

Penalized maximum likelihood estimation and effective dimension

Vladimir Spokoiny¹

Weierstrass-Institute and Humboldt University Berlin, Higher School of Economics, IITP RAS, MIPT, Mohrenstr. 39, 10117 Berlin, Germany.
E-mail: spokoiny@wias-berlin.de

Received 29 July 2014; revised 24 August 2015; accepted 5 October 2015

Abstract. This paper extends some prominent statistical results including *Fisher Theorem* and *Wilks phenomenon* to the penalized maximum likelihood estimation with a quadratic penalization. It appears that sharp expansions for the penalized MLE $\tilde{\theta}_G$ and for the penalized maximum likelihood can be obtained without involving any asymptotic arguments, the results only rely on smoothness and regularity properties of the of the considered log-likelihood function. The error of estimation is specified in terms of the effective dimension p_G of the parameter set which can be much smaller than the true parameter dimension and even allows an infinite dimensional functional parameter. In the i.i.d. case, the Fisher expansion for the penalized MLE can be established under the constraint “ p_G^2/n is small” while the remainder in the Wilks result is of order $\sqrt{p_G^3/n}$.

Résumé. Cet article généralise certains résultats statistiques importants dont le *Théorème de Fisher* et le *phénomène de Wilks* à l'estimation du maximum de vraisemblance pénalisée de façon quadratique. Il apparaît que des développements précis pour l'EMV pénalisée $\tilde{\theta}_G$ et le maximum de vraisemblance pénalisé peuvent être obtenus sans arguments asymptotiques, les résultats reposent alors seulement sur la régularité et les propriétés de la fonction de log-vraisemblance. L'erreur d'estimation est spécifiée en fonction de la dimension effective p_G de l'ensemble des paramètres qui peut être beaucoup plus petite que la véritable dimension et permet ainsi de considérer un cas infini dimensionnel. Dans le cas i.i.d., le développement de Fisher pour l'EMV pénalisée peut être établi sous la contrainte « p_G^2/n est petit » tandis que le reste dans le résultat de Wilks est d'ordre $\sqrt{p_G^3/n}$.

MSC: Primary 62F10; secondary 62J12; 62F25; 62H12

Keywords: Penalty; Wilks and Fisher expansions

1. Introduction

The Fisher and Wilks theorems belong to the short list of most fascinating results in the statistical theory. In particular, the Wilks result in its simple form claims that the likelihood ratio test statistic is close in distribution to the χ_p^2 distribution as the sample size increases, where p means the parameter dimension. So, the limiting distribution of this test statistic only depends on the dimension of the parameter space whatever the parametric model is. This explains why this result is sometimes called the *Wilks phenomenon*. This paper aims at reconsidering the mentioned results from different viewpoints. One important issue is that the presented results are stated for *finite samples*. There are only few general finite-sample results in statistical inference; see Boucheron and Massart [7] and references therein in context of i.i.d. modeling. The novel approach from Spokoiny [25] offered a general framework for a *finite sample*

¹The research was partly supported by the Russian Science Foundation grant (project 14-50-00150). Financial support by the German Research Foundation (DFG) through the Collaborative Research Center 649 “Economic Risk” is gratefully acknowledged.

theory, and the present paper makes a further step in this direction: the classical large sample results are extended to the finite sample case with *explicit and sharp* error bounds.

Another important point is a possible *model misspecification*. The classical parametric theory requires the parametric assumption to be exactly fulfilled. Any violation of the parametric specification may destroy the Fisher and Wilks results; cf. Huber [13]. This study admits from the very beginning that the parametric specification is probably wrong. This automatically extends the applicability of the proposed approach.

The further issue is the use of *penalization* for reducing the *model complexity*. If the parameter dimension is too large, the classical statistical results become almost intractable because the corresponding error is proportional to the dimension of parameter space. Sieve parametric approach is often used to replace the an infinite dimensional problem with a finite dimensional one; see e.g. Shen and Wong [24], Shen [23], Van de Geer [29], Birgé and Massart [4], Barron et al. [2], and references therein. Some asymptotic results for generalized regression models are available in Fan et al. [8].

Another standard way of reducing the complexity of the model is by introducing some penalty in the likelihood function. In this paper we focus on quadratic-type penalization. Roughness penalty approach provides a popular example; cf. Green and Silverman [12]. Koenker et al. [16] explained how roughness penalty works in context of quantile regression. Tikhonov regularization and ridge regression are the other examples which are often used in linear inverse problems. It is well known that the use of a penalization in context of an inverse problem provides regularization and uncertainty reduction at the same time. Our results show that the use of penalization indeed leads to some improvement in the obtained error bounds. Namely, one can replace the original parameter dimension p by the so called *effective dimension* p_G which can be much smaller than p . Even the case of a functional parameter θ with $p = \infty$ can be included. In this paper the penalty term is supposed to be given in advance. In general, a model selection procedure based on a proper choice of penalization is a high topic, one of the central in nonparametric statistics. We refer to Shen [23], Birgé and Massart [4], van de Geer [28] for the general models and to Birgé and Massart [5,6] for Gaussian model selection where one can find an extensive overview of the vast literature on this problem.

The final issue is the *critical parameter dimension* which is measured by the effective dimension p_G . The problem of statistical inference for models with growing parameter dimension is quite involved. There are some specific issues even if a simple linear or exponential model is considered, the results from Portnoy [19,20] requires “ p^2/n small” for asymptotic normality of the MLE. Depending on the considered problem and the model at hand, the conditions on the critical parameter dimension p may differ. For instance, Portnoy [22] obtained the Fisher and Wilks results for a generalized linear model under $p^{3/2}/n \rightarrow 0$, Mammen [18] established similar results for high-dimensional linear models. A general Wilks result can be stated under the condition that p^3/n is small; see e.g. Belloni and Chernozhukov [3]. Below we show that the conditions on the critical dimension in penalized ML estimation can be given in terms of the effective dimension p_G rather than the parameter dimension p . In particular, in the i.i.d. case, the Fisher expansion can be stated under “ p_G^2/n small” and “ p_G^3/n small” is sufficient for the Wilks result.

First we specify our set-up. Let \mathbf{Y} denote the observed data and \mathbb{P} mean their distribution. A general parametric assumption (PA) means that \mathbb{P} belongs to p -dimensional family $(\mathbb{P}_\theta, \theta \in \Theta \subseteq \mathbb{R}^p)$ dominated by a measure μ_0 . This family yields the log-likelihood function $L(\theta) = L(\mathbf{Y}, \theta) \stackrel{\text{def}}{=} \log \frac{d\mathbb{P}_\theta}{d\mu_0}(\mathbf{Y})$. The PA can be misspecified, so, in general, $L(\theta)$ is a *quasi log-likelihood*. The classical likelihood principle suggests to estimate θ by maximizing the function $L(\theta)$:

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

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

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

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 classical *Fisher Theorem* claims the expansion for the MLE $\tilde{\theta}$:

$$D(\tilde{\theta} - \theta^*) - \xi \xrightarrow{\mathbb{P}} 0,$$

where $D^2 = -\nabla^2 \mathbb{E}L(\theta^*)$ and $\xi \stackrel{\text{def}}{=} D^{-1} \nabla L(\theta^*)$. Under the correct model specification, D^2 is the total Fisher information matrix and the vector ξ is centered and standardized. So, it is asymptotically standard normal under general CLT conditions.

It is well known that many important properties of the quasi MLE $\tilde{\theta}$ like concentration or coverage probability can be described in terms of the *excess* or *quasi maximum likelihood* $L(\tilde{\theta}, \theta^*) \stackrel{\text{def}}{=} L(\tilde{\theta}) - L(\theta^*) = \max_{\theta \in \Theta} L(\theta) - L(\theta^*)$, which is the difference between the maximum of the process $L(\theta)$ and its value at the “true” point θ^* . The *Wilks phenomenon* claims that the distribution of the twice excess $2L(\tilde{\theta}, \theta^*)$ can be approximated by $\|\xi\|^2$ which is asymptotically χ_p^2 , where p is the dimension of the parameter space:

$$2L(\tilde{\theta}, \theta^*) - \|\xi\|^2 \xrightarrow{\mathbb{P}} 0, \quad \|\xi\|^2 \xrightarrow{w} \chi_p^2.$$

This fact is very attractive and yields asymptotic confidence and concentration sets as well as the limiting critical values for the likelihood ratio tests. However, practical applications of all mentioned results are limited: they require true parametric distribution, large samples and a fixed parameter dimension.

Modern applications stimulate a further extension of the classical theory beyond the classical parametric assumptions. Spokoiny [25] offers a general approach which appears to be very useful for such an extension. The whole approach is based on the following local bracketing result:

$$\mathbb{L}_\epsilon(\theta, \theta^*) - \diamond_\epsilon \leq L(\theta) - L(\theta^*) \leq \mathbb{L}_\epsilon(\theta, \theta^*) + \diamond_\epsilon, \quad \theta \in \Theta_0. \quad (1.2)$$

Here $\mathbb{L}_\epsilon(\theta, \theta^*)$ and $\mathbb{L}_\epsilon(\theta, \theta^*)$ are quadratic in $\theta - \theta^*$ expressions and Θ_0 is a local vicinity of the central point θ^* . This result can be viewed as an extension of the famous Le Cam *local asymptotic normality* (LAN) condition. The LAN condition considers just one quadratic process for approximating the log-likelihood $L(\theta)$. The use of bracketing with two different quadratic expressions allows one to keep control of the error terms \diamond_ϵ , \diamond_ϵ even for relatively large neighborhoods Θ_0 of θ^* while the LAN approach is essentially restricted to a root- n vicinity of θ^* . It also allows to incorporate a large parameter dimension and a model misspecification. However, the approach from Spokoiny [25] has natural limitations: the parameter dimension p cannot be too large. For instance, in the i.i.d. case, the error terms \diamond_ϵ and \diamond_ϵ are of order $\sqrt{p^3/n}$ which destroys the Wilks result if $p > n^{1/3}$.

A standard way of overcoming this difficulty is to impose a kind of smoothness assumption on the unknown parameter value θ^* . Here we discuss one general way to deal with such smoothness assumptions using a quadratic penalization. Section 2 offers a new approach to studying the properties of the penalized MLE which is based on a linear approximation of the gradient of the log-likelihood process. Compared to the bracketing approach (1.2), it allows to establish a Fisher type expansion for the penalized MLE under weaker conditions on the critical dimension of the problem. Another important novelty of the approach is the systematic use of the *effective dimension* \mathfrak{p}_G in place of the original dimension p of the parameter space. Usually \mathfrak{p}_G is much smaller than p . It is even possible to treat the case of a functional parameter if the effective dimension of the parameter set remains finite. Our main results include the Fisher and Wilks expansions for the penalized MLE. In the important special case of an i.i.d. model, the error term in the Wilks expansion is of order \mathfrak{p}_G^3/n , while the Fisher expansion requires \mathfrak{p}_G^2/n small.

Also we discuss an implication of these results to the bias-variance decomposition of the squared risk of the penalized MLE. In all our results, the error terms only depend on the effective dimension \mathfrak{p}_G .

The paper is organized as follows. Section 2 states the analog of Fisher and Wilks results for the penalized MLE procedure. Section 2.5 explains how the Fisher expansion can be used for obtaining the bias-variance decomposition of the quadratic risk of the penalized MLE. Section 3 specifies the general conditions and results to i.i.d. and generalized linear models. Appendix A collects some deviation bounds for general quadratic forms. Appendix B contains technical results on the maximum of a stochastic process under exponential moment conditions.

2. Fisher and Wilks theorems under quadratic penalization

Let $\text{pen}(\theta)$ be a penalty function on Θ . A big value of $\text{pen}(\theta)$ corresponds to a large degree of roughness or a small amount of smoothness of θ . The underlying assumption on the model is that the true value θ^* is smooth in the sense

that $\text{pen}(\theta^*)$ is relatively small. A penalized (quasi) MLE approach leads to maximizing the penalized log-likelihood:

$$\tilde{\theta} = \underset{\theta \in \Theta}{\operatorname{argmax}} \{L(\theta) - \text{pen}(\theta)\}.$$

Below we discuss an important special case of a quadratic penalty $\text{pen}(\theta) = \|G\theta\|^2/2$ for a given symmetric matrix G ; see e.g. Green and Silverman [12] or Koenker et al. [16] for particular examples. Denote

$$L_G(\theta) \stackrel{\text{def}}{=} L(\theta) - \|G\theta\|^2/2,$$

$$\tilde{\theta}_G \stackrel{\text{def}}{=} \underset{\theta \in \Theta}{\operatorname{argmax}} L_G(\theta).$$

The use of a penalty changes the target of estimation which is now defined as

$$\theta_G^* \stackrel{\text{def}}{=} \underset{\theta \in \Theta}{\operatorname{argmax}} \mathbb{E}L_G(\theta). \quad (2.1)$$

So, introducing a penalty leads to some estimation bias: the new target θ_G^* may be different from θ^* . At the same time, similarly to linear modeling, the use of penalization reduces the variability of the estimate $\tilde{\theta}_G$ and improves its concentration properties. An interesting question is the total impact and a possible gain of using the penalized procedure. A preliminary answer is that the penalty term $\|G\theta^*\|^2$ at the true point should not be too large relative to the squared error of estimation for the penalized model. This rule is known under the name “bias-variance trade-off.”

Another important message of this study is that the use of penalization allows to reduce the parameter dimension to the *effective dimension* which characterizes the entropy of the penalized parameter space. The resulting confidence and concentration sets depend on the effective dimension rather than on the real parameter dimension and they can be much more narrow than in the non-penalized case.

The principle steps of the study are as follows. The *concentration* step claims that the penalized MLE $\tilde{\theta}_G$ is concentrated in a local vicinity $\Theta_{0,G}(r_G)$ of the point θ_G^* . It is based on the upper function method which bounds the penalized log-likelihood $L_G(\theta)$ from above by a deterministic function. Theorem 2.1 states that $\tilde{\theta}_G$ belongs to the local set $\Theta_{0,G}(r_G)$ with a dominating probability, and this local set can be much smaller than the similar set for the non-penalized results. As the next step, Spokoiny [25] applied the *bracketing* approach to bound from above and from below the log-likelihood process $L(\theta)$ by two quadratic in $\theta - \theta^*$ expressions. Here the bracketing step is changed essentially by using a local linear approximation of the vector gradient process $\nabla L(\theta)$. This helps to get a sharper bound on the error of approximation and improve the quality of the Fisher expansion. Similarly to Spokoiny [25], the obtained results are stated for finite samples and do not involve any asymptotic arguments. An advantage of the proposed approach is that it combines an accurate local approximation with rather rough large deviation arguments and allows one to obtain usual asymptotic statements including asymptotic normality of the penalized MLE; see Section 3.1 for the i.i.d. case.

2.1. Effective dimension

Let V^2 be the matrix shown in condition (E_0G) in Section 2.2. Typically $V^2 = \text{Var}\{\nabla L(\theta_G^*)\}$ and this matrix measures the local variability of the process $L_G(\cdot)$. Let also D_G^2 be a penalized information matrix defined as

$$D_G^2 = -\nabla^2 \mathbb{E}L_G(\theta_G^*) = D^2 + G^2$$

with $D^2 = -\nabla^2 \mathbb{E}L(\theta_G^*)$. One can redefine $D^2 = -\nabla^2 \mathbb{E}L(\theta^*)$ under condition (\mathcal{L}_0G) below and the so called small modeling bias condition; see Section 2.5. The *effective dimension* \mathfrak{p}_G is defined as the trace of the matrix $B_G \stackrel{\text{def}}{=} D_G^{-1} V^2 D_G^{-1}$:

$$\mathfrak{p}_G \stackrel{\text{def}}{=} \text{tr}(B_G). \quad (2.2)$$

Below we show that the use of penalization enables us to replace the original dimension p in our risk bounds with the effective dimension \mathfrak{p}_G which can be much smaller than p depending on relations between the matrices D^2 , V^2 , and G^2 .

In our results the value \mathfrak{p}_G will be used via another quantity $z(B_G, \mathfrak{x})$ which also depends on a fixed constant \mathfrak{x} and for moderate values of \mathfrak{x} can be defined as

$$z(B_G, \mathfrak{x}) = \sqrt{\mathfrak{p}_G} + \sqrt{2\mathfrak{x}\lambda_G}, \tag{2.3}$$

where $\lambda_G \stackrel{\text{def}}{=} \lambda_{\max}(B_G)$ is the largest eigenvalue of B_G ; see (A.4) for a precise definition.

Now we present a couple of typical examples of using the quadratic penalty: blockwise penalization and estimation under a Sobolev smoothness constraint. For simplicity of presentation we assume that $V^2 = D^2 = n\mathbb{I}_p$, while G^2 is diagonal with non-decreasing eigenvalues g_j^2 . Then $D_G^2 = D^2 + G^2 = \text{diag}\{n + g_1^2, \dots, n + g_p^2\}$. It holds that $B_G = \text{diag}\{(1 + n^{-1}g_1^2)^{-1}, \dots, (1 + n^{-1}g_p^2)^{-1}\}$, and we apply (2.2) for computing the effective dimension \mathfrak{p}_G .

Block penalization

Consider the case when G is of a simple two-block structure: $G = \text{diag}\{0, G_1\}$. Many blocks can be considered in the similar way. The first block of dimension p_0 corresponds to the unconstrained part of the parameter vector while the second block of dimension p_1 corresponds to the low energy component. An interesting question is the minimal penalization G_1 making the impact of the low energy part inessential. Assume for simplicity that $G_1 = g\mathbb{I}_{p_1}$. Then

$$\mathfrak{p}_G = \text{tr } B_G = p_0 + p_1/(1 + n^{-1}g^2).$$

One can see that the impact of the second block G_1 in the effective dimension is inessential if $g^2/n \gg p_1/p_0$.

Sobolev smoothness constraint

Consider the case with $D^2 = V^2 = n\mathbb{I}_p$ and $G^2 = \text{diag}\{g_1^2, \dots, g_p^2\}$ with $g_j = Lj^\beta$ for $\beta > 1/2$. The value β is usually considered as the Sobolev smoothness parameter. It holds

$$\mathfrak{p}_G = \sum_{j=1}^p \frac{1}{1 + L^2 j^{2\beta}/n}.$$

Define also the index \mathfrak{p}_e as the largest j satisfying $L^2 j^{2\beta} \leq n$. It is straightforward to see that $\beta > 1/2$ yields $\mathfrak{p}_G \leq C(\beta, L)\mathfrak{p}_e$ for a constant $C(\beta, L)$ depending on β, L only.

Linear inverse problem

The next example corresponds to the case of a linear inverse problem. Assume for simplicity of notation the sequence space representation, the noise is inhomogeneous with increasing eigenvalues $V^2 = \text{diag}\{v_1^2, \dots, v_p^2\}$ and the information matrix D^2 is proportional to identity, that is, $D^2 = n\mathbb{I}_p$. Then the effective dimension is given by the sum

$$\mathfrak{p}_G = \sum_{j=1}^p \frac{v_j^2}{n + g_j^2}.$$

To keep the effective dimension small, one has to compensate the increase of the eigenvalues v_j^2 by the penalization g_j^2 .

2.2. Conditions

This section presents the list of conditions which are similar to ones from the non-penalized case in Spokoiny [25]. However, the use of penalization leads to some change in each condition. Most important fact is that the use of penalization helps to state the large deviation (LD) result for much smaller local neighborhoods than in the non-penalized case. Spokoiny [25] presented the LD result for local sets of the form $\Theta_0(\mathfrak{x}) = \{\theta : \|V(\theta - \theta^*)\| \leq \mathfrak{x}\}$ with

a proper $\varkappa \asymp p^{1/2}$. Now we redefine this set by using D_G^2 in place of V^2 and θ_G^* in place of θ^* :

$$\Theta_{0,G}(\varkappa) \stackrel{\text{def}}{=} \{\theta: \|D_G(\theta - \theta_G^*)\| \leq \varkappa\}.$$

Moreover, the radius \varkappa can be selected of order $p_G^{1/2}$, which can be very useful for large or infinite p . Our conditions mainly assume some regularity and smoothness of the penalized log-likelihood process $L_G(\theta)$. The first condition states some smoothness properties of the expected log-likelihood $\mathbb{E}L_G(\theta)$ as a function of θ in a vicinity $\Theta_{0,G}(\varkappa)$ of θ_G^* . More precisely, it effectively means that the expected log-likelihood $\mathbb{E}L(\theta)$ is twice continuously differentiable on the local set $\Theta_{0,G}(\varkappa)$.

Below each condition is given in penalized and non-penalized form for the sake of comparison. Already now it is worth saying that the use of penalization helps to relax most of conditions. Define

$$\mathbb{F}_G(\theta) \stackrel{\text{def}}{=} -\nabla^2 \mathbb{E}L_G(\theta) = -\nabla^2 \mathbb{E}L(\theta) + G^2.$$

Then $D_G^2 = \mathbb{F}_G(\theta_G^*)$. The conditions involve a radius \varkappa_G which separates the local zone and the zone of large deviations. This value will be made precise in Theorem 2.1.

First we consider the stochastic component of the log-likelihood process $L_G(\theta)$ which is the same as in the non-penalized case:

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

We assume that it is twice differentiable and denote by $\nabla \zeta(\theta)$ its gradient and by $\nabla^2 \zeta(\theta)$ its Hessian matrix. The next two conditions are to ensure that the random vector $\nabla \zeta(\theta_G^*)$ and the random processes $\nabla^2 \zeta(\theta)$ are stochastically bounded with exponential moments. The conditions involve a $p \times p$ -matrix V which normalizes the vector $\nabla \zeta(\theta_G^*)$, and a similar matrix V_2 normalizing $\nabla^2 \zeta(\theta)$.

(E_0G) There exist a positively semi-definite symmetric matrix V^2 , and constants $\mathfrak{g} > 0$, $\nu_0 \geq 1$ such that $\text{Var}\{\nabla \zeta(\theta_G^*)\} \leq V^2$ and

$$\sup_{\boldsymbol{y} \in \mathbb{R}^p} \log \mathbb{E} \exp \left\{ \lambda \frac{\boldsymbol{y}^\top \nabla \zeta(\theta_G^*)}{\|V \boldsymbol{y}\|} \right\} \leq \frac{\nu_0^2 \lambda^2}{2}, \quad |\lambda| \leq \mathfrak{g}.$$

(E_2G) There exist a positively semi-definite symmetric matrix V_2^2 and a value $\omega > 0$ such that it holds for any $\theta \in \Theta_{0,G}(\varkappa_0)$:

$$\sup_{\boldsymbol{y}_1, \boldsymbol{y}_2 \in \mathbb{R}^p} \log \mathbb{E} \exp \left\{ \frac{\lambda}{\omega} \frac{\boldsymbol{y}_1^\top \nabla^2 \zeta(\theta) \boldsymbol{y}_2}{\|V_2 \boldsymbol{y}_1\| \cdot \|V_2 \boldsymbol{y}_2\|} \right\} \leq \frac{\nu_0^2 \lambda^2}{2}, \quad |\lambda| \leq \mathfrak{g}.$$

Below we only need that the constant $\mathfrak{g}(\varkappa)$ is larger than C_{pG} for a fixed constant C . This allows to reduce the condition to the case with a fixed \mathfrak{g} which does not depend on the distance \varkappa .

Their non-penalized versions are almost identical: one has to replace θ_G^* with θ^* and $\Theta_{0,G}(\varkappa)$ with $\Theta_0(\varkappa)$.

$$(E_0) \quad \sup_{\boldsymbol{y} \in \mathbb{R}^p} \log \mathbb{E} \exp \left\{ \lambda \frac{\boldsymbol{y}^\top \nabla \zeta(\theta^*)}{\|V \boldsymbol{y}\|} \right\} \leq \frac{\nu_0^2 \lambda^2}{2}, \quad |\lambda| \leq \mathfrak{g}.$$

$$(E_2) \quad \sup_{\boldsymbol{y}_1, \boldsymbol{y}_2 \in \mathbb{R}^p} \log \mathbb{E} \exp \left\{ \frac{\lambda}{\omega} \frac{\boldsymbol{y}_1^\top \nabla^2 \zeta(\theta) \boldsymbol{y}_2}{\|V_2 \boldsymbol{y}_1\| \cdot \|V_2 \boldsymbol{y}_2\|} \right\} \leq \frac{\nu_0^2 \lambda^2}{2}, \quad |\lambda| \leq \mathfrak{g}(\varkappa).$$

The conditions (E_0) and (E_0G) are very similar while (E_2G) is restricted to the vicinity $\Theta_{0,G}(\varkappa)$ which can be much smaller than $\Theta_0(\varkappa)$.

