

## PARAMETRIC ESTIMATION. FINITE SAMPLE THEORY

BY VLADIMIR SPOKOINY<sup>1</sup>

*Weierstrass-Institute, Humboldt University Berlin and Moscow Institute of  
Physics and Technology*

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 as follows: (1) the study is nonasymptotic, 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. These two features enable to bridge the gap between parametric and nonparametric theory and to build a unified framework for statistical estimation. The main results include large deviation bounds for the (quasi) maximum likelihood and the local quadratic bracketing 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 nonclassical 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 nonasymptotic statements. At the same time, the new bracketing device works well in the situations with large or growing parameter dimension in which the classical parametric theory fails. The general results are illustrated for the i.i.d. setup 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.

**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 ( $\mathbb{P}_\theta, \theta \in \Theta \subseteq \mathbb{R}^p$ ), where  $p$  stands for the number of parameters. This assumption allows to reduce the problem of statistical inference about  $\mathbb{P}$  to recovering the parameter  $\theta$ . The theory of parameter estimation and inference is nicely developed in a quite general setup. There is a vast literature on this issue. We only mention the book by [Ibragimov and Khas'minskiĭ \(1981\)](#), 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

---

Received November 2011; revised August 2012.

<sup>1</sup>Supported by Predictive Modeling Laboratory, MIPT, RF government grant, ag. 11.G34.31.0073. *MSC2010 subject classifications.* Primary 62F10; secondary 62J12, 62F25, 62H12.

*Key words and phrases.* Maximum likelihood, local quadratic bracketing, deficiency, concentration.

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 limit the scope of applicability for the whole approach. Indeed, the PA is usually only an approximation of real data distribution and in most 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 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'minskiĭ \(1981\)](#)]: first one localizes the problem to a neighborhood of the target parameter. Then one uses a local quadratic expansion of the log-likelihood to solve the corresponding estimation problem. There is, however, one feature of the proposed approach which makes it essentially different from the classical scheme. Namely, the use of the bracketing device instead of classical Taylor expansion allows to consider much larger local neighborhoods than in the LAN theory. More specifically, the classical LAN theory effectively requires a strict localization to a root- $n$  vicinity of the true point. At this point, the LAN theory fails in extending to the nonparametric situation. Our approach works for any local vicinity of the true point. This opens the door to building a unified theory including most of the classical parametric and nonparametric results.

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  $\mathbb{P}_\theta$  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  $\mu_0$  is any dominating measure for the family  $(\mathbb{P}_\theta)$ . We focus on the properties of the process  $L(\mathbf{Y}, \theta)$  as a function of the parameter  $\theta$ . 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  $\mathbf{Y}$ . By  $L(\theta, \theta^*) \stackrel{\text{def}}{=} L(\theta) - L(\theta^*)$  we denote the log-likelihood ratio. The classical likelihood principle suggests to estimate  $\theta$  by maximizing the corresponding log-

likelihood function  $L(\boldsymbol{\theta})$ :

$$(1.1) \quad \tilde{\boldsymbol{\theta}} \stackrel{\text{def}}{=} \arg \max_{\boldsymbol{\theta} \in \Theta} L(\boldsymbol{\theta}).$$

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

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

$$(1.2) \quad \boldsymbol{\theta}^* \stackrel{\text{def}}{=} \arg \max_{\boldsymbol{\theta} \in \Theta} \mathbb{E}L(\boldsymbol{\theta}),$$

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.

The results below show that the main properties of the quasi-MLE  $\tilde{\boldsymbol{\theta}}$  like concentration or coverage probability can be described in terms of the *excess* which is the difference between the maximum of the process  $L(\boldsymbol{\theta})$  and its value at the “true” point  $\boldsymbol{\theta}^*$ :

$$L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) \stackrel{\text{def}}{=} L(\tilde{\boldsymbol{\theta}}) - L(\boldsymbol{\theta}^*) = \max_{\boldsymbol{\theta} \in \Theta} L(\boldsymbol{\theta}) - L(\boldsymbol{\theta}^*).$$

The established results can be split into two big groups. A large deviation bound states some concentration properties of the estimator  $\tilde{\boldsymbol{\theta}}$ . For specific local sets  $\Theta_0(\varkappa)$  with elliptic shape, the deviation probability  $\mathbb{P}(\tilde{\boldsymbol{\theta}} \notin \Theta_0(\varkappa))$  is exponentially small in  $\varkappa$ . This concentration bound allows to restrict the parameter space to a properly selected vicinity  $\Theta_0(\varkappa)$ . Our main results concern the local properties of the process  $L(\boldsymbol{\theta})$  within  $\Theta_0(\varkappa)$  including a bracketing bound and its corollaries.

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-log-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  $\Theta_0(\varkappa)$ . The idea is to sandwich the underlying (quasi) log-likelihood process  $L(\boldsymbol{\theta})$  for  $\boldsymbol{\theta} \in \Theta_0(\varkappa)$  between two quadratic (in parameter) expressions. Then the maximum of  $L(\boldsymbol{\theta})$  over  $\Theta_0(\varkappa)$  will be sandwiched as well by the maxima of the lower and upper processes. The quadratic structure of these processes helps 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 estimator  $\tilde{\boldsymbol{\theta}}$ , polynomial risk

bounds, etc. In contrary to the classical theory, all the results are nonasymptotic 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  $\tilde{\theta}$  deviates from the local set  $\Theta_0(x)$ . In the modern statistical literature there are a number of studies considering maximum likelihood or, more generally, minimum contrast estimators in a general i.i.d. situation, when the parameter set  $\Theta$  is a subset of some functional space. We mention the papers of van de Geer (1993), Birgé and Massart (1993, 1998), Birgé (2006) and the references therein. The established results are based on deep probabilistic facts from empirical process theory; see, for example, Talagrand (1996, 2001, 2005), van der Vaart and Wellner (1996) and Boucheron, Lugosi and Massart (2003). The general result presented in Section 2 of the supplement [Spokoiny (2012b)] 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 per 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 nonasymptotic 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 2 of the supplement [Spokoiny (2012b)] presents an exponential bound for a general quadratic form which is very important 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 1 of the supplement [Spokoiny (2012b)] explains how the 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.

The proposed approach can be useful in many further research directions including penalized maximum likelihood and semiparametric estimation [Andresen and Spokoiny (2012)], contraction rate and asymptotic normality of the posterior within the Bayes approach [Spokoiny (2012a)] and local adaptive quantile estimation [Spokoiny, Wang and Härdle (2012)].

**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. Once checked, one obtains automatically all the established results. 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. It is important to stress that all the conditions only concern the properties of the quasi-likelihood process  $L(\boldsymbol{\theta})$ . Even if the process  $L(\cdot)$  is not a sufficient statistic, the whole analysis is entirely based on its geometric structure and probabilistic properties. The conditions are not restrictive and can be effectively checked in many particular situations. Some examples are given in Section 5 for i.i.d. setup, generalized linear models and for median regression.

The imposed conditions 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  $\boldsymbol{\theta}^*$ . Below we suppose that degree of locality is described by a number  $r$ . The local zone corresponds to  $r \leq r_0$  for a fixed  $r_0$ . The global conditions concern  $r > 0$ .

2.1. *Local conditions.* Local conditions describe the properties of  $L(\boldsymbol{\theta})$  in a vicinity of the central point  $\boldsymbol{\theta}^*$  from (1.2).

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

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

Below we suppose that the random function  $\zeta(\boldsymbol{\theta})$  is differentiable in  $\boldsymbol{\theta}$  and its gradient  $\nabla\zeta(\boldsymbol{\theta}) = \partial\zeta(\boldsymbol{\theta})/\partial\boldsymbol{\theta} \in \mathbb{R}^p$  has some exponential moments. Our first condition describes the property of the gradient  $\nabla\zeta(\boldsymbol{\theta}^*)$  at the central point  $\boldsymbol{\theta}^*$ .

(ED<sub>0</sub>) *There exist a positive symmetric matrix  $V_0^2$  and constants  $g > 0$ ,  $v_0 \geq 1$  such that  $\text{Var}\{\nabla\zeta(\boldsymbol{\theta}^*)\} \leq V_0^2$  and for all  $|\lambda| \leq g$*

$$\sup_{\boldsymbol{y} \in \mathbb{R}^p} \log \mathbb{E} \exp \left\{ \lambda \frac{\boldsymbol{y}^\top \nabla\zeta(\boldsymbol{\theta}^*)}{\|V_0 \boldsymbol{y}\|} \right\} \leq v_0^2 \lambda^2 / 2.$$

In a 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  $\theta^*$ . In particular, define the local elliptic neighborhoods of  $\theta^*$  as

$$(2.1) \quad \Theta_0(r) \stackrel{\text{def}}{=} \{\theta \in \Theta : \|V_0(\theta - \theta^*)\| \leq r\}.$$

The further conditions are restricted to such defined neighborhoods  $\Theta_0(r)$ .

(ED<sub>1</sub>) For each  $r \leq r_0$ , there exists a constant  $\omega(r) \leq 1/2$  such that it holds for all  $\theta \in \Theta_0(r)$

$$\sup_{\gamma \in \mathbb{R}^p} \log \mathbb{E} \exp \left\{ \lambda \frac{\gamma^\top \{\nabla \zeta(\theta) - \nabla \zeta(\theta^*)\}}{\omega(r) \|V_0 \gamma\|} \right\} \leq v_0^2 \lambda^2 / 2, \quad |\lambda| \leq g.$$

Here the constant  $g$  is the same as in (ED<sub>0</sub>).

The main bracketing result also requires second order smoothness of the expected log-likelihood  $\mathbb{E}L(\theta)$ . By definition,  $L(\theta^*, \theta^*) \equiv 0$  and  $\nabla \mathbb{E}L(\theta^*) = 0$  because  $\theta^*$  is the extreme point of  $\mathbb{E}L(\theta)$ . Therefore,  $-\mathbb{E}L(\theta, \theta^*)$  can be approximated by a quadratic function of  $\theta - \theta^*$  in the neighborhood of  $\theta^*$ . The *local identifiability* condition quantifies this quadratic approximation from above and from below on the set  $\Theta_0(r)$  from (2.1).

(L<sub>0</sub>) There is a symmetric strictly positive-definite matrix  $D_0^2$  and for each  $r \leq r_0$  and a constant  $\delta(r) \leq 1/2$ , such that it holds on the set  $\Theta_0(r) = \{\theta : \|V_0(\theta - \theta^*)\| \leq r\}$ ,

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

Usually  $D_0^2$  is defined as the negative Hessian of  $\mathbb{E}L(\theta^*)$ :  $D_0^2 = -\nabla^2 \mathbb{E}L(\theta^*)$ . If  $L(\theta, \theta^*)$  is the log-likelihood ratio and  $\mathbb{P} = \mathbb{P}_{\theta^*}$ , then  $-\mathbb{E}L(\theta, \theta^*) = \mathbb{E}_{\theta^*} \log(d\mathbb{P}_{\theta^*}/d\mathbb{P}_\theta) = \mathcal{K}(\mathbb{P}_{\theta^*}, \mathbb{P}_\theta)$ , the Kullback–Leibler divergence between  $\mathbb{P}_{\theta^*}$  and  $\mathbb{P}_\theta$ . Then condition (L<sub>0</sub>) with  $D_0 = V_0$  follows from the usual regularity conditions on the family  $(\mathbb{P}_\theta)$ ; cf. Ibragimov and Khas'minskiĭ (1981). If the log-likelihood process  $L(\theta)$  is sufficiently smooth in  $\theta$ , for example, three times stochastically differentiable, then the quantities  $\omega(r)$  and  $\delta(r)$  can be taken proportional to the value  $\varrho(r)$  defined as

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

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

The *identifiability condition* relates the matrices  $D_0^2$  and  $V_0^2$ .

