proof. Denote $\delta = \|\mathbf{v} - \mathbf{v}^\circ\|$, $\boldsymbol{\gamma} = (\mathbf{v} - \mathbf{v}^\circ)/\delta$. Then

$$\mathcal{U}(\mathbf{v}, \mathbf{v}^\circ) = \delta \boldsymbol{\gamma}^\top \int_0^1 \nabla \mathcal{U}(\mathbf{v}^\circ + t\delta \boldsymbol{\gamma}) dt$$

and $\|H^*(\mathbf{v} - \mathbf{v}^\circ)\| = \delta\|H^*\boldsymbol{\gamma}\|$. Now the Hölder inequality and $(\mathcal{E}D)$ yield

$$\begin{aligned} & \mathbb{E} \exp \left\{ \lambda \frac{\mathcal{U}(\mathbf{v}, \mathbf{v}^\circ)}{\|H^*(\mathbf{v} - \mathbf{v}^\circ)\|} - \frac{\nu_0^2 \lambda^2}{2} \right\} \\ &= \mathbb{E} \exp \left\{ \int_0^1 \left[\lambda \frac{\boldsymbol{\gamma}^\top \nabla \mathcal{U}(\mathbf{v}^\circ + t\delta\boldsymbol{\gamma})}{\|H^*\boldsymbol{\gamma}\|} - \frac{\nu_0^2 \lambda^2}{2} \right] dt \right\} \\ &\leq \int_0^1 \mathbb{E} \exp \left\{ \lambda \frac{\boldsymbol{\gamma}^\top \nabla \mathcal{U}(\mathbf{v}^\circ + t\delta\boldsymbol{\gamma})}{\|H^*\boldsymbol{\gamma}\|} - \frac{\nu_0^2 \lambda^2}{2} \right\} dt \leq 1 \end{aligned}$$

as required. \square

The result of Lemma 7.6 enables us to define $d(\mathbf{v}, \mathbf{v}') = \|H^*(\mathbf{v} - \mathbf{v}^\circ)\|$ so that the corresponding ball coincides with the ellipsoid $B(\mathbf{r}, \mathbf{v}^\circ)$. Now we bound the value $\mathbb{Q}(\mathcal{Y}^\circ)$ for $\mathcal{Y}^\circ = B(\mathbf{r}_1, \mathbf{v}^\circ)$.

Lemma 7.7. *Let $\mathcal{Y}^\circ = B(\mathbf{r}_1, \mathbf{v}^\circ)$. Under the conditions of Lemma 7.6, it holds $\mathbb{Q}(\mathcal{Y}^\circ) \leq \mathbf{c}_1 p$, where $\mathbf{c}_1 = 2$ for $p \geq 2$, and $\mathbf{c}_1 = 2.4$ for $p = 1$.*

Proof. The set \mathcal{Y}° coincides with the ellipsoid $B(\mathbf{r}_1, \mathbf{v}^\circ)$ while the d -ball $\mathcal{B}_k(\mathbf{v})$ coincides with the ellipsoid $B(\mathbf{r}_k, \mathbf{v}^\circ)$ for each $k \geq 2$. By change of variables, the study can be reduced to the case with $\mathbf{v}^\circ = 0$, $H^* \equiv I_p$, $\mathbf{r}_1 = 1$, so that $B(\mathbf{r}, \mathbf{v})$ is the usual Euclidean ball in \mathbb{R}^p of radius \mathbf{r} . It is obvious that the measure of the overlap of two balls $B(1, 0)$ and $B(2^{-k+1}, \mathbf{v})$ for $\|\mathbf{v}\| \leq 1$ is minimized when $\|\mathbf{v}\| = 1$, and this value is the same for all such \mathbf{v} . Define the number $a_{k,p}$ by

$$a_{k,p}^p \stackrel{\text{def}}{=} \frac{\pi(B(1, 0))}{\pi(B(1, 0) \cap B(2^{-k+1}, \mathbf{v}))}$$