The *identifiability condition* relates the matrices V^2 and V_2^2 and to D_G^2 .

($\mathcal{I}G$) There is a constant $\alpha_G > 0$ such that

$$\alpha_G^2 D_G^2 \geq V^2, \quad \alpha_G^2 D_G^2 \geq V_2^2.$$

In the non-penalized case of Spokoiny [25], this condition reads as

$$(\mathcal{I}) \quad \alpha^2 D^2 \geq V^2 \text{ with } D^2 = -\nabla^2 \mathbb{E}L(\boldsymbol{\theta}^*).$$

Therefore, the use of regularization helps to improve the identifiability in the regularized problem relative to the non-penalized one as $D^2 \leq D_G^2$.

Finally, we consider the expected log-likelihood $\mathbb{E}L_G(\boldsymbol{\theta})$. The local condition requires that it is nearly quadratic in the vicinity $\Theta_{0,G}(r_G)$ of $\boldsymbol{\theta}_G^*$ while the global condition assumes a linear growth in the complement of this vicinity. Here and below $\|A\|_{\text{op}}$ means the operator norm of a matrix A .

(\mathcal{L}_0G) For each $r \leq r_G$, there is a constant $\delta_G(r) \leq 1/2$ such that

$$\|D_G^{-1} \mathbb{F}_G(\boldsymbol{\theta}) D_G^{-1} - \mathbb{I}_p\|_{\text{op}} \leq \delta_G(r), \quad \boldsymbol{\theta} \in \Theta_{0,G}(r). \tag{2.4}$$

Under condition (\mathcal{L}_0G), it follows from the second order Taylor expansion at $\boldsymbol{\theta}_G^*$:

$$| -2\mathbb{E}L_G(\boldsymbol{\theta}, \boldsymbol{\theta}_G^*) - \|D_G(\boldsymbol{\theta} - \boldsymbol{\theta}_G^*)\|^2 | \leq \delta_G(r) \|D_G(\boldsymbol{\theta} - \boldsymbol{\theta}_G^*)\|^2, \quad \boldsymbol{\theta} \in \Theta_{0,G}(r).$$

A non-penalized version of (2.4) claims a similar approximation of $\mathbb{F}(\boldsymbol{\theta}) = -\nabla^2 \mathbb{E}L(\boldsymbol{\theta})$ by $D^2 \stackrel{\text{def}}{=} \mathbb{F}(\boldsymbol{\theta}^*)$ in the vicinity $\Theta_0(r_0)$ centered at $\boldsymbol{\theta}^*$ instead of $\boldsymbol{\theta}_G^*$:

$$(\mathcal{L}_0) \quad \|D^{-1} \mathbb{F}(\boldsymbol{\theta}) D^{-1} - \mathbb{I}_p\|_{\text{op}} \leq \delta(r_0), \quad \boldsymbol{\theta} \in \Theta_0(r_0) = \{\boldsymbol{\theta} : \|D(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\| \leq r_0\}.$$

As the quadratic penalty $\|G\boldsymbol{\theta}\|^2$ does not change the smoothness properties of the expected contrast $\mathbb{E}L_G(\boldsymbol{\theta})$, the conditions (\mathcal{L}_0G) and (\mathcal{L}_0) are essentially equivalent provided that the points $\boldsymbol{\theta}^*$ and $\boldsymbol{\theta}_G^*$ are not too far from each others.

The local condition (\mathcal{L}_0G) describes the behavior of $\mathbb{E}L_G(\boldsymbol{\theta})$ within $\Theta_{0,G}(r_G)$. In particular, $\mathbb{E}L_G(\boldsymbol{\theta}_G^*) - \mathbb{E}L_G(\boldsymbol{\theta}) \approx r_G^2/2$ on the boundary of this local set. The global condition means that $\mathbb{E}L_G(\boldsymbol{\theta}_G^*) - \mathbb{E}L_G(\boldsymbol{\theta})$ can be lower bounded by a linear function on the complement of this set.

($\mathcal{L}G$) For each $\boldsymbol{\theta}$ with $r = \|D_G(\boldsymbol{\theta} - \boldsymbol{\theta}_G^*)\| \geq r_G$

$$\mathbb{E}L_G(\boldsymbol{\theta}_G^*) - \mathbb{E}L_G(\boldsymbol{\theta}) \geq (1 - \delta(r_G)) \left(r_G r - \frac{r_G^2}{2} \right) + C_1 r^2, \tag{2.5}$$

for a small constant C_1 ; see Theorem 2.1 below for a precise bound.

A non-penalized version of this condition is obtained by letting $G^2 = 0$.

$$(\mathcal{L}) \quad \mathbb{E}L(\boldsymbol{\theta}^*) - \mathbb{E}L(\boldsymbol{\theta}) \geq (1 - \delta(r_0)) (r_0 r - r_0^2/2) + C_1 r^2 \quad \text{for } r = \|D(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|.$$

Remark 2.1. Conditions (\mathcal{L}_0G) and ($\mathcal{L}G$) can be effectively checked if the function $f(\boldsymbol{\theta}) \stackrel{\text{def}}{=} -\mathbb{E}L_G(\boldsymbol{\theta})$ is smooth and convex in $\boldsymbol{\theta}$. Continuity of the second derivative $\nabla^2 f(\boldsymbol{\theta})$ in $\Theta_{0,G}(r_G)$ implies (\mathcal{L}_0G). Convexity of f implies for any $\boldsymbol{\theta}^\circ = \boldsymbol{\theta}_G^* + \rho(\boldsymbol{\theta} - \boldsymbol{\theta}_G^*)$ with $\rho = r_G / \|D_G(\boldsymbol{\theta} - \boldsymbol{\theta}_G^*)\| \leq 1$

$$f(\boldsymbol{\theta}) \geq f(\boldsymbol{\theta}^\circ) + (\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)^\top \nabla f(\boldsymbol{\theta}^\circ) \geq f(\boldsymbol{\theta}^\circ) + \{D_G(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)\}^\top D_G^{-1} \nabla f(\boldsymbol{\theta}^\circ).$$

Condition (\mathcal{L}_0G) implies in view of $\|D_G(\boldsymbol{\theta}^\circ - \boldsymbol{\theta}_G^*)\| = r_G$ that $f(\boldsymbol{\theta}^\circ) \geq (1 - \delta)r_G^2/2$ and $\nabla^2 f(\boldsymbol{\theta}^\circ) \geq (1 - \delta)D_G^2$ for $\delta = \delta(r_G)$. As $\|D_G(\boldsymbol{\theta}^\circ - \boldsymbol{\theta})\| = r - r_G$, we conclude that

$$f(\boldsymbol{\theta}) \geq (1 - \delta) \left\{ \frac{r_G^2}{2} + r_G(r - r_G) \right\}$$

and (2.5) follows for $C_1 = 0$. The case $C_1 > 0$ can be similarly checked under a strong convexity of $-\mathbb{E}L_G(\boldsymbol{\theta})$. In the case of linear or generalized linear models, one can use $C_1 = 0$, in regular situations, it suffices that ($\mathcal{L}G$) holds with C_1 of order $n^{-1/2}$.

We briefly comment on examples for which the conditions can be easily verified. Conditions (E_0G) and (E_2G) require some exponential moments of log-likelihood ratio and its derivatives. Usually one assumes some finite moments of the normalized increments of the likelihood function; cf. Ibragimov and Khas'minskij [14], Chapter 2. Our conditions (E_0G) and (E_2G) a bit more restrictive but it allows one to obtain some finite sample bounds. Note that majority of finite samples results are stated under Gaussian or sub-Gaussian stochastic errors. The sub-Gaussian case corresponds to $g = \infty$ in (E_0G) and (E_2G) . Our results apply for sub-exponential errors with $g < \infty$ as well. Condition (\mathcal{L}_0G) only requires some regularity of the considered parametric family and is not restrictive. As already mentioned, condition $(\mathcal{L}G)$ can be easily checked if $\mathbb{E}L(\theta)$ is smooth and concave in θ . It is also easy to verify if $\nabla^2 \mathbb{E}L_G(\theta)$ is bounded from below by a positive matrix.

The i.i.d. case will be considered in details below in Section 3.1. Section 3.6 demonstrates how the conditions can be checked for generalized linear models in terms of design regularity, smoothness of the link function, and exponential moment conditions on the errors. The regression and generalized regression models are included as well; cf. Ghosal [9,10] or Kim [15]. Spokoiny [25], Section 5.2, argued that (E_0G) and (E_2G) are fulfilled when regression errors fulfill some exponential moments condition. If this condition is too restrictive and a more stable (robust) estimation procedure is desirable, one can apply the LAD-type contrast leading to median regression. Spokoiny [25], Section 5.3, showed for the case of linear median regression that all the required conditions are fulfilled automatically if the sample size n exceeds Cp for a fixed constant C . Spokoiny et al. [26] applied this approach for local polynomial quantile regression. Zaitsev et al. [31] applied the approach to the problem of regression with Gaussian process where the unknown parameters enter in the likelihood function in a rather complicated way. We conclude that the imposed conditions are quite general and can be verified for many classical examples met in the statistical literature.

2.3. Concentration and a large deviation bound

This section demonstrates that the use of the penalty term helps to strengthen the concentration properties of the penalized quasi maximum likelihood estimator (qMLE) $\tilde{\theta}_G$. Namely, we show that $\tilde{\theta}_G$ belongs with a dominating probability to a set $\Theta_{0,G}(r_G)$ which can be much smaller than a similar set from the non-penalized case; see Remark 2.2. All our results involve a value x . We say that a generic random set $\Omega(x)$ is of a *dominating probability* if $\mathbb{P}(\Omega(x)) \geq 1 - Ce^{-x}$ for a fixed constant C like 1 or 2. We also use two growing functions $z(B_G, x)$ and $\mathfrak{z}_{\mathbb{H}}(x)$ of the argument x . The functions $z(B_G, x)$ already mentioned in (2.3) and it describes the quantiles of the norm of the normalized score vector ξ_G ; see (2.8) below. The formal definition of $z(B_G, x)$ is given in (A.4). The function $\mathfrak{z}_{\mathbb{H}}(x)$ is related to the penalized entropy of the parameter space and it is given by (B.3). In typical situations one can use the upper bounds $z^2(B_G, x) \leq C(p_G + x)$ and $\mathfrak{z}_{\mathbb{H}}^2(x) \leq C(p_G + x)$ for both functions.

Theorem 2.1. *Let (E_0G) , (E_2G) , $(\mathcal{I}G)$, (\mathcal{L}_0G) , and $(\mathcal{L}G)$ hold with*

$$\{1 - \delta(r_G)\}r_G \geq 2z(B_G, x), \quad (2.6)$$

where $z(B_G, x)$ is from (A.4), and let the constant C_1 in $(\mathcal{L}G)$ satisfy

$$C_1 \geq \sup_{r > r_G} \varrho_G(r, x), \quad \text{with } \varrho_G(r, x) \stackrel{\text{def}}{=} \sqrt{8}v_0 a_G \mathfrak{z}_{\mathbb{H}}(x + \log(2r/r_G))\omega$$

with the function $\mathfrak{z}_{\mathbb{H}}(x)$ given by (B.3). Then

$$\mathbb{P}(\tilde{\theta}_G \notin \Theta_{0,G}(r_G)) \leq 3e^{-x}. \quad (2.7)$$

Remark 2.2. *This result explains a proper r_G ensuring (2.7). In the non-penalized case of Spokoiny [25], a similar condition reads as $r_0 \geq C(\sqrt{p} + \sqrt{2x})$, so the use of penalization helps to improve the concentration properties of the penalized MLE. We conclude that the use of penalization leads to weaker conditions and to a stronger concentration property. The only problem is that the corresponding estimate $\tilde{\theta}_G$ concentrates around θ_G^* instead of θ^* . This can yield a bias effect; see Section 2.5 below.*

Proof of Theorem 2.1. By definition $\sup_{\boldsymbol{\theta} \in \Theta_{0,G}(r_G)} L_G(\boldsymbol{\theta}, \boldsymbol{\theta}_G^*) \geq 0$. So, it suffices to check that $L_G(\boldsymbol{\theta}, \boldsymbol{\theta}_G^*) < 0$ for all $\boldsymbol{\theta} \in \Theta \setminus \Theta_{0,G}(r_G)$. The proof is based on the following bound: for each r

$$\mathbb{P}\left(\sup_{\boldsymbol{\theta} \in \Theta_{0,G}(r)} |\zeta(\boldsymbol{\theta}, \boldsymbol{\theta}_G^*) - (\boldsymbol{\theta} - \boldsymbol{\theta}_G^*)^\top \nabla \zeta(\boldsymbol{\theta}_G^*)| \geq \sqrt{8} v_0 \alpha_G \mathfrak{J}_{\mathbb{H}}(\mathbf{x}) \omega r^2\right) \leq e^{-x}.$$

This bound follows from Theorem B.15; see (2.18) for more details. It implies by Theorem B.3 with $\rho = 1/2$ on a set of dominating probability at least $1 - e^{-x}$ that for all $r \geq r_G$ and all $\boldsymbol{\theta}$ with $\|D_G(\boldsymbol{\theta} - \boldsymbol{\theta}_G^*)\| \leq r$

$$|\zeta(\boldsymbol{\theta}, \boldsymbol{\theta}_G^*) - (\boldsymbol{\theta} - \boldsymbol{\theta}_G^*)^\top \nabla \zeta(\boldsymbol{\theta}_G^*)| \leq \varrho_G(r, \mathbf{x}) r^2,$$

where $\varrho_G(r, \mathbf{x}) = v_0 \alpha_G \mathfrak{J}_{\mathbb{H}}(\mathbf{x} + \log(2r/r_G)) \omega$. The use of $\nabla \mathbb{E} L_G(\boldsymbol{\theta}_G^*) = 0$ yields

$$\sup_{\boldsymbol{\theta} \in \Theta_{0,G}(r)} |L_G(\boldsymbol{\theta}, \boldsymbol{\theta}_G^*) - \mathbb{E} L_G(\boldsymbol{\theta}, \boldsymbol{\theta}_G^*) - (\boldsymbol{\theta} - \boldsymbol{\theta}_G^*)^\top \nabla L_G(\boldsymbol{\theta}_G^*)| \leq \varrho_G(r, \mathbf{x}) r^2.$$

Also the vector $\boldsymbol{\xi}_G = D_G^{-1} \nabla L_G(\boldsymbol{\theta}_G^*) = D_G^{-1} \nabla \zeta(\boldsymbol{\theta}_G^*)$ can be bounded with a dominating probability: by Theorem A.1 $\mathbb{P}(\|\boldsymbol{\xi}_G\| \geq z(B_G, \mathbf{x})) \leq 2e^{-x}$. We ignore here the negligible term Ce^{-x^c} . The condition $\|\boldsymbol{\xi}_G\| \leq z(B_G, \mathbf{x})$ implies for each $r \geq r_G$

$$\begin{aligned} & \sup_{\boldsymbol{\theta} \in \Theta_{0,G}(r)} |(\boldsymbol{\theta} - \boldsymbol{\theta}_G^*)^\top \nabla L_G(\boldsymbol{\theta}_G^*)| \\ & \leq \sup_{\boldsymbol{\theta} \in \Theta_{0,G}(r)} \|D_G(\boldsymbol{\theta} - \boldsymbol{\theta}_G^*)\| \times \|D_G^{-1} \nabla \zeta(\boldsymbol{\theta}_G^*)\| = r \|\boldsymbol{\xi}_G\| \leq z(B_G, \mathbf{x}) r. \end{aligned}$$

Condition $(\mathcal{L}G)$ implies for each $\boldsymbol{\theta}$ with $\|D_G(\boldsymbol{\theta} - \boldsymbol{\theta}_G^*)\| = r > r_0$ that

$$-\mathbb{E} L_G(\boldsymbol{\theta}, \boldsymbol{\theta}_G^*) > (1 - \delta) \left(r_0 r - \frac{1}{2} r^2 \right) + C_1 r^2 \geq z(B_G, \mathbf{x}) r + \varrho_G(r, \mathbf{x}) r^2$$

provided that $r_0 \geq 2(1 - \delta)^{-1} z(B_G, \mathbf{x})$ and $C_1 \geq \varrho_G(r, \mathbf{x})$. This ensures that $L_G(\boldsymbol{\theta}, \boldsymbol{\theta}_G^*) < 0$ for all $\boldsymbol{\theta} \notin \Theta_{0,G}(r_G)$ with a dominating probability. \square

2.4. Wilks and Fisher expansions

This section collects the main results of the paper. Let $\boldsymbol{\theta}_G^*$ be the point of concentration from (2.1) and let $\zeta(\boldsymbol{\theta}) = L_G(\boldsymbol{\theta}) - \mathbb{E} L_G(\boldsymbol{\theta}) = L(\boldsymbol{\theta}) - \mathbb{E} L(\boldsymbol{\theta})$. Define a random p -vector

$$\boldsymbol{\xi}_G \stackrel{\text{def}}{=} D_G^{-1} \nabla \zeta(\boldsymbol{\theta}_G^*) = D_G^{-1} \{ \nabla L(\boldsymbol{\theta}_G^*) - G^2 \boldsymbol{\theta}_G^* \}. \quad (2.8)$$

Theorem 2.2. Suppose that r_G is selected to ensure (2.6). Suppose also that the conditions (E_0G) , (E_2G) , $(\mathcal{I}G)$ hold. On a random set $\Omega(\mathbf{x})$ of a dominating probability at least $1 - 4e^{-x}$, it holds

$$\|D_G(\tilde{\boldsymbol{\theta}}_G - \boldsymbol{\theta}_G^*) - \boldsymbol{\xi}_G\| \leq \diamond_G(\mathbf{x}), \quad (2.9)$$

where $\diamond_G(\mathbf{x})$ is given by

$$\diamond_G(\mathbf{x}) \stackrel{\text{def}}{=} \{ \delta_G(r_G) + \sqrt{8} v_0 \alpha_G \mathfrak{J}_{\mathbb{H}}(\mathbf{x}) \omega \} r_G \quad (2.10)$$

for $\mathfrak{J}_{\mathbb{H}}(\mathbf{x})$ given by (B.3).

The proof of this and the next result is based on a linear expansion of the gradient $\nabla L_G(\boldsymbol{\theta})$ and will be given in Section 2.6.

Now we present a result on the excess $L_G(\tilde{\boldsymbol{\theta}}_G, \boldsymbol{\theta}_G^*) = L_G(\tilde{\boldsymbol{\theta}}_G) - L_G(\boldsymbol{\theta}_G^*)$. The classical Wilks result claims that the twice excess is nearly χ_p^2 . Our result describes the quality of its approximation by a quadratic form $\|\boldsymbol{\xi}_G\|^2$.

Theorem 2.3. *Suppose that (\mathcal{L}_0G) , (E_0G) , and (E_2G) hold. Suppose also that \varkappa_G is selected to ensure (2.6). On a random set $\Omega(\mathbf{x})$ of a dominating probability at least $1 - 5e^{-\varkappa}$, it holds with $\diamond_G(\mathbf{x})$ from (2.10)*

$$|2L_G(\tilde{\theta}_G, \theta_G^*) - \|\xi_G\|^2| \leq 2\varkappa_G \diamond_G(\mathbf{x}) + \diamond_G^2(\mathbf{x}), \tag{2.11}$$

$$\left| \sqrt{2L_G(\tilde{\theta}_G, \theta_G^*)} - \|\xi_G\| \right| \leq 3\diamond_G(\mathbf{x}). \tag{2.12}$$

One can see that the Fisher expansion (2.9) and the square root Wilks expansion (2.12) require $\diamond_G(\mathbf{x})$ small, while the standard Wilks expansion (2.11) is accurate if $\varkappa_G \diamond_G(\mathbf{x})$ is small. This makes some difference if the parameter dimension is large. Below we address this question for the important special case of an i.i.d. likelihood.

The classical Fisher and Wilks results include some statements about the limiting behavior of the vector ξ_G and of the quadratic form $\|\xi_G\|^2$. In the i.i.d. case, one can easily show that the vector ξ_G is asymptotically standard normal as $n \rightarrow \infty$; see Section 3.5 below. However, it is well known that the convergence of $\|\xi_G\|^2$ to the χ^2 -distribution is quite slow even in the case of a fixed dimension p . For finite sample inference, we recommend to combine the approximations (2.9) to (2.12) with any resampling technique which mimics the specific behavior of the quadratic form $\|\xi_G\|^2$; see Spokoiny and Zhilova [27].

2.5. Quadratic risk bound and modeling bias

This section demonstrates the applicability of the obtained general results to bounding the quadratic risk of estimation. For the penalized MLE $\tilde{\theta}_G$ of the parameter θ , consider the quadratic loss of estimation $\|W(\tilde{\theta}_G - \theta^*)\|^2$ for a given non-negative symmetric matrix W . A special case includes the usual quadratic loss $\|\tilde{\theta}_G - \theta^*\|^2$. Here the point $\theta^* \in \Theta$ is a proxy for the true parameter value which describes the best parametric fit of the true measure \mathbb{P} by the family (\mathbb{P}_θ) :

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

The use of penalization $\|G\theta\|^2/2$ introduces some estimation bias: the penalized MLE $\tilde{\theta}_G$ estimates θ_G^* from (2.1) rather than θ^* . The value $\|W(\theta^* - \theta_G^*)\|^2$ is called the modeling bias and it describes the modeling error caused by using the penalization. The variance term $\|W(\tilde{\theta}_G - \theta_G^*)\|^2$ describes the error *within the penalized model*, and it can be studied with the help of the Fisher expansion of Theorem 2.2: $\|D_G(\tilde{\theta}_G - \theta_G^*) - \xi_G\| \leq \diamond_G(\mathbf{x})$ on a set $\Omega(\mathbf{x})$ of dominating probability for $\xi_G = D_G^{-1} \nabla \zeta(\theta_G^*)$. This yields the following result on $\Omega(\mathbf{x})$:

$$\|D_G(\tilde{\theta}_G - \theta^* - b_G) - \xi_G\| \leq \diamond_G(\mathbf{x})$$

with the *bias* $b_G = \theta_G^* - \theta^*$. For any positive symmetric $p \times p$ matrix W satisfying $W^2 \leq D_G^2$, it implies the probability bound for the squared loss

$$\|W(\tilde{\theta}_G - \theta^*)\| = \|Wb_G + WD_G^{-1}\xi_G\| \pm \diamond_G(\mathbf{x}).$$

One can see that analysis of the quadratic risk of the penalized MLE $\tilde{\theta}_G$ can be reduced to the analysis of $\|Wb_G + WD_G^{-1}\xi_G\|^2$. Now we consider an implication of this bound to the squared risk $\mathbb{E}\|W(\tilde{\theta}_G - \theta^*)\|^2$. The use of the identity $\mathbb{E}\nabla \zeta(\theta_G^*) = 0$ and $\operatorname{Var}(\nabla \zeta(\theta_G^*)) \leq V^2$ yields

$$\begin{aligned} \mathbb{E}\|Wb_G + WD_G^{-1}\xi_G\|^2 &= \|Wb_G\|^2 + \mathbb{E}\|WD_G^{-2}\nabla \zeta(\theta_G^*)\|^2 \\ &= \|Wb_G\|^2 + \operatorname{tr}(WD_G^{-2}\operatorname{Var}\{\nabla \zeta(\theta_G^*)\}D_G^{-2}W) \\ &\leq \|Wb_G\|^2 + \operatorname{tr}(WD_G^{-2}V^2D_G^{-2}W). \end{aligned}$$

Denote $\mathcal{X}_G \stackrel{\text{def}}{=} \operatorname{tr}(WD_G^{-2}V^2D_G^{-2}W)$ and

$$\mathcal{R}_G \stackrel{\text{def}}{=} \|Wb_G\|^2 + \mathcal{X}_G = \|Wb_G\|^2 + \operatorname{tr}(WD_G^{-2}V^2D_G^{-2}W). \tag{2.13}$$

Theorem 2.4. Let (E_0G) , (E_2G) , (\mathcal{L}_0G) , $(\mathcal{I}G)$, and $(\mathcal{L}G)$ hold. If $W^2 \leq D_G^2$, then it holds with \mathcal{R}_G from (2.13)

$$\mathbb{E} \|W(\tilde{\theta}_G - \theta^*)\|^2 \leq \{\mathcal{R}_G^{1/2} + \diamond_G^*\}^2, \quad (2.14)$$

where

$$\diamond_G^* = 4\{\delta_G(r_G)r_G + 2\nu_0\mathbf{a}_G r_G(\mathbb{H}_1 + \mathbb{H}_2/g + 4)\omega\}.$$

Remark 2.3. If the error term \diamond_G^* in (2.14) is relatively small, this result implies $\mathbb{E} \|W(\tilde{\theta}_G - \theta^*)\|^2 \approx \mathcal{R}_G = \|D_G b_G\|^2 + \mathcal{X}_G$. This is the usual decomposition of the quadratic risk in term of the squared bias $\|W(\theta_G^* - \theta^*)\|^2$ and the variance term \mathcal{X}_G . The condition “ $\|Wb_G\|^2/\mathcal{X}_G$ is small” yields $\mathcal{R}_G \approx \mathcal{X}_G$. This condition can be naturally called the small modeling bias (SMB) condition, often it is referred to as undersmoothing. The bias-variance trade-off corresponds to the situation with $\|Wb_G\|^2 \asymp \mathcal{X}_G$. Oversmoothing means that the bias terms $\|Wb_G\|^2$ dominates.