(I) There is a constant  $\alpha > 0$  such that  $\alpha^2 D_0^2 \geq V_0^2$ .

2.2. *Global conditions.* The global conditions have to be fulfilled for all  $\theta$  lying beyond  $\Theta_0(r_0)$ . We only impose one condition on the smoothness of the stochastic component of the process  $L(\theta)$  in term of its gradient and one identifiability condition in terms of the expectation  $\mathbb{E}L(\theta, \theta^*)$ .

The first condition is similar to the local condition  $(ED_0)$  and it requires some exponential moment of the gradient  $\nabla\zeta(\theta)$  for all  $\theta \in \Theta$ . However, the constant  $g$  may be dependent of the radius  $r = \|V_0(\theta - \theta^*)\|$ .

$(Er)$  For any  $r$ , there exists a value  $g(r) > 0$  such that for all  $\lambda \leq g(r)$

$$\sup_{\theta \in \Theta_0(r)} \sup_{\gamma \in \mathbb{R}^p} \log \mathbb{E} \exp \left\{ \lambda \frac{\gamma^\top \nabla \zeta(\theta)}{\|V_0 \gamma\|} \right\} \leq v_0^2 \lambda^2 / 2.$$

The global identification property means that the deterministic component  $\mathbb{E}L(\theta, \theta^*)$  of the log-likelihood is competitive with its variance  $\text{Var} L(\theta, \theta^*)$ .

$(Lr)$  There is a function  $b(r)$  such that  $rb(r)$  monotonously increases in  $r$  and for each  $r \geq r_0$

$$\inf_{\theta : \|V_0(\theta - \theta^*)\| = r} |\mathbb{E}L(\theta, \theta^*)| \geq b(r)r^2.$$

**3. Local inference.** The *local asymptotic normality* (LAN) condition since introduced 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'minskiĭ (1981) for a nice presentation of the LAN theory including asymptotic efficiency of MLE and Bayes estimators. The LAN property was extended to *mixed LAN* or *local asymptotic quadraticity* (LAQ); see, for example, Le Cam and Yang (2000). All these notions are very much asymptotic and very much local. The LAN theory also requires that  $L(\theta)$  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 for free: 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 nonasymptotic 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  $L(\theta)$  can be approximated by a quadratic in  $\theta$  expression in a vicinity of the central point  $\theta^*$ . 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 nonasymptotic version of the LAN property in which the local quadratic approximation of  $L(\theta)$  is replaced by bounding this process from above and from below by two different quadratic in  $\theta$  processes. More precisely, we apply the *bracketing* idea: the difference  $L(\theta, \theta^*) = L(\theta) - L(\theta^*)$  is put between two quadratic processes  $\mathbb{L}_{\underline{\epsilon}}(\theta, \theta^*)$  and  $\mathbb{L}_{\epsilon}(\theta, \theta^*)$ :

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

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

The sandwiching result (3.1) naturally leads to two important notions: the value of the problem and the spread. It turns out that most of the statements like confidence and concentration probability rely upon the maximum of  $L(\theta, \theta^*)$  over  $\theta$  which we call *the excess*. Its expectation will be referred to as *the value of the problem*. Due to (3.1), the excess can be bounded from above and from below using the similar quantities  $\max_{\theta} \mathbb{L}_{\underline{\epsilon}}(\theta, \theta^*)$  and  $\max_{\theta} \mathbb{L}_{\epsilon}(\theta, \theta^*)$  which can be called *the lower and upper excess*, while their expectations are *the values of the lower and upper problems*. Note that  $\max_{\theta} \{\mathbb{L}_{\epsilon}(\theta, \theta^*) - \mathbb{L}_{\underline{\epsilon}}(\theta, \theta^*)\}$  can be very large or even infinite. However, this is not crucial. What really matters is the difference between the upper and the lower excess. The *spread*  $\Delta_{\epsilon}$  can be defined as the width of the interval bounding the excess due to (3.1), that is, as the sum of the approximation errors and of the difference between the upper and the lower excess:

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

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 spread.” The further sections explain in detail 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 nonasymptotic bounds.

3.1. *Local quadratic bracketing.* 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  $r$  of the local neighborhood  $\Theta_0(r)$  be fixed in a way that the deviation probability  $\mathbb{P}(\tilde{\theta} \notin \Theta_0(r))$  is sufficiently small. Precise results about the

choice of  $r$  which ensures this property are postponed until Section 4. In this neighborhood  $\Theta_0(r)$  we aim at building some quadratic lower and upper bounds for the process  $L(\theta)$ . The first step is the usual decomposition of this process into deterministic and stochastic components:

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

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

$$(3.2) \quad L(\theta, \theta^*) \approx \mathbb{L}(\theta, \theta^*) \stackrel{\text{def}}{=} (\theta - \theta^*)^\top \nabla \zeta(\theta^*) - \|D_0(\theta - \theta^*)\|^2/2.$$

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

Let  $\delta, \varrho$  be nonnegative constants. Introduce for a vector  $\epsilon = (\delta, \varrho)$  the following notation:

$$(3.3) \quad \begin{aligned} \mathbb{L}_\epsilon(\theta, \theta^*) &\stackrel{\text{def}}{=} (\theta - \theta^*)^\top \nabla L(\theta^*) - \|D_\epsilon(\theta - \theta^*)\|^2/2 \\ &= \xi_\epsilon^\top D_\epsilon(\theta - \theta^*) - \|D_\epsilon(\theta - \theta^*)\|^2/2, \end{aligned}$$

where  $\nabla L(\theta^*) = \nabla \zeta(\theta^*)$  by  $\nabla \mathbb{E}L(\theta^*) = 0$  and

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

Here we implicitly assume that with the proposed choice of the constants  $\delta$  and  $\varrho$ , the matrix  $D_\epsilon^2$  is nonnegative:  $D_\epsilon^2 \geq 0$ . The representation (3.3) indicates that the process  $\mathbb{L}_\epsilon(\theta, \theta^*)$  has the geometric structure of log-likelihood of a linear Gaussian model. We do not require that the vector  $\xi_\epsilon$  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  $\delta, \varrho$  are positive, the quadratic drift component of the process  $\mathbb{L}_\epsilon(\theta, \theta^*)$  is shrunk relative to  $\mathbb{L}(\theta, \theta^*)$  in (3.2) for  $\epsilon$  positive and it is stretched if  $\delta, \varrho$  are negative. Now, given  $r$ , fix some  $\delta \geq \delta(r)$  and  $\varrho \geq 3\nu_0\omega(r)$  with the value  $\delta(r)$  from condition  $(\mathcal{L}_0)$  and  $\omega(r)$  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  $x$  the values  $\varrho \geq 3\nu_0\omega(x)$  and  $\delta \geq \delta(x)$  be such that  $D_0^2(1 - \delta) - \varrho V_0^2 \geq 0$ . Then*

$$(3.4) \quad \mathbb{L}_\epsilon(\boldsymbol{\theta}, \boldsymbol{\theta}^*) - \diamond_\epsilon(x) \leq L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \leq \mathbb{L}_\epsilon(\boldsymbol{\theta}, \boldsymbol{\theta}^*) + \diamond_\epsilon(x), \quad \boldsymbol{\theta} \in \Theta_0(x),$$

with  $\mathbb{L}_\epsilon(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ ,  $\mathbb{L}_\epsilon(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$  defined by (3.3). The error terms  $\diamond_\epsilon(x)$  and  $\diamond_\epsilon(x)$  satisfy the bound (3.11) from Proposition 3.7.

The proof of this theorem is given in Proposition 3.7.

**REMARK 3.1.** This bracketing bound (3.4) describes some properties of the log-likelihood process and the estimator  $\tilde{\boldsymbol{\theta}}$  is not shown there. However, it directly implies most of our inference results. We therefore formulate (3.4) as a separate statement. Section 3.3 below presents some exponential bounds on the error terms  $\diamond_\epsilon(x)$  and  $\diamond_\epsilon(x)$ . The main message is that under rather broad conditions, these errors are small and have only minor impact on the inference for the quasi-MLE  $\tilde{\boldsymbol{\theta}}$ .

**3.2. Local inference.** 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 *bracketing result* (3.4) and the *geometric structure* of the processes  $\mathbb{L}_\epsilon$  and  $\mathbb{L}_\epsilon$ . Define the *spread*  $\Delta_\epsilon(x)$  by

$$(3.5) \quad \Delta_\epsilon(x) \stackrel{\text{def}}{=} \diamond_\epsilon(x) + \diamond_\epsilon(x) + (\|\boldsymbol{\xi}_\epsilon\|^2 - \|\boldsymbol{\xi}_\epsilon\|^2)/2.$$

Here  $\boldsymbol{\xi}_\epsilon = D_\epsilon^{-1}\nabla L(\boldsymbol{\theta}^*)$  and  $\boldsymbol{\xi}_\epsilon = D_\epsilon^{-1}\nabla L(\boldsymbol{\theta}^*)$ . The quantity  $\Delta_\epsilon(x)$  appears to be the price induced by our bracketing device. Section 3.3 below presents some probabilistic bounds on the spread showing that it is small relative to the other terms. All our corollaries below are stated under conditions of Theorem 3.1 and implicitly assume that the spread can be nearly ignored.

**3.2.1. Local coverage probability.** Our first result describes the probability of covering  $\boldsymbol{\theta}^*$  by the random set

$$(3.6) \quad \mathcal{E}(z) = \{\boldsymbol{\theta} : 2L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}) \leq z\}.$$

**COROLLARY 3.2.** *For any  $z > 0$*

$$(3.7) \quad \mathbb{P}\{\mathcal{E}(z) \not\ni \boldsymbol{\theta}^*, \tilde{\boldsymbol{\theta}} \in \Theta_0(x)\} \leq \mathbb{P}\{\|\boldsymbol{\xi}_\epsilon\|^2 \geq z - \diamond_\epsilon(x)\}.$$

**PROOF.** The bound (3.7) follows from the upper bound of Theorem 3.1 and the statement (3.12) of Lemma 3.8 below.  $\square$

Below [see (3.14)] we also present an exponential bound which helps to answer a very important question about a proper choice of the critical value  $z$  ensuring a prescribed covering probability.

3.2.2. *Local expansion, Wilks theorem and local concentration.* Now we show how the bound (3.4) can be used for obtaining a local expansion of the quasi-MLE  $\tilde{\theta}$ . All our results will be conditioned to the random set  $C_\epsilon(x)$  defined as

$$(3.8) \quad C_\epsilon(x) \stackrel{\text{def}}{=} \{\tilde{\theta} \in \Theta_0(x), \|V_0 D_\epsilon^{-1} \xi_\epsilon\| \leq x\}.$$

The second inequality in the definition of  $C_\epsilon(x)$  is related to the solution of the upper and lower problems (cf. Lemma 3.8):  $\|V_0 D_\epsilon^{-1} \xi_\epsilon\| \leq x$  means  $\tilde{\theta}_\epsilon \notin \Theta_0(x)$ , where  $\tilde{\theta}_\epsilon = \arg \min_{\theta} \mathbb{L}_\epsilon(\theta, \theta^*)$ .

Below in Section 3.3 we present some upper bounds on the value  $x$  ensuring a dominating probability of this random set. The first result can be viewed as a finite sample version of the famous Wilks theorem.

COROLLARY 3.3. *On the random set  $C_\epsilon(x)$  from (3.8), it holds*

$$(3.9) \quad \|\xi_\epsilon\|^2/2 - \diamond_\epsilon(x) \leq L(\tilde{\theta}, \theta^*) \leq \|\xi_\epsilon\|^2/2 + \diamond_\epsilon(x).$$

The next result is an extension of another prominent asymptotic result, namely, the Fisher expansion of the MLE.