for any \mathbf{v} with $\|\mathbf{v}\| = 1$. It is easy to see that $a_{k,1} = 2^k$ and $a_{k,p}$ decreases with p . So, $M_k \leq 2^{pk}$ and

$$\begin{aligned} \mathbb{Q}(\mathcal{Y}^\circ) &\leq \frac{1}{3} \log(2^{1+p}) + \frac{4}{3} \sum_{k=2}^{\infty} 2^{-k} \log(2^{1+kp}) \\ &= \frac{\log 2}{3} \left[3 + p + 2p \sum_{k=1}^{\infty} (k+1) 2^{-k} \right] = (3 + 7p) \frac{\log 2}{3} \leq \mathbf{c}_1 p, \end{aligned}$$

where $\mathbf{c}_1 = 2$ for $p \geq 2$, and $\mathbf{c}_1 = 2.4$ for $p = 1$, and the result follows. \square

7.4.1 Local central bound

Here we specify the local bounds of Theorem 7.1 and the central result of Corollary 7.3 to the smooth case.

Theorem 7.8. *Suppose $(\mathcal{E}d)$. For any $\lambda \leq \nu_0 \mathbf{g}$, $\mathbf{r}_1 > 0$, and $\mathbf{v}^\circ \in \mathcal{T}$*

$$\log \mathbb{E} \exp \left\{ \frac{\lambda}{3\nu_0 \mathbf{r}_1} \sup_{\mathbf{v} \in B(\mathbf{r}_1, \mathbf{v}^\circ)} |\mathcal{U}(\mathbf{v}) - \mathcal{U}(\mathbf{v}^\circ)| \right\} \leq \lambda^2/2 + \mathbb{Q},$$

where $\mathbb{Q} = \mathbf{c}_1 p$.

We consider the local sets of the elliptic form $\mathcal{T}_\circ(\mathbf{r}) \stackrel{\text{def}}{=} \{\mathbf{v} : \|H_0(\mathbf{v} - \mathbf{v}_0)\| \leq \mathbf{r}\}$, where H_0 dominates $H(\mathbf{v})$ on this set: $H(\mathbf{v}) \preceq H_0$.

Theorem 7.9. *Let $(\mathcal{E}D)$ hold with some \mathbf{g} and a matrix $H(\mathbf{v})$. Suppose that $H(\mathbf{v}) \preceq H_0$ for all $\mathbf{v} \in \mathcal{T}_\circ(\mathbf{r})$. Then*

$$\mathbb{P} \left(\sup_{\mathbf{v} \in \mathcal{T}_\circ(\mathbf{r})} \left\{ \frac{1}{3\nu_0} \mathcal{U}(\mathbf{v}, \mathbf{v}_0) - \frac{1}{2} \|H_0(\mathbf{v} - \mathbf{v}_0)\|^2 \right\} \geq \mathfrak{z}_0(\mathbf{x}, p) \right) \leq \exp(-\mathbf{x}), \quad (7.15)$$

where $\mathfrak{z}_0(\mathbf{x}, p)$ coincides with $\mathfrak{z}_0(\mathbf{x}, \mathbb{Q})$ from (7.7) with $\mathbb{Q} = \mathbf{c}_1 p$.

Remark 7.3. An important feature of the established result is that the bound in the right hand-side of (7.15) does not depend on the value \mathbf{r} describing the radius of the local vicinity around the central point \mathbf{v}_0 . In the ideal case one would apply this result with $\mathbf{r} = \infty$ provided that the conditions $H(\mathbf{v}) \leq H_0$ is fulfilled uniformly over \mathcal{T} .