Remark 2.4. As already mentioned, the result (2.14) is informative if the remainder \diamond_G^* is relatively small and can be ignored. For the special case $W^2 = D_G^2$, it holds $\mathcal{X}_G = \mathbb{P}_G \asymp r_G^2$. In the i.i.d. situation (see Section 3.5 below)

$$r_G^{-1} \diamond_G^* \leq C\sqrt{\mathbb{P}_G/n}$$

which yields a sharp risk bound $\mathbb{E} \|W(\tilde{\theta}_G - \theta^*)\|^2 = \mathcal{R}_G(1 + o(1))$ under “ \mathbb{P}_G/n small.”

Remark 2.5. The bias induced by penalization can be measured in terms of the value $\|G\theta^*\|^2$. To be more precise, consider the case with $W^2 = D^2$, where $D^2 = -\nabla^2 \mathbb{E}L(\theta^*)$ is the non-penalized Fisher information matrix. The definition of θ^* and θ_G^* implies

$$\mathbb{E}L(\theta^*) - \|G\theta^*\|^2/2 \leq \mathbb{E}L(\theta_G^*) - \|G\theta_G^*\|^2/2 \leq \mathbb{E}L(\theta_G^*).$$

Condition (\mathcal{L}_0G) implies $\mathbb{E}L(\theta^*) - \mathbb{E}L(\theta_G^*) \approx \|D(\theta^* - \theta_G^*)\|^2/2$ and

$$\|D(\theta^* - \theta_G^*)\|^2 \leq \|G\theta^*\|^2 - \|G\theta_G^*\|^2 \leq \|G\theta^*\|^2.$$

So, if the true point is “smooth” in there sense that $\|G\theta^*\|^2$ is small, then the squared bias $\|D(\theta^* - \theta_G^*)\|^2$ caused by penalization is small as well.

Proof of Theorem 2.4. The Fisher expansion from Theorem 2.2 can be written as

$$\mathbb{P}(\|D_G(\tilde{\theta}_G - \theta^*) - D_G b_G - \xi_G\| \geq \diamond_G(x)) \leq 4e^{-x}.$$

The definition (2.10) of $\diamond_G(x)$ and (B.4) of Theorem B.1 imply

$$\mathbb{E}^{1/2} \|D_G(\tilde{\theta}_G - \theta^*) - D_G b_G - \xi_G\|^2 \leq 4\{\delta_G(r_G)r_G + 2\nu_0\mathbf{a}_G r_G(\mathbb{H}_1 + \mathbb{H}_2/g + 4)\omega\}.$$

By the result follows by the triangle inequality

$$\mathbb{E}^{1/2} \|D_G(\tilde{\theta}_G - \theta^*)\|^2 \leq \mathbb{E}^{1/2} \|D_G(\tilde{\theta}_G - \theta^*) - D_G b_G - \xi_G\|^2 + \mathbb{E}^{1/2} \|D_G b_G + \xi_G\|^2.$$

This yields the assertion of the theorem. □

2.6. Proofs of the Fisher and Wilks expansions

This section presents the proofs of the main results and some additional statements which can be of independent interest. The principle step of the proof is a bound on the local linear approximation of the gradient $\nabla L_G(\theta)$. Below we

study separately its stochastic and deterministic components coming from the decomposition $L(\boldsymbol{\theta}) = \mathbb{E}L(\boldsymbol{\theta}) + \zeta(\boldsymbol{\theta})$. With $D_G^2 = -\nabla^2 \mathbb{E}L_G(\boldsymbol{\theta}_G^*)$, this leads to the decomposition

$$\begin{aligned} \chi(\boldsymbol{\theta}, \boldsymbol{\theta}_G^*) &\stackrel{\text{def}}{=} D_G^{-1} \{ \nabla L_G(\boldsymbol{\theta}) - \nabla L_G(\boldsymbol{\theta}_G^*) \} + D_G(\boldsymbol{\theta} - \boldsymbol{\theta}_G^*) \\ &= D_G^{-1} \{ \nabla \zeta(\boldsymbol{\theta}) - \nabla \zeta(\boldsymbol{\theta}_G^*) \} \\ &\quad + D_G^{-1} \{ \nabla \mathbb{E}L_G(\boldsymbol{\theta}) - \nabla \mathbb{E}L_G(\boldsymbol{\theta}_G^*) \} + D_G(\boldsymbol{\theta} - \boldsymbol{\theta}_G^*). \end{aligned}$$

First we check the deterministic part. For any $\boldsymbol{\theta}$ with $\|D_G(\boldsymbol{\theta} - \boldsymbol{\theta}_G^*)\| \leq r$ and any unit vector $\mathbf{u} \in \mathbb{R}^p$, it holds

$$\begin{aligned} \mathbf{u}^\top \mathbb{E} \chi(\boldsymbol{\theta}, \boldsymbol{\theta}_G^*) &= \mathbf{u}^\top D_G^{-1} \{ \nabla \mathbb{E}L_G(\boldsymbol{\theta}) - \nabla \mathbb{E}L_G(\boldsymbol{\theta}_G^*) + D_G^2(\boldsymbol{\theta} - \boldsymbol{\theta}_G^*) \} \\ &= \mathbf{u}^\top \{ \mathbb{I}_p - D_G^{-1} \mathbb{F}_G(\boldsymbol{\theta}^\circ) D_G^{-1} \} D_G(\boldsymbol{\theta} - \boldsymbol{\theta}_G^*), \end{aligned}$$

where $\boldsymbol{\theta}^\circ = \boldsymbol{\theta}^\circ(\mathbf{u})$ is a point on the line connecting $\boldsymbol{\theta}_G^*$ and $\boldsymbol{\theta}$. This implies by (\mathcal{L}_0G)

$$\|\mathbb{E} \chi(\boldsymbol{\theta}, \boldsymbol{\theta}_G^*)\| \leq \|\mathbb{I}_p - D_G^{-1} \mathbb{F}_G(\boldsymbol{\theta}^\circ) D_G^{-1}\|_{\text{op}} r \leq \delta_G(r) r. \quad (2.15)$$

Now we study the stochastic part. Consider the vector process

$$\boldsymbol{\mathcal{U}}(\boldsymbol{\theta}, \boldsymbol{\theta}_G^*) \stackrel{\text{def}}{=} D_G^{-1} \{ \nabla \zeta(\boldsymbol{\theta}) - \nabla \zeta(\boldsymbol{\theta}_G^*) \}. \quad (2.16)$$

Further, define $\mathbf{v} = V_2(\boldsymbol{\theta} - \boldsymbol{\theta}_G^*)$ and introduce a vector process $\boldsymbol{\mathcal{Y}}(\mathbf{v})$ with

$$\boldsymbol{\mathcal{Y}}(\mathbf{v}) \stackrel{\text{def}}{=} V_2^{-1} [\nabla \zeta(\boldsymbol{\theta}) - \nabla \zeta(\boldsymbol{\theta}_G^*)].$$

It obviously holds $\nabla \boldsymbol{\mathcal{Y}}(\mathbf{v}) = V_2^{-1} \nabla^2 \zeta(\boldsymbol{\theta}) V_2^{-1}$. Moreover, for any $\boldsymbol{\gamma}_1, \boldsymbol{\gamma}_2 \in \mathbb{R}^p$ with $\|\boldsymbol{\gamma}_1\| = \|\boldsymbol{\gamma}_2\| = 1$, condition (E_2G) implies for $|\lambda| \leq g(r)$

$$\log \mathbb{E} \exp \left\{ \frac{\lambda}{\omega} \boldsymbol{\gamma}_1^\top \nabla \boldsymbol{\mathcal{Y}}(\mathbf{v}) \boldsymbol{\gamma}_2 \right\} = \log \mathbb{E} \exp \left\{ \frac{\lambda}{\omega} \boldsymbol{\gamma}_1^\top V_2^{-1} \nabla^2 \zeta(\boldsymbol{\theta}) V_2^{-1} \boldsymbol{\gamma}_2 \right\} \leq \frac{v_0^2 \lambda^2}{2}.$$

Define $\mathcal{Y}_o(r) \stackrel{\text{def}}{=} \{ \mathbf{v} : \|\mathbf{v}\| \leq r, \|\mathcal{S}\mathbf{v}\| \leq r \}$ for $S^{-2} = \mathbf{a}_G^{-2} D_G^{-1} V_2^2 D_G^{-1}$. Then

$$\sup_{\boldsymbol{\theta} \in \Theta_{0,G}(r)} \|\boldsymbol{\mathcal{U}}(\boldsymbol{\theta}, \boldsymbol{\theta}_G^*)\| \leq \sup_{\mathbf{v} \in \mathcal{Y}_o(r)} \|A \boldsymbol{\mathcal{Y}}(\mathbf{v})\| \quad (2.17)$$

for $A = \mathbf{a}_G^{-1} D_G^{-1} V_2$. Theorem B.15 yields

$$\sup_{\mathbf{v} \in \mathcal{Y}_o(r)} \|A \boldsymbol{\mathcal{Y}}(\mathbf{v})\| \leq \sqrt{8} v_0 \mathfrak{J}_{\mathbb{H}}(\mathbf{x}) \mathbf{a}_G \omega r \quad (2.18)$$

on a set of a dominating probability at least $1 - e^{-x}$, where the function $\mathfrak{J}_{\mathbb{H}}(\mathbf{x})$ is given by (B.3). Putting together the bounds (2.15) and (2.17) imply the following result.

Theorem 2.5. *Suppose that the matrix $\mathbb{F}_G(\boldsymbol{\theta}) \stackrel{\text{def}}{=} -\nabla^2 \mathbb{E}L_G(\boldsymbol{\theta})$ fulfills the condition (\mathcal{L}_0G) and let (E_0G) and (E_2G) be fulfilled on $\Theta_{0,G}(r)$ for any fixed $r \leq r^*$. Then*

$$\mathbb{P} \left\{ \sup_{\boldsymbol{\theta} \in \Theta_{0,G}(r)} \|D_G^{-1} \{ \nabla L_G(\boldsymbol{\theta}) - \nabla L_G(\boldsymbol{\theta}_G^*) \} + D_G(\boldsymbol{\theta} - \boldsymbol{\theta}_G^*)\| \geq \diamond_G(r, \mathbf{x}) \right\} \leq e^{-x},$$

where

$$\diamond_G(r, \mathbf{x}) \stackrel{\text{def}}{=} \{ \delta_G(r) + \sqrt{8} v_0 \mathfrak{J}_{\mathbb{H}}(\mathbf{x}) \mathbf{a}_G \omega \} r. \quad (2.19)$$

The result of Theorem 2.5 can be extended to the increments of the process $\mathcal{U}(\boldsymbol{\theta})$: on a random set of probability at least $1 - e^{-x}$, it holds for any $\boldsymbol{\theta}, \boldsymbol{\theta}^\circ \in \Theta_{0,G}(x)$ and $\chi(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ) = D_G^{-1}\{\nabla L_G(\boldsymbol{\theta}) - \nabla L_G(\boldsymbol{\theta}^\circ)\} + D_G(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)$

$$\begin{aligned} \mathbb{E}[\chi(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ)] &\leq \delta_G(x) \|D_G(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)\| \leq 2x\delta_G(x), \\ \|\chi(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ)\| &\leq 2\diamond_G(x, \mathbf{x}). \end{aligned} \quad (2.20)$$

Now we present the proof of Theorem 2.2 about the Fisher expansion for the qMLE $\tilde{\boldsymbol{\theta}}_G$ defined by maximization of $L_G(\boldsymbol{\theta})$. Let r_G be selected to ensure that $\mathbb{P}\{\tilde{\boldsymbol{\theta}}_G \notin \Theta_{0,G}(r_G)\} \leq e^{-x}$. Furthermore, the definition of $\tilde{\boldsymbol{\theta}}_G$ yields $\nabla L_G(\tilde{\boldsymbol{\theta}}_G) = 0$ and

$$\chi(\tilde{\boldsymbol{\theta}}_G, \boldsymbol{\theta}_G^*) = -D_G^{-1}\nabla L_G(\boldsymbol{\theta}_G^*) + D_G(\tilde{\boldsymbol{\theta}}_G - \boldsymbol{\theta}_G^*).$$

By Theorem 2.5, it holds on a set of a dominating probability

$$\|D_G(\tilde{\boldsymbol{\theta}}_G - \boldsymbol{\theta}_G^*) - \boldsymbol{\xi}_G\| \leq \diamond_G(x) \quad (2.21)$$

as required.

As the next step, we apply the obtained results to evaluate the quality of the Wilks expansion $2L_G(\tilde{\boldsymbol{\theta}}_G, \boldsymbol{\theta}_G^*) \approx \|\boldsymbol{\xi}_G\|^2$. For this we derive a uniform deviation bound on the error of a quadratic approximation

$$\alpha(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ) \stackrel{\text{def}}{=} L_G(\boldsymbol{\theta}) - L_G(\boldsymbol{\theta}^\circ) - (\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)^\top \nabla L_G(\boldsymbol{\theta}^\circ) + \frac{1}{2} \|D_G(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)\|^2$$

in all $\boldsymbol{\theta}, \boldsymbol{\theta}^\circ \in \Theta_0$, where Θ_0 is some vicinity of a fixed point $\boldsymbol{\theta}_G^*$. With $\boldsymbol{\theta}^\circ$ fixed, the gradient $\nabla \alpha(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ) \stackrel{\text{def}}{=} \frac{d}{d\boldsymbol{\theta}} \alpha(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ)$ fulfills

$$\nabla \alpha(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ) = \nabla L_G(\boldsymbol{\theta}) - \nabla L_G(\boldsymbol{\theta}^\circ) + D_G^2(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ) = D_G \chi(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ);$$

cf. (2.16). This implies

$$\alpha(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ) = (\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)^\top \nabla \alpha(\boldsymbol{\theta}', \boldsymbol{\theta}^\circ),$$

where $\boldsymbol{\theta}'$ is a point on the line connecting $\boldsymbol{\theta}$ and $\boldsymbol{\theta}^\circ$. Further,

$$|\alpha(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ)| = |(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)^\top D_G D_G^{-1} \nabla \alpha(\boldsymbol{\theta}', \boldsymbol{\theta}^\circ)| \leq \|D_G(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)\| \sup_{\boldsymbol{\theta}' \in \Theta_{0,G}(x)} |\chi(\boldsymbol{\theta}', \boldsymbol{\theta}^\circ)|,$$

and one can apply (2.20). This yields the following result.

Theorem 2.6. *Suppose (\mathcal{L}_0G) , (E_0G) , and (E_2G) . For each x , it holds on a random set $\Omega(x)$ of a dominating probability at least $1 - e^{-x}$, it holds with any $\boldsymbol{\theta}, \boldsymbol{\theta}^\circ \in \Theta_{0,G}(x)$ and $\diamond_G(x, \mathbf{x})$ is from (2.19)*

$$\begin{aligned} \frac{|\alpha(\boldsymbol{\theta}, \boldsymbol{\theta}_G^*)|}{\|D_G(\boldsymbol{\theta} - \boldsymbol{\theta}_G^*)\|} &\leq \diamond_G(x, \mathbf{x}), \quad |\alpha(\boldsymbol{\theta}, \boldsymbol{\theta}_G^*)| \leq x \diamond_G(x, \mathbf{x}), \\ \frac{|\alpha(\boldsymbol{\theta}_G^*, \boldsymbol{\theta})|}{\|D_G(\boldsymbol{\theta} - \boldsymbol{\theta}_G^*)\|} &\leq 2\diamond_G(x, \mathbf{x}), \quad |\alpha(\boldsymbol{\theta}_G^*, \boldsymbol{\theta})| \leq 2x \diamond_G(x, \mathbf{x}), \\ \frac{|\alpha(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ)|}{\|D_G(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)\|} &\leq 2\diamond_G(x, \mathbf{x}), \quad |\alpha(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ)| \leq 4x \diamond_G(x, \mathbf{x}). \end{aligned}$$

The result of Theorem 2.6 for the special case with $\boldsymbol{\theta} = \boldsymbol{\theta}_G^*$ and $\boldsymbol{\theta}^\circ = \tilde{\boldsymbol{\theta}}_G$ yields in view of $\nabla L_G(\tilde{\boldsymbol{\theta}}_G) = 0$ for $r = r_G$ and $\diamond_G(x) = \diamond_G(r_G, \mathbf{x})$ under the condition $\tilde{\boldsymbol{\theta}}_G \in \Theta_{0,G}(r_G)$

$$|L_G(\tilde{\boldsymbol{\theta}}_G, \boldsymbol{\theta}_G^*) - \|D_G(\tilde{\boldsymbol{\theta}}_G - \boldsymbol{\theta}_G^*)\|^2/2| = |\alpha(\boldsymbol{\theta}_G^*, \tilde{\boldsymbol{\theta}}_G)| \leq 2r_G \diamond_G(x).$$

Furthermore, with $\boldsymbol{\theta} = \tilde{\boldsymbol{\theta}}_G$ and $\boldsymbol{\theta}^\circ = \boldsymbol{\theta}_G^*$

$$|L_G(\tilde{\boldsymbol{\theta}}_G, \boldsymbol{\theta}_G^*) - \boldsymbol{\xi}_G^\top D_G(\tilde{\boldsymbol{\theta}}_G - \boldsymbol{\theta}_G^*) + \|D_G(\tilde{\boldsymbol{\theta}}_G - \boldsymbol{\theta}_G^*)\|^2/2| = |\alpha(\tilde{\boldsymbol{\theta}}_G, \boldsymbol{\theta}_G^*)| \leq r_G \diamond_G(\mathbf{x})$$

which implies

$$|L(\tilde{\boldsymbol{\theta}}_G, \boldsymbol{\theta}_G^*) - \|\boldsymbol{\xi}_G\|^2 + \|D_G(\tilde{\boldsymbol{\theta}}_G - \boldsymbol{\theta}_G^*) - \boldsymbol{\xi}_G\|^2| \leq 2r_G \diamond_G(\mathbf{x}).$$

Now it follows by (2.21) that

$$|L(\tilde{\boldsymbol{\theta}}_G, \boldsymbol{\theta}_G^*) - \|\boldsymbol{\xi}_G\|^2/2| \leq r_G \diamond_G(\mathbf{x}) + \diamond_G^2(\mathbf{x})/2.$$

The error term can be improved if the squared root of the excess is considered. Indeed, if $\tilde{\boldsymbol{\theta}}_G \in \Theta_{0,G}(r_G)$

$$\begin{aligned} |\{2L_G(\tilde{\boldsymbol{\theta}}_G, \boldsymbol{\theta}_G^*)\}^{1/2} - \|D_G(\tilde{\boldsymbol{\theta}}_G - \boldsymbol{\theta}_G^*)\|| &\leq \frac{|2L_G(\tilde{\boldsymbol{\theta}}_G, \boldsymbol{\theta}_G^*) - \|D_G(\tilde{\boldsymbol{\theta}}_G - \boldsymbol{\theta}_G^*)\|^2|}{\|D_G(\tilde{\boldsymbol{\theta}}_G - \boldsymbol{\theta}_G^*)\|} \\ &\leq \frac{2|\alpha(\tilde{\boldsymbol{\theta}}_G, \boldsymbol{\theta}_G^*)|}{\|D_G(\tilde{\boldsymbol{\theta}}_G - \boldsymbol{\theta}_G^*)\|} \leq \sup_{\boldsymbol{\theta} \in \Theta_{0,G}(r_G)} \frac{2|\alpha(\boldsymbol{\theta}, \boldsymbol{\theta}_G^*)|}{\|D_G(\boldsymbol{\theta} - \boldsymbol{\theta}_G^*)\|} \leq 2\diamond_G(\mathbf{x}). \end{aligned}$$

The Fisher expansion (2.21) allows to replace here the norm of the standardized error $D_G(\tilde{\boldsymbol{\theta}}_G - \boldsymbol{\theta}_G^*)$ with the norm of the normalized score $\boldsymbol{\xi}_G$. This completes the proof of Theorem 2.3.

3. Examples

This section illustrates the general results for two particularly important cases of i.i.d. and generalized linear models. The primary focus of the study is to compare the penalized and non-penalized cases and to quantify the impact of penalization.

3.1. *I.i.d. case*

The model with independent identically distributed (i.i.d.) observations is one of the most popular setups in statistical literature and in statistical applications. The essential and the most developed part of the statistical theory is designed for the i.i.d. modeling. Especially, the classical asymptotic parametric theory is almost complete including asymptotic root- n normality and efficiency of the MLE and Bayes estimators under rather mild assumptions; see e.g. Chapter 2 and 3 in Ibragimov and Khas'minskij [14]. So, the i.i.d. model can naturally serve as a benchmark for any extension of the statistical theory: being applied to the i.i.d. setup, the new approach should lead to essentially the same conclusions as in the classical theory. Similar reasons apply to the regression model and its extensions. Below we try demonstrate that the proposed non-asymptotic viewpoint is able to reproduce the existing brilliant and well established results of the classical parametric theory. Surprisingly, the majority of classical efficiency results can be easily derived from the obtained general non-asymptotic bounds.

3.2. *Quasi MLE in an i.i.d. model*

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

parametric family (P_θ) . The case of non-identically distributed observations can be done similarly at cost of more complicated notation.

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

$$L(\theta) = \sum \ell(Y_i, \theta), \quad (3.1)$$

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

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

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

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

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

$$\nabla \zeta(\theta^*) = \sum \nabla \zeta_i(\theta^*) = \sum \nabla \ell_i(\theta^*). \quad (3.2)$$

3.3. Conditions in the i.i.d. case

I.i.d. structure of the Y_i 's allows for rewriting the conditions (E_0) , (E_2) , (\mathcal{I}) , (\mathcal{L}_0) , and (\mathcal{L}) in terms of the marginal distribution. In the following conditions the index i runs from 1 to n .

(ed_0) There exists a positive symmetric matrix v_0 , such that for all $|\lambda| \leq g_1$

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

A natural candidate on v_0^2 is given by the variance of the gradient $\nabla \ell(Y_1, \theta^*)$, that is, $v_0^2 = \operatorname{Var} \nabla \ell(Y_1, \theta) = \operatorname{Var} \nabla \zeta_1(\theta)$. Note that (ed_0) is automatically fulfilled if the model is correctly specified and $P = P_{\theta^*}$ because $E_{\theta^*} \exp\{\ell(Y_1, \theta) - \ell(Y_1, \theta^*)\} \equiv 1$.

Next consider the local sets

$$\Theta_0(r) = \{\theta : \|v_0(\theta - \theta^*)\| \leq r/n^{1/2}\}.$$

The local smoothness conditions (E_2) and (\mathcal{L}_0) require to specify the functions $\delta(r)$ and $\varrho(r)$. If the log-likelihood function $\ell(y, \theta)$ is sufficiently smooth in θ , these functions can be selected proportional to r .