COROLLARY 3.4. *On the random set  $C_\epsilon(x)$  from (3.8), it holds*

$$(3.10) \quad \|D_\epsilon(\tilde{\theta} - \theta^*) - \xi_\epsilon\|^2 \leq 2\Delta_\epsilon(x).$$

The proof of Corollaries 3.3 and 3.4 relies on the solution of the upper and lower problems and it is given below at the end of this section.

Now we describe concentration properties of  $\tilde{\theta}$  assuming that  $\tilde{\theta}$  is restricted to  $\Theta_0(x)$ . More precisely, we bound the probability that  $\|D_\epsilon(\tilde{\theta} - \theta^*)\| > z$  for a given  $z > 0$ .

COROLLARY 3.5. *For any  $z > 0$ , it holds*

$$\mathbb{P}\{\|D_\epsilon(\tilde{\theta} - \theta^*)\| > z, C_\epsilon(x)\} \leq \mathbb{P}\{\|\xi_\epsilon\| > z - \sqrt{2\Delta_\epsilon(x)}\}.$$

An interesting and important question is for which  $z$  in (3.6) the coverage probability of the event  $\{\mathcal{E}(z) \ni \theta^*\}$  or for which  $z$  the concentration probability of the event  $\{\|D_\epsilon(\tilde{\theta} - \theta^*)\| \leq z\}$  becomes close to one. It will be addressed in Section 3.3.

3.2.3. *A local risk bound.* Below we also bound the moments of the excess  $L(\tilde{\theta}, \theta^*)$  and of the normalized loss  $D_\epsilon(\tilde{\theta} - \theta^*)$  when  $\tilde{\theta}$  is restricted to  $\Theta_0(x)$ . The result follows directly from Corollaries 3.3 and 3.4.

COROLLARY 3.6. *For  $u > 0$*

$$\mathbb{E}\{L^u(\tilde{\theta}, \theta^*) \mathbb{1}(\tilde{\theta} \in \Theta_0(x))\} \leq \mathbb{E}[\{\|\xi_\epsilon\|^2/2 + \diamond_\epsilon(x)\}^u].$$

Moreover, it holds

$$\mathbb{E}\{\|D_\epsilon(\tilde{\theta} - \theta^*)\|^u \mathbb{1}(C_\epsilon(x))\} \leq \mathbb{E}[\{\|\xi_\epsilon\| + \sqrt{2\Delta_\epsilon(x)}\}^u].$$

3.2.4. *Comparing with the asymptotic theory.* This section briefly discusses the relation between the established nonasymptotic bounds and the classical asymptotic results in parametric estimation. This comparison is not straightforward because the asymptotic theory involves the sample size or noise level as the asymptotic parameter, while our setup is very general and works even for a “single” observation. Here we simply treat  $\epsilon = (\delta, \varrho)$  as a small parameter. This is well justified by the i.i.d. case with  $n$  observations, where it holds  $\delta = \delta(x) \asymp \sqrt{x/n}$  and similarly for  $\varrho$ ; see Section 5 for more details. The bounds below in Section 3.3 show that the spread  $\Delta_\epsilon(x)$  from (3.5) is small and can be ignored in the asymptotic calculations. The results of Corollary 3.2 through 3.6 represent the desired bounds in terms of deviation bounds for the quadratic form  $\|\xi_\epsilon\|^2$ .

For better understanding the essence of the presented results, consider first the “true” parametric model with the correctly specified log-likelihood  $L(\theta)$ . Then  $D_0^2 = V_0^2$  is the total Fisher information matrix. In the i.i.d. case it becomes  $n\mathbf{f}_0$  where  $\mathbf{f}_0$  is the usual Fisher information matrix of the considered parametric family at  $\theta^*$ . In particular,  $\text{Var}\{\nabla L(\theta^*)\} = n\mathbf{f}_0$ . So, if  $D_\epsilon$  is close to  $D_0$ , then  $\xi_\epsilon$  can be treated as the normalized score. Under usual assumptions,  $\xi \stackrel{\text{def}}{=} D_0^{-1}\nabla L(\theta^*)$  is the asymptotically standard normal  $p$ -vector. The same applies to  $\xi_\epsilon$ . Now one can observe that Corollaries 3.2 through 3.6 directly imply most of the classical asymptotic statements. In particular, Corollary 3.3 shows that the twice excess  $2L(\tilde{\theta}, \theta^*)$  is nearly  $\|\xi_\epsilon\|^2$  and thus nearly  $\chi_p^2$  (Wilks’ theorem). Corollary 3.4 yields the expansion  $D_\epsilon(\tilde{\theta} - \theta^*) \approx \xi_\epsilon$  (the Fisher expansion) and, hence,  $D_\epsilon(\tilde{\theta} - \theta^*)$  is asymptotically standard normal. Asymptotic variance of  $D_\epsilon(\tilde{\theta} - \theta^*)$  is nearly one, so  $\tilde{\theta}$  achieves the Cramér–Rao efficiency bound in the asymptotic setup.

3.3. *Spread.* This section presents some bounds on the spread  $\Delta_\epsilon(x)$  from (3.5). This quantity is random but it can be easily evaluated under the conditions made. We present two different results: one bounds the errors  $\diamond_\epsilon(x)$ ,  $\diamond_\epsilon(x)$ , while the other presents a deviation bound on quadratic forms like  $\|\xi_\epsilon\|^2$ . The results are stated under conditions  $(ED_0)$  and  $(ED_1)$  in a nonasymptotic way, so the formulation is quite technical. An informal discussion at the end of this section explains the typical behavior of the spread. The first result accomplishes the bracketing bound (3.4).

PROPOSITION 3.7. *Assume  $(ED_1)$ . The error  $\diamond_\epsilon(x)$  in (3.4) fulfills*

$$(3.11) \quad \mathbb{P}\{\varrho^{-1}\diamond_\epsilon(x) \geq \mathfrak{z}_0(x, \mathbb{Q})\} \leq \exp(-x)$$

with  $\mathfrak{z}_0(x, \mathbb{Q})$  given for  $\mathfrak{g}_0 = \mathfrak{g}v_0 \geq 3$  by

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

where  $\mathbb{Q} = c_1 p$  with  $c_1 = 2$  for  $p \geq 2$  and  $c_1 = 2.7$  for  $p = 1$ . Similarly for  $\diamond_\epsilon(x)$ .

REMARK 3.2. The bound (3.11) essentially depends on the value  $\mathfrak{g}$  from condition  $(ED_1)$ . The result requires that  $\mathfrak{g}v_0 \geq 3$ . However, this constant can usually be taken of order  $n^{1/2}$ ; see Section 5 for examples. If  $\mathfrak{g}^2$  is larger in order than  $p + x$ , then  $\mathfrak{J}_0(x, \mathbb{Q}) \approx c_1 p + x$ .

PROOF. Consider for fixed  $x$  and  $\epsilon = (\delta, \varrho)$  the quantity

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

As  $\delta \geq \delta(x)$ , it holds  $-\mathbb{E}L(\theta, \theta^*) \geq (1 - \delta)D_0^2$  and  $L(\theta, \theta^*) - \mathbb{E}L(\theta, \theta^*) \leq \diamond_{\epsilon}(x)$ . Moreover, in view of  $\nabla \mathbb{E}L(\theta^*) = 0$ , the definition of  $\diamond_{\epsilon}(x)$  can be rewritten as

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

Now the claim of the theorem can be easily reduced to an exponential bound for the quantity  $\diamond_{\epsilon}(x)$ . We apply Theorem 2.11 of the supplement [Spokoiny (2012b)] to the process

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

and  $H_0 = V_0$ . Condition  $(\mathcal{E}D)$  follows from  $(ED_1)$  with the same  $v_0$  and  $\mathfrak{g}$  in view of  $\nabla \mathcal{U}(\theta, \theta^*) = \{ \nabla \zeta(\theta) - \nabla \zeta(\theta^*) \} / \omega(x)$ . So, the conditions of Theorem 2.11 in the supplement [Spokoiny (2012b)] are fulfilled, yielding (3.11) in view of  $\varrho \geq 3v_0\omega(x)$ .  $\square$

Due to the main bracketing result, the local excess  $\sup_{\theta \in \Theta_0(x)} L(\theta, \theta^*)$  can be put between similar quantities for the upper and lower approximating processes up to the error terms  $\diamond_{\epsilon}(x), \underline{\diamond}_{\epsilon}(x)$ . The random quantity  $\sup_{\theta \in \mathbb{R}^p} \mathbb{L}_{\epsilon}(\theta, \theta^*)$  can be called the *upper excess* while  $\sup_{\theta \in \Theta_0(x_0)} \underline{\mathbb{L}}_{\epsilon}(\theta, \theta^*)$  is the *lower excess*. The quadratic (in  $\theta$ ) structure of the functions  $\mathbb{L}_{\epsilon}(\theta, \theta^*)$  and  $\underline{\mathbb{L}}_{\epsilon}(\theta, \theta^*)$  enables us to explicitly solve the problem of maximizing the corresponding function w.r.t.  $\theta$ .

LEMMA 3.8. *It holds*

$$(3.12) \quad \sup_{\theta \in \mathbb{R}^p} \mathbb{L}_{\epsilon}(\theta, \theta^*) = \|\xi_{\epsilon}\|^2 / 2.$$

On the random set  $\{ \|V_0 D_{\epsilon}^{-1} \xi_{\epsilon}\| \leq x \}$ , it also holds

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

PROOF. The unconstrained maximum of the quadratic form  $\mathbb{L}_\epsilon(\theta, \theta^*)$  w.r.t.  $\theta$  is attained at  $\tilde{\theta}_\epsilon = D_\epsilon^{-1}\xi_\epsilon = D_\epsilon^{-2}\nabla L(\theta^*)$ , yielding the expression (3.12). The lower excess is computed similarly.  $\square$

Our next step is in bounding the difference  $\|\xi_\epsilon\|^2 - \|\xi_\epsilon\|^2$ . It can be decomposed as

$$\|\xi_\epsilon\|^2 - \|\xi_\epsilon\|^2 = \|\xi_\epsilon\|^2 - \|\xi\|^2 + \|\xi\|^2 - \|\xi_\epsilon\|^2$$

with  $\xi = D_0^{-1}\nabla L(\theta^*)$ . If the values  $\delta, \varrho$  are small, then the difference  $\|\xi_\epsilon\|^2 - \|\xi_\epsilon\|^2$  is automatically smaller than  $\|\xi\|^2$ .