Proof. Lemma 7.7 implies $(\mathcal{E}d)$ with $d(\mathbf{v}, \mathbf{v}_0) = \|H_0(\mathbf{v} - \mathbf{v}_0)\|^2/2$. Now the result follows from Corollary 7.3. \square

7.4.2 A global upper function

Now we specify the general result of Theorem 7.4 to the smooth case. To make the formulation more transparent, the matrix $H(\mathbf{v})$ from condition $(\mathcal{E}D)$ is assumed to be uniformly bounded by a fixed matrix H^* . Let \mathbf{r} be fixed with $\mathbf{r}^2 \geq p/2$. We aim to build an upper function for the process $\mathcal{U}(\cdot)$ on the complement of the central set $\mathcal{T}_\circ(\mathbf{r}) \stackrel{\text{def}}{=} \{\mathbf{v} : \|H_0 \mathbf{v}\| \leq \mathbf{r}\}$. The penalty function $\mathbf{t}(\mathbf{v})$ is taken independent of μ as a logarithmic function of $\|H_0 \mathbf{v}\|$.

Theorem 7.10. *Assume (\mathcal{E}) and $(\mathcal{E}D)$ with $H(\mathbf{v}) \preceq H^*$ for all $\mathbf{v} \in \mathcal{T}$. Let $\mathbf{r}_1 \geq 1$ be such that $\det(\mathbf{r}_1^{-1} H^*) \leq \det(H_0)$, and s be such that $3\nu_0 \mathbf{r}_1/s \leq \mathbf{g} \wedge \sqrt{2\mathbf{c}_1 p}$. Given $\mathbf{r} \geq \sqrt{p/2}$, define*

$$\mathbf{t}(\mathbf{v}) \stackrel{\text{def}}{=} (p+2) \log(\|H_0 \mathbf{v}\|), \quad \mathbf{v} \in \mathcal{T}.$$

Then for any $\mathbf{x} > 0$, it holds with probability at least $1 - 2\text{Me}^{-\mathbf{x}}$ for $\mathbf{M} = \sum_{\mu \in \mathbb{M}} \mu^p$

$$\mu \mathcal{U}(\mathbf{v}) - \mathbb{S}_\mu \mathfrak{N}(\mu, \mathbf{v}) - (1+s)\mathbf{t}(\mathbf{v}) \leq (1+s)\mathbf{x} + 2s\mathbf{c}_1 p, \quad \mathbf{v} \in \mathcal{T} \setminus \mathcal{T}_\circ(\mathbf{r}). \quad (7.16)$$

Proof. For the measure Lebesgue measure π on \mathbb{R}^p , it holds for $\mathcal{B}_\mu(\mathbf{v}^\circ) = B(\mathbf{r}_1/\mu, \mathbf{v}^\circ)$:

$$\frac{1}{\pi(\mathcal{B}_\mu(\mathbf{v}^\circ))} = \frac{\det(\mathbf{r}_1^{-1}\mu H^*)}{\omega_p} \leq \frac{\mu^p \det(H_0)}{\omega_p},$$

where ω_p is the measure of the unit ball in \mathbb{R}^p . In particular, this measure does not depend on the location \mathbf{v}° and thus, $\nu_1 = 1$; see (7.10). The change of variables yields

$$\begin{aligned} \mathcal{T} &= \sum_{\mu \in \mathbb{M}} \int_{\mathcal{V}^\circ} \frac{1}{\pi(\mathcal{B}_\mu(\mathbf{v}))} \exp\{-\mathbf{t}(\mathbf{v})\} d\pi(\mathbf{v}) \\ &\leq \sum_{\mu \in \mathbb{M}} \mu^p \int_{\mathbb{R}^p} \frac{\det(H_0)}{\omega_p} \|H_0 \mathbf{v}\|^{-p-2} \mathbb{I}(\|H_0 \mathbf{v}\| \geq \mathbf{r}) d\pi(\mathbf{v}) \\ &= \frac{\mathbf{M}}{\omega_p} \int_{\|\mathbf{u}\| \geq \mathbf{r}} \|\mathbf{u}\|^{-p-2} d\mathbf{u} \leq \mathbf{M}p/(2\mathbf{r}^2) \leq \mathbf{M}. \end{aligned}$$