(ed_2) There exist a value $\omega^* > 0$ and for each $r > 0$, a constant $g(r) > 0$ such that

$$\sup_{\boldsymbol{y}_1, \boldsymbol{y}_2 \in \mathbb{R}^p} \log \mathbb{E} \exp \left\{ \frac{\lambda}{\omega^*} \frac{\boldsymbol{y}_1^\top \nabla^2 \zeta_i(\theta) \boldsymbol{y}_2}{\|v_0 \boldsymbol{y}_1\| \cdot \|v_0 \boldsymbol{y}_2\|} \right\} \leq \frac{v_0^2 \lambda^2}{2}, \quad |\lambda| \leq g(r).$$

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

(ℓ_0) The function $\bar{\ell}(\theta)$ is two times differentiable and the matrix function $H(\theta) = -\nabla^2 \mathbb{E}\ell(Y_1, \theta)$ fulfills with $H_0 \stackrel{\text{def}}{=} H(\theta^*)$ for some constant δ^* :

$$\sup_{\theta \in \Theta_0(r)} \|H_0^{-1/2} H(\theta) H_0^{-1/2} - \mathbb{I}_p\|_{\text{op}} \leq \frac{\delta^* r}{\sqrt{n}}.$$

In the regular parametric case with $P \in (P_\theta)$, the matrices v_0^2 and H_0 coincide with the Fisher information matrix $H_0 = H(\theta^*)$ of the family (P_θ) at the point θ^* .

The consistency result for $\tilde{\theta}$ requires certain growth of the value $\bar{\ell}(\theta^*) - \bar{\ell}(\theta)$ as $\|\theta - \theta^*\|$ grows. The marginal version of the global condition (\mathcal{L}) reads as follows:

$(\bar{\ell})$ There exists $C_1 \geq 0$ such that with $\delta = \delta^*/\sqrt{n}$ and $r = \sqrt{n}\|H_0^{1/2}(\theta - \theta^*)\|$

$$n\{\bar{\ell}(\theta^*) - \bar{\ell}(\theta)\} \geq (1 - \delta)(r r_0 - r_0^2/2) - C_1 r^2.$$

Remark 3.1. If the parametric i.i.d. model is correct, then

$$\bar{\ell}(\theta^*) - \bar{\ell}(\theta) = \mathcal{K}(\theta^*, \theta) = E_{\theta^*} \log \frac{dP_{\theta^*}}{dP_\theta}(Y_1)$$

is the Kullback–Leibler divergence for the family (P_θ) . Condition $(\bar{\ell})$ is fulfilled automatically if $\bar{\ell}(\theta^*, \theta) > 0$ for $\theta \neq \theta^*$ and Θ is a compact set. Then

$$\inf_{\theta \in \Theta} \frac{\bar{\ell}(\theta^*) - \bar{\ell}(\theta)}{\|H_0^{1/2}(\theta - \theta^*)\|^2} \geq b > 0.$$

This and (ℓ_0) imply $(\bar{\ell})$.

The *identifiability condition* relates the matrices v_0^2 and H_0 .

(i) There is a constant $\alpha > 0$ such that $\alpha^2 H_0 \geq v_0^2$.

Lemma 3.1. Let Y_1, \dots, Y_n be i.i.d. Then (ed_0) , (ed_2) , (ℓ_0) , $(\bar{\ell})$, and (i) imply (E_0) , (E_2) , (\mathcal{L}_0) , (\mathcal{L}) , (\mathcal{I}) , with $V^2 = n v_0^2$, $D^2 = n H_0$, $\omega = \omega^*/\sqrt{n}$, $\delta(r) = \delta^* r/\sqrt{n}$, and the same constants v_0 , α , $g \stackrel{\text{def}}{=} g_1 \sqrt{n}$.

Proof. The identities $V^2 = n v_0^2$, $D^2 = n H_0$ follow from the i.i.d. structure of the observations Y_i . We briefly comment on condition (E_0) . The use once again of the i.i.d. structure yields by (3.2) in view of $V^2 = n v_0^2$

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

as long as $\lambda \leq n^{1/2} g_1 \leq g$. Similarly one can check (E_2) . The conditions (\mathcal{L}_0) , (\mathcal{L}) , and (\mathcal{I}) follow from (ℓ_0) and $(\bar{\ell})$, and (i) due to $D^2 = n H_0$ and $\mathbb{E} L(\theta) = n \bar{\ell}(\theta)$. \square

Below we specify the obtained general results to the i.i.d. setup.

3.4. Results in the non-penalized i.i.d. case

Here we specify the general results of previous chapters to the i.i.d. case. In particular, we explicitly state the large deviation bound and show that it yields a root- n consistency of the qMLE $\tilde{\theta}$. Then we comment on the Fisher and Wilks theorems.

First we describe the large deviation probability for the event $\{\tilde{\theta} \notin \Theta_0(r_0)\}$ for a fixed r_0 . The next result specifies the general large deviation statement of Theorem 2.1 to the finite dimensional non-penalized i.i.d. case and states the inference results.

Theorem 3.2. Let (ed_0) , (ed_2) , (ℓ_0) , (i), and $(\bar{\ell})$ hold with

$$\{1 - \delta(r_0)\} r_0 \geq 2z(B, \mathbf{x}), \quad C_1 \geq \varrho(r, \mathbf{x}), \quad r > r_0, \quad (3.3)$$

where $B = H_0^{-1/2} v_0^2 H_0^{-1/2} = D^{-1} V^2 D^{-1}$, $z(B, \mathbf{x})$ is given by (A.4), and

$$\varrho(r, \mathbf{x}) \stackrel{\text{def}}{=} v_0 \mathfrak{z}_{\mathbb{H}}(\mathbf{x} + \log(2r/r_0)) \omega^* / \sqrt{n}$$

with $\mathfrak{z}_{\mathbb{H}}(\mathbf{x}) \leq C\sqrt{p + \mathbf{x}}$. Then it holds on a set $\Omega(\mathbf{x})$ with $\mathbb{P}(\Omega(\mathbf{x})) \geq 1 - 5e^{-\mathbf{x}}$

$$\sqrt{n} \|H_0^{1/2}(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\| \leq r_0. \tag{3.4}$$

Furthermore, on this set $\Omega(\mathbf{x})$, it holds

$$\begin{aligned} \|\sqrt{nH_0}(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) - \boldsymbol{\xi}\| &\leq C\sqrt{(p + \mathbf{x})^2/n}, \\ |\sqrt{2L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*)} - \|\boldsymbol{\xi}\|| &\leq C\sqrt{(p + \mathbf{x})^2/n}, \\ |2L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) - \|\boldsymbol{\xi}\|^2| &\leq C\sqrt{(p + \mathbf{x})^3/n}. \end{aligned}$$

The constant C here depends in an explicit way on the constants α_G, \mathfrak{g}_1 , and v_0 from our conditions, and

$$\boldsymbol{\xi} \stackrel{\text{def}}{=} (nH_0)^{-1/2} \sum_{i=1}^n \nabla \ell(Y_i, \boldsymbol{\theta}^*). \tag{3.5}$$

Proof. Condition (t) implies $B = H_0^{-1/2} v_0^2 H_0^{-1/2} \leq \alpha^2 \mathbb{I}_p$ and thus, $\text{tr}(B) \leq \alpha^2 p$. Therefore, the value $z(B, \mathbf{x})$ fulfills $z^2(B, \mathbf{x}) \leq C(p + \mathbf{x})$. The same bound holds for $\mathfrak{z}_{\mathbb{H}}^2(\mathbf{x})$. Condition (3.3) with $b(r_0) \approx 1$ yields $r_0^2 \approx 4z^2(B, \mathbf{x}) \approx C(p + \mathbf{x})$. This yields in view of $\delta(r_0) \leq \delta^* r_0 / \sqrt{n}$ and $\omega = \omega^* n^{-1/2}$

$$\diamond(r_0, \mathbf{x}) \leq \{\delta(r_0) + v_0 \mathfrak{z}_{\mathbb{H}}(\mathbf{x}) \omega\} r_0 \leq C(p + \mathbf{x}) / \sqrt{n}.$$

Similarly

$$\Delta(r_0, \mathbf{x}) \leq \{\delta(r_0) + v_0 \mathfrak{z}_{\mathbb{H}}(\mathbf{x}) \omega\} r_0^2 \leq C\sqrt{(p + \mathbf{x})^3/n}.$$

The results follow now from general theorems of Section 2. □

For the classical asymptotic setup when n tends to infinity, the random vector $\boldsymbol{\xi}$ from (3.5) fulfills $\text{Var}(\boldsymbol{\xi}) \leq H_0^{-1/2} v_0^2 H_0^{-1/2} = B$ and by the central limit theorem $\boldsymbol{\xi}$ is asymptotically normal $\mathcal{N}(0, B)$. This yields by Theorem 3.2 that $\sqrt{nH_0}(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)$ is asymptotically normal $\mathcal{N}(0, B)$ as well. The correct model specification implies $B \equiv \mathbb{I}_p$ and hence $\tilde{\boldsymbol{\theta}}$ is asymptotically efficient; see Ibragimov and Khas'minskij [14]. Also $2L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) \approx \|\boldsymbol{\xi}\|^2$ which is nearly χ^2 r.v. with p degrees of freedom. This result is known as asymptotic Wilks theorem.

In the non-asymptotic framework of this paper, the error terms still depend on n and they can only be small if n is large. However, we show in explicit way how these error terms depend on the parameter dimension. It appears that the root- n consistency result (3.4) requires “ p/n small.” The Fisher and square root Wilks results apply if “ p^2/n is small.” Finally, the Wilks expansion is valid under “ p^3/n small.” Existing statistical literature addresses the issue of a growing parameter dimension in different set-ups. The classical results by Portnoy [19–21] provide some constraints on parameter dimension for consistency and asymptotic normality of the M-estimator for regression models. Our results are consistent with the conclusion of that papers. We refer to Andresen and Spokoiny [1] for a version of such result in context of semiparametric profile estimation. That paper also provides an example of an i.i.d. model in which the Fisher expansion of Theorem 3.2 fails for $p^2 \geq n$. The next section demonstrates how these constraints on the parameter dimension can be relaxed by using a penalization.

3.5. Roughness penalization for an i.i.d. sample

This section discusses the impact of penalization in the case of an i.i.d. model with n observations. For penalty term $\text{pen}(\boldsymbol{\theta}) = \|G\boldsymbol{\theta}\|^2/2$, the penalized log-likelihood is given by $L_G(\boldsymbol{\theta}) = L(\boldsymbol{\theta}) + \|G\boldsymbol{\theta}\|^2/2$, where $L(\boldsymbol{\theta})$ is from (3.1). With $\boldsymbol{\theta}_G^* = \arg\max_{\boldsymbol{\theta} \in \Theta} \mathbb{E}L_G(\boldsymbol{\theta})$, define

$$D_G^2 = n\mathbb{H}(\boldsymbol{\theta}_G^*) + G^2, \quad V^2 = nv_0^2, \quad \boldsymbol{\xi}_G = D_G^{-1} \sum_{i=1}^n \nabla \ell(Y_i, \boldsymbol{\theta}_G^*),$$

where $\mathbb{H}(\boldsymbol{\theta}) = -\nabla^2 \mathbb{E}\ell(Y_1, \boldsymbol{\theta})$, $v_0^2 = \text{Var}\{\ell(Y_1, \boldsymbol{\theta}_G^*)\}$. The value \mathfrak{p}_G is defined as previously by (2.2).

Note that all the introduced quantities including the parameter set Θ , the parameter dimension p , and the effective dimension \mathfrak{p}_G , may depend on n . Here we also allow a functional parameter $\boldsymbol{\theta}$ with $p = \infty$. The main goal is to show that the presented general approach yields sharp results in this special case.

Suppose that the conditions of Section 3.3 are fulfilled. One can easily check the conditions from Section 2.2 with $\delta_G(\mathbf{x}) = C\mathbf{x}/\sqrt{n}$ and $\omega = C/\sqrt{n}$; cf. Lemma 3.1. The large deviation bound of Theorem 2.1 applies for $\mathbf{r}_G \approx 2z(B_G, \mathbf{x}) \asymp \sqrt{\mathfrak{p}_G + \mathbf{x}}$. The general statements of Theorems 2.2 and 2.3 apply with $\diamond_G(\mathbf{x}) \leq C(\mathfrak{p}_G + \mathbf{x})/\sqrt{n}$ yielding the following expansions.

Theorem 3.3. *Suppose also that the conditions (ed_0) , (ed_2) , (ℓ_0) , $(\bar{\ell})$, and (ι) are fulfilled, and \mathbf{r}_G and C_1 fulfill*

$$\{1 - \delta_G(\mathbf{r}_G)\}\mathbf{r}_G \geq 2z(B_G, \mathbf{x}), \quad C_1 \geq \varrho_G(\mathbf{r}, \mathbf{x}), \quad \mathbf{r} > \mathbf{r}_G,$$

with $B_G = D_G^{-1}V^2D_G^{-1/2}$, then on a set of dominating probability $1 - 5e^{-\mathbf{x}}$, it holds

$$\begin{aligned} \|D_G(\tilde{\boldsymbol{\theta}}_G - \boldsymbol{\theta}_G^*) - \boldsymbol{\xi}_G\| &\leq C\sqrt{(\mathfrak{p}_G + \mathbf{x})^2/n}, \\ \left| \sqrt{2L_G(\tilde{\boldsymbol{\theta}}_G, \boldsymbol{\theta}_G^*)} - \|\boldsymbol{\xi}_G\| \right| &\leq C\sqrt{(\mathfrak{p}_G + \mathbf{x})^2/n}, \\ |2L_G(\tilde{\boldsymbol{\theta}}_G, \boldsymbol{\theta}_G^*) - \|\boldsymbol{\xi}_G\|^2| &\leq C\sqrt{(\mathfrak{p}_G + \mathbf{x})^3/n}. \end{aligned}$$

The constant C here depends in an explicit way on the constants α_G , \mathfrak{g}_1 , and v_0 from our conditions.

A short look at the results for non-penalized and penalized estimates indicates that the quality of the penalized MLE $\tilde{\boldsymbol{\theta}}_G$ improves relative to the non-penalized case because the matrix D_G^2 can be much larger than D^2 , the variance of the stochastic term $\boldsymbol{\xi}_G$ is of order \mathfrak{p}_G instead of p for the variance of $\boldsymbol{\xi}$, and, simultaneously, the error terms in the Fisher and Wilks expansions become smaller due to reduction of the effective dimension \mathfrak{p}_G in place of the full dimension p .

Andresen and Spokoiny [1] provides a simple example of estimating the squared norm $\|\boldsymbol{\theta}\|^2$ which shows that the Fisher expansion fails if p/\sqrt{n} is not small. The result can be easily adjusted to the penalized case.

3.6. Generalized linear models (GLM)

Generalized linear models (GLM) are frequently used for modeling the data with special structure: categorical data, binary data, Poissonian and exponential data, volatility models, etc. All these examples can be treated in a unified way by a GLM approach. This section specifies the results and conditions to this case.

Let $\mathbf{Y} = (Y_1, \dots, Y_n)^\top \sim \mathbb{P}$ be a sample of independent r.v.'s. The parametric GLM is given by $Y_i \sim P_{\Psi_i^\top \boldsymbol{\theta}} \in (P_\nu)$, where Ψ_i are given factors in \mathbb{R}^p , $\boldsymbol{\theta} \in \mathbb{R}^p$ is the unknown parameter in \mathbb{R}^p , and (P_ν) is an exponential family with canonical parametrization yielding the log-density $\ell(\mathbf{y}, \boldsymbol{\nu}) = \mathbf{y}\boldsymbol{\nu} - g(\boldsymbol{\nu})$ for a convex function $g(\boldsymbol{\nu})$. Below we suppose that the function $g(\boldsymbol{\nu})$ is sufficiently smooth, in particular, three times differentiable.

The (quasi) log-likelihood $L(\boldsymbol{\theta})$ can be represented in the form

$$L(\boldsymbol{\theta}) = \sum_{i=1}^n \{Y_i \boldsymbol{\psi}_i^\top \boldsymbol{\theta} - g(\boldsymbol{\psi}_i^\top \boldsymbol{\theta})\} = S^\top \boldsymbol{\theta} - A(\boldsymbol{\theta}) \quad (3.6)$$

with a random p -vector S and a function $A(\boldsymbol{\theta})$ given by

$$S \stackrel{\text{def}}{=} \sum_{i=1}^n Y_i \boldsymbol{\psi}_i, \quad A(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \sum_i g(\boldsymbol{\psi}_i^\top \boldsymbol{\theta}).$$

The MLE $\tilde{\boldsymbol{\theta}}$ and the target $\boldsymbol{\theta}^*$ for this GLM read as

$$\begin{aligned} \tilde{\boldsymbol{\theta}} &= \underset{\boldsymbol{\theta}}{\operatorname{argmax}} L(\boldsymbol{\theta}) = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \{S^\top \boldsymbol{\theta} - A(\boldsymbol{\theta})\}, \\ \boldsymbol{\theta}^* &= \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \mathbb{E}L(\boldsymbol{\theta}) = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \{\mathbb{E}S^\top \boldsymbol{\theta} - A(\boldsymbol{\theta})\}, \end{aligned} \quad (3.7)$$

where

$$\mathbb{E}S = \sum_{i=1}^n \mathbb{E}Y_i \boldsymbol{\psi}_i.$$

The definition of $\boldsymbol{\theta}^*$ implies the identity $\nabla \mathbb{E}L(\boldsymbol{\theta}^*) = 0$ which yields

$$\mathbb{E}S = \nabla A(\boldsymbol{\theta}^*).$$

An important feature of a GLM is that the stochastic component $\zeta(\boldsymbol{\theta})$ of $L(\boldsymbol{\theta})$ is *linear in $\boldsymbol{\theta}$* : with $\varepsilon_i = Y_i - \mathbb{E}Y_i$

$$\begin{aligned} \zeta(\boldsymbol{\theta}) &= L(\boldsymbol{\theta}) - \mathbb{E}L(\boldsymbol{\theta}) = \sum_{i=1}^n \varepsilon_i \boldsymbol{\psi}_i^\top \boldsymbol{\theta}, \\ \nabla \zeta(\boldsymbol{\theta}) &= S - \mathbb{E}S = \sum_{i=1}^n \varepsilon_i \boldsymbol{\psi}_i. \end{aligned} \quad (3.8)$$

In the contrary to the linear case, the Fisher information matrix $D^2 = \mathbb{F}(\boldsymbol{\theta}^*)$ for

$$\mathbb{F}(\boldsymbol{\theta}) \stackrel{\text{def}}{=} -\nabla^2 \mathbb{E}L(\boldsymbol{\theta}) = \sum_{i=1}^n \boldsymbol{\psi}_i \boldsymbol{\psi}_i^\top g''(\boldsymbol{\psi}_i^\top \boldsymbol{\theta}) \quad (3.9)$$

depends on the true data distribution via the target $\boldsymbol{\theta}^*$. As $g(\cdot)$ is convex, it holds $g''(u) \geq 0$ for any u and thus $\mathbb{F}(\boldsymbol{\theta}) \geq 0$.

Linearity in $\boldsymbol{\theta}$ of the stochastic component $\zeta(\boldsymbol{\theta})$ and concavity of the deterministic part $\mathbb{E}L(\boldsymbol{\theta})$ allow for a simple and straightforward proof of the result about concentration of the MLE $\tilde{\boldsymbol{\theta}}$. Recall the definition of the local vicinity $\Theta_0(r)$ of $\boldsymbol{\theta}^*$:

$$\Theta_0(r) \stackrel{\text{def}}{=} \{\boldsymbol{\theta}: \|D(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\| \leq r\}.$$

Also define

$$B \stackrel{\text{def}}{=} D^{-1} \operatorname{Var}(S) D^{-1}.$$

Theorem 3.4. *If for some $r_0 > 0$, $\mathbb{F}(\boldsymbol{\theta})$ from (3.9) fulfill for $D^2 = \mathbb{F}(\boldsymbol{\theta}^*)$*

$$\sup_{\boldsymbol{\theta} \in \Theta_0(r_0)} \|D^{-1}\mathbb{F}(\boldsymbol{\theta})D^{-1} - \mathbb{I}_p\|_{\text{op}} \leq \delta(r_0) \quad (3.10)$$

with $\delta(r_0) < 1$, and if S from (3.8) follows for $\mathbf{x} > 0$ the probability bound

$$\mathbb{P}(\|D^{-1}(S - \mathbb{E}S)\| > z(B, \mathbf{x})) \leq 2e^{-\mathbf{x}}, \quad (3.11)$$

then the solution $\tilde{\boldsymbol{\theta}}$ of (3.7) satisfies

$$\mathbb{P}(\tilde{\boldsymbol{\theta}} \notin \Theta_0(r_0)) \leq 2e^{-\mathbf{x}}$$

provided that

$$r_0\{1 - \delta(r_0)\} \geq 2z(B, \mathbf{x}).$$

Proof. The function $L(\boldsymbol{\theta})$ is concave in $\boldsymbol{\theta}$ because

$$-\nabla^2 L(\boldsymbol{\theta}) = \mathbb{F}(\boldsymbol{\theta}) \geq 0. \quad (3.12)$$

If $\tilde{\boldsymbol{\theta}} \notin \Theta_0(r_0)$, denote by $\check{\boldsymbol{\theta}}$ the point at which the line connecting $\boldsymbol{\theta}^*$ and $\tilde{\boldsymbol{\theta}}$ crosses the boundary of $\Theta_0(r_0)$. It is easy to see that

$$\check{\boldsymbol{\theta}} - \boldsymbol{\theta}^* = \frac{\|D(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\|}{\|D(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\|} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) = \frac{r_0}{\|D(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\|} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*).$$

Concavity of $L(\boldsymbol{\theta})$ implies for the point of maximum $\tilde{\boldsymbol{\theta}}$ that

$$L(\tilde{\boldsymbol{\theta}}) - L(\boldsymbol{\theta}^*) \geq L(\check{\boldsymbol{\theta}}) - L(\boldsymbol{\theta}^*).$$

Therefore, it suffices to check that for each $\boldsymbol{\theta}$ with $\|D(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\| = r_0$ that

$$L(\boldsymbol{\theta}^*) - L(\boldsymbol{\theta}) > 0$$

on a set $\Omega(\mathbf{x})$ of probability $1 - 2e^{-\mathbf{x}}$. Then the event $\tilde{\boldsymbol{\theta}} \notin \Theta_0(r_0)$ is impossible on $\Omega(\mathbf{x})$. For any such $\boldsymbol{\theta}$, we apply the second order Taylor expansion of $L(\boldsymbol{\theta})$ at $\boldsymbol{\theta}^*$. By definition of $\boldsymbol{\theta}^*$, it holds $\nabla \mathbb{E}L(\boldsymbol{\theta}^*) = 0$ and thus $\nabla L(\boldsymbol{\theta}^*) = \nabla \zeta(\boldsymbol{\theta}^*) = (S - \mathbb{E}S)$. The use of (3.12), (3.10) yields now for $\boldsymbol{\xi} = D^{-1}(S - \mathbb{E}S)$ and for $\boldsymbol{\theta}$ with $\|D(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\| = r_0$

$$\begin{aligned} L(\boldsymbol{\theta}^*) - L(\boldsymbol{\theta}) &= (\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top \nabla L(\boldsymbol{\theta}^*) + \frac{1}{2} \|\sqrt{\mathbb{F}(\boldsymbol{\theta}^\circ)}(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2 \\ &\geq (S - \mathbb{E}S)^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^*) + \frac{1 - \delta(r_0)}{2} \|D(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2 \\ &= \boldsymbol{\xi}^\top D(\boldsymbol{\theta} - \boldsymbol{\theta}^*) + \frac{1 - \delta(r_0)}{2} r_0^2 \geq -\|\boldsymbol{\xi}\| r_0 + \frac{1 - \delta(r_0)}{2} r_0^2. \end{aligned}$$

Here $\boldsymbol{\theta}^\circ$ is a point from $\Omega(\mathbf{x})$ on the interval connecting $\boldsymbol{\theta}$ and $\boldsymbol{\theta}^*$. If $\|\boldsymbol{\xi}\| \leq r_0\{1 - \delta(r_0)\}/2$, then this implies $L(\boldsymbol{\theta}^*) - L(\boldsymbol{\theta}) > 0$, and the result follows. \square

As a corollary, we obtain Fisher and Wilks expansions for the quasi MLE $\tilde{\boldsymbol{\theta}}$ in a generalized linear model.