LEMMA 3.9. Suppose (I) and let  $\tau_\epsilon \stackrel{\text{def}}{=} \delta + \varrho\alpha^2 < 1$ . Then

$$(3.13) \quad \begin{aligned} D_\epsilon^2 &\geq (1 - \tau_\epsilon)D_0^2, & D_\epsilon^2 &\leq (1 + \tau_\epsilon)D_0^2, \\ \|\mathbf{I}_p - D_\epsilon D_\epsilon^{-2} D_\epsilon\|_\infty &\leq \alpha_\epsilon \stackrel{\text{def}}{=} \frac{2\tau_\epsilon}{1 - \tau_\epsilon^2}. \end{aligned}$$

Moreover,

$$\begin{aligned} \|\xi_\epsilon\|^2 - \|\xi\|^2 &\leq \frac{\tau_\epsilon}{1 - \tau_\epsilon} \|\xi\|^2, & \|\xi\|^2 - \|\xi_\epsilon\|^2 &\leq \frac{\tau_\epsilon}{1 + \tau_\epsilon} \|\xi\|^2, \\ \|\xi_\epsilon\|^2 - \|\xi_\epsilon\|^2 &\leq \alpha_\epsilon \|\xi\|^2. \end{aligned}$$

Our final step is in showing that under  $(ED_0)$ , the norm  $\|\xi\|$  behaves essentially as a norm of a Gaussian vector with the same covariance matrix. Define for  $\mathbb{B} \stackrel{\text{def}}{=} D_0^{-1}V_0^2D_0^{-1}$

$$\mathfrak{p} \stackrel{\text{def}}{=} \text{tr}(\mathbb{B}), \quad \mathfrak{v}^2 \stackrel{\text{def}}{=} 2 \text{tr}(\mathbb{B}^2), \quad \lambda_0 \stackrel{\text{def}}{=} \|\mathbb{B}\|_\infty = \lambda_{\max}(\mathbb{B}).$$

Under the identifiability condition (I), one can bound

$$\mathbb{B}^2 \leq \alpha^2 \mathbf{I}_p, \quad \mathfrak{p} \leq \alpha^2 p, \quad \mathfrak{v}^2 \leq 2\alpha^4 p, \quad \lambda_0 \leq \alpha^2.$$

Similarly to the previous result, we assume that the constant  $g$  from condition  $(ED_0)$  is sufficiently large, namely,  $g^2 \geq 2\mathfrak{p}$ . Define  $\mu_c = 2/3$  and

$$\begin{aligned} Y_c^2 &\stackrel{\text{def}}{=} g^2/\mu_c^2 - \mathfrak{p}/\mu_c, \\ g_c &\stackrel{\text{def}}{=} \mu_c Y_c = \sqrt{g^2 - \mu_c \mathfrak{p}}, \\ 2x_c &\stackrel{\text{def}}{=} \mu_c Y_c^2 + \log \det(\mathbf{I}_p - \mu_c \mathbb{B}^2/\lambda_0). \end{aligned}$$

It is easy to see that  $Y_c^2 \geq 3g^2/2$  and  $g_c \geq \sqrt{2/3}g$ .

**THEOREM 3.10.** *Let  $(ED_0)$  hold with  $v_0 = 1$  and  $\sigma^2 \geq 2p$ . Then  $\mathbb{E}\|\xi\|^2 \leq p$ , and for each  $x \leq x_c$*

$$(3.14) \quad \mathbb{P}(\|\xi\|^2/\lambda_0 \geq \mathfrak{z}(x, \mathbb{B})) \leq 2e^{-x} + 8.4e^{-x_c},$$

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

$$\mathfrak{z}(x, \mathbb{B}) \stackrel{\text{def}}{=} \begin{cases} p + 2\sqrt{vx}^{1/2}, & x \leq v/18, \\ p + 6x, & v/18 < x \leq x_c. \end{cases}$$

Moreover, for  $x > x_c$ , it holds with  $\mathfrak{z}(x, \mathbb{B}) = |y_c + 2(x - x_c)/g_c|^2$

$$\mathbb{P}(\|\xi\|^2/\lambda_0 \geq \mathfrak{z}(x, \mathbb{B})) \leq 8.4e^{-x}.$$

**PROOF.** It follows from condition  $(ED_0)$  that

$$\begin{aligned} \mathbb{E}\|\xi\|^2 &= \mathbb{E} \operatorname{tr} \xi \xi^\top \\ &= \operatorname{tr} D_0^{-1} [\mathbb{E} \nabla L(\theta^*) \{ \nabla L(\theta^*) \}^\top] D_0^{-1} = \operatorname{tr} [D_0^{-2} \operatorname{Var} \{ \nabla L(\theta^*) \}] \end{aligned}$$

and  $(ED_0)$  implies  $\boldsymbol{y}^\top \operatorname{Var} \{ \nabla L(\theta^*) \} \boldsymbol{y} \leq \boldsymbol{y}^\top V_0^2 \boldsymbol{y}$  and, thus,  $\mathbb{E}\|\xi\|^2 \leq p$ . The deviation bound (3.14) is proved in Corollary 2.5 of the supplement [Spokoiny (2012b)].  $\square$

**REMARK 3.3.** This small remark concerns the term  $8.4e^{-x_c}$  in the probability bound (3.14). As already mentioned, this bound implicitly assumes that the constant  $g$  is large (usually  $g \asymp n^{1/2}$ ). Then  $x_c \asymp \sigma^2 \asymp n$  is large as well. So,  $e^{-x_c}$  is very small and asymptotically negligible. Below we often ignore this term. For  $x \leq x_c$ , we can use  $\mathfrak{z}(x, \mathbb{B}) = p + 6x$ .

**REMARK 3.4.** The exponential bound of Theorem 3.10 helps to describe the critical value of  $\mathfrak{z}$  ensuring a prescribed deviation probability  $\mathbb{P}(\|\xi\|^2 \geq \mathfrak{z})$ . Namely, this probability starts to gradually decrease when  $\mathfrak{z}$  grows over  $\lambda_0 p$ . In particular, this helps to answer a very important question about a proper choice of the critical value  $\mathfrak{z}$  providing the prescribed covering probability, or of the value  $z$  ensuring the dominating concentration probability  $\mathbb{P}(\|D_\epsilon(\tilde{\theta} - \theta^*)\| \leq z)$ .

The definition of the set  $C_\epsilon(r)$  from (3.8) involves the event  $\{\|V_0 D_\epsilon^{-1} \xi_\epsilon\| > r\}$ . Under  $(\mathcal{I})$ , it is included in the set  $\{\|\xi_\epsilon\| > (1 + \alpha_\epsilon)^{-1} a^{-1} r\}$  [see (3.13)], and its probability is of order  $e^{-x}$  for  $r^2 \geq C(x + p)$  with a fixed  $C > 0$ .

By Theorem 3.7, one can use  $\max\{\diamond_\epsilon(r), \diamond_\epsilon(r)\} \leq \varrho \mathfrak{z}_0(x, \mathbb{Q})$  on a set of probability at least  $1 - e^{-x}$ . Further,  $\|\xi\|^2/\lambda_0 \leq \mathfrak{z}(x, \mathbb{B})$  with a probability of order  $1 - e^{-x}$ ; see (3.14). Putting together the obtained bounds yields for the spread  $\Delta_\epsilon(r)$  with a probability about  $1 - 4e^{-x}$

$$\Delta_\epsilon(r) \leq 2\varrho \mathfrak{z}_0(x, \mathbb{Q}) + \alpha_\epsilon \lambda_0 \mathfrak{z}(x, \mathbb{B}).$$

The results obtained in Section 3.2 are sharp and meaningful if the spread  $\Delta_\epsilon(r)$  is smaller in order than the value  $\mathbb{E}\|\xi\|^2$ . Theorem 3.10 states that  $\|\xi\|^2$  does not significantly deviate over its expected value  $\underline{p} \stackrel{\text{def}}{=} \mathbb{E}\|\xi\|^2$  which is our leading term. We know that  $\mathfrak{z}_0(x, \mathbb{Q}) \approx \mathbb{Q} + x = c_1 p + x$  if  $x$  is not too large. Also,  $\mathfrak{z}(x, \mathbb{B}) \leq p + 6x$ , where  $p$  is of order  $p$  due to  $(\mathcal{I})$ . Summarizing the above discussion yields that the local results apply if the regularity condition  $(\mathcal{I})$  holds and the values  $\varrho$  and  $\alpha_\epsilon$  or, equivalently,  $\omega(r), \delta(r)$  are small. In Section 5 we show for the i.i.d. example that  $\omega(r) \asymp \sqrt{r^2/n}$  and similarly for  $\delta(r)$ .

3.4. *Proof of Corollaries 3.3 and 3.4.* The bound (3.4) together with Lemma 3.8 yield on  $C_\epsilon(r)$

$$\begin{aligned}
 L(\tilde{\theta}, \theta^*) &= \sup_{\theta \in \Theta_0(r)} L(\theta, \theta^*) \\
 &\geq \sup_{\theta \in \Theta_0(r)} \mathbb{L}_\epsilon(\theta, \theta^*) - \diamond_\epsilon(r) = \|\xi_\epsilon\|^2/2 - \diamond_\epsilon(r).
 \end{aligned}
 \tag{3.15}$$

Similarly,

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

yielding (3.9). For getting (3.10), we again apply the inequality  $L(\theta, \theta^*) \leq \mathbb{L}_\epsilon(\theta, \theta^*) + \diamond_\epsilon(r)$  from Theorem 3.1 for  $\theta$  equal to  $\tilde{\theta}$ . With  $\xi_\epsilon = D_\epsilon^{-1} \nabla L(\theta^*)$  and  $\mathbf{u}_\epsilon \stackrel{\text{def}}{=} D_\epsilon(\tilde{\theta} - \theta^*)$ , this gives

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

Therefore, by (3.15),

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

or, equivalently,

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

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

**4. Upper function approach and concentration of the qMLE.** A very important step in the analysis of the qMLE  $\tilde{\theta}$  is *localization*. This property means that  $\tilde{\theta}$  concentrates in a small vicinity of the central point  $\theta^*$ . This section states such a concentration bound under the global conditions of Section 2. Given  $r_0$ , the deviation bound describes the probability  $\mathbb{P}(\tilde{\theta} \notin \Theta_0(r_0))$  that  $\tilde{\theta}$  does not belong to the local vicinity  $\Theta_0(r_0)$  of  $\theta$ . The question of interest is to check a possibility of selecting  $r_0$  in a way that the local bracketing result and the deviation bound apply simultaneously; see the discussion at the end of the section.

Below we suppose that a sufficiently large constant  $x$  is fixed to specify the accepted level be of order  $e^{-x}$  for this deviation probability. All the constructions below depend upon this constant. We do not indicate it explicitly for ease of notation.

The key step in this large deviation bound is made in terms of an *upper function* for the process  $L(\theta, \theta^*) \stackrel{\text{def}}{=} L(\theta) - L(\theta^*)$ . Namely,  $u(\theta)$  is a deterministic *upper function* if it holds with a high probability:

$$(4.1) \quad \sup_{\theta \in \Theta} \{L(\theta, \theta^*) + u(\theta)\} \leq 0.$$

Such bounds are usually called for in the analysis of the posterior measure in the Bayes approach. Below we present sufficient conditions ensuring (4.1). Now we explain how (4.1) can be used for describing *concentration sets* for  $\theta$ .

LEMMA 4.1. *Let  $u(\theta)$  be an upper function in the sense*

$$(4.2) \quad \mathbb{P}\left(\sup_{\theta \in \Theta} \{L(\theta, \theta^*) + u(\theta)\} \geq 0\right) \leq e^{-x}$$

for  $x > 0$ . Given a subset  $\Theta_0 \subset \Theta$  with  $\theta^* \in \Theta_0$ , the condition  $u(\theta) \geq 0$  for  $\theta \notin \Theta_0$  ensures

$$\mathbb{P}(\tilde{\theta} \notin \Theta_0) \leq e^{-x}.$$

PROOF. If  $\Theta^\circ$  is a subset of  $\Theta$  not containing  $\theta^*$ , then the event  $\tilde{\theta} \in \Theta^\circ$  is only possible if  $\sup_{\theta \in \Theta^\circ} L(\theta, \theta^*) \geq 0$ , because  $L(\theta^*, \theta^*) \equiv 0$ .  $\square$

A possible way of checking the condition (4.2) is based on a lower quadratic bound for the negative expectation  $-\mathbb{E}L(\theta, \theta^*) \geq b(r)\|V_0(\theta - \theta^*)\|^2/2$  in the sense of condition  $(\mathcal{L}_r)$  from Section 2.2. We present two different results. The first one assumes that the values  $b(r)$  can be fixed universally for all  $r \geq r_0$ .