Symmetricity arguments imply $\mathbb{S}_\mu \|H_0 \mathbf{v}\| = \|H_0 \mathbf{v}\|$ and concavity of the log-function yields $\mathbb{S}_\mu \mathbf{t}(\mathbf{v}) \leq \mathbf{t}(\mathbf{v})$. Now the result (7.16) follows from Theorem 7.4 in view of $\mathbb{Q} \leq \mathbf{c}_1 p$. \square

7.5 Auxiliary facts

Lemma 7.11. *For any r.v.'s ξ_k and $\lambda_k \geq 0$ such that $\Lambda = \sum_k \lambda_k \leq 1$*

$$\log \mathbb{E} \exp\left(\sum_k \lambda_k \xi_k\right) \leq \sum_k \lambda_k \log \mathbb{E} e^{\xi_k}.$$

Proof. Convexity of e^x and concavity of x^Λ imply

$$\begin{aligned} \mathbb{E} \exp\left\{\frac{\Lambda}{\Lambda} \sum_k \lambda_k (\xi_k - \log \mathbb{E} e^{\xi_k})\right\} &\leq \mathbb{E}^\Lambda \exp\left\{\frac{1}{\Lambda} \sum_k \lambda_k (\xi_k - \log \mathbb{E} e^{\xi_k})\right\} \\ &\leq \left\{\frac{1}{\Lambda} \sum_k \lambda_k \mathbb{E} \exp(\xi_k - \log \mathbb{E} e^{\xi_k})\right\}^\Lambda = 1. \end{aligned}$$

\square

Lemma 7.12. *Let a r.v. ξ fulfill $\mathbb{E}\xi = 0$, $\mathbb{E}\xi^2 = 1$ and $\mathbb{E} \exp(\lambda_1 |\xi|) = \varkappa < \infty$ for some $\lambda_1 > 0$. Then for any $\varrho < 1$ there is a constant C_1 depending on \varkappa , λ_1 and ϱ only such that for $\lambda < \varrho \lambda_1$*

$$\log \mathbb{E} e^{\lambda \xi} \leq C_1 \lambda^2 / 2.$$

Moreover, there is a constant $\lambda_2 > 0$ such that for all $\lambda \leq \lambda_2$

$$\log \mathbb{E} e^{\lambda \xi} \geq \varrho \lambda^2 / 2.$$

Proof. Define $h(x) = (\lambda - \lambda_1)x + m \log(x)$ for $m \geq 0$ and $\lambda < \lambda_1$. It is easy to see by a simple algebra that

$$\max_{x \geq 0} h(x) = -m + m \log \frac{m}{\lambda_1 - \lambda}.$$

Therefore for any $x \geq 0$

$$\lambda x + m \log(x) \leq \lambda_1 x + \log \left(\frac{m}{e(\lambda_1 - \lambda)} \right)^m.$$

This implies for all $\lambda < \lambda_1$

$$\mathbb{E}|\xi|^m \exp(\lambda|\xi|) \leq \left(\frac{m}{e(\lambda_1 - \lambda)} \right)^m \mathbb{E} \exp(\lambda_1|\xi|).$$

Suppose now that for some $\lambda_1 > 0$, it holds $\mathbb{E} \exp(\lambda_1|\xi|) = \varkappa(\lambda_1) < \infty$. Then the function $h_0(\lambda) = \mathbb{E} \exp(\lambda\xi)$ fulfills $h_0(0) = 1$, $h'_0(0) = \mathbb{E}\xi = 0$, $h''_0(0) = 1$ and for $\lambda < \lambda_1$,

$$h''_0(\lambda) = \mathbb{E}\xi^2 e^{\lambda\xi} \leq \mathbb{E}\xi^2 e^{\lambda|\xi|} \leq \frac{1}{(\lambda_1 - \lambda)^2} \mathbb{E} \exp(\lambda_1|\xi|).$$

This implies by the Taylor expansion for $\lambda < \varrho\lambda_1$ that

$$h_0(\lambda) \leq 1 + C_1 \lambda^2 / 2$$

with $C_1 = \varkappa(\lambda_1) / \{\lambda_1^2(1 - \varrho)^2\}$, and hence, $\log h_0(\lambda) \leq C_1 \lambda^2 / 2$. □