Theorem 3.5. *Suppose the conditions of Theorem 3.4 for some r_0 . Then it holds on a set $\Omega(\mathbf{x})$ with $\mathbb{P}(\Omega(\mathbf{x})) \geq 1 - 2e^{-\mathbf{x}}$*

$$\|D(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) - \boldsymbol{\xi}\| \leq r_0\delta(r_0),$$

$$|2L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) - \|\boldsymbol{\xi}\|^2| \leq 2r_0^2\delta(r_0) + r_0^2\delta^2(r_0),$$

$$|\sqrt{2L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*)} - \|\boldsymbol{\xi}\|| \leq 3r_0\delta(r_0).$$

Proof. The large deviation bound of Theorem 3.4 allows to restrict the whole parameter space to the local vicinity $\Theta_0(r_0)$. In this vicinity, the log-likelihood ratio $L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) = L(\boldsymbol{\theta}) - L(\boldsymbol{\theta}^*)$ can be well approximated by the quadratic expansion $\mathbb{L}(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$:

$$L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) = (S - \mathbb{E}S)^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^*) + \mathbb{E}S^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^*) - A(\boldsymbol{\theta}) + A(\boldsymbol{\theta}^*),$$

$$\mathbb{L}(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \stackrel{\text{def}}{=} (S - \mathbb{E}S)^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^*) - \frac{1}{2} \|D(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2.$$

Lemma 3.6. Suppose (3.10) for some r_0 . The difference $L(\boldsymbol{\theta}) - \mathbb{L}(\boldsymbol{\theta})$ is deterministic and it holds for each $\boldsymbol{\theta} \in \Theta_0(r_0)$

$$|L(\boldsymbol{\theta}) - L(\boldsymbol{\theta}^*) - \mathbb{L}(\boldsymbol{\theta}, \boldsymbol{\theta}^*)| \leq \frac{\delta(r_0)}{2} \|D(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2 \leq \frac{\delta(r_0)}{2} r_0^2,$$

$$\|D^{-1}\{\nabla L(\boldsymbol{\theta}) - \nabla \mathbb{L}(\boldsymbol{\theta})\}\| \leq r_0\delta(r_0). \quad (3.13)$$

Proof. The linear stochastic terms $(S - \mathbb{E}S)^\top \boldsymbol{\theta}$ are the same for $L(\boldsymbol{\theta})$ and $\mathbb{L}(\boldsymbol{\theta})$. For the deterministic terms $\mathbb{E}S^\top \boldsymbol{\theta} - A(\boldsymbol{\theta})$ we use the Taylor formula of the second order at $\boldsymbol{\theta}^*$, the extreme point equation $\nabla A(\boldsymbol{\theta}^*) = \mathbb{E}S$, and the definition $D^2 = \mathbb{F}(\boldsymbol{\theta}^*)$:

$$|\mathbb{E}L(\boldsymbol{\theta}) - \mathbb{E}\mathbb{L}(\boldsymbol{\theta})| = |A(\boldsymbol{\theta}) - A(\boldsymbol{\theta}^*) - (\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top \nabla A(\boldsymbol{\theta}^*) - \|D(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2/2|$$

$$= \frac{1}{2} |(\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top \{\mathbb{F}(\boldsymbol{\theta}^*) - \mathbb{F}(\boldsymbol{\theta}^c)\}(\boldsymbol{\theta} - \boldsymbol{\theta}^*)|,$$

where $\boldsymbol{\theta}^c$ is a point on the interval between $\boldsymbol{\theta}$ and $\boldsymbol{\theta}^*$. Now the condition (3.10) implies

$$|\mathbb{E}L(\boldsymbol{\theta}) - \mathbb{E}\mathbb{L}(\boldsymbol{\theta})| \leq \frac{\delta(r_0)}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top D^2 (\boldsymbol{\theta} - \boldsymbol{\theta}^*) = \frac{\delta(r_0)}{2} \|D(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2 \leq \frac{\delta(r_0)}{2} r_0^2$$

and the first assertion follows. The second one can be proved similarly. \square

With the approximation (3.13), all the statements of the theorem follow from the general results of Theorem 2.6. \square

To complete the study of a generalized linear model, we translate the general conditions of Theorem 3.4 into conditions on the design Ψ and on individual errors ε_i .

- *Design regularity* is measured by the value

$$\delta_\Psi \stackrel{\text{def}}{=} \max_i \|D^{-1}\boldsymbol{\psi}_i\|.$$

In the case of a regular or random design, the Fisher design matrix $D^2 = \mathbb{F}(\boldsymbol{\theta}^*)$ is proportional to the sample size and thus, the value δ_Ψ is of order $n^{-1/2}$. Our results only apply if this value is small, in particular, the condition $\delta_\Psi < 1/2$ has to be fulfilled.

- *Exponential moments of the errors* $\varepsilon_i = Y_i - \mathbb{E}Y_i$. Suppose that for some values s_i and fixed constants $C_0, \lambda_0 > 0$

$$\mathbb{E} \exp(\lambda_0 s_i^{-1} \varepsilon_i) \leq C_0, \quad i = 1, \dots, n. \quad (3.14)$$

This condition means that the errors ε_i have exponential moments. In most of cases one can use $s_i^2 = \text{Var}(Y_i)$. Condition (3.14) implies that there are another constants $\mathfrak{g}_1 \leq \lambda_0$ and ν_0 such that the following condition is fulfilled:

$$\mathbb{E} \exp(\lambda_0 s_i^{-1} \varepsilon_i) \leq \frac{1}{2} \nu_0^2 \lambda^2, \quad i = 1, \dots, n, |\lambda| \leq \mathfrak{g}_1. \quad (3.15)$$

This follows from the fact that each function $\log \mathbb{E} \exp(\lambda_0 s_i^{-1} \varepsilon_i)$ analytic in λ in a vicinity of the point zero and can be well approximated by $\lambda^2/2$; see Golubev and Spokoiny [11] for more details.

- *Noise homogeneity* is measured by the variability of the values s_i :

$$a_s \stackrel{\text{def}}{=} \max_{i,j=1,\dots,n} s_i/s_j. \quad (3.16)$$

- *Smoothness of the link function* $g(v)$ can be measured by its third derivative. It will be assumed that given \mathbf{x} , there is a constant $a_g(\mathbf{x})$

$$\frac{|g'''(\Psi_i^\top \boldsymbol{\theta})|}{g''(\Psi_i^\top \boldsymbol{\theta}^*)} \leq a_g(\mathbf{x}), \quad \boldsymbol{\theta} \in \Theta_0(\mathbf{x}), i = 1, \dots, n. \quad (3.17)$$

- *Signal-to-noise ratio* is measured by relationship between the matrices D^2 and V^2 , where the matrix V^2 defined as

$$V^2 \stackrel{\text{def}}{=} \sum_{i=1}^n s_i^2 \Psi_i \Psi_i^\top. \quad (3.18)$$

If the observation Y_i follow the GLM assumption P_{v_i} for $v_i = \Psi_i^\top \boldsymbol{\theta}^*$, that is, the model is correctly specified, then $\text{Var}(Y_i) = g''(v_i)$ and the matrices V^2 and D^2 coincide. In the general case under a possible model misspecification, the matrices V^2 and D^2 may be different. In this case we need an identifiability condition

$$V^2 \leq \alpha^2 D^2 \quad (3.19)$$

for some $\alpha > 0$. This condition can be spelled out as

$$\sum_{i=1}^n s_i^2 \Psi_i \Psi_i^\top \leq \alpha^2 \sum_{i=1}^n g''(\Psi_i^\top \boldsymbol{\theta}^*) \Psi_i \Psi_i^\top.$$

First we discuss a deviation bound for the norm of the vector $\boldsymbol{\xi}$ given by

$$\boldsymbol{\xi} = D^{-1}(S - \mathbb{E}S) = D^{-1} \sum_{i=1}^n \varepsilon_i \Psi_i.$$

The squared norm $\|\boldsymbol{\xi}\|^2$ is a quadratic form of the ε_i 's and one can directly apply general results for quadratic forms from Appendix A.

Theorem 3.7. Suppose (3.15), (3.16), (3.17), and (3.19). For $z(p, \mathbf{x})$ from (A.3) with $z(p, \mathbf{x}) \leq \sqrt{p} + \sqrt{2\mathbf{x}}$, fix

$$\mathbf{r}_0 = 4\nu_0 z(p, \mathbf{x}), \quad (3.20)$$

and suppose that δ_Ψ is small enough to ensure

$$a_g(\mathbf{r}_0) \delta_\Psi \mathbf{r}_0 < 1/2. \quad (3.21)$$

Then the conditions of Theorem 3.4 are fulfilled with $\delta(\mathbf{r}_0) \leq a_g(\mathbf{r}_0) \delta_\Psi \mathbf{r}_0$ and the results of this theorem continue to apply.

Proof. Let \mathbf{r}_0 be fixed by (3.20). First we bound the value $\delta(\mathbf{r}_0)$.

Lemma 3.8. The condition (3.10) is fulfilled with $\delta(\mathbf{r}_0) = a_g(\mathbf{r}_0) \delta_\Psi \mathbf{r}_0$.

Proof. For each $\boldsymbol{\theta} \in \Theta_0(\mathbf{r}_0)$ and $i \leq n$, it holds by (3.23)

$$|\boldsymbol{\Psi}_i^\top \boldsymbol{\theta} - \boldsymbol{\Psi}_i^\top \boldsymbol{\theta}^*| = |(D^{-1} \boldsymbol{\Psi}_i)^\top D(\boldsymbol{\theta} - \boldsymbol{\theta}^*)| \leq \|D^{-1} \boldsymbol{\Psi}_i\| \mathbf{r}_0 \leq \delta_\Psi \mathbf{r}_0. \quad (3.22)$$

This implies for the difference $\mathbb{F}(\boldsymbol{\theta}) - \mathbb{F}(\boldsymbol{\theta}^*)$

$$\mathbb{F}(\boldsymbol{\theta}) - \mathbb{F}(\boldsymbol{\theta}^*) = \sum_{i=1}^n \{g''(\boldsymbol{\Psi}_i^\top \boldsymbol{\theta}) - g''(\boldsymbol{\Psi}_i^\top \boldsymbol{\theta}^*)\} \boldsymbol{\Psi}_i \boldsymbol{\Psi}_i^\top.$$

Next, for each $i \leq n$, there exists a point $\boldsymbol{\theta}^\circ$ on the interval between $\boldsymbol{\theta}^*$ and $\boldsymbol{\theta}$ (possibly depending on i) such that

$$g''(\boldsymbol{\Psi}_i^\top \boldsymbol{\theta}) - g''(\boldsymbol{\Psi}_i^\top \boldsymbol{\theta}^*) = \frac{g'''(\boldsymbol{\Psi}_i^\top \boldsymbol{\theta}^\circ)}{g''(\boldsymbol{\Psi}_i^\top \boldsymbol{\theta}^*)} (\boldsymbol{\Psi}_i^\top \boldsymbol{\theta} - \boldsymbol{\Psi}_i^\top \boldsymbol{\theta}^*) g''(\boldsymbol{\Psi}_i^\top \boldsymbol{\theta}^*).$$

Now (3.17) and (3.22) imply

$$\max_{i \leq n} \left| \frac{g'''(\boldsymbol{\Psi}_i^\top \boldsymbol{\theta}^\circ)}{g''(\boldsymbol{\Psi}_i^\top \boldsymbol{\theta}^*)} (\boldsymbol{\Psi}_i^\top \boldsymbol{\theta} - \boldsymbol{\Psi}_i^\top \boldsymbol{\theta}^*) \right| \leq a_g(\mathbf{r}_0) \delta_\Psi \mathbf{r}_0$$

and

$$\pm \{\mathbb{F}(\boldsymbol{\theta}) - \mathbb{F}(\boldsymbol{\theta}^*)\} \leq a_g(\mathbf{r}_0) \delta_\Psi \mathbf{r}_0 \sum_{i=1}^n g''(\boldsymbol{\Psi}_i^\top \boldsymbol{\theta}^*) \boldsymbol{\Psi}_i \boldsymbol{\Psi}_i^\top = a_g(\mathbf{r}_0) \delta_\Psi \mathbf{r}_0 D^2$$

yielding (3.10) in an obvious way. □

This lemma and (3.21) imply $\delta(\mathbf{r}_0) < 1/2$. Now we show that (3.15) implies (3.11).

Lemma 3.9. *Let the errors $\varepsilon_i = Y_i - \mathbb{E}Y_i$ be independent and follow (3.15). Then*

$$\log \mathbb{E} \exp\{\mathbf{u}^\top V^{-1}(S - \mathbb{E}S)\} \leq \frac{\nu_0^2}{2} \|\mathbf{u}\|^2, \quad \|\mathbf{u}\| \leq \mathfrak{g} \stackrel{\text{def}}{=} \frac{\mathfrak{g}_1}{\delta_\Psi a_S}, \quad (3.23)$$

where V^2 is from (3.18) and a_S from (3.16).

Proof. The formula (3.8) and independence of the ε_i 's imply for any vector $\mathbf{u} \in \mathbb{R}^p$ with $\|\mathbf{u}\| \leq \mathfrak{g}$

$$\log \mathbb{E} \exp\{\mathbf{u}^\top V^{-1}(S - \mathbb{E}S)\} = \sum_{i=1}^n \log \mathbb{E} \exp(\lambda_i s_i^{-1} \varepsilon_i),$$

where the definition (3.23) of \mathfrak{g} and condition (3.16) imply for $\lambda_i = s_i \mathbf{u}^\top V^{-1} \boldsymbol{\Psi}_i$

$$|\lambda_i| = |\mathbf{u}^\top V^{-1} \boldsymbol{\Psi}_i| s_i \leq \mathfrak{g} \|V^{-1} \boldsymbol{\Psi}_i\| s_i \leq \mathfrak{g}_1.$$

Therefore, by (3.15) and the definition of V^2

$$\log \mathbb{E} \exp\{\mathbf{u}^\top V^{-1}(S - \mathbb{E}S)\} \leq \frac{\nu_0^2}{2} \sum_{i=1}^n \lambda_i^2 = \frac{\nu_0^2}{2} \sum_{i=1}^n \mathbf{u}^\top V^{-1} (\boldsymbol{\Psi}_i \boldsymbol{\Psi}_i^\top s_i^2) V^{-1} \mathbf{u} = \frac{\nu_0^2}{2} \|\mathbf{u}\|^2,$$

and the assertion follows. □

The result of Lemma 3.9 provides exponential moments of ξ and one can apply Theorem A.1 from Appendix A yielding the bound (3.11) under the condition

$$\frac{1 - \delta(\varkappa_0)}{2} \varkappa_0 \geq \nu_0 z(p, \mathbf{x})$$

which is obviously fulfilled for our choice of $\varkappa_0 = 4\nu_0 z(p, \mathbf{x})$ in view of $\delta(\varkappa_0) < 1/2$. This will also provide (3.11). All the conditions of Theorem 3.4 have been checked. \square

3.7. Estimation for a penalized GLM

This section briefly discusses what will be changed if the GLM (3.6) is penalized by a roughness penalty term $\|G\boldsymbol{\theta}\|^2/2$. The corresponding penalized log-likelihood $L_G(\boldsymbol{\theta})$ reads as

$$L_G(\boldsymbol{\theta}) = S^\top \boldsymbol{\theta} - A(\boldsymbol{\theta}) - \|G\boldsymbol{\theta}\|^2/2.$$

The penalized MLE and its target are defined by maximizing $L_G(\boldsymbol{\theta})$ and its expectation:

$$\begin{aligned} \tilde{\boldsymbol{\theta}}_G &\stackrel{\text{def}}{=} \operatorname{argmax}_{\boldsymbol{\theta} \in \Theta} \{S^\top \boldsymbol{\theta} - A(\boldsymbol{\theta}) - \|G\boldsymbol{\theta}\|^2/2\}, \\ \boldsymbol{\theta}_G^* &\stackrel{\text{def}}{=} \operatorname{argmax}_{\boldsymbol{\theta} \in \Theta} \{\mathbb{E}S^\top \boldsymbol{\theta} - A(\boldsymbol{\theta}) - \|G\boldsymbol{\theta}\|^2/2\}. \end{aligned} \tag{3.24}$$

Further, define the matrix D_G by $D_G^2 = \mathbb{F}_G(\boldsymbol{\theta}_G^*)$ for

$$\mathbb{F}_G(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \mathbb{F}(\boldsymbol{\theta}) + G^2 = \sum_{i=1}^n \boldsymbol{\psi}_i \boldsymbol{\psi}_i^\top g''(\boldsymbol{\psi}_i^\top \boldsymbol{\theta}) + G^2. \tag{3.25}$$

One can see that the use of penalization leads to a growth of the ‘‘information matrix’’ D_G^2 relative to the non-penalized case. The stochastic term $(S - \mathbb{E}S)^\top \boldsymbol{\theta}$ of $L_G(\boldsymbol{\theta})$ remains the same as in the non-penalized case, thus, the matrix V^2 from (3.18) can be used here as well and the identifiability condition (3.19) continues to hold.

The local vicinity $\Theta_{0,G}(\varkappa)$ of $\boldsymbol{\theta}_G^*$ is now defined as

$$\Theta_{0,G}(\varkappa) \stackrel{\text{def}}{=} \{\boldsymbol{\theta}: \|D_G(\boldsymbol{\theta} - \boldsymbol{\theta}_G^*)\| \leq \varkappa\}.$$

The concentration result for $\tilde{\boldsymbol{\theta}}_G$ can be easily extended to the penalized case.

Theorem 3.10. *Let, for some $\varkappa_G > 0$, the matrix function $\mathbb{F}_G(\boldsymbol{\theta})$ from (3.25) fulfill with $D_G^2 = \mathbb{F}_G(\boldsymbol{\theta}_G^*)$*

$$\sup_{\boldsymbol{\theta} \in \Theta_{0,G}(\varkappa_0)} \|D_G^{-1} \mathbb{F}_G(\boldsymbol{\theta}) D_G^{-1} - \mathbb{I}_p\|_{\text{op}} \leq \delta(\varkappa_G)$$

for $\delta(\varkappa_G) < 1$. Let also S from (3.8) follow for $\mathbf{x} > 0$ the probability bound

$$\mathbb{P}(\|D_G^{-1}(S - \mathbb{E}S)\| > z(B_G, \mathbf{x})) \leq 2e^{-\mathbf{x}}.$$

If

$$\varkappa_G \{1 - \delta(\varkappa_G)\} \geq 2z(B_G, \mathbf{x}),$$

with $z(B_G, \mathbf{x}) \leq \sqrt{pG} + \sqrt{2\mathbf{x}}$ from (A.4), then the solution $\tilde{\boldsymbol{\theta}}_G$ of (3.25) satisfies

$$\mathbb{P}(\tilde{\boldsymbol{\theta}} \notin \Theta_{0,G}(\varkappa_G)) \leq 2e^{-\mathbf{x}}.$$

Then all the statements of Theorem 3.5 hold for the pair $\tilde{\boldsymbol{\theta}}_G, \boldsymbol{\theta}_G^*$ with $\boldsymbol{\xi}_G \stackrel{\text{def}}{=} D_G^{-1}(S - \mathbb{E}S)$. Then it holds on a set $\Omega_G(\mathbf{x})$ with $\mathbb{P}(\Omega_G(\mathbf{x})) \geq 1 - 2e^{-\mathbf{x}}$

$$\begin{aligned} \|D_G(\tilde{\boldsymbol{\theta}}_G - \boldsymbol{\theta}_G^*) - \boldsymbol{\xi}_G\| &\leq r_G \delta(r_G), \\ |2L_G(\tilde{\boldsymbol{\theta}}_G, \boldsymbol{\theta}_G^*) - \|\boldsymbol{\xi}_G\|^2| &\leq 2r_G^2 \delta(r_G) + r_G^2 \delta^2(r_G), \\ \left| \sqrt{2L_G(\tilde{\boldsymbol{\theta}}_G, \boldsymbol{\theta}_G^*)} - \|\boldsymbol{\xi}_G\| \right| &\leq 3r_G \delta(r_G). \end{aligned}$$

The proof of the non-penalized case applies here with obvious changes in notation. However, at one place the difference is essential. Namely, the radius r_G can be much smaller and it depends on the effective dimension $\mathfrak{p}_G = \text{tr}(B_G) = \text{tr}(D_G^{-1}V^2D_G^{-1})$ rather than on the total dimension p .

Appendix A: Deviation bounds for quadratic forms

This section collects some probability bounds for non-Gaussian quadratic forms. The presented results can be viewed as a slight improvement of the bounds from Spokoiny [25] using the deviation bound from Laurent and Massart [17]. The proofs are very similar to ones from Spokoiny [25] and are omitted by the space reasons.

Let a random vector $\boldsymbol{\xi} \in \mathbb{R}^p$ has some exponential moments. More exactly, suppose for some fixed $g > 0$ that

$$\log \mathbb{E} \exp(\boldsymbol{\gamma}^\top \boldsymbol{\xi}) \leq \|\boldsymbol{\gamma}\|^2/2, \quad \boldsymbol{\gamma} \in \mathbb{R}^p, \|\boldsymbol{\gamma}\| \leq g. \quad (\text{A.1})$$

For ease of presentation, assume below that g is sufficiently large, namely, $0.3g \geq \sqrt{p}$. In typical examples of an i.i.d. sample, $g \asymp \sqrt{n}$. Define

$$\begin{aligned} x_c &\stackrel{\text{def}}{=} g^2/4, \\ z_c^2 &\stackrel{\text{def}}{=} p + \sqrt{pg^2} + g^2/2 = g^2(1/2 + \sqrt{p/g^2} + p/g^2), \\ g_c &\stackrel{\text{def}}{=} \frac{g(1/2 + \sqrt{p/g^2} + p/g^2)^{1/2}}{1 + \sqrt{p/g^2}}. \end{aligned}$$

Note that with $\alpha = \sqrt{p/g^2} \leq 0.3$, one has

$$\begin{aligned} z_c^2 &= g^2(1/2 + \alpha + \alpha^2), \\ g_c &= g \frac{(1/2 + \alpha + \alpha^2)^{1/2}}{1 + \alpha} \end{aligned}$$

so that $z_c^2/g^2 \in [1/2, 1]$ and $g_c^2/g^2 \in [1/2, 1]$.

Theorem A.1. *Let (A.1) hold and $0.3g \geq \sqrt{p}$. Then for each $x > 0$*

$$\mathbb{P}(\|\boldsymbol{\xi}\| \geq z(p, x)) \leq 2e^{-x} + 8.4e^{-x_c} \mathbb{1}(x < x_c), \quad (\text{A.2})$$

where $z(p, x)$ is defined by

$$z(p, x) \stackrel{\text{def}}{=} \begin{cases} (p + 2\sqrt{px} + 2x)^{1/2}, & x \leq x_c, \\ z_c + 2g_c^{-1}(x - x_c), & x > x_c. \end{cases} \quad (\text{A.3})$$

Depending on the value x , we have two types of tail behavior of the quadratic form $\|\boldsymbol{\xi}\|^2$. For $x \leq x_c = g^2/4$, we have the same deviation bounds as in the Gaussian case with the extra-factor two in the deviation probability. Remind that one can use a simplified expression $(p + 2\sqrt{px} + 2x)^{1/2} \leq \sqrt{p} + \sqrt{2x}$. For $x > x_c$, we switch to the special

regime driven by the exponential moment condition (A.1). Usually g^2 is a large number (of order n in the i.i.d. setup) and the second term in (A.2) can be simply ignored.