THEOREM 4.2. *Suppose  $(E_r)$  and  $(\mathcal{L}_r)$  with  $b(r) \equiv b$ . Let, for  $r \geq r_0$ ,*

$$(4.3) \quad 1 + \sqrt{x + Q} \leq 3v_0^2g(r)/b,$$

$$(4.4) \quad 6v_0\sqrt{x + Q} \leq rb,$$

with  $x + Q \geq 2.5$  and  $Q = c_1 p$ . Then

$$(4.5) \quad \mathbb{P}(\tilde{\theta} \notin \Theta_0(r_0)) \leq e^{-x}.$$

PROOF. The result follows from Theorem 2.8 of the supplement [Spokoiny (2012b)] with  $\mu = \frac{b}{3v_0}$ ,  $t(\mu) \equiv 0$ ,  $U(\theta) = L(\theta) - \mathbb{E}L(\theta)$  and  $M(\theta, \theta^*) = -\mathbb{E}L(\theta, \theta^*) \geq \frac{b}{2}\|V_0(\theta - \theta^*)\|^2$ .  $\square$

REMARK 4.1. The bound (4.5) requires only two conditions. Condition (4.3) means that the value  $g(r)$  from condition  $(Er)$  fulfills  $g^2(r) \geq C(x + p)$ , that is, we need a qualified rate in the exponential moment conditions. This is similar to requiring finite polynomial moments for the score function. Condition (4.4) requires that  $r$  exceeds some fixed value, namely,  $r^2 \geq C(x + p)$ . This bound is helpful for fixing the value  $r_0$  providing a sensible deviation probability bound.

If  $b(r)$  decreases with  $r$ , the result is a bit more involved. The key requirement is that  $b(r)$  decreases not too fast, so that the product  $rb(r)$  grows to infinity with  $r$ . The idea is to include the complement of the central set  $\Theta_0$  in  $\Theta$  in the union of the growing sets  $\Theta_0(r_k)$  with  $b(r_k) \geq b(r_0)2^{-k}$ , and then apply Theorem 4.2 for each  $\Theta_0(r_k)$ .

THEOREM 4.3. *Suppose  $(Er)$  and  $(Lr)$ . Let  $r_k$  be such that  $b(r_k) \geq b(r_0)2^{-k}$  for  $k \geq 1$ . If the conditions*

$$1 + \sqrt{x + Q + ck} \leq 3v_0^2 g(r_k)/b(r_k),$$

$$6v_0 \sqrt{x + Q + ck} \leq r_k b(r_k),$$

are fulfilled for  $c = \log(2)$ , then it holds

$$\mathbb{P}(\tilde{\theta} \notin \Theta_0(r_0)) \leq e^{-x}.$$

PROOF. The result (4.5) is applied to each set  $\Theta_0(r_k)$  and  $x_k = x + ck$ . This yields

$$\mathbb{P}(\tilde{\theta} \notin \Theta_0(r_0)) \leq \sum_{k \geq 1} \mathbb{P}(\tilde{\theta} \notin \Theta_0(r_k)) \leq \sum_{k \geq 1} e^{-x-ck} = e^{-x}$$

as required.  $\square$

REMARK 4.2. Here we briefly discuss the very important question: how one can fix the value  $r_0$  ensuring the bracketing result in the local set  $\Theta_0(r_0)$  and a small probability of the related set  $C_\epsilon(r)$  from (3.8)? The event  $\{\|V_0 D_\epsilon^{-1} \xi_\epsilon\| > r\}$  requires  $r^2 \geq C(x + p)$ . Further, we inspect the deviation bound for the complement  $\Theta \setminus \Theta_0(r_0)$ . For simplicity, assume  $(Lr)$  with  $b(r) \equiv b$ . Then the condition (4.4) of Theorem 4.2 requires that

$$(4.6) \quad r_0^2 \geq Cb^{-2}(x + p).$$

In words, the squared radius  $r_0^2$  should be at least of order  $p$ . The other condition (4.3) of Theorem 4.2 is technical and only requires that  $g(r)$  is sufficiently large, while the local results only require that  $\delta(r)$  and  $\varrho(r)$  are small for such  $r$ . In the asymptotic setup one can typically bring these conditions together. Section 5 provides further discussion for the i.i.d. setup.

**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, for example, Chapters 2 and 3 in [Ibragimov and Khas'minskiĭ \(1981\)](#). 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 to demonstrate that the proposed nonasymptotic viewpoint is able to reproduce the existing brilliant and well-established results of the classical parametric theory. Surprisingly, the majority of classical efficiency results can be easily derived from the obtained general nonasymptotic 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 challenges in 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 of 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.* An i.i.d. parametric model means that the observations  $\mathbf{Y} = (Y_1, \dots, Y_n)$  are independent identically distributed from a distribution  $P$  which belongs to 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 and 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  $(P_\theta)$ . The case of nonidentically distributed observations can be done similarly at the cost of more complicated notation.

In what follows the parametric family  $(P_\theta)$  is supposed to be dominated by a measure  $\mu_0$ , and each density  $p(y, \theta) = dP_\theta/d\mu_0(y)$  is two times continuously

differentiable in  $\theta$  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  $\theta \in \Theta$ :

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

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

$$\theta^* \stackrel{\text{def}}{=} \arg \max_{\theta \in \Theta} \mathbb{E}L(\theta) = \arg \max_{\theta \in \Theta} \mathbb{E}\ell(Y_1, \theta).$$

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

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

I.i.d. structure of the  $Y_i$ 's allows to rewrite the local conditions  $(Er)$ ,  $(ED_0)$ ,  $(ED_1)$ , and  $(\mathcal{L}_0)$ , and  $(\mathcal{I})$  in terms of the marginal distribution.

*(ed<sub>0</sub>) There exists a positively definite symmetric matrix  $\mathbf{v}_0$ , such that for all  $|\lambda| \leq \mathfrak{g}_1$*

$$\sup_{\boldsymbol{\gamma} \in \mathbb{R}^p} \log \mathbb{E} \exp \left\{ \lambda \frac{\boldsymbol{\gamma}^\top \nabla \zeta_1(\theta^*)}{\|\mathbf{v}_0 \boldsymbol{\gamma}\|} \right\} \leq v_0^2 \lambda^2 / 2.$$

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

Next consider the local sets

$$\Theta_{\text{loc}}(u) = \{\theta : \|\mathbf{v}_0(\theta - \theta^*)\| \leq u\}.$$

In view of  $V_0^2 = n\mathbf{v}_0^2$ , it holds  $\Theta_0(r) = \Theta_{\text{loc}}(u)$  with  $r^2 = nu^2$ .

Below we distinguish between local conditions for  $u \leq u_0$  and the global conditions for all  $u > 0$ , where  $u_0$  is some fixed value.

The local smoothness conditions  $(ED_1)$  and  $(\mathcal{L}_0)$  require to specify the functions  $\delta(r)$  and  $\varrho(r)$  for  $r \leq r_0$  where  $r_0^2 = nu_0^2$ . If the log-likelihood function  $\ell(y, \theta)$  is sufficiently smooth in  $\theta$ , these functions can be selected proportional to  $u = r/n^{1/2}$ .

*(ed<sub>1</sub>) There are constants  $\omega^* > 0$  and  $\mathfrak{g}_1 > 0$  such that for each  $u \leq u_0$  and  $|\lambda| \leq \mathfrak{g}_1$*

$$\sup_{\boldsymbol{\gamma} \in \mathbb{R}^p} \sup_{\theta \in \Theta_{\text{loc}}(u)} \log \mathbb{E} \exp \left\{ \lambda \frac{\boldsymbol{\gamma}^\top [\nabla \zeta_1(\theta) - \nabla \zeta_1(\theta^*)]}{\omega^* u \|\mathbf{v}_0 \boldsymbol{\gamma}\|} \right\} \leq v_0^2 \lambda^2 / 2.$$

Further, we restate the local identifiability condition  $(\mathcal{L}_0)$  in terms of the expected value  $\mathbb{k}(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \stackrel{\text{def}}{=} -\mathbb{E}\{\ell(Y_i, \boldsymbol{\theta}) - \ell(Y_i, \boldsymbol{\theta}^*)\}$  for each  $i$ . We suppose that  $\mathbb{k}(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$  is two times differentiable w.r.t.  $\boldsymbol{\theta}$ . The definition of  $\boldsymbol{\theta}^*$  implies  $\nabla \mathbb{E}\ell(Y_i, \boldsymbol{\theta}^*) = 0$ . Define also the matrix  $\mathbf{f}_0 = -\nabla^2 \mathbb{E}\ell(Y_i, \boldsymbol{\theta}^*)$ . In the parametric case  $P = P_{\boldsymbol{\theta}^*}$ ,  $\mathbb{k}(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$  is the *Kullback–Leibler divergence* between  $P_{\boldsymbol{\theta}^*}$  and  $P_{\boldsymbol{\theta}}$ , while the matrices  $\mathbf{v}_0^2 = \mathbf{f}_0$  are equal to each other and coincide with the *Fisher information matrix* of the family  $(P_{\boldsymbol{\theta}})$  at  $\boldsymbol{\theta}^*$ .

$(\ell_0)$  There is a constant  $\delta^*$  such that it holds for each  $u \leq u_0$

$$\sup_{\boldsymbol{\theta} \in \Theta_{\text{loc}}(u)} \left| \frac{2\mathbb{k}(\boldsymbol{\theta}, \boldsymbol{\theta}^*)}{(\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top \mathbf{f}_0 (\boldsymbol{\theta} - \boldsymbol{\theta}^*)} - 1 \right| \leq \delta^* u.$$

(i) There is a constant  $\alpha > 0$  such that  $\alpha^2 \mathbf{f}_0^2 \geq \mathbf{v}_0^2$ .

(eu) For each  $u > 0$ , there exists  $g_1(u) > 0$ , such that for all  $|\lambda| \leq g_1(u)$

$$\sup_{\boldsymbol{y} \in \mathbb{R}^p} \sup_{\boldsymbol{\theta} \in \Theta_{\text{loc}}(u)} \log \mathbb{E} \exp \left\{ \lambda \frac{\boldsymbol{y}^\top \nabla \zeta_1(\boldsymbol{\theta})}{\|\mathbf{v}_0 \boldsymbol{y}\|} \right\} \leq v_0^2 \lambda^2 / 2.$$

$(\ell u)$  For each  $u > 0$ , there exists  $b(u) > 0$  such that

$$\sup_{\boldsymbol{\theta} \in \Theta : \|\mathbf{v}_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\| = u} \frac{\mathbb{k}(\boldsymbol{\theta}, \boldsymbol{\theta}^*)}{\|\mathbf{v}_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2} \geq b(u),$$

LEMMA 5.1. Let  $Y_1, \dots, Y_n$  be i.i.d. Then  $(eu)$ ,  $(ed_0)$ ,  $(ed_1)$ , (i) and  $(\ell_0)$  imply  $(Er)$ ,  $(ED_0)$ ,  $(ED_1)$ ,  $(I)$  and  $(\mathcal{L}_0)$  with  $V_0^2 = n\mathbf{v}_0^2$ ,  $D_0^2 = n\mathbf{f}_0$ ,  $\omega(r) = \omega^* r / n^{1/2}$ ,  $\delta(r) = \delta^* r / n^{1/2}$ , and  $g = g_1 \sqrt{n}$ .