References

- Bednorz, W. (2006). A theorem on majorizing measures. *Ann. Probab.*, 34(5):1771–1781.
- Birgé, L. (2006). Model selection via testing: an alternative to (penalized) maximum likelihood estimators. *Annales de l'Institut Henri Poincaré (B) Probability and Statistics*.
- Birgé, L. and Massart, P. (1993). Rates of convergence for minimum contrast estimators. *Probab. Theory Relat. Fields*, 97(1-2):113–150.
- Birgé, L. and Massart, P. (1998). Minimum contrast estimators on sieves: Exponential bounds and rates of convergence. *Bernoulli*, 4(3):329–375.
- Boucheron, S., Lugosi, G., and Massart, P. (2003). Concentration inequalities using the entropy method. *Ann. Probab.*, 31(3):1583–1614.

- Ibragimov, I. and Khas'minskij, R. (1981). *Statistical estimation. Asymptotic theory. Transl. from the Russian by Samuel Kotz.* New York - Heidelberg -Berlin: Springer-Verlag .
- Le Cam, L. (1960). Locally asymptotically normal families of distributions. Certain approximations to families of distributions and their use in the theory of estimation and testing hypotheses. *Univ. California Publ. Stat.*, 3:37–98.
- Le Cam, L. and Yang, G. L. (2000). *Asymptotics in Statistics: Some Basic Concepts.* Springer-Verlag, New York.
- McCullagh, P. and Nelder, J. (1989). *Generalized linear models. 2nd ed.* Monographs on Statistics and Applied Probability. 37. London etc.: Chapman and Hall. xix, 511 p. .
- Talagrand, M. (1996). Majorizing measures: The generic chaining. *Ann. Probab.*, 24(3):1049–1103.
- Talagrand, M. (2001). Majorizing measures without measures. *Ann. Probab.*, 29(1):411–417.
- Talagrand, M. (2005). *The generic chaining.* Springer Monographs in Mathematics. Springer-Verlag, Berlin. Upper and lower bounds of stochastic processes.
- Van de Geer, S. (1993). Hellinger-consistency of certain nonparametric maximum likelihood estimators. *Ann. Stat.*, 21(1):14–44.
- van der Vaart, A. and Wellner, J. A. (1996). *Weak convergence and empirical processes. With applications to statistics.* Springer Series in Statistics. New York, Springer.

SFB 649 Discussion Paper Series 2011

For a complete list of Discussion Papers published by the SFB 649, please visit <http://sfb649.wiwi.hu-berlin.de>.