Next we present a bound for a quadratic form $\xi^\top B \xi$, where ξ satisfies (A.1) and B is a given symmetric positive $p \times p$ matrix. Define

$$p \stackrel{\text{def}}{=} \text{tr}(B), \quad v^2 \stackrel{\text{def}}{=} \text{tr}(B^2), \quad \lambda \stackrel{\text{def}}{=} \lambda_{\max}(B).$$

For ease of presentation, suppose that $0.3g \geq \sqrt{p}$ so that $\alpha = \sqrt{p/g^2} \leq 0.3$. The other case only changes the constants in the inequalities. Define also

$$\begin{aligned} x_c &\stackrel{\text{def}}{=} g^2/4, \\ z_c^2 &\stackrel{\text{def}}{=} p + vg + \lambda g^2/2, \\ g_c &\stackrel{\text{def}}{=} \frac{\sqrt{p/\lambda + gv/\lambda + g^2/2}}{1 + v/(\lambda g)}. \end{aligned}$$

Theorem A.2. *Let (A.1) hold and $0.3g \geq \sqrt{p/\lambda}$. Then for each $x > 0$*

$$\mathbb{P}(\|B^{1/2}\xi\| \geq z(B, x)) \leq 2e^{-x} + 8.4e^{-x_c} \mathbb{1}(x < x_c),$$

where $z(B, x)$ is defined by

$$z(B, x) \stackrel{\text{def}}{=} \begin{cases} \sqrt{p + 2vx^{1/2} + 2\lambda x}, & x \leq x_c, \\ z_c + 2\lambda(x - x_c)/g_c, & x > x_c. \end{cases} \quad (\text{A.4})$$

Similarly to the case $B = \mathbb{I}_p$, the upper quantile $z(B, x) = \sqrt{p + 2vx^{1/2} + 2\lambda x}$ can be upper bounded by $\sqrt{p} + \sqrt{2\lambda x}$:

$$z(B, x) \leq \begin{cases} \sqrt{p} + \sqrt{2\lambda x}, & x \leq x_c, \\ z_c + 2\lambda(x - x_c)/g_c, & x > x_c. \end{cases}$$

Appendix B: Deviation bounds for random processes

This chapter presents some general results of the theory of empirical processes. We assume some exponential moment conditions on the increments of the process which allow to apply the well developed chaining arguments in Orlicz spaces; see e.g. van der Vaart and Wellner [30], Chapter 2.2. We state the results in a slightly different form and present an independent and self-contained proof.

B.1. Chaining and covering numbers

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

($\mathcal{E}d$) There exist $g > 0$, $r_0 > 0$, $\nu_0 \geq 1$, such that for any $\lambda \leq g$ and $\mathbf{v}, \mathbf{v}' \in \mathcal{Y}$ with $d(\mathbf{v}, \mathbf{v}') \leq r_0$

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

By $\mathcal{B}_r(\mathbf{v})$ we denote the d -ball centered at \mathbf{v} of radius r :

$$\mathcal{B}_r(\mathbf{v}) \stackrel{\text{def}}{=} \{\mathbf{v}' \in \mathcal{Y} : d(\mathbf{v}, \mathbf{v}') \leq r\}.$$

Let Υ° be a subset of a ball in Υ with center at \mathbf{v}^* and radius r_0 , and let a sequence r_k be fixed with $r_k = r_0 2^{-k}$. For each k , by \mathcal{M}_k we denote a r_k -net in Υ° , so that

$$\Upsilon^\circ \subseteq \bigcup_{\mathbf{v} \in \mathcal{M}_k} \mathcal{B}_{r_k}(\mathbf{v}).$$

Let also $\Pi_k \mathbf{v}$ be the closest to \mathbf{v} point from \mathcal{M}_k , so that $d(\mathbf{v}, \Pi_k \mathbf{v}) \leq r_k$. We assume that \mathcal{M}_0 consists of one point \mathbf{v}^* , that is, $\Pi_0 \mathbf{v} = \mathbf{v}^*$. Let $\mathbb{N}_k \stackrel{\text{def}}{=} |\mathcal{M}_k|$ denote the cardinality of \mathcal{M}_k . Finally set $c_k = 2^{-k}$ for $k \geq 1$, and define the values $\mathbb{Q}_1(\Upsilon^\circ)$ and $\mathbb{Q}_2(\Upsilon^\circ)$ by

$$\begin{aligned} \mathbb{Q}_1(\Upsilon^\circ) &\stackrel{\text{def}}{=} \sum_{k=1}^{\infty} c_k \sqrt{2 \log(2\mathbb{N}_k)} = \sum_{k=1}^{\infty} 2^{-k} \sqrt{2 \log(2\mathbb{N}_k)}, \\ \mathbb{Q}_2(\Upsilon^\circ) &\stackrel{\text{def}}{=} \sum_{k=1}^{\infty} 2c_k \log(2\mathbb{N}_k) = \sum_{k=1}^{\infty} 2^{-k+1} \log(2\mathbb{N}_k). \end{aligned} \tag{B.1}$$

By the Cauchy–Schwarz inequality $\mathbb{Q}_1^2(\Upsilon^\circ) \leq \mathbb{Q}_2(\Upsilon^\circ)$. The inverse relation is not generally true and one can build some examples with $\mathbb{Q}_1(\Upsilon^\circ)$ finite and $\mathbb{Q}_2(\Upsilon^\circ)$ infinite. If the process $\mathcal{U}(\mathbf{v})$ is sub-Gaussian and $(\mathcal{E}d)$ is fulfilled with $g = \infty$, then one can only operate with $\mathbb{Q}_1(\Upsilon^\circ)$ which is equivalent to the Dudley integral.

Theorem B.1. *Let \mathcal{U} be a separable process and Υ° be a ball in Υ with center \mathbf{v}° and radius r_0 for the distance $d(\cdot, \cdot)$, i.e. $d(\mathbf{v}, \mathbf{v}^\circ) \leq r_0$ for all $\mathbf{v} \in \Upsilon^\circ$. If $(\mathcal{E}d)$ holds with $g = \infty$ then for any $x \geq 1/2$, it holds with $\mathbb{Q}_1 = \mathbb{Q}_1(\Upsilon^\circ)$ and $\mathbb{Q}_2 = \mathbb{Q}_2(\Upsilon^\circ)$*

$$\mathbb{P}\left(\frac{1}{v_0 r_0} \sup_{\mathbf{v} \in \Upsilon^\circ} \mathcal{U}(\mathbf{v}, \mathbf{v}^*) \geq \mathfrak{J}_{\mathbb{H}}(x)\right) \leq e^{-x} \tag{B.2}$$

with

$$\mathfrak{J}_{\mathbb{H}}(x) \stackrel{\text{def}}{=} 2\mathbb{Q}_1 + \sqrt{8x}.$$

If $g < \infty$ in $(\mathcal{E}d)$ then (B.2) holds with $\mathfrak{J}_{\mathbb{H}}(x)$ given by one of the following rules:

$$\begin{aligned} \mathfrak{J}_{\mathbb{H}}(x) &= 2\mathbb{Q}_1 + \sqrt{8x} + 2g^{-1}(g^{-2}x + 1)\mathbb{Q}_2, \\ \mathfrak{J}_{\mathbb{H}}(x) &= \begin{cases} 2\sqrt{\mathbb{Q}_2 + 2x}, & \text{if } \mathbb{Q}_2 + 2x \leq g^2, \\ 2g^{-1}x + g^{-1}\mathbb{Q}_2 + g, & \text{if } \mathbb{Q}_2 + 2x > g^2. \end{cases} \end{aligned} \tag{B.3}$$

Moreover, the r.v. $\mathcal{U}^*(r_0) \stackrel{\text{def}}{=} \sup_{\mathbf{v} \in \Upsilon^\circ} \mathcal{U}(\mathbf{v}, \mathbf{v}^*)$ fulfills

$$\begin{aligned} \mathbb{E}\mathcal{U}^*(r_0) &\leq 2v_0 r_0 (\mathbb{Q}_1 + \mathbb{Q}_2/g + 3), \\ \{\mathbb{E}|\mathcal{U}^*(r_0)|^2\}^{1/2} &\leq 2v_0 r_0 (\mathbb{Q}_1 + \mathbb{Q}_2/g + 4). \end{aligned} \tag{B.4}$$

Proof. We start the proof by stating some general facts for a convex combinations of sub-exponential r.v.’s ζ_k such that

$$\log \mathbb{E} \exp(\lambda \zeta_k) \leq (q_k^2 + \lambda^2)/2, \quad |\lambda| \leq g, k = 0, 1, 2, \dots, \tag{B.5}$$

where $q_k \geq 1$ are fixed numbers, and g is some positive value or infinity. We aim at bounding a sum S of the form $S = \sum_k c_k \zeta_k$ for a sequence of positive weights c_k satisfying $\sum_k c_k = 1$. We implicitly assume that the numbers q_k grow with k in a way that $\sum_k \exp(-q_k) \leq 1$. Define

$$\mathbb{H}_1 \stackrel{\text{def}}{=} \sum_k c_k q_k, \quad \mathbb{H}_2 \stackrel{\text{def}}{=} \sum_k c_k q_k^2.$$

Lemma B.2. Suppose that random variables ζ_k follow (B.5) with $g = \infty$ and $\sum_k \exp(-q_k) \leq 1$. Let also $\sum_k c_k = 1$. Then it holds for the sum $S = \sum_k c_k \zeta_k$ and $x \geq 1/2$

$$\begin{aligned} \log \mathbb{E} \exp(S) &\leq \mathbb{H}_1, \\ \mathbb{P}(S \geq \mathbb{H}_1 + \sqrt{2x}) &\leq e^{-x}. \end{aligned} \tag{B.6}$$

If (B.5) holds for $g < \infty$, then for each $\lambda > 0$ with $|\lambda| \leq g$ and $x \geq 1/2$

$$\log \mathbb{E} \exp(\lambda S) \leq (\mathbb{H}_2 + \lambda^2)/2, \tag{B.7}$$

$$\mathbb{P}\{S \geq \mathfrak{z}_{\mathbb{H}}(x)\} \leq e^{-x}, \tag{B.8}$$

where $\mathfrak{z}_{\mathbb{H}}(x)$ is given by (B.3). Moreover, if $g^2 \geq \mathbb{H}_2 + 1$, then

$$\mathbb{E} S \leq \mathbb{H}_1 + \mathbb{H}_2/g + 3, \quad \{\mathbb{E} S^2\}^{1/2} \leq \mathbb{H}_1 + \mathbb{H}_2/g + 4.$$

Proof. Consider first the sub-Gaussian case with $g = \infty$. Define $\alpha_k = c_k/q_k$. Obviously $\sum_k \alpha_k \leq \sum_k c_k = 1$. By the Hölder inequality and (B.5), it holds

$$\begin{aligned} \log \mathbb{E} \exp\left(\sum_k c_k \zeta_k\right) &= \log \mathbb{E} \exp\left(\sum_k \alpha_k q_k \zeta_k\right) \leq \sum_k \alpha_k \log \mathbb{E} \exp(q_k \zeta_k) \\ &\leq \frac{1}{2} \sum_k \alpha_k (q_k^2 + q_k^2) \leq \sum_k c_k q_k. \end{aligned}$$

Further, by the same arguments, it holds

$$\log \mathbb{E} \exp(\lambda S) \leq \sum_k c_k \log \mathbb{E} \exp(\lambda \zeta_k) \leq \frac{1}{2} \sum_k c_k (q_k^2 + \lambda^2)$$

and the assertion (B.7) follows as well.

Let $x \geq 1/2$ be fixed. With $z_k = q_k + \sqrt{2x}$, it follows by (B.5) for $\lambda_k = z_k$ in view of $\sum_k e^{-q_k} \leq 1$

$$\begin{aligned} &\mathbb{P}\left(\sum_k c_k (\zeta_k - z_k) \geq 0\right) \\ &\leq \sum_k \mathbb{P}(\zeta_k - z_k \geq 0) \leq \sum_k \mathbb{E} \exp\{\lambda_k (\zeta_k - z_k)\} \\ &\leq \sum_k \exp\left(-\lambda_k z_k + \frac{\lambda_k^2}{2} + \frac{q_k^2}{2}\right) = \sum_k \exp\left(-\frac{z_k^2}{2} + \frac{q_k^2}{2}\right) = \sum_k e^{-x - q_k \sqrt{2x}} \leq e^{-x}. \end{aligned} \tag{B.9}$$

This implies (B.6) by

$$\sum_k c_k z_k = \sum_k c_k (q_k + \sqrt{2x}) = \mathbb{H}_1 + \sqrt{2x}. \tag{B.10}$$

Now we briefly discuss how the condition (B.5) can be relaxed to the case of a finite g . Suppose that (B.5) holds for all $\lambda \leq g < \infty$. Define $k(x)$ as the largest index k , for which $\lambda_k = q_k + \sqrt{2x} \leq g$. For $k > k(x)$, define $\lambda_k = g$ and

$$z_k = \frac{x + q_k}{g} + \frac{g}{2} + \frac{q_k^2}{2g}. \tag{B.11}$$

The above arguments yield for $k > k(x)$

$$\mathbb{P}(\zeta_k \geq z_k) \leq \exp\{-g z_k + (q_k^2 + g^2)/2\} = \exp(-x - q_k).$$

This and (B.9) yield

$$\sum_k \mathbb{P}(\zeta_k \geq z_k) \leq \sum_{k \leq k(x)} e^{-x - q_k \sqrt{2x}} + \sum_{k > k(x)} e^{-x - q_k} \leq \sum_k e^{-x - q_k} \leq e^{-x}.$$

Further, as $q_k > g$ for $k > k(x)$, it follows from the definition (B.11)

$$\begin{aligned} \sum_{k > k(x)} c_k z_k &= \frac{1}{g} \sum_{k > k(x)} c_k (x + q_k) + \frac{g}{2} \sum_{k > k(x)} c_k + \frac{1}{2g} \sum_{k > k(x)} c_k q_k^2 \\ &\leq \frac{1}{g} \sum_{k > k(x)} c_k q_k + \left(\frac{x}{g^3} + \frac{1}{g}\right) \sum_{k > k(x)} c_k q_k^2. \end{aligned}$$

This and (B.10) imply due to $g \geq 1$

$$\sum_k c_k z_k \leq \sum_k c_k q_k + \left(\frac{x}{g^3} + \frac{1}{g}\right) \sum_k c_k q_k^2 + \sqrt{2x} \leq \mathbb{H}_1 + \left(\frac{x}{g^3} + \frac{1}{g}\right) \mathbb{H}_2 + \sqrt{2x}.$$

In particular, if $x \leq g^2$, then

$$\sum_k c_k z_k \leq \mathbb{H}_1 + \frac{2}{g} \mathbb{H}_2 + \sqrt{2x}.$$

Now (B.8) with $\mathfrak{z}(x) = \mathbb{H}_1 + \sqrt{2x} + g^{-1}(g^{-2}x + 1)\mathbb{H}_2$ follows similarly to (B.6). Further, if $\mathfrak{z}(x) = \sqrt{\mathbb{H}_2 + 2x} \leq g$, then (B.7) with $\lambda = \mathfrak{z}(x)$ and the exponential Chebyshev inequality implies again

$$\mathbb{P}(S \geq \mathfrak{z}(x)) \leq \exp\left(-\lambda \mathfrak{z}(x) + \frac{\mathbb{H}_2 + \lambda^2}{2}\right) = \exp\left(\frac{-\mathfrak{z}^2(x) + \mathbb{H}_2}{2}\right) = \exp(-x).$$

Similarly one can check the case with $\lambda = g$ and $\mathfrak{z}(x) = x/g + (\mathbb{H}_2/g + g)/2 > g$.

To bound the moments of S , we apply the following technical result: if $\mathbb{P}(S \geq \mathfrak{z}(x)) \leq e^{-x}$ for all $x \geq x_0$ and if $\mathfrak{z}(\cdot)$ is absolutely continuous, then

$$\mathbb{E}S \leq \mathfrak{z}(x_0) + \int_{x_0}^{\infty} \mathfrak{z}'(x) e^{-x} dx, \quad \mathbb{E}S^2 \leq \mathfrak{z}^2(x_0) + 2 \int_{x_0}^{\infty} \mathfrak{z}(x) \mathfrak{z}'(x) e^{-x} dx.$$

For $\mathfrak{z}(x) = \mathbb{H}_1 + \sqrt{2x} + g^{-1}(g^{-2}x + 1)\mathbb{H}_2$, it holds $\mathfrak{z}'(x) \leq 1 + g^{-3}$. In view of $g^2 \geq \mathbb{H}_2 + 1$

$$\mathbb{E}S \leq \mathbb{H}_1 + 1 + (\mathbb{H}_2 + 1/2)/g + \int_{1/2}^{\infty} (1 + g^{-3}) e^{-x} dx \leq \mathbb{H}_1 + \mathbb{H}_2/g + 3.$$

Similarly one can bound

$$\mathbb{E}S^2 \leq (\mathbb{H}_1 + \mathbb{H}_2/g + 3/2)^2 + 2 \int_{1/2}^{\infty} \left(\frac{1}{\sqrt{2x}} + g^{-3}\right) \mathfrak{z}(x) e^{-x} dx \leq (\mathbb{H}_1 + \mathbb{H}_2/g + 4)^2$$

as required. □

Now we show how the statement of the theorem can be reduced to the bounds of Lemma B.2. Denote for $i < k$ by Π_i^k the product $\Pi_i^k = \Pi_i \Pi_{i+1} \cdots \Pi_k$. As $\Pi_0 \mathbf{v} \equiv \mathbf{v}^*$, the telescopic sum devices yields

$$|\mathcal{U}(\Pi_k \mathbf{v}) - \mathcal{U}(\mathbf{v}^*)| \leq \sum_{i=1}^k |\mathcal{U}(\Pi_{i-1}^k \mathbf{v}) - \mathcal{U}(\Pi_i^k \mathbf{v})|.$$

Separability of $\mathcal{U}(\cdot)$ implies that $\lim_{k \rightarrow \infty} \mathcal{U}(\Pi_k \mathbf{v}) = \mathcal{U}(\mathbf{v})$. Therefore, for any $\mathbf{v} \in \mathcal{Y}^\circ$

$$|\mathcal{U}(\mathbf{v}) - \mathcal{U}(\mathbf{v}^*)| = \lim_{k \rightarrow \infty} |\mathcal{U}(\Pi_k \mathbf{v}) - \mathcal{U}(\mathbf{v}^*)| \leq \sum_{k=1}^{\infty} \xi_k^*,$$

where

$$\xi_k^* \stackrel{\text{def}}{=} \max_{\mathbf{v} \in \mathcal{M}_k} |\mathcal{U}(\mathbf{v}) - \mathcal{U}(\Pi_{k-1} \mathbf{v})|.$$

For each $\mathbf{v} \in \mathcal{M}_k$, it holds $d(\mathbf{v}, \Pi_{k-1} \mathbf{v}) \leq r_{k-1}$ and

$$|\mathcal{U}(\mathbf{v}) - \mathcal{U}(\Pi_{k-1} \mathbf{v})| \leq r_{k-1} \frac{|\mathcal{U}(\mathbf{v}) - \mathcal{U}(\Pi_{k-1} \mathbf{v})|}{d(\mathbf{v}, \Pi_{k-1} \mathbf{v})}.$$

This implies by the Jensen inequality and (E*d*) in view of $e^{|x|} \leq e^x + e^{-x}$ for each $k \geq 1$ and $|\lambda| \leq \mathfrak{g}$

$$\mathbb{E} \exp\left(\frac{\lambda}{r_{k-1}} \xi_k^*\right) \leq 2 \sum_{\mathbf{v} \in \mathcal{M}_k} \mathbb{E} \exp\left(\lambda \frac{|\mathcal{U}(\mathbf{v}) - \mathcal{U}(\Pi_{k-1} \mathbf{v})|}{d(\mathbf{v}, \Pi_{k-1} \mathbf{v})}\right) \leq 2\mathbb{N}_k \exp(\lambda^2/2). \quad (\text{B.12})$$

For $k \geq 1$, define $q_k^2/2 = \log(2\mathbb{N}_k)$, $c_k = 2^{-k}$, and $\zeta_k = \xi_k^*/r_{k-1} = c_k^{-1} \xi_k^*/(2r_0)$. Then (B.12) implies by $r_{k-1} = 2^{-k+1} r_0$

$$\log \mathbb{E} \exp(\lambda \zeta_k) \leq \log(2\mathbb{N}_k) + \lambda^2/2 = (q_k^2 + \lambda^2)/2.$$

Now we apply Lemma B.2 with $c_k = 2^{-k}$. By construction

$$\sum_{k=1}^{\infty} c_k \zeta_k = \frac{1}{2r_0} \sum_{k=1}^{\infty} \xi_k^*$$

and the results follow with $\mathbb{H}_1 = \mathbb{Q}_1(\mathcal{Y}^\circ)$, $\mathbb{H}_2 = \mathbb{Q}_2(\mathcal{Y}^\circ)$. □

B.2. A large deviation bound

Due to the result of Theorem B.1, the bound for the maximum of $\mathcal{U}(\mathbf{v}, \mathbf{v}^*)$ over $\mathbf{v} \in \mathcal{B}_r(\mathbf{v}^*)$ grows linearly in r . So, its applications to situations with $r \gg \mathbb{Q}_1(\mathcal{Y}^\circ)$ are limited. The next result shows that introducing a negative drift helps to state a uniform in r local probability bound. Namely, the bound for the process $\mathcal{U}(\mathbf{v}, \mathbf{v}^*) - f(d(\mathbf{v}, \mathbf{v}^*))$ for some function $f(r)$ over a ball $\mathcal{B}_r(\mathbf{v}^*)$ around the point \mathbf{v}^* does not depend on r . Here the generic chaining arguments are accomplished with the slicing technique. The idea is for a given $r^* > 1$ to split the ball $\mathcal{B}_{r^*}(\mathbf{v}^*)$ into the slices $\mathcal{B}_{r_k}(\mathbf{v}^*) \setminus \mathcal{B}_{r_{k-1}}(\mathbf{v}^*)$ and to apply Theorem B.1 to each slice separately.

Theorem B.3. *Let r^* be such that (E*d*) holds on $\mathcal{B}_{r^*}(\mathbf{v}^*)$. Let also $\mathbb{Q}_1(\mathcal{B}_{r^*}(\mathbf{v}^*)) \leq \mathbb{H}_1$ and $\mathbb{Q}_2(\mathcal{B}_{r^*}(\mathbf{v}^*)) \leq \mathbb{H}_2$ for $r \leq r^*$. Given $r_0 < r^*$, let a monotonous function $f(r, r_0)$ fulfill for some $\rho < 1$*

$$f(r, r_0) \geq \nu_0 r \mathfrak{J}_{\mathbb{H}}(x + \log(r/r_0)), \quad r_0 \leq r \leq r^*, \quad (\text{B.13})$$

where the function $\mathfrak{z}_{\mathbb{H}}(\cdot)$ is given by (B.3). Then it holds

$$\mathbb{P}\left(\sup_{r_0 \leq r \leq r^*} \sup_{\mathbf{v} \in \mathcal{B}_r(\mathbf{v}^*)} \{\mathcal{U}(\mathbf{v}, \mathbf{v}^*) - f(\rho^{-1}r, r_0)\} \geq 0\right) \leq \frac{\rho}{1-\rho} e^{-x}.$$

Remark B.1. Formally the bound applies even with $r^* = \infty$ provided that $(\mathcal{E}d)$ is fulfilled on the whole set Υ° .

Remark B.2. If $\mathfrak{g} = \infty$, then $\mathfrak{z}_{\mathbb{H}}(x) = 2\mathbb{H}_1 + \sqrt{8x}$ and the condition (B.13) on the drift simplifies to $(2v_0r)^{-1}f(r, r_0) \geq \mathbb{H}_1 + \sqrt{2x + 2\log(r/r_0)}$.