PROOF. The identities  $V_0^2 = n\mathbf{v}_0^2$ ,  $D_0^2 = n\mathbf{f}_0$  follow from the i.i.d. structure of the observations  $Y_i$ . We briefly comment on condition  $(Er)$ . The use of the i.i.d. structure once again yields by (5.1) in view of  $V_0^2 = n\mathbf{v}_0^2$

$$\log \mathbb{E} \exp \left\{ \lambda \frac{\boldsymbol{y}^\top \nabla \zeta(\boldsymbol{\theta})}{\|V_0 \boldsymbol{y}\|} \right\} = n \mathbb{E} \exp \left\{ \frac{\lambda}{n^{1/2}} \frac{\boldsymbol{y}^\top \nabla \zeta_1(\boldsymbol{\theta})}{\|\mathbf{v}_0 \boldsymbol{y}\|} \right\} \leq v_0^2 \lambda^2 / 2$$

as long as  $\lambda \leq n^{1/2} g_1(u) \leq g(r)$ . Similarly for  $(ED_0)$  and  $(ED_1)$ .  $\square$

REMARK 5.1. This remark discusses how the presented conditions relate to what is usually assumed in statistical literature. One general remark concerns the choice of the parametric family  $(P_{\boldsymbol{\theta}})$ . The point of the classical theory is that the true measure is in this family, so the conditions should be as weak as possible. The viewpoint of this paper is slightly different: whatever family  $(P_{\boldsymbol{\theta}})$  is taken, the true measure is never included, any model is only an approximation of reality. From the other side, the choice of the parametric model  $(P_{\boldsymbol{\theta}})$  is always done by a statistician. Sometimes some special stylized features of the model force to include

an irregularity in this family. Otherwise any smoothness condition on the density  $\ell(y, \theta)$  can be secured by a proper choice of the family  $(P_\theta)$ .

The presented list also includes the exponential moment conditions  $(ed_0)$  and  $(ed_1)$  on the gradient  $\nabla \ell(Y_1, \theta)$ . We need exponential moments for establishing some nonasymptotic risk bounds; the classical concentration bounds require even stronger conditions that the considered random variables are bounded.

The identifiability condition  $(\ell_u)$  is very easy to check in the usual asymptotic setup. Indeed, if the parameter set  $\Theta$  is compact, the Kullback–Leibler divergence  $\mathbb{K}(\theta, \theta^*)$  is continuous and positive for all  $\theta \neq \theta^*$ , then  $(\ell_u)$  is fulfilled automatically with a universal constant  $b$ . If  $\Theta$  is not compact, the condition is still fulfilled but the function  $b(u)$  may depend on  $u$ .

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

5.1.1. *A large deviation bound.* This section presents some sufficient conditions ensuring a small deviation probability for the event  $\{\tilde{\theta} \notin \Theta_{\text{loc}}(u_0)\}$  for a fixed  $u_0$ . Below  $\mathbb{Q} = c_1 p$ . We only discuss the case  $b(u) \equiv b$ . The general case only requires more complicated notation. The next result follows from Theorem 4.2 with the obvious changes.

THEOREM 5.2. *Suppose  $(eu)$  and  $(\ell_u)$  with  $b(u) \equiv b$ . If, for  $u_0 > 0$ ,*

$$(5.2) \quad \begin{aligned} n^{1/2} u_0 b &\geq 6v_0 \sqrt{x + \mathbb{Q}}, \\ 1 + \sqrt{x + \mathbb{Q}} &\leq 3b^{-1} v_0^2 g_1(u_0) n^{1/2}, \end{aligned}$$

then

$$\mathbb{P}(\tilde{\theta} \notin \Theta_{\text{loc}}(u_0)) = \mathbb{P}(\|\mathbf{v}_0(\tilde{\theta} - \theta^*)\| > u_0) \leq e^{-x}.$$

REMARK 5.2. The presented result helps to qualify two important values  $u_0$  and  $n$  providing a sensible deviation probability bound. For simplicity suppose that  $g_1(u) \equiv g_1 > 0$ . Then the condition (5.2) can be written as  $nu_0^2 \gg x + \mathbb{Q}$ . In other words, the result of the theorem claims a large deviation bound for the vicinity  $\Theta_{\text{loc}}(u_0)$  with  $u_0^2$  of order  $p/n$ . In classical asymptotic statistics this result is usually referred to as *root- $n$  consistency*. Our approach yields this result in a very strong form and for finite samples.

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  $\mathbf{f}_0$  from conditions  $(ed_0)$ ,  $(ed_1)$  and  $(\ell_0)$  determine their drift and variance components. Define

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

If  $\tau_\epsilon \stackrel{\text{def}}{=} \delta + \alpha^2 \varrho < 1$ , then

$$\mathbf{f}_\epsilon \geq (1 - \tau_\epsilon) \mathbf{f}_0 > 0.$$

Further,  $D_\epsilon^2 = n \mathbf{f}_\epsilon$  and

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

The upper bracketing 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 can be viewed 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<sub>0</sub>). Given  $\mathfrak{u}_0$ , assume (ed<sub>1</sub>), ( $\ell_0$ ) and (i) on  $\Theta_{\text{loc}}(\mathfrak{u}_0)$ , and let  $\varrho = 3\nu_0 \omega^* \mathfrak{u}_0$ ,  $\delta = \delta^* \mathfrak{u}_0$ , and  $\tau_\epsilon \stackrel{\text{def}}{=} \delta + \alpha^2 \varrho < 1$ . Then the results of Theorem 3.1 and all its corollaries apply to the case of i.i.d. modeling with  $\mathfrak{r}_0^2 = n \mathfrak{u}_0^2$ . In particular, on the random set  $C_\epsilon(\mathfrak{r}_0) = \{\tilde{\boldsymbol{\theta}} \in \Theta_{\text{loc}}(\mathfrak{u}_0), \|\boldsymbol{\xi}_\epsilon\| \leq \mathfrak{r}_0\}$ , it holds*

$$\begin{aligned} \|\boldsymbol{\xi}_\epsilon\|^2/2 - \diamond_\epsilon(\mathfrak{r}_0) &\leq L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) \leq \|\boldsymbol{\xi}_\epsilon\|^2/2 + \diamond_\epsilon(\mathfrak{r}_0), \\ \|\sqrt{n \mathbf{f}_\epsilon}(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) - \boldsymbol{\xi}_\epsilon\|^2 &\leq 2\Delta_\epsilon(\mathfrak{r}_0). \end{aligned}$$

The random quantities  $\diamond_\epsilon(\mathfrak{r}_0)$ ,  $\diamond_\epsilon(\mathfrak{r}_0)$  and  $\Delta_\epsilon(\mathfrak{r}_0)$  follow the probability bounds of Theorems 3.7 and 3.10.

Now we briefly discuss the implications of Theorem 5.2 and 5.3 to the classical asymptotic setup with  $n \rightarrow \infty$ . We fix  $\mathfrak{u}_0^2 = Cp/n$  for a constant  $C$  ensuring the deviation bound of Theorem 5.2. Then  $\delta$  is of order  $\mathfrak{u}_0$  and the same for  $\varrho$ . For a sufficiently large  $n$ , both quantities are small and, thus, the spread  $\Delta_\epsilon(\mathfrak{r}_0)$  is small as well; see Section 3.3.

Further, under (ed<sub>0</sub>) condition, the normalized score

$$\boldsymbol{\xi} \stackrel{\text{def}}{=} (n \mathbf{f}_0)^{-1/2} \sum \nabla \ell(Y_i, \boldsymbol{\theta}^*)$$

is zero mean asymptotically normal by the central limit theorem. Moreover, if  $\mathbf{f}_0 = \mathbf{v}_0^2$ , then  $\boldsymbol{\xi}$  is asymptotically standard normal. The same holds for  $\boldsymbol{\xi}_\epsilon$ . This immediately yields all classical asymptotic results like Wilks theorem or the Fisher expansion for MLE in the i.i.d. setup as well as the asymptotic efficiency of the MLE. Moreover, our results' bounds yield the asymptotic result for the case when the parameter dimension  $p = p_n$  grows linearly with  $n$ . Below  $u_n = o_n(p_n)$  means that  $u_n/p_n \rightarrow 0$  as  $n \rightarrow \infty$ .

**THEOREM 5.4.** *Let  $Y_1, \dots, Y_n$  be i.i.d.  $\mathbb{P}_{\theta^*}$  and let  $(ed_0), (ed_1), (\ell_0), (\iota), (e\iota)$  and  $(\ell\iota)$  with  $\mathfrak{b}(\iota) \equiv \mathfrak{b}$  hold. If  $n > Cp_n$  for a fixed constant  $C$  depending on constants in the above conditions only, then*

$$\|\sqrt{n\mathbf{f}_0}(\tilde{\theta} - \theta^*) - \xi\|^2 = o_n(p_n), \quad 2L(\tilde{\theta}, \theta^*) - \|\xi\|^2 = o_n(p_n).$$

This result particularly yields that  $\sqrt{n\mathbf{f}_0}(\tilde{\theta} - \theta^*)$  is nearly standard normal and  $2L(\tilde{\theta}, \theta^*)$  is nearly  $\chi_p^2$ .

**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, for example, [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  $d(w) = e^w$  and 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  $\psi_i \in \mathbb{R}^p$ :

$$(5.3) \quad Y_i \sim P_{\psi_i^\top \theta^*}.$$

The corresponding log-density of a GLM reads as

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

Under  $\mathbb{P}_{\theta^*}$  each observation  $Y_i$  follows (5.3), in particular,  $\mathbb{E}Y_i = d'(\psi_i^\top \theta^*)$ . However, similarly to the previous sections, it is accepted that the parametric model (5.3) is misspecified. Response misspecification means that the vector  $\mathbf{f} \stackrel{\text{def}}{=} \mathbb{E}\mathbf{Y}$  cannot be represented in the form  $d'(\Psi^\top \theta)$  whatever  $\theta$  is. The other sort of misspecification concerns the data distribution. The model (5.3) 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 target of estimation  $\theta^*$  is defined by

$$\theta^* \stackrel{\text{def}}{=} \arg \max_{\theta} \mathbb{E}L(\theta).$$

The quasi-MLE  $\tilde{\theta}$  is defined by maximization of  $L(\theta)$ :

$$\tilde{\theta} = \arg \max_{\theta} L(\theta) = \arg \max_{\theta} \sum \{Y_i \psi_i^\top \theta - d(\psi_i^\top \theta)\}.$$

Convexity of  $d(\cdot)$  implies that  $L(\theta)$  is a concave function of  $\theta$ , 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  $\theta^*$  is the maximizer of the expected log-likelihood:

$$\theta^* = \arg \max_{\theta} \mathbb{E}L(\theta) = \arg \max_{\theta} \sum \{f_i \Psi_i^\top \theta - d(\Psi_i^\top \theta)\}$$

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

(e<sub>1</sub>) *There exist some constants  $v_0$  and  $g_1 > 0$ , and for every  $i$  a constant  $s_i$  such that  $\mathbb{E}(\varepsilon_i/s_i)^2 \leq 1$  and*

$$(5.4) \quad \log \mathbb{E} \exp(\lambda \varepsilon_i / s_i) \leq v_0^2 \lambda^2 / 2, \quad |\lambda| \leq g_1.$$

A natural candidate for  $s_i$  is  $\sigma_i$  where  $\sigma_i^2 = \mathbb{E}\varepsilon_i^2$  is the variance of  $\varepsilon_i$ ; see Lemma 2.13 of the supplement [Spokoiny (2012b)]. Under (5.4), introduce a  $p \times p$  matrix  $V_0$  defined by

$$(5.5) \quad V_0^2 \stackrel{\text{def}}{=} \sum s_i^2 \Psi_i \Psi_i^\top.$$

Condition (e<sub>1</sub>) effectively means that each error term  $\varepsilon_i = Y_i - \mathbb{E}Y_i$  has some bounded exponential moments: for  $|\lambda| \leq g_1$ , it holds  $f(\lambda) \stackrel{\text{def}}{=} \log \mathbb{E} \exp(\lambda \varepsilon_i / s_i) < \infty$ . This implies the quadratic upper bound for the function  $f(\lambda)$  for  $|\lambda| \leq g_1$ ; see Lemma 2.13 of the supplement [Spokoiny (2012b)]. In words, condition (e<sub>1</sub>) requires a light (exponentially decreasing) tail for the marginal distribution of each  $\varepsilon_i$ .