- 001 "Localising temperature risk" by Wolfgang Karl Härdle, Brenda López Cabrera, Ostap Okhrin and Weining Wang, January 2011.
- 002 "A Confidence Corridor for Sparse Longitudinal Data Curves" by Shuzhuan Zheng, Lijian Yang and Wolfgang Karl Härdle, January 2011.
- 003 "Mean Volatility Regressions" by Lu Lin, Feng Li, Lixing Zhu and Wolfgang Karl Härdle, January 2011.
- 004 "A Confidence Corridor for Expectile Functions" by Esra Akdeniz Duran, Mengmeng Guo and Wolfgang Karl Härdle, January 2011.
- 005 "Local Quantile Regression" by Wolfgang Karl Härdle, Vladimir Spokoiny and Weining Wang, January 2011.
- 006 "Sticky Information and Determinacy" by Alexander Meyer-Gohde, January 2011.
- 007 "Mean-Variance Cointegration and the Expectations Hypothesis" by Till Strohsal and Enzo Weber, February 2011.
- 008 "Monetary Policy, Trend Inflation and Inflation Persistence" by Fang Yao, February 2011.
- 009 "Exclusion in the All-Pay Auction: An Experimental Investigation" by Dietmar Fehr and Julia Schmid, February 2011.
- 010 "Unwillingness to Pay for Privacy: A Field Experiment" by Alastair R. Beresford, Dorothea Kübler and Sören Preibusch, February 2011.
- 011 "Human Capital Formation on Skill-Specific Labor Markets" by Runli Xie, February 2011.
- 012 "A strategic mediator who is biased into the same direction as the expert can improve information transmission" by Lydia Mechtenberg and Johannes Münster, March 2011.
- 013 "Spatial Risk Premium on Weather Derivatives and Hedging Weather Exposure in Electricity" by Wolfgang Karl Härdle and Maria Osipenko, March 2011.
- 014 "Difference based Ridge and Liu type Estimators in Semiparametric Regression Models" by Esra Akdeniz Duran, Wolfgang Karl Härdle and Maria Osipenko, March 2011.
- 015 "Short-Term Herding of Institutional Traders: New Evidence from the German Stock Market" by Stephanie Kremer and Dieter Nautz, March 2011.
- 016 "Oracally Efficient Two-Step Estimation of Generalized Additive Model" by Rong Liu, Lijian Yang and Wolfgang Karl Härdle, March 2011.
- 017 "The Law of Attraction: Bilateral Search and Horizontal Heterogeneity" by Dirk Hofmann and Salmai Qari, March 2011.
- 018 "Can crop yield risk be globally diversified?" by Xiaoliang Liu, Wei Xu and Martin Odening, March 2011.
- 019 "What Drives the Relationship Between Inflation and Price Dispersion? Market Power vs. Price Rigidity" by Sascha Becker, March 2011.
- 020 "How Computational Statistics Became the Backbone of Modern Data Science" by James E. Gentle, Wolfgang Härdle and Yuichi Mori, May 2011.
- 021 "Customer Reactions in Out-of-Stock Situations – Do promotion-induced phantom positions alleviate the similarity substitution hypothesis?" by Jana Luisa Diels and Nicole Wiebach, May 2011.

SFB 649, Spandauer Str. 1, D-10178 Berlin
<http://sfb649.wiwi.hu-berlin.de>

This research was supported by the Deutsche
Forschungsgemeinschaft through the SFB 649 "Economic Risk".



SFB 649 Discussion Paper Series 2011

For a complete list of Discussion Papers published by the SFB 649, please visit <http://sfb649.wiwi.hu-berlin.de>.

- 022 "Extreme value models in a conditional duration intensity framework" by Rodrigo Herrera and Bernhard Schipp, May 2011.
- 023 "Forecasting Corporate Distress in the Asian and Pacific Region" by Russ Moro, Wolfgang Härdle, Saeideh Aliakbari and Linda Hoffmann, May 2011.
- 024 "Identifying the Effect of Temporal Work Flexibility on Parental Time with Children" by Juliane Scheffel, May 2011.
- 025 "How do Unusual Working Schedules Affect Social Life?" by Juliane Scheffel, May 2011.
- 026 "Compensation of Unusual Working Schedules" by Juliane Scheffel, May 2011.
- 027 "Estimation of the characteristics of a Lévy process observed at arbitrary frequency" by Johanna Kappus and Markus Reiß, May 2011.
- 028 "Asymptotic equivalence and sufficiency for volatility estimation under microstructure noise" by Markus Reiß, May 2011.
- 029 "Pointwise adaptive estimation for quantile regression" by Markus Reiß, Yves Rozenholc and Charles A. Cuenod, May 2011.
- 030 "Developing web-based tools for the teaching of statistics: Our Wikis and the German Wikipedia" by Sigbert Klinke, May 2011.
- 031 "What Explains the German Labor Market Miracle in the Great Recession?" by Michael C. Burda and Jennifer Hunt, June 2011.
- 032 "The information content of central bank interest rate projections: Evidence from New Zealand" by Gunda-Alexandra Detmers and Dieter Nautz, June 2011.
- 033 "Asymptotics of Asynchronicity" by Markus Bibinger, June 2011.
- 034 "An estimator for the quadratic covariation of asynchronously observed Itô processes with noise: Asymptotic distribution theory" by Markus Bibinger, June 2011.
- 035 "The economics of TARGET2 balances" by Ulrich Bindseil and Philipp Johann König, June 2011.
- 036 "An Indicator for National Systems of Innovation - Methodology and Application to 17 Industrialized Countries" by Heike Belitz, Marius Clemens, Christian von Hirschhausen, Jens Schmidt-Ehmcke, Axel Werwatz and Petra Zloczynski, June 2011.
- 037 "Neurobiology of value integration: When value impacts valuation" by Soyoung Q. Park, Thorsten Kahnt, Jörg Rieskamp and Hauke R. Heekeren, June 2011.
- 038 "The Neural Basis of Following Advice" by Guido Biele, Jörg Rieskamp, Lea K. Krugel and Hauke R. Heekeren, June 2011.
- 039 "The Persistence of "Bad" Precedents and the Need for Communication: A Coordination Experiment" by Dietmar Fehr, June 2011.
- 040 "News-driven Business Cycles in SVARs" by Patrick Bunk, July 2011.
- 041 "The Basel III framework for liquidity standards and monetary policy implementation" by Ulrich Bindseil and Jeroen Lamoot, July 2011.
- 042 "Pollution permits, Strategic Trading and Dynamic Technology Adoption" by Santiago Moreno-Bromberg and Luca Taschini, July 2011.
- 043 "CRRA Utility Maximization under Risk Constraints" by Santiago Moreno-Bromberg, Traian A. Pirvu and Anthony Réveillac, July 2011.