Proof of Theorem B.3. By (B.13) and Theorem B.1 for any $r > r_0$

$$\begin{aligned} & \mathbb{P}\left(\sup_{\mathbf{v} \in \mathcal{B}_r(\mathbf{v}^*) \setminus \mathcal{B}_{\rho r}(\mathbf{v}^*)} \{\mathcal{U}(\mathbf{v}, \mathbf{v}^*) - f(r, r_0)\} \geq 0\right) \\ & \leq \mathbb{P}\left(\frac{1}{v_0 r} \sup_{\mathbf{v} \in \mathcal{B}_r(\mathbf{v}^*)} \mathcal{U}(\mathbf{v}, \mathbf{v}^*) \geq \mathfrak{z}(x + \log(r/r_0))\right) \leq \frac{r_0}{r} e^{-x}. \end{aligned} \tag{B.14}$$

Now defined $r_k = r_0 \rho^{-k}$ for $k = 0, 1, 2, \dots$. Define also $k^* \stackrel{\text{def}}{=} \log(r^*/r_0) + 1$. It follows from (B.14) that

$$\begin{aligned} & \mathbb{P}\left(\sup_{\mathbf{v} \in \mathcal{B}_{r^*}(\mathbf{v}^*) \setminus \mathcal{B}_{r_0}(\mathbf{v}^*)} \{\mathcal{U}(\mathbf{v}, \mathbf{v}^*) - f(\rho^{-1}d(\mathbf{v}, \mathbf{v}^*), r_0)\} \geq 0\right) \\ & \leq \sum_{k=1}^{k^*} \mathbb{P}\left(\frac{1}{r_k} \sup_{\mathbf{v} \in \mathcal{B}_{r_k}(\mathbf{v}^*) \setminus \mathcal{B}_{r_{k-1}}(\mathbf{v}^*)} \{\mathcal{U}(\mathbf{v}, \mathbf{v}^*) - f(r_k, r_0)\} \geq 0\right) \leq e^{-x} \sum_{k=1}^{k^*} \rho^k \leq \frac{\rho}{1-\rho} e^{-x} \end{aligned}$$

as required. □

B.3. Finite-dimensional smooth case

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

$(\mathcal{E}D)$ There exist $\mathfrak{g} > 0$, $v_0 \geq 1$, and for each $\mathbf{v} \in \Upsilon$, a symmetric non-negative matrix $\mathbb{V}(\mathbf{v})$ such that for any $\lambda \leq \mathfrak{g}$ and any unit vector $\boldsymbol{\gamma} \in \mathbb{R}^p$, it holds

$$\log \mathbb{E} \exp \left\{ \lambda \frac{\boldsymbol{\gamma}^\top \nabla \mathcal{U}(\mathbf{v})}{\|\mathbb{V}(\mathbf{v})\boldsymbol{\gamma}\|} \right\} \leq \frac{v_0^2 \lambda^2}{2}.$$

A natural candidate for $\mathbb{V}^2(\mathbf{v})$ is the covariance matrix $\text{Var}(\nabla \mathcal{U}(\mathbf{v}))$ provided that this matrix is well posed. Then the constant v_0 can be taken close to one by reducing the value \mathfrak{g} .

In what follows we fix a subset Υ° of Υ and establish a bound for the maximum of the process $\mathcal{U}(\mathbf{v}, \mathbf{v}^\circ) = \mathcal{U}(\mathbf{v}) - \mathcal{U}(\mathbf{v}^\circ)$ on Υ° for a fixed point \mathbf{v}° . We assume existence of a matrix $\mathbb{V} = \mathbb{V}(\Upsilon^\circ)$ such that $\mathbb{V}(\mathbf{v}) \preceq \mathbb{V}$ for all $\mathbf{v} \in \Upsilon^\circ$. We also assume that π is the Lebesgue measure on Υ . First we show that $(\mathcal{E}D)$ implies $(\mathcal{E}d)$.

Lemma B.4. Assume that $(\mathcal{E}D)$ holds with some \mathfrak{g} and $\mathbb{V}(\mathbf{v}) \preceq \mathbb{V}$ for $\mathbf{v} \in \Upsilon^\circ$. Consider any $\mathbf{v}, \mathbf{v}^\circ \in \Upsilon^\circ$. Then it holds for $|\lambda| \leq \mathfrak{g}$

$$\log \mathbb{E} \exp \left\{ \lambda \frac{\mathcal{U}(\mathbf{v}, \mathbf{v}^\circ)}{\|\mathbb{V}(\mathbf{v} - \mathbf{v}^\circ)\|} \right\} \leq \frac{v_0^2 \lambda^2}{2}.$$

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

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

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

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

as required. \square

The result of Lemma B.4 enables us to define $d(\mathbf{v}, \mathbf{v}') = \|\nabla(\mathbf{v} - \mathbf{v}')\|$ so that the corresponding d -ball coincides with the following ellipsoidal set $\mathcal{B}(\mathfrak{r}, \mathbf{v}^\circ)$:

$$\mathcal{B}(\mathfrak{r}, \mathbf{v}^\circ) \stackrel{\text{def}}{=} \{\mathbf{v}: \|\nabla(\mathbf{v} - \mathbf{v}^\circ)\| \leq \mathfrak{r}\}.$$

Now we bound the value $\mathbb{Q}(\mathcal{Y}^\circ)$ for $\mathcal{Y}^\circ = \mathcal{B}(\mathfrak{r}, \mathbf{v}^\circ)$. Note that by change of variable one can reduce the study to the case $\nabla = \mathbb{I}_p$ and consider the entropy of the unit ball in \mathbb{R}^p w.r.t. the Euclidean distance. We use the following general result which allows to upperbound the covering number of a convex set in \mathbb{R}^p for the Euclidean metric.

Lemma B.5. *Let \mathcal{Y}° be a convex set in \mathbb{R}^p , $\delta > 0$, and B be the unit ball in \mathbb{R}^p . Then the covering number $\mathbb{N}(\mathcal{Y}^\circ, \delta)$ fulfills*

$$\mathbb{N}(\mathcal{Y}^\circ, \delta) \leq \frac{\text{vol}(\mathcal{Y}^\circ + (\delta/2)B)}{\text{vol}(B)} (2/\delta)^p.$$

Proof. Let $(\mathbf{v}^{(i)}, i = 1, \dots, \mathbb{N})$ be a maximal subset of \mathcal{Y}° such that $\|\mathbf{v}^{(i)} - \mathbf{v}^{(j)}\| \geq \delta$ for all $i \neq j$. By maximality, $(\mathbf{v}^{(i)})$ is a δ -net of \mathcal{Y}° . Let also B be the unit ball in \mathbb{R}^p . Note that the balls $\mathbf{v}^{(i)} + (\delta/2)B$ are disjoint and included in $\mathcal{Y}^\circ + (\delta/2)B$. Therefore,

$$\sum_{i \leq \mathbb{N}} \text{vol} \left(\mathbf{v}^{(i)} + \frac{\delta}{2} B \right) \leq \text{vol} \left(\mathcal{Y}^\circ + \frac{\delta}{2} B \right),$$

where $\text{vol}(A)$ means the Lebesgue measure of the set A . This yields

$$\mathbb{N}(\delta/2)^p \text{vol}(B) \leq \text{vol}(\mathcal{Y}^\circ + (\delta/2)B)$$

and the claim of the lemma follows. \square

Lemma B.6 (Entropy of a ball). *Let $\mathcal{Y}^\circ = \mathcal{B}(\mathfrak{r}_\circ, \mathbf{v}^*)$ and $\mathfrak{r}_k = 2^{-k} \mathfrak{r}_\circ$. Then the covering numbers \mathbb{N}_k fulfill with $\delta = \mathfrak{r}_k / \mathfrak{r}_\circ = 2^{-k}$*

$$\mathbb{N}_k \leq (1 + 2/\delta)^p = (1 + 2^{k+1})^p.$$

Moreover, with $\mathfrak{c}_2 = 4.67$,

$$\mathbb{Q}_2(\mathcal{Y}^\circ) \leq 2 \log 2 + \mathfrak{c}_2 p \leq 6p,$$

$$\mathbb{Q}_1(\mathcal{Y}^\circ) \leq \sqrt{2 \log 2 + \mathfrak{c}_2 p} \leq \sqrt{6p}.$$

(B.15)

Proof. A change of variable reduces the statement to the case $\mathbb{V} = \mathbb{I}_p$ and $r_o = 1$. For $\delta = 2^{-k}$, this implies by Lemma B.5 in view of $\mathcal{Y}^\circ = B$

$$\text{vol}\left(\mathcal{Y}^\circ + \frac{\delta}{2}B\right) = (1 + \delta/2)^p \text{vol}(B),$$

that $\mathbb{N}_k \leq (1 + 2/\delta)^p$ as claimed. Now we derive

$$\begin{aligned} \mathbb{Q}_2(\mathcal{Y}^\circ) &\leq \sum_{k=1}^{\infty} 2^{-k+1} \log(2\mathbb{N}_k) \leq \sum_{k=1}^{\infty} 2^{-k+1} \{\log 2 + 2p \log(1 + 2^{k+1})\} \\ &\leq 2\log 2 + p \sum_{k=0}^{\infty} 2^{-k+1} \log(1 + 2^k) \leq 2\log 2 + c_2 p \end{aligned}$$

as required. □

Now we specify the local bounds of Theorem B.1 to the smooth case. We consider the local sets of the elliptic form $\mathcal{Y}_o(x) \stackrel{\text{def}}{=} \{\mathbf{v} : \|\mathbb{V}(\mathbf{v} - \mathbf{v}^*)\| \leq x\}$, where \mathbb{V} dominates $\mathbb{V}(\mathbf{v})$ on this set: $\mathbb{V}(\mathbf{v}) \leq \mathbb{V}$.

Theorem B.7. Let (E D) hold with some $\varrho > 0$, and matrices $\mathbb{V}(\mathbf{v})$ such that $\mathbb{V}(\mathbf{v}) \leq \mathbb{V}$ for all $\mathbf{v} \in \mathcal{Y}_o(x)$ and a fixed x . For any $x \geq 1/2$

$$\mathbb{P}\left\{\frac{1}{v_0 x} \sup_{\mathbf{v} \in \mathcal{Y}_o(x)} |\mathcal{U}(\mathbf{v}, \mathbf{v}^*)| \geq \mathfrak{z}_{\mathbb{H}}(x)\right\} \leq e^{-x},$$

where $\mathfrak{z}_{\mathbb{H}}(x)$ is given by (B.3) with $\mathbb{Q}_1(\mathcal{Y}^\circ)$ and $\mathbb{Q}_2(\mathcal{Y}^\circ)$ from (B.15).

Proof. Lemma B.6 implies (E d) with $d(\mathbf{v}, \mathbf{v}^*) = \|\mathbb{V}(\mathbf{v} - \mathbf{v}^*)\|$. Now the result follows from Theorem B.1. □

B.4. Entropy of an ellipsoid

Let H be a positive self adjoint operator in \mathbb{R}^∞ . We are interested to describe the entropy of the elliptic set

$$\mathcal{E}_H(r_o) \stackrel{\text{def}}{=} \{\mathbf{v} : \|H(\mathbf{v} - \mathbf{v}^\circ)\| \leq r_o\} \tag{B.16}$$

for given $\mathbf{v}^\circ \in \mathbb{R}^\infty$ and $r_o > 0$ with respect to the usual Euclidean distance in \mathbb{R}^∞ . Below we evaluate the entropy of this set assuming that $\|H^{-1}\|_{\text{op}} = 1$ and H^{-2} is a trace operator, i.e., $h_1 = 1$ and

$$\mathfrak{p}_H \stackrel{\text{def}}{=} \text{tr}(H^{-2}) = \sum_{j=1}^{\infty} h_j^{-2} < \infty, \tag{B.17}$$

where $h_1 \leq h_2 \leq \dots$ are the ordered eigenvalues of H .

Theorem B.8. Suppose that for some $\alpha > 1$

$$\mathfrak{p}_H(\alpha) \stackrel{\text{def}}{=} \sum_{j=1}^{\infty} h_j^{-2} \log^\alpha(h_j^2) < \infty. \tag{B.18}$$

Then for $\mathcal{E} = \mathcal{E}_H(r_o)$

$$\mathbb{Q}_1(\mathcal{E}) \leq C(\alpha - 1)^{-1/2} \sqrt{\mathfrak{p}_H(\alpha)}, \tag{B.19}$$

where C is an absolute constant. Furthermore,

$$\mathbb{Q}_2(\mathcal{E}) \leq C \mathfrak{P}_H^* = C \sum_{j=1}^{\infty} h_j^{-1}.$$

Remark B.3. The log-factor in the definition of $\mathfrak{P}_H(\alpha)$ can be removed by using a more advanced generic chaining and majorising measure technique. However, in most of situations, the bound in terms of $\mathfrak{P}_H(\alpha)$ is also sharp.

The term \mathfrak{P}_H^* only appears in the sub-exponential case when $\mathfrak{g} < \infty$. In this case we need the condition $\mathfrak{P}_H^* < \infty$ which requires $\sum_j h_j^{-1} < \infty$, that is, a more rapid growth of the values h_j is necessary than in (B.18).

Proof of Theorem B.8. We begin by a general lemma which bounds the covering numbers for the elliptic set \mathcal{E} for the Euclidean distance.

Lemma B.9 (Entropy of the ellipsoid). Let $\mathcal{E} = \mathcal{E}_H(r_\circ)$ be an elliptic set from (B.16) with $\|H^{-1}\|_{\text{op}} = 1$ and $\text{tr}(H^{-2}) < \infty$. Let also $d(\mathbf{v}, \mathbf{v}') = \|\mathbf{v} - \mathbf{v}'\|$. Then for $r_k = 2^{-k}r_\circ$, the value $\mathbb{Q}_1(\mathcal{E})$ from (B.1) satisfies

$$\mathbb{Q}_1(\mathcal{E}) \leq \sum_{k=1}^{\infty} 2^{-k} \sqrt{\log 2 + 2L_H(m_k)}, \tag{B.20}$$

where m_k is the index j for which $h_{m_k}^2 = 2^{2k+1}$ and hence,

$$h_j^2 \leq 2^{2k+1}, \quad j \leq m_k, \tag{B.21}$$

and

$$L_H(m) \stackrel{\text{def}}{=} \sum_{j=1}^m \log(3h_m/h_j).$$

Remark B.4. For the ease of presentation, we supposed in the lemma that for each $k \geq 1$, there exists some m_k with $h_{m_k} = 2^{k+1/2}$. The results easily extend to the case when this equality is approximate.

Proof of Lemma B.9. Without loss of generality assume $\mathbf{v}^\circ = 0$. A basis transform reduces the study to the case when H is diagonal: $H = \text{diag}\{h_1, h_2, \dots\}$. We only have to evaluate the covering numbers \mathbb{N}_k . Let us fix $k \geq 1$ and let m_k be given by (B.21). For any point $\mathbf{v} = (\mathbf{v}_1, \mathbf{v}_2, \dots)^\top$ in \mathcal{E} , it holds

$$\begin{aligned} \sum_{j=m_k+1}^{\infty} \mathbf{v}_j^2 &= \sum_{j=m_k+1}^{\infty} h_j^{-2} h_j^2 \mathbf{v}_j^2 \leq h_{m_k+1}^{-2} \sum_{j=m_k+1}^{\infty} h_j^2 \mathbf{v}_j^2 \\ &\leq h_{m_k+1}^{-2} \sum_{j=1}^{\infty} h_j^2 \mathbf{v}_j^2 \leq 2^{-2k-1} r_\circ^2 \leq r_k^2/2. \end{aligned} \tag{B.22}$$

Consider the elliptic set \mathcal{E}_k in \mathbb{R}^{m_k} obtained by projection Π_k of \mathcal{E} on the first m_k coordinates:

$$\mathcal{E}_k \stackrel{\text{def}}{=} \left\{ (\mathbf{v}_1, \dots, \mathbf{v}_{m_k})^\top : \sum_{j=1}^{m_k} h_j^2 \mathbf{v}_j^2 \leq r_\circ^2 \right\}.$$

Let \mathcal{M}_k be a ϵ_k -net in \mathcal{E}_k for $\epsilon_k^2 = r_k^2/2$. A r_k -net in \mathcal{E} can be constructed from \mathcal{M}_k in a simple way: just fix to zero the remaining coordinates $\mathbf{v}_j = 0$ for $j > m_k$. If \mathbf{v}° is constructed in this way, then $\|H\mathbf{v}^\circ\| = \|H\Pi_k\mathbf{v}^\circ\| \leq 1$, that is, $\mathbf{v}^\circ \in \mathcal{E}$. Moreover, for any other point $\mathbf{v} \in \mathcal{E}$, take \mathbf{v}° such that their projections satisfy $\|\Pi_k(\mathbf{v} - \mathbf{v}^\circ)\| \leq \epsilon_k$. Then

by (B.22)

$$\|\mathbf{v} - \mathbf{v}^\circ\|^2 = \|\Pi_k(\mathbf{v} - \mathbf{v}^\circ)\|^2 + \|(\mathbb{I} - \Pi_k)\mathbf{v}\|^2 \leq r_k^2/2 + r_k^2/2 = r_k^2.$$

Therefore, the covering number $\mathbb{N}(\mathcal{E}, r_k)$ of the infinite dimensional elliptic set \mathcal{E} does not exceed the covering number $\mathbb{N}(\mathcal{E}_k, \epsilon_k)$ for the m_k -dimensional ellipsoid \mathcal{E}_k . By Lemma B.5 with $\delta = \epsilon_k$,

$$\mathbb{N}(\mathcal{E}_k, \epsilon_k) \leq \frac{\text{vol}(\mathcal{E}_k + (\epsilon_k/2)B_k)}{\text{vol}(B_k)} (2/\epsilon_k)^{m_k},$$

where B_k is the unit ball in \mathbb{R}^{m_k} . The bound $h_j^{-2} \geq 2^{-2k-1}$ for $j \leq m_k$ implies that $\mathcal{E}_k + (\epsilon_k/2)B_k$ is contained in the elliptic set $(3/2)\mathcal{E}_k$.

The definition implies due to $h_{m_k}^2 = 2^{2k+1}$

$$\mathbb{N}(\mathcal{E}, r_k) \leq \log \frac{\text{vol}((3/2)\mathcal{E}_k)}{(\epsilon_k/2)^{m_k} \text{vol}(B_k)} \leq \sum_{j=1}^{m_k} \log \frac{3h_j^{-1}}{\epsilon_k} \leq \sum_{j=1}^{m_k} \log \left(\frac{3h_{m_k}}{h_j} \right) = L_H(m_k). \tag{B.23}$$

Now the result (B.20) follows by the definition of $\mathbb{Q}_1(\mathcal{E})$. □

Denote $\mathbb{N}_k = \mathbb{N}(\mathcal{E}, r_k)$. By the Cauchy–Schwarz inequality for $\alpha > 1$

$$\mathbb{Q}_1(\mathcal{E}) = \sum_{k=1}^{\infty} 2^{-k} \sqrt{2 \log(2\mathbb{N}_k)} \leq \left\{ \sum_{k=1}^{\infty} k^{-\alpha} \sum_{k=1}^{\infty} k^\alpha 2^{-2k} 2 \log(2\mathbb{N}_k) \right\}^{1/2}. \tag{B.24}$$

The use of $h_{m_\ell}^{-2} = 2^{2\ell+1}$ and $h_j^2 \geq 4h_{m_{\ell-1}}^2$ for $j \in (m_{\ell-1}, m_\ell]$ yields by (B.23) with $n_\ell \stackrel{\text{def}}{=} m_\ell - m_{\ell-1}$

$$2 \log(2\mathbb{N}_k) = \sum_{\ell=1}^k \sum_{j=m_{\ell-1}+1}^{m_\ell} \log \frac{9h_{m_k}^2}{h_j^2} \leq \sum_{\ell=1}^k \{k - \ell + \log(36)\} n_\ell.$$

Further, in view of $h_{m_k} = 2^k$

$$\begin{aligned} \sum_{k=1}^{\infty} k^\alpha 2^{-2k} 2\mathbb{N}_k &\leq \sum_{k=1}^{\infty} k^\alpha 2^{-2k} \sum_{\ell=1}^k \{k - \ell + \log(36)\} n_\ell = \sum_{\ell=1}^{\infty} \sum_{k \geq \ell} k^\alpha 2^{-2k} \{k - \ell + \log(36)\} n_\ell \\ &= \sum_{\ell=1}^{\infty} n_\ell 2^{-2\ell} \sum_{k \geq \ell} k^\alpha 2^{-2(k-\ell)} \{k - \ell + \log(36)\} = C \sum_{\ell=1}^{\infty} n_\ell 2^{-2\ell} \ell^\alpha. \end{aligned}$$

It remains to note that $2^{2\ell-1} \leq h_j^2 \leq 2^{2\ell+1}$ for $m_{\ell-1} < j \leq m_\ell$ and

$$\sum_{\ell=1}^{\infty} n_\ell 2^{-2\ell} \ell^\alpha \leq \sum_{\ell=1}^{\infty} \sum_{j=m_{\ell-1}+1}^{m_\ell} h_j^{-2} \log^\alpha(h_j^2) = \sum_{j=1}^{\infty} h_j^{-2} \log^\alpha(h_j^2) = \mathfrak{p}_H(\alpha). \tag{B.25}$$

The assertion (B.19) now follows from (B.24) in view of $\sum_{k \geq 1} k^{-\alpha} \leq C(\alpha - 1)^{-1}$.

The result on $\mathbb{Q}_2(\mathcal{E})$ requires to bound the sum of $2^{-k} \log \mathbb{N}_k$. Similarly to the above, one easily derives

$$\begin{aligned} \sum_{k=1}^{\infty} 2^{-k} \mathbb{N}_k &\leq \sum_{k=1}^{\infty} 2^{-k} \sum_{\ell=1}^k \{k - \ell + \log(36)\} n_\ell = \sum_{\ell=1}^{\infty} \sum_{k \geq \ell} 2^{-k} \{k - \ell + \log(36)\} n_\ell \\ &= \sum_{\ell=1}^{\infty} n_\ell 2^{-\ell} \sum_{k \geq \ell} 2^{-(k-\ell)} \{k - \ell + \log(36)\} = C \sum_{\ell=1}^{\infty} n_\ell 2^{-\ell} \leq C \sum_{j=1}^{\infty} h_j^{-1} = C \mathfrak{p}_H^*. \end{aligned}$$

Theorem is proved. □

Now we present a special case for which the entropy can be bounded via the effective dimension \mathfrak{p}_H of Υ° defined in (B.17).

Theorem B.10. *Let $h_j^2 = f(j)$ for a monotonously increasing smooth function $f(x) > 0$. If $xf'(x)/f(x) \leq \beta$, then*

$$\begin{aligned} \mathbb{Q}_2(\mathcal{E}) &\leq C\beta\mathfrak{p}_H, \\ \mathbb{Q}_1(\mathcal{E}) &\leq C\sqrt{\beta\mathfrak{p}_H}, \end{aligned} \tag{B.26}$$

where the effective dimension \mathfrak{p}_H is defined in (B.17).

Proof. Obviously

$$\sum_{j=1}^m \log\left(\frac{h_m^2}{h_j^2}\right) \leq \int_0^m \log\left(\frac{f(m)}{f(t)}\right) dt.$$

Now we note that the function

$$F(x) \stackrel{\text{def}}{=} \int_0^x \log\left(\frac{f(x)}{f(t)}\right) dt$$

fulfills $F(0) = 0$ and $F'(x) = xf'(x)/f(x)$ yielding

$$\sum_{j=1}^m \log\left(\frac{h_m^2}{h_j^2}\right) \leq \int_0^m \log\left(\frac{f(m)}{f(t)}\right) dt = \int_0^m \frac{xf'(x)}{f(x)} dx.$$

Moreover, in particular, if $F'(x) \leq \beta$, then $F(x) \leq \beta x$ and thus, $L_H(m_k) \leq \beta m_k$. Now it holds similarly to (B.25)

$$\sum_{k=1}^{\infty} 2^{-k} m_k = \sum_{k=1}^{\infty} 2^{-k} \sum_{\ell=1}^k n_\ell \leq \sum_{\ell=1}^{\infty} n_\ell \sum_{k \geq \ell} 2^{-k} = \sum_{\ell=1}^{\infty} n_\ell 2^{-\ell} \leq 2 \sum_{j=1}^{\infty} h_j^{-2} = 2\mathfrak{p}_H,$$

and the statement (B.19) follows. □

Now we evaluate the entropy for the cases when h_j grow polynomially.

Theorem B.11. *Let $h_j^2 = 1 + \varkappa^2 j^{2\beta}$ for $\beta > 1/2$ and some small value \varkappa . Then*

$$\begin{aligned} \mathbb{Q}_1(\mathcal{E}) &\leq C(2\beta - 1)^{-1/2} \varkappa^{-1/(2\beta)}, \\ \mathbb{Q}_2(\mathcal{E}) &\leq C(2\beta - 1)^{-1} \varkappa^{-1/\beta}, \end{aligned}$$

where C is an absolute constant.