Define also

$$(5.6) \quad N^{-1/2} \stackrel{\text{def}}{=} \max_i \sup_{\mathcal{Y} \in \mathbb{R}^p} \frac{s_i |\Psi_i^\top \mathcal{Y}|}{\|V_0 \mathcal{Y}\|}.$$

LEMMA 5.5. *Assume (e<sub>1</sub>) and let  $V_0$  be defined by (5.5) and  $N$  by (5.6). Then conditions (ED<sub>0</sub>) and (E $\mathcal{X}$ ) follow from (e<sub>1</sub>) with the matrix  $V_0$  due to (5.5) and  $g = g_1 N^{1/2}$ . Moreover, the stochastic component  $\zeta(\theta)$  is linear in  $\theta$  and the condition (ED<sub>1</sub>) is fulfilled with  $\omega(x) \equiv 0$ .*

PROOF. The gradient of the stochastic component  $\zeta(\theta)$  of  $L(\theta)$  does not depend on  $\theta$ , namely,  $\nabla \zeta(\theta) = \sum \Psi_i \varepsilon_i$  with  $\varepsilon_i = Y_i - \mathbb{E}Y_i$ . Now, for any unit vector  $\mathcal{Y} \in \mathbb{R}^p$  and  $\lambda \leq g$ , independence of the  $\varepsilon_i$ 's implies that

$$\log \mathbb{E} \exp \left\{ \frac{\lambda}{\|V_0 \mathcal{Y}\|} \mathcal{Y}^\top \sum \Psi_i \varepsilon_i \right\} = \sum \log \mathbb{E} \exp \left\{ \frac{\lambda s_i \Psi_i^\top \mathcal{Y}}{\|V_0 \mathcal{Y}\|} \varepsilon_i / s_i \right\}.$$

By definition,  $\mathfrak{s}_i |\Psi_i^\top \boldsymbol{\gamma}| / \|V_0 \boldsymbol{\gamma}\| \leq N^{-1/2}$  and, therefore,  $\lambda \mathfrak{s}_i |\Psi_i^\top \boldsymbol{\gamma}| / \|V_0 \boldsymbol{\gamma}\| \leq \mathfrak{g}_1$ . Hence, (5.4) implies

$$(5.7) \quad \log \mathbb{E} \exp \left\{ \frac{\lambda}{\|V_0 \boldsymbol{\gamma}\|} \boldsymbol{\gamma}^\top \sum \Psi_i \varepsilon_i \right\} \leq \frac{v_0^2 \lambda^2}{2 \|V_0 \boldsymbol{\gamma}\|^2} \sum \mathfrak{s}_i^2 |\Psi_i^\top \boldsymbol{\gamma}|^2 = \frac{v_0^2 \lambda^2}{2},$$

and  $(ED_0)$  follows.  $\square$

It only remains 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 5.6. *It holds*

$$(5.8) \quad \begin{aligned} -\mathbb{E}L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) &= \sum \{d(\Psi_i^\top \boldsymbol{\theta}) - d(\Psi_i^\top \boldsymbol{\theta}^*) - d'(\Psi_i^\top \boldsymbol{\theta}^*) \Psi_i^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^*)\} \\ &= \mathcal{K}(\mathbb{P}_{\boldsymbol{\theta}^*}, \mathbb{P}_{\boldsymbol{\theta}}), \end{aligned}$$

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,

$$(5.9) \quad -\mathbb{E}L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) = \|D(\boldsymbol{\theta}^\circ)(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2 / 2,$$

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 (5.8) follows. The Taylor expansion of the second order around  $\boldsymbol{\theta}^*$  yields the expansion (5.9).  $\square$

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 (5.5). 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  $\mathfrak{s}_i^2 = d''(\Psi_i^\top \boldsymbol{\theta}^*)$ . The matrix  $V_0$  describes a local elliptic neighborhood of the central point  $\boldsymbol{\theta}^*$  in the form  $\Theta_0(r) = \{\boldsymbol{\theta} : \|V_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\| \leq r\}$ . If the matrix function  $D^2(\boldsymbol{\theta})$  is continuous in this vicinity  $\Theta_0(r)$ , then the value  $\delta(r)$  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  $\Theta_0(r)$ .

LEMMA 5.7. *Suppose that*

$$(5.10) \quad \|\mathbf{I}_p - D_0^{-1} D^2(\boldsymbol{\theta}) D_0^{-1}\|_\infty \leq \delta(x), \quad \boldsymbol{\theta} \in \Theta_0(x).$$

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

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

with

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

Linearity of the stochastic component  $\zeta(\boldsymbol{\theta})$  in the considered GLM implies the important fact that the quantities  $\diamond_\epsilon(x)$ ,  $\diamond_{\underline{\epsilon}}(x)$  in the general bracketing bound (3.4) vanish for any  $x$ . Therefore, in the GLM case, the deficiency can be defined as the difference between upper and lower excess and it can be easily evaluated:

$$\Delta(x) = \|\xi_\epsilon\|^2/2 - \|\xi_{\underline{\epsilon}}\|^2/2 = \delta \|\xi\|^2.$$

Our result assumes some concentration properties of the squared norm  $\|\xi\|^2$  of the vector  $\xi$ . These properties can be established by general results of Section 1 of the complement under the regularity condition: for some  $\alpha$

$$(5.11) \quad V_0 \leq \alpha D_0.$$

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

THEOREM 5.8. *Let  $(e_1)$  hold. Then for  $\delta \geq \delta(x)$  any  $z > 0$  and  $\mathfrak{z} > 0$ , it holds*

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

Moreover, on the set  $C_\epsilon(x) = \{\|V_0(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\| \leq r, \|\xi_\epsilon\| \leq r\}$ , it holds

$$(5.12) \quad \|D_0(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) - \xi\|^2 \leq \frac{2\delta}{1 - \delta^2} \|\xi\|^2.$$

If the function  $d(w)$  is quadratic, then the approximation error  $\delta$  vanishes as well and the expansion (5.12) becomes equality which is also fulfilled globally, a localization step is not required. However, if  $d(w)$  is not quadratic, the result applies only locally and it has to be accomplished with a large deviation bound. The GLM structure is helpful in the large deviation zone as well. Indeed, the gradient  $\nabla \zeta(\boldsymbol{\theta})$  does not depend on  $\boldsymbol{\theta}$  and, hence, the most delicate condition  $(E_r)$  is fulfilled automatically with  $\mathfrak{g} = \mathfrak{g}_1 N^{1/2}$  for all local sets  $\Theta_0(x)$ . Further, the

identifiability condition  $(\mathcal{L}_r)$  easily follows from Lemma 5.6: it suffices to bound from below the matrix  $D(\theta)$  for  $\theta \in \Theta_0(r)$ :

$$D(\theta) \geq b(r)V_0, \quad \theta \in \Theta_0(r).$$

An interesting question, similarly to the i.i.d. case, is the minimal radius  $r_0$  of the local vicinity  $\Theta_0(r_0)$  ensuring the desirable concentration property. Suppose for the moment that the constants  $b(r)$  are all the same for different  $r$ :  $b(r) \equiv b$ . Under the regularity condition (5.11), a sufficient lower bound for  $r_0$  can be based on Corollary 4.3. The required condition can be restated as

$$1 + \sqrt{x + Q} \leq 3v_0^2g/b, \quad 6v_0\sqrt{x + Q} \leq rb.$$

It remains to note that  $Q = c_1p$  and  $g = g_1N^{1/2}$ . So, the required conditions are fulfilled for  $r^2 \geq r_0^2 = C(x + p)$ , where  $C$  only depends on  $v_0, b$ , and  $g$ .

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  $\Psi_i \in \mathbb{R}^p$ :

$$(5.13) \quad Y_i = \Psi_i^\top \theta + \varepsilon_i,$$

where  $\varepsilon_i$  means the  $i$ th individual error. As usual, the true data distribution can deviate from the linear model. In addition, we admit contaminated data which naturally leads to the idea of robust estimation. This section offers a qMLE view on the robust estimation problem. Our parametric family assumes the linear dependence (5.13) with i.i.d. errors  $\varepsilon_i$  which follow the double exponential (Laplace) distribution with the density  $(1/2)e^{-|y|}$ . Then the corresponding log-likelihood reads as

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

and  $\tilde{\theta} \stackrel{\text{def}}{=} \arg \max_{\theta} L(\theta)$  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  $\theta^*$  is usually defined by the equation  $\theta^* = \arg \max_{\theta} \mathbb{E}L(\theta)$ .

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

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

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

Note that the difference  $L(\theta) - L(\theta^*)$  is bounded by  $\frac{1}{2} \sum |\Psi_i^\top (\theta - \theta^*)|$ . Next we check conditions  $(ED_0)$  and  $(ED_1)$ . Denote  $\xi_i(\theta) = \mathbb{1}(Y_i - \Psi_i^\top \theta \leq 0) - q_i(\theta)$

for  $q_i(\boldsymbol{\theta}) = \mathbb{P}(Y_i - \boldsymbol{\Psi}_i^\top \boldsymbol{\theta} \leq 0)$ . This is a centered Bernoulli random variable, and it is easy to check that

$$(5.14) \quad \nabla \zeta(\boldsymbol{\theta}) = - \sum \xi_i(\boldsymbol{\theta}) \boldsymbol{\Psi}_i.$$

This expression differs from the similar ones from the linear and generalized linear regression because the stochastic terms  $\xi_i$  now depend on  $\boldsymbol{\theta}$ . First we check the global condition (Er). Fix any  $g_1 < 1$ . Then it holds for a Bernoulli r.v.  $Z$  with  $\mathbb{P}(Z = 1) = q$ ,  $\xi = Z - q$ , and  $|\lambda| \leq g_1$

$$(5.15) \quad \begin{aligned} \log \mathbb{E} \exp(\lambda \xi) &= \log[q \exp\{\lambda(1 - q)\} + (1 - q) \exp(-\lambda q)] \\ &\leq v_0^2 q(1 - q) \lambda^2 / 2, \end{aligned}$$

where  $v_0 \geq 1$  depends on  $g_1$  only. Let now a vector  $\boldsymbol{\gamma} \in \mathbb{R}^p$  and  $\rho > 0$  be such that  $\rho |\boldsymbol{\Psi}_i^\top \boldsymbol{\gamma}| \leq g_1$  for all  $i = 1, \dots, n$ . Then

$$(5.16) \quad \begin{aligned} \log \mathbb{E} \exp\{\rho \boldsymbol{\gamma}^\top \nabla \zeta(\boldsymbol{\theta})\} &\leq \frac{v_0^2 \rho^2}{2} \sum q_i(\boldsymbol{\theta}) \{1 - q_i(\boldsymbol{\theta})\} |\boldsymbol{\Psi}_i^\top \boldsymbol{\gamma}|^2 \\ &\leq \frac{v_0^2 \rho^2}{2} \|V(\boldsymbol{\theta}) \boldsymbol{\gamma}\|^2, \end{aligned}$$

where

$$V^2(\boldsymbol{\theta}) = \sum q_i(\boldsymbol{\theta}) \{1 - q_i(\boldsymbol{\theta})\} \boldsymbol{\Psi}_i \boldsymbol{\Psi}_i^\top.$$

Denote also

$$(5.17) \quad V_0^2 = \frac{1}{4} \sum \boldsymbol{\Psi}_i \boldsymbol{\Psi}_i^\top.$$

Clearly,  $V(\boldsymbol{\theta}) \leq V_0$  for all  $\boldsymbol{\theta}$  and condition (Er) is fulfilled with the matrix  $V_0$  and  $g(x) \equiv g = g_1 N^{1/2}$  for  $N$  defined by

$$(5.18) \quad N^{-1/2} \stackrel{\text{def}}{=} \max_i \sup_{\boldsymbol{\gamma} \in \mathbb{R}^p} \frac{\boldsymbol{\Psi}_i^\top \boldsymbol{\gamma}}{2 \|V_0 \boldsymbol{\gamma}\|};$$

cf. (5.7).