SFB 649, Spandauer Str. 1, D-10178 Berlin
<http://sfb649.wiwi.hu-berlin.de>

This research was supported by the Deutsche
Forschungsgemeinschaft through the SFB 649 "Economic Risk".



SFB 649 Discussion Paper Series 2011

For a complete list of Discussion Papers published by the SFB 649, please visit <http://sfb649.wiwi.hu-berlin.de>.

- 044 "Predicting Bid-Ask Spreads Using Long Memory Autoregressive Conditional Poisson Models" by Axel Groß-Klußmann and Nikolaus Hautsch, July 2011.
- 045 "Bayesian Networks and Sex-related Homicides" by Stephan Stahlschmidt, Helmut Tausendteufel and Wolfgang K. Härdle, July 2011.
- 046 "The Regulation of Interdependent Markets", by Raffaele Fiocco and Carlo Scarpa, July 2011.
- 047 "Bargaining and Collusion in a Regulatory Model", by Raffaele Fiocco and Mario Gilli, July 2011.
- 048 "Large Vector Auto Regressions", by Song Song and Peter J. Bickel, August 2011.
- 049 "Monetary Policy, Determinacy, and the Natural Rate Hypothesis", by Alexander Meyer-Gohde, August 2011.
- 050 "The impact of context and promotion on consumer responses and preferences in out-of-stock situations", by Nicole Wiebach and Jana L. Diels, August 2011.
- 051 "A Network Model of Financial System Resilience", by Kartik Anand, Prasanna Gai, Sujit Kapadia, Simon Brennan and Matthew Willison, August 2011.
- 052 "Rollover risk, network structure and systemic financial crises", by Kartik Anand, Prasanna Gai and Matteo Marsili, August 2011.
- 053 "When to Cross the Spread: Curve Following with Singular Control" by Felix Naujokat and Ulrich Horst, August 2011.
- 054 "TVICA - Time Varying Independent Component Analysis and Its Application to Financial Data" by Ray-Bing Chen, Ying Chen and Wolfgang K. Härdle, August 2011.
- 055 "Pricing Chinese rain: a multi-site multi-period equilibrium pricing model for rainfall derivatives" by Wolfgang K. Härdle and Maria Osipenko, August 2011.
- 056 "Limit Order Flow, Market Impact and Optimal Order Sizes: Evidence from NASDAQ TotalView-ITCH Data" by Nikolaus Hautsch and Ruihong Huang, August 2011.
- 057 "Optimal Display of Iceberg Orders" by Gökhan Cebiroğlu and Ulrich Horst, August 2011.
- 058 "Optimal liquidation in dark pools" by Peter Kratz and Torsten Schöneborn, September 2011.
- 059 "The Merit of High-Frequency Data in Portfolio Allocation" by Nikolaus Hautsch, Lada M. Kyj and Peter Malec, September 2011.
- 060 "On the Continuation of the Great Moderation: New evidence from G7 Countries" by Wenjuan Chen, September 2011.
- 061 "Forward-backward systems for expected utility maximization" by Ulrich Horst, Ying Hu, Peter Imkeller, Anthony Réveillac and Jianing Zhang.
- 062 "On heterogeneous latent class models with applications to the analysis of rating scores" by Aurélie Bertrand and Christian M. Hafner, October 2011.
- 063 "Multivariate Volatility Modeling of Electricity Futures" by Luc Bauwens, Christian Hafner and Diane Pierret, October 2011.