Let some  $r_0 > 0$  be fixed. We will specify this choice later. Now we check the local conditions within the elliptic vicinity  $\Theta_0(r_0) = \{\boldsymbol{\theta} : \|V_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\| \leq r_0\}$  of the central point  $\boldsymbol{\theta}^*$  for  $V_0$  from (5.17). Then condition (ED<sub>0</sub>) with the matrix  $V_0$  and  $g = N^{1/2} g_1$  is fulfilled on  $\Theta_0(r_0)$  due to (5.16). Next, in view of (5.18), it holds  $|\boldsymbol{\Psi}_i^\top \boldsymbol{\gamma}| \leq 2N^{-1/2} \|V_0 \boldsymbol{\gamma}\|$  for any vector  $\boldsymbol{\gamma} \in \mathbb{R}^p$ . By (5.14),

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

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

$$\xi_i(\boldsymbol{\theta}) - \xi_i(\boldsymbol{\theta}^*) = \mathbb{1}(\boldsymbol{\Psi}_i^\top \boldsymbol{\theta}^* \leq Y_i < \boldsymbol{\Psi}_i^\top \boldsymbol{\theta}) - \mathbb{P}(\boldsymbol{\Psi}_i^\top \boldsymbol{\theta}^* \leq Y_i < \boldsymbol{\Psi}_i^\top \boldsymbol{\theta}).$$

Similarly, for  $\Psi_i^\top \boldsymbol{\theta} < \Psi_i^\top \boldsymbol{\theta}^*$

$$\xi_i(\boldsymbol{\theta}) - \xi_i(\boldsymbol{\theta}^*) = -\mathbb{1}(\Psi_i^\top \boldsymbol{\theta} \leq Y_i < \Psi_i^\top \boldsymbol{\theta}^*) + \mathbb{P}(\Psi_i^\top \boldsymbol{\theta} \leq Y_i < \Psi_i^\top \boldsymbol{\theta}^*).$$

Define  $q_i(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \stackrel{\text{def}}{=} |q_i(\boldsymbol{\theta}) - q_i(\boldsymbol{\theta}^*)|$ . Now (5.15) yields similarly to (5.16)

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

with

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

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

$$|q_i(\boldsymbol{\theta}) - q_i(\boldsymbol{\theta}^*)| \leq C |\Psi_i^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^*)| \leq CN^{-1/2} \|V_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\| \leq CN^{-1/2} r.$$

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

LEMMA 5.9. *It holds for any  $\boldsymbol{\theta}$*

$$(5.19) \quad -\frac{\partial^2}{\partial^2 \boldsymbol{\theta}} \mathbb{E}L(\boldsymbol{\theta}) = D^2(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \sum p_i(\Psi_i^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^*)) \Psi_i \Psi_i^\top,$$

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

$$(5.20) \quad \begin{aligned} -\mathbb{E}L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) &= \frac{1}{2} \sum |\Psi_i^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^*)|^2 f_i(\Psi_i^\top (\boldsymbol{\theta}^\circ - \boldsymbol{\theta}^*)) \\ &= (\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top D^2(\boldsymbol{\theta}^\circ) (\boldsymbol{\theta} - \boldsymbol{\theta}^*) / 2. \end{aligned}$$

PROOF. Obviously

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

The identity (5.19) 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

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

Now (5.20) follows by the Taylor expansion of the second order at  $\boldsymbol{\theta}^*$ .  $\square$

Define

$$(5.21) \quad D_0^2 \stackrel{\text{def}}{=} \sum |\Psi_i^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^*)|^2 f_i(0).$$

Due to this lemma, condition  $(\mathcal{L}_0)$  is fulfilled in  $\Theta_0(r)$  with this choice  $D_0$  for  $\delta(r)$  from (5.10); see Lemma 5.7. Moreover, if  $f_i(0) \geq \alpha^2/4$  for  $\alpha > 0$ , then the identifiability condition  $(\mathcal{I})$  is also satisfied. Now all the local conditions are fulfilled, yielding the general bracketing bound of Theorem 3.1 and all its corollaries.

It only remains to accomplish them by a large deviation bound, that is, to specify the local vicinity  $\Theta_0(r_0)$  providing the prescribed deviation bound. A sufficient condition for the concentration property is that the expectation  $\mathbb{E}L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$  grows in absolute value with the distance  $\|V_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|$ . We use the representation (5.19). Suppose that for some fixed  $\delta < 1/2$  and  $\rho > 0$

$$(5.22) \quad |f_i(u)/f_i(0) - 1| \leq \delta, \quad |u| \leq \rho.$$

For any  $\boldsymbol{\theta}$  with  $\|V_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\| = r \geq r_0$ , and for any  $i = 1, \dots, n$ , it holds

$$|\Psi_i^\top(\boldsymbol{\theta} - \boldsymbol{\theta}^*)| \leq N^{-1/2} \|V_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\| = N^{-1/2} r.$$

Therefore, for  $r \leq \rho N^{1/2}$  and any  $\boldsymbol{\theta} \in \Theta_0(r)$  with  $\|V_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\| = r$ , it holds  $f_i(\Psi_i^\top(\boldsymbol{\theta} - \boldsymbol{\theta}^*)) \geq (1 - \delta)f_i(0)$ . Now Lemma 5.9 implies

$$-\mathbb{E}L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \geq \frac{1 - \delta}{2} \|D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2 \geq \frac{1 - \delta}{2\alpha^2} \|V_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2 = \frac{1 - \delta}{2\alpha^2} r^2.$$

By Lemma 5.9 the function  $-\mathbb{E}L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$  is convex. This easily yields

$$-\mathbb{E}L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \geq \frac{1 - \delta}{2\alpha^2} \rho N^{1/2} r$$

for all  $r \geq \rho N^{1/2}$ . Thus,

$$\text{rb}(r) \geq \begin{cases} (1 - \delta)(2\alpha^2)^{-1} r, & \text{if } r \leq \rho N^{1/2}, \\ (1 - \delta)(2\alpha^2)^{-1} \rho N^{1/2}, & \text{if } r > \rho N^{1/2}. \end{cases}$$

So, the global identifiability condition  $(\mathcal{L}_1)$  is fulfilled if  $r_0^2 \geq C_1 \alpha^2(x + \mathbb{Q})$  and if  $\rho^2 N \geq C_2 \alpha^2(x + \mathbb{Q})$  for some fixed constants  $C_1$  and  $C_2$ .

Putting this all together yields the following result.

**THEOREM 5.10.** *Let  $Y_i$  be independent,  $\boldsymbol{\theta}^* = \arg \max_{\boldsymbol{\theta}} \mathbb{E}L(\boldsymbol{\theta})$ ,  $D_0^2$  be given by (5.21), and  $V_0^2$  by (5.17). Let also the densities  $f_i(\cdot)$  of  $Y_i - \Psi_i^\top \boldsymbol{\theta}^*$  be uniformly bounded by a constant  $C$ , fulfill (5.22) for some  $\rho > 0$  and  $\delta > 0$ , and  $f_i(0) \geq \alpha^2/4$  for all  $i$ . Finally, let  $N \geq C_2 \rho^{-2} \alpha^2(x + p)$  for some fixed  $x > 0$  and  $C_2$ . Then on the random set of probability at least  $1 - e^{-x}$ , one obtains for  $\boldsymbol{\xi} \stackrel{\text{def}}{=} D_0^{-1} \nabla L(\boldsymbol{\theta}^*)$  the bounds*

$$\|\sqrt{D_0}(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) - \boldsymbol{\xi}\|^2 = o(p), \quad 2L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) - \|\boldsymbol{\xi}\|^2 = o(p).$$

**Acknowledgments.** Financial support by the German Research Foundation (DFG) through the Collaborative Research Center 649 “Economic Risk” is gratefully acknowledged. Critics and suggestions of two anonymous referees helped a lot in improving the paper.

## SUPPLEMENTARY MATERIAL

**Some results from the theory of empirical processes** (DOI: [10.1214/12-AOS1054SUPP](https://doi.org/10.1214/12-AOS1054SUPP); .pdf). This part collects some general deviation bounds for non-Gaussian quadratic forms and for general centered random processes used in the text.

## REFERENCES

- ANDRESEN, A. and SPOKOINY, V. (2012). Wilks theorem for a quasi profile maximum likelihood. Unpublished manuscript.
- BEDNORZ, W. (2006). A theorem on majorizing measures. *Ann. Probab.* **34** 1771–1781. [MR2271481](#)
- BIRGÉ, L. (2006). Model selection via testing: An alternative to (penalized) maximum likelihood estimators. *Ann. Inst. Henri Poincaré Probab. Stat.* **42** 273–325. [MR2219712](#)
- BIRGÉ, L. and MASSART, P. (1993). Rates of convergence for minimum contrast estimators. *Probab. Theory Related Fields* **97** 113–150. [MR1240719](#)
- BIRGÉ, L. and MASSART, P. (1998). Minimum contrast estimators on sieves: Exponential bounds and rates of convergence. *Bernoulli* **4** 329–375. [MR1653272](#)
- BOUCHERON, S., LUGOSI, G. and MASSART, P. (2003). Concentration inequalities using the entropy method. *Ann. Probab.* **31** 1583–1614. [MR1989444](#)
- IBRAGIMOV, I. A. and KHAS’MINSKIĬ, R. Z. (1981). *Statistical Estimation: Asymptotic Theory. Applications of Mathematics* **16**. Springer-Verlag, New York-Berlin. Translated from the Russian by Samuel Kotz. [MR0620321](#)
- 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. Statist.* **3** 37–98. [MR0126903](#)
- LE CAM, L. and YANG, G. L. (2000). *Asymptotics in Statistics: Some Basic Concepts*, 2nd ed. Springer, New York. [MR1784901](#)
- MCCULLAGH, P. and NELDER, J. A. (1989). *Generalized Linear Models*, 2nd ed. Chapman & Hall, London.
- SPOKOINY, V. (2012a). Roughness penalty, Wilks phenomenon, and Bernstein–von Mises theorem. Unpublished manuscript. Available at [arXiv:1205.0498](https://arxiv.org/abs/1205.0498) [stat.ME].
- SPOKOINY, V. (2012b). Supplement to “Parametric estimation. Finite sample theory.” DOI: [10.1214/12-AOS1054SUPP](https://doi.org/10.1214/12-AOS1054SUPP).
- SPOKOINY, V., WANG, W. and HÄRDLE, W. (2012). Local quantile regression. Unpublished manuscript. Available at [arXiv:1208.5384](https://arxiv.org/abs/1208.5384) [math.ST].
- TALAGRAND, M. (1996). Majorizing measures: The generic chaining. *Ann. Probab.* **24** 1049–1103. [MR1411488](#)
- TALAGRAND, M. (2001). Majorizing measures without measures. *Ann. Probab.* **29** 411–417. [MR1825156](#)
- TALAGRAND, M. (2005). *The Generic Chaining: Upper and Lower Bounds of Stochastic Processes*. Springer, Berlin. [MR2133757](#)

VAN DE GEER, S. (1993). Hellinger-consistency of certain nonparametric maximum likelihood estimators. *Ann. Statist.* **21** 14–44. [MR1212164](#)

VAN DER VAART, A. W. and WELLNER, J. A. (1996). *Weak Convergence and Empirical Processes: With Applications to Statistics*. Springer, New York. [MR1385671](#)

WEIERSTRASS-INSTITUTE

MOHRENSTR. 39, 10117 BERLIN

GERMANY

AND

HUMBOLDT UNIVERSITY BERLIN

GERMANY

AND

MOSCOW INSTITUTE OF PHYSICS AND TECHNOLOGY

E-MAIL: [spokoiny@wias-berlin.de](mailto:spokoiny@wias-berlin.de)