SFB 649, Spandauer Str. 1, D-10178 Berlin
<http://sfb649.wiwi.hu-berlin.de>

This research was supported by the Deutsche
Forschungsgemeinschaft through the SFB 649 "Economic Risk".



SFB 649 Discussion Paper Series 2011

For a complete list of Discussion Papers published by the SFB 649, please visit <http://sfb649.wiwi.hu-berlin.de>.

- 064 "Semiparametric Estimation with Generated Covariates" by Enno Mammen, Christoph Rothe and Melanie Schienle, October 2011.
- 065 "Linking corporate reputation and shareholder value using the publication of reputation rankings" by Sven Tischer and Lutz Hildebrandt, October 2011.
- 066 "Monitoring, Information Technology and the Labor Share" by Dorothee Schneider, October 2011.
- 067 "Minimal Supersolutions of BSDEs with Lower Semicontinuous Generators" by Gregor Heyne, Michael Kupper and Christoph Mainberger, October 2011.
- 068 "Bargaining, Openness, and the Labor Share" by Dorothee Schneider, October 2011.
- 069 "The Labor Share: A Review of Theory and Evidence" by Dorothee Schneider, October 2011.
- 070 "The Power of Sunspots: An Experimental Analysis" by Dietmar Fehr, Frank Heinemann and Aniol Llorente-Saguer, October 2011.
- 071 "Econometric analysis of volatile art markets" by Fabian Y. R. P. Bocart and Christian M. Hafner, October 2011.
- 072 "Financial Network Systemic Risk Contributions" by Nikolaus Hautsch, Julia Schaumburg and Melanie Schienle, October 2011.
- 073 "Calibration of self-decomposable Lévy models" by Mathias Trabs, November 2011.
- 074 "Time-Varying Occupational Contents: An Additional Link between Occupational Task Profiles and Individual Wages" by Alexandra Fedorets, November 2011.
- 075 "Changes in Occupational Demand Structure and their Impact on Individual Wages" by Alexandra Fedorets, November 2011.
- 076 "Nonparametric Nonstationary Regression with Many Covariates" by Melanie Schienle, November 2011.
- 077 "Increasing Weather Risk: Fact or Fiction?" by Weining Wang, Ihtiyor Bobojonov, Wolfgang Karl Härdle and Martin Odening, November 2011.
- 078 "Spatially Adaptive Density Estimation by Localised Haar Projections" by Florian Gach, Richard Nickl and Vladimir Spokoiny, November 2011.
- 079 "Martingale approach in pricing and hedging European options under regime-switching" by Grigori N. Milstein and Vladimir Spokoiny, November 2011.
- 080 "Sparse Non Gaussian Component Analysis by Semidefinite Programming" by Elmar Diederichs, Anatoli Juditsky, Arkadi Nemirovski and Vladimir Spokoiny, November 2011.
- 081 "Parametric estimation. Finite sample theory" by Vladimir Spokoiny, November 2011.

SFB 649, Spandauer Str. 1, D-10178 Berlin
<http://sfb649.wiwi.hu-berlin.de>

This research was supported by the Deutsche
Forschungsgemeinschaft through the SFB 649 "Economic Risk".