Proof. For $f(x) = 1 + \varkappa^2 x^{2\beta}$, it holds $xf'(x)/f(x) \leq 2\beta$ and we can apply the result of Theorem B.10. With $\beta > 1/2$, the effective dimension \mathfrak{p}_H from (B.17) fulfills

$$\mathfrak{p}_H \leq \sum_{j=1}^{\infty} h_j^{-2} = \sum_{j=1}^{\infty} \frac{1}{1 + \varkappa^2 j^{2\beta}} \leq \int_0^{\infty} \frac{1}{1 + \varkappa^2 x^{2\beta}} dx = C\varkappa^{-1/\beta} \frac{1}{2\beta - 1}$$

and the result follows by (B.26). □

B.5. Roughness constraints for dimension reduction

The local bounds of Theorems B.1 and B.3 can be extended in several directions. Here we briefly discuss one extension related to the use of a smoothness condition on the parameter \mathbf{v} . Let $\text{pen}(\mathbf{v})$ be a non-negative *penalty* function on \mathcal{Y} . A particular example of such penalty function is the *roughness penalty* $\text{pen}(\mathbf{v}) = \|G\mathbf{v}\|^2$ for a given p -matrix G^2 . Let r be fixed. Consider the intersection of the ball $\mathcal{B}_r(\mathbf{v}^\circ)$ with the set \mathcal{Y} given by the constraint $\text{pen}(\mathbf{v}) \leq 1$:

$$\mathcal{Y}_{\text{pen}}(r) = \{\mathbf{v} \in \mathcal{Y} : d(\mathbf{v}, \mathbf{v}^\circ) \leq r; \text{pen}(\mathbf{v}) \leq 1\},$$

for a fixed central point \mathbf{v}° and the radius r . Here and below we assume that the central point \mathbf{v}° is “smooth” in the sense that $\text{pen}(\mathbf{v}^\circ) < 1$. One can easily check that the results of Theorems B.1 and B.3 and their corollaries extend to this situation without any change. The only difference is in the definition of the values $\mathbb{Q}_1(\mathcal{Y}_\circ)$ and $\mathbb{Q}_2(\mathcal{Y}_\circ)$ for $\mathcal{Y}_\circ = \mathcal{Y}_{\text{pen}}(r)$. Examples below show that the use of the penalization can substantially reduce these values relative to the non-penalized case.

We consider the case of a smooth process \mathcal{U} given on a local set $\mathcal{Y}_G(r)$ of the form

$$\mathcal{Y}_G(r) = \{\mathbf{v} \in \mathcal{Y} : \|\nabla(\mathbf{v} - \mathbf{v}^\circ)\| \leq r; \|G\mathbf{v}\| \leq 1\}, \tag{B.27}$$

with the distance $d(\mathbf{v}, \mathbf{v}^\circ) = \|\nabla(\mathbf{v} - \mathbf{v}^\circ)\|$ and a smoothness constraint $\|G\mathbf{v}\|^2 \leq 1$. Then the set $\mathcal{Y}_G(r)$ is contained in an elliptic set

$$\mathcal{Y}_\circ \stackrel{\text{def}}{=} \{\boldsymbol{\theta} : \|G\mathbf{v}\|^2 + \|\nabla(\mathbf{v} - \mathbf{v}^\circ)\|^2 \leq 1 + r^2\}. \tag{B.28}$$

Define

$$\mathbb{V}_G^2 = \mathbb{V}^2 + G^2, \quad \mathbf{v}_G = \mathbb{V}_G^{-2} \mathbb{V}^2 \mathbf{v}^\circ.$$

Then $\mathbf{v}^\circ - \mathbf{v}_G = (\mathbb{I}_p - \mathbb{V}_G^{-2} \mathbb{V}^2) \mathbf{v}^\circ = \mathbb{V}_G^{-2} G^2 \mathbf{v}^\circ$ and one can get by simple algebra

$$\begin{aligned} \|G\mathbf{v}\|^2 + \|\nabla(\mathbf{v} - \mathbf{v}^\circ)\|^2 &= \|\mathbb{V}_G(\mathbf{v} - \mathbf{v}_G)\|^2 + \|G\mathbf{v}_G\|^2 + \|\nabla(\mathbf{v}_G - \mathbf{v}^\circ)\|^2 \\ &= \|\mathbb{V}_G(\mathbf{v} - \mathbf{v}_G)\|^2 + \mathbf{v}^{\circ\top} G^2 \mathbb{V}_G^{-2} \mathbb{V}^2 \mathbf{v}^\circ = \|\mathbb{V}_G(\mathbf{v} - \mathbf{v}^\circ)\|^2 + d_G \end{aligned}$$

with $d_G = \mathbf{v}^{\circ\top} G^2 \mathbb{V}_G^{-2} \mathbb{V}^2 \mathbf{v}^\circ \leq \|G\mathbf{v}^\circ\|^2 < 1$. A change of variables $\mathbf{v} \rightarrow \mathbb{V}(\mathbf{v} - \mathbf{v}_G)$ allows us to reduce the study to the case of an ellipsoid considered in Section B.4. For H defined by $H^{-2} = \mathbb{V} \mathbb{V}_G^{-2} \mathbb{V}$, the set \mathcal{Y}_\circ from (B.28) is transferred into the elliptic set

$$\mathcal{Y}_H(r) = \{\mathbf{v} : \|H\mathbf{v}\|^2 \leq 1 + r^2 - d_G\},$$

whose entropy for the Euclidean distance is given via the trace $\mathfrak{p}_H = \text{tr}(H^{-2})$.

Now we are prepared to state the penalized bound for the process $\mathcal{U}(\cdot)$ over \mathcal{Y}_\circ which naturally generalizes the result of Theorem B.7 to the non-penalized case.

Theorem B.12. *Let $\mathcal{Y}_\circ = \mathcal{Y}_{\text{pen}}(r)$ be given by (B.27) and $\|G\mathbf{v}^\circ\| \leq 1$. Let also (E D) hold with some \mathfrak{g} and a matrix $\mathbb{V}(\mathbf{v}) \leq \mathbb{V}$ for all $\mathbf{v} \in \mathcal{Y}_\circ$. For H defined by $H^{-2} = \mathbb{V} \mathbb{V}_G^{-2} \mathbb{V}$, let the entropy values $\mathbb{Q}_1(\mathcal{Y}^\circ)$ and $\mathbb{Q}_2(\mathcal{Y}^\circ)$ for the elliptic set $\mathcal{Y}_H(r)$ from (B.28) be given in Section B.4. Then for any $x \geq 1/2$*

$$\mathbb{P} \left\{ \frac{1}{\nu_{0x}} \sup_{\mathbf{v} \in \mathcal{Y}_{\text{pen}}(r)} |\mathcal{U}(\mathbf{v}, \mathbf{v}^\circ)| \geq \mathfrak{z}_{\mathbb{H}}(x) \right\} \leq e^{-x},$$

where $\mathfrak{z}_{\mathbb{H}}(x)$ is from (B.3) with these values $\mathbb{Q}_1(\mathcal{Y}^\circ)$ and $\mathbb{Q}_2(\mathcal{Y}^\circ)$.

B.6. Bound for a bivariate process

Consider a smooth bivariate process $\mathcal{U}(\mathbf{v}) = \mathcal{U}(\mathbf{v}_1, \mathbf{v}_2)$ over a product set $\mathcal{Y} = \mathcal{Y}_1 \times \mathcal{Y}_2$, where $\mathcal{Y}_j \subseteq \mathbb{R}^{p_j}$ for $j = 1, 2$. We suppose that partial derivatives of \mathcal{U} have uniform exponential moments.

($\mathcal{E}D_p$) There exist $\mathfrak{g} > 0$, $\nu_0 \geq 1$, and for each $\mathbf{v} = (\mathbf{v}_1, \mathbf{v}_2) \in \mathcal{Y} = \mathcal{Y}_1 \times \mathcal{Y}_2$, symmetric non-negative $p_j \times p_j$ matrices \mathbb{V}_j , $j = 1, 2$, such that for any $\lambda \leq \mathfrak{g}$ and any unit vector $\boldsymbol{\gamma} \in \mathbb{R}^p$, it holds

$$\log \mathbb{E} \exp \left\{ \lambda \frac{\boldsymbol{\gamma}^\top \nabla_j \mathcal{U}(\mathbf{v})}{\|\mathbb{V}_j \boldsymbol{\gamma}\|} \right\} \leq \frac{\nu_0^2 \lambda^2}{2}, \quad j = 1, 2.$$

Here $\nabla_j \mathcal{U}$ denotes the partial derivative $\partial \mathcal{U} / \partial \mathbf{v}_j$ for $j = 1, 2$.

This allows to establish an exponential bound for the process $\mathcal{U}(\mathbf{v})$. Let us fix the central point $\mathbf{v}^\circ = (\mathbf{v}_1^\circ, \mathbf{v}_2^\circ)$ and a radius r . As usual,

$$\mathcal{Y}_j(r) = \{ \mathbf{v}_j \in \mathcal{Y}_j : \|\mathbb{V}_j(\mathbf{v}_j - \mathbf{v}_j^\circ)\| \leq r \}$$

denotes the ball in \mathcal{Y}_j with this radius.

Theorem B.13. *Let a bivariate random process $\mathcal{U}(\mathbf{v})$ on $\mathcal{Y} = \mathcal{Y}_1 \times \mathcal{Y}_2$ satisfy ($\mathcal{E}D_p$). Then for any r_\circ and $x \geq 1/2$, it holds on the product set $\mathcal{Y}_\circ = \mathcal{Y}_1(r_\circ) \times \mathcal{Y}_2(r_\circ)$*

$$\mathbb{P} \left\{ \frac{1}{\sqrt{8\nu_0 r_\circ}} \sup_{\mathbf{v} \in \mathcal{Y}_\circ} |\mathcal{U}(\mathbf{v}, \mathbf{v}^\circ)| \geq \mathfrak{z}_{\mathbb{H}}(x) \right\} \leq e^{-x},$$

with $\mathfrak{z}_{\mathbb{H}}(x)$ from (B.3) for $\mathbb{Q}_1(\mathcal{Y}^\circ) = \mathbb{Q}_1(\mathcal{Y}_1) + \mathbb{Q}_1(\mathcal{Y}_2)$ and $\mathbb{Q}_2(\mathcal{Y}^\circ) = \mathbb{Q}_2(\mathcal{Y}_1) + \mathbb{Q}_2(\mathcal{Y}_2)$.

Proof. By the Hölder inequality, (B.31), and (B.30), it holds for $\|\boldsymbol{\gamma}_1\| = \|\boldsymbol{\gamma}_2\| = 1$ and $\mathbf{v} \in \mathcal{Y}$.

$$\begin{aligned} & \log \mathbb{E} \exp \left\{ \frac{\lambda}{2} (\boldsymbol{\gamma}_1, \boldsymbol{\gamma}_2)^\top \nabla \mathcal{U}(\mathbf{v}) \right\} \\ & \leq \frac{1}{2} \log \mathbb{E} \exp \{ \lambda \boldsymbol{\gamma}_1^\top \nabla_1 \mathcal{U}(\mathbf{v}) \} + \frac{1}{2} \log \mathbb{E} \exp \{ \lambda \boldsymbol{\gamma}_2^\top \nabla_2 \mathcal{U}(\mathbf{v}) \} \\ & \leq \frac{1}{2} \log \mathbb{E} \exp \{ \lambda \boldsymbol{\gamma}_1^\top \nabla_1 \mathcal{U}(\mathbf{v}) \} + \frac{1}{2} \log \mathbb{E} \exp \{ \lambda \boldsymbol{\gamma}_2^\top \nabla_2 \mathcal{U}(\mathbf{v}) \} \leq \frac{\nu_0^2 \lambda^2}{2}, \quad |\lambda| \leq \mathfrak{g}. \end{aligned}$$

This means that the bivariate process $\mathcal{U}(\mathbf{v})/2$ fulfills the full dimensional condition ($\mathcal{E}D$) with $\mathbb{V} = \text{block}(\mathbb{V}_1, \mathbb{V}_2)$. Let $\mathbf{v} = (\mathbf{v}_1, \mathbf{v}_2)$ and $\mathbf{v}^\circ = (\mathbf{v}_1^\circ, \mathbf{v}_2^\circ)$ be a couple of points in \mathcal{Y} such that $\|\mathbb{V}_j(\mathbf{v}_j - \mathbf{v}_j^\circ)\| \leq \varepsilon$ for $j = 1, 2$. Then obviously

$$\|\mathbb{V}(\mathbf{v} - \mathbf{v}^\circ)\|^2 \leq 2\varepsilon^2. \tag{B.29}$$

Therefore, the direct product of two ε -nets $\mathcal{M}_j(\varepsilon)$ in \mathcal{Y}_j for $j = 1, 2$ yield a $\sqrt{2}\varepsilon$ -net $\mathcal{M}(\varepsilon) = \mathcal{M}_1(\varepsilon) \times \mathcal{M}_2(\varepsilon)$ in the product space \mathcal{Y} .

Due to (B.29), the product set $\mathcal{Y}_\circ \stackrel{\text{def}}{=} \mathcal{Y}_1(r_\circ) \times \mathcal{Y}_2(r_\circ)$ has the radius r_\circ . Now we can easily bound the entropy of the product set \mathcal{Y}_\circ via the entropy of \mathcal{Y}_1 and \mathcal{Y}_2 . Indeed, with $r_k = 2^{-k} r_\circ$, the cardinality \mathbb{N}_k of $\mathcal{M}_k = \mathcal{M}(r_k)$ fulfills $\mathbb{N}_k = \mathbb{N}_k(\mathcal{Y}_1) \mathbb{N}_k(\mathcal{Y}_2)$ and

$$\begin{aligned} \mathbb{Q}_2(\mathcal{Y}_\circ) & \leq \sum_{k=1}^{\infty} 2^{-k+1} \log(2\mathbb{N}_k) \\ & \leq \sum_{k=1}^{\infty} 2^{-k+1} \log(2\mathbb{N}_k(\mathcal{Y}_1)) + \sum_{k=1}^{\infty} 2^{-k+1} \log(2\mathbb{N}_k(\mathcal{Y}_2)) \leq \mathbb{Q}_2(\mathcal{Y}_1) + \mathbb{Q}_2(\mathcal{Y}_2). \end{aligned}$$

Similarly

$$\begin{aligned} \mathbb{Q}_1(\mathcal{Y}_o) &\leq \sum_{k=1}^{\infty} 2^{-k} \sqrt{2 \log(2\mathbb{N}_k)} \\ &\leq \sum_{k=1}^{\infty} 2^{-k} \sqrt{2 \log(2\mathbb{N}_k(\mathcal{Y}_1)) + 2 \log(2\mathbb{N}_k(\mathcal{Y}_2))} \leq \mathbb{Q}_1(\mathcal{Y}_1) + \mathbb{Q}_1(\mathcal{Y}_2). \end{aligned}$$

Now we just apply the assertion of Theorem B.7 to the process $\mathcal{U}(\mathbf{v})/2$ and account for the fact that by (B.29) the radius of \mathcal{Y}_o is $\sqrt{2}r_o$. \square

B.7. A bound for the norm of a vector random process

Let $\mathcal{Y}(\mathbf{v})$, $\mathbf{v} \in \mathcal{Y}$, be a smooth centered random vector process with values in \mathbb{R}^q , where $\mathcal{Y} \subseteq \mathbb{R}^p$. Let also $\mathcal{Y}(\mathbf{v}^*) = 0$ for a fixed point $\mathbf{v}^* \in \mathcal{Y}$. Without loss of generality assume $\mathbf{v}^* = 0$. We aim to bound the maximum of the norm $\|\mathcal{Y}(\mathbf{v})\|$ over a vicinity \mathcal{Y}_o of \mathbf{v}^* . By $\nabla \mathcal{U}(\mathbf{v})$ we denote the $p \times q$ matrix with entries $\nabla_{v_i} \mathcal{U}_j$, $i \leq p$, $j \leq q$. Suppose that $\mathcal{Y}(\mathbf{v})$ satisfies for each $\boldsymbol{\gamma}_1 \in \mathbb{R}^p$ and $\boldsymbol{\gamma}_2 \in \mathbb{R}^q$ with $\|\boldsymbol{\gamma}_1\| = \|\boldsymbol{\gamma}_2\| = 1$

$$\sup_{\mathbf{v} \in \mathcal{Y}} \log \mathbb{E} \exp \left\{ \lambda \boldsymbol{\gamma}_1^\top \nabla \mathcal{Y}(\mathbf{v}) \boldsymbol{\gamma}_2 \right\} \leq \frac{v_0^2 \lambda^2}{2}, \quad |\lambda| \leq g. \quad (\text{B.30})$$

Condition (B.30) implies for any $\mathbf{v} \in \mathcal{Y}_o$ with $\|\mathbf{v}\| \leq r$ and $\boldsymbol{\gamma} \in \mathbb{R}^q$ with $\|\boldsymbol{\gamma}\| = 1$ in view of $\mathcal{Y}(\mathbf{v}^*) = 0$ by Lemma B.4

$$\log \mathbb{E} \exp \left\{ \frac{\lambda}{r} \mathcal{Y}(\mathbf{v})^\top \boldsymbol{\gamma} \right\} \leq \frac{v_0^2 \lambda^2 \|\mathbf{v}\|^2}{2r^2}, \quad |\lambda| \leq g. \quad (\text{B.31})$$

In what follows, we use the representation

$$\|\mathcal{Y}(\mathbf{v})\| = \sup_{\|\mathbf{u}\| \leq r} \frac{1}{r} \mathbf{u}^\top \mathcal{Y}(\mathbf{v}). \quad (\text{B.32})$$

This implies for $\mathcal{Y}_o(r) = \{\mathbf{v} \in \mathcal{Y} : \|\mathbf{v} - \mathbf{v}^*\| \leq r\}$

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

Consider a bivariate process $\mathbf{u}^\top \mathcal{Y}(\mathbf{v})$ of $\mathbf{u} \in \mathbb{R}^q$ and $\mathbf{v} \in \mathcal{Y} \subset \mathbb{R}^p$. By definition $\mathbb{E} \mathbf{u}^\top \mathcal{Y}(\mathbf{v}) = 0$. Further, $\nabla_{\mathbf{u}}[\mathbf{u}^\top \mathcal{Y}(\mathbf{v})] = \mathcal{Y}(\mathbf{v})$ while $\nabla_{\mathbf{v}}[\mathbf{u}^\top \mathcal{Y}(\mathbf{v})] = \mathbf{u}^\top \nabla \mathcal{Y}(\mathbf{v}) = \|\mathbf{u}\| \boldsymbol{\gamma}^\top \nabla \mathcal{Y}(\mathbf{v})$ for $\boldsymbol{\gamma} = \mathbf{u}/\|\mathbf{u}\|$. Suppose that $\mathbf{u} \in \mathbb{R}^q$ and $\mathbf{v} \in \mathcal{Y}$ are such that $\|\mathbf{u}\| \leq r$ and $\|\mathbf{v}\| \leq r$. By (B.30), it holds for $\boldsymbol{\gamma} \in \mathbb{R}^p$ with $\|\boldsymbol{\gamma}\| = 1$ and $\mathbf{v} \in \mathcal{Y}_o(r)$

$$\log \mathbb{E} \exp \left\{ \frac{\lambda}{r} \nabla_{\mathbf{v}}[\mathbf{u}^\top \mathcal{Y}(\mathbf{v})] \boldsymbol{\gamma} \right\} \leq \log \mathbb{E} \exp \left\{ \frac{\lambda}{r} \mathbf{u}^\top \nabla \mathcal{Y}(\mathbf{v}) \boldsymbol{\gamma} \right\} \leq \frac{v_0^2 \lambda^2}{2},$$

and by (B.31) for a unit vector $\boldsymbol{\gamma} \in \mathbb{R}^q$

$$\log \mathbb{E} \exp \left\{ \frac{\lambda}{r} \nabla_{\mathbf{u}}[\mathbf{u}^\top \mathcal{Y}(\mathbf{v})] \boldsymbol{\gamma} \right\} \leq \log \mathbb{E} \exp \left\{ \frac{\lambda}{r} \mathcal{Y}(\mathbf{v}) \boldsymbol{\gamma} \right\} \leq \frac{v_0^2 \lambda^2}{2}.$$

Therefore, $(\mathcal{E}D_p)$ is fulfilled for $\mathbf{u}^\top \mathcal{Y}(\mathbf{v})$ and Theorem B.7 applies. We summarize our findings in the following theorem.

Theorem B.14. Let a random p -vector process $\mathcal{Y}(\mathbf{v})$ for $\mathbf{v} \in \mathcal{Y} \subseteq \mathbb{R}^p$ fulfill $\mathcal{Y}(\mathbf{v}^*) = 0$, $\mathbb{E}\mathcal{Y}(\mathbf{v}) \equiv 0$, and the condition (B.30) be satisfied. Then for each r and any $x \geq 1/2$, it holds for $\mathcal{Y}_\circ = \mathcal{Y}_\circ(r)$

$$\mathbb{P}\left\{\sup_{\mathbf{v} \in \mathcal{Y}_\circ(r)} \|\mathcal{Y}(\mathbf{v})\| \geq \sqrt{8}v_0 r \mathfrak{z}_{\mathbb{H}}(x)\right\} \leq e^{-x}, \quad (\text{B.33})$$

where $\mathfrak{z}_{\mathbb{H}}(x)$ is given by (B.3) with $\mathbb{Q}_1 = \mathbb{Q}_1(\mathcal{Y}_\circ) + \sqrt{6q}$ and $\mathbb{Q}_2 = \mathbb{Q}_2(\mathcal{Y}_\circ) + 6q$.

B.8. A bound for a family of quadratic forms

Now we consider an extension of the previous result with a quadratic form $\|A\mathcal{Y}(\mathbf{v})\|^2$ to be bounded under the conditions (B.30) and (B.31) on $\mathcal{Y}(\mathbf{v})$ for $\mathbf{v} \in \mathcal{Y} \subseteq \mathbb{R}^p$. Here $\mathcal{Y}(\cdot)$ is a vector process with values in \mathbb{R}^q and A is a $q \times q$ matrix with $\|A^\top A\|_{\text{op}} \leq 1$. The idea is to use the representation (B.32) in which we replace \mathbf{u} with $A\mathbf{u}$. The bound (B.33) implies for any r

$$\mathbb{P}\left\{\sup_{\mathbf{v} \in \mathcal{Y}_\circ(r), \|A\mathbf{u}\| \leq r} \mathbf{u}^\top A \mathcal{Y}(\mathbf{v}) > \sqrt{8}v_0 r \mathfrak{z}_{\mathbb{H}}(x)\right\} \leq e^{-x},$$

where $\mathfrak{z}_{\mathbb{H}}(x)$ corresponds to $\mathbb{Q}_1 = \sqrt{\mathbb{Q}_2} = \sqrt{6p + \mathbb{Q}_2(\mathcal{Y}_\circ)}$.

Now we discuss how this bound can be refined if A is a smoothing operator. For simplicity assume that A fulfills the condition of Theorem B.10. One can expect that the dimension q can be replaced by the effective dimension p_A . The arguments similar to the above yield

$$\|A\mathcal{Y}(\mathbf{v})\| = \sup_{\mathbf{u} \in \mathbb{R}^q: \|\mathbf{u}\| \leq r} \frac{1}{r} \mathbf{u}^\top A \mathcal{Y}(\mathbf{v}),$$

and we again consider a bivariate process $\mathbf{u}^\top A \mathcal{Y}(\mathbf{v})$ of $\mathbf{u} \in \mathbb{R}^q$ and $\mathbf{v} \in \mathcal{Y} \subseteq \mathbb{R}^p$. The conditions (B.30) and (B.31) imply for any two unit vectors $\boldsymbol{\gamma}_1 \in \mathbb{R}^q$ and $\boldsymbol{\gamma}_2 \in \mathbb{R}^p$ and any points $\mathbf{u} \in \mathbb{R}^q$ with $\|A\mathbf{u}\| \leq r$ and $\mathbf{v} \in \mathcal{Y}_\circ(r)$, it holds

$$\log \mathbb{E} \exp\left\{\frac{\lambda}{r} \nabla_{\mathbf{v}}[\mathbf{u}^\top A \mathcal{Y}(\mathbf{v})] \boldsymbol{\gamma}_2\right\} = \log \mathbb{E} \exp\left\{\frac{\lambda}{r} \mathbf{u}^\top A \nabla \mathcal{Y}(\mathbf{v}) \boldsymbol{\gamma}_2\right\} \leq \frac{v_0^2 \lambda^2}{2},$$

and by (B.31) with $\nabla_{\mathbf{v}}^2 = A^\top A$

$$\log \mathbb{E} \exp\left\{\frac{\lambda}{\|\nabla_{\mathbf{v}} \boldsymbol{\gamma}_1\|} \boldsymbol{\gamma}_1^\top \nabla_{\mathbf{u}}[\mathbf{u}^\top A \mathcal{Y}(\mathbf{v})]\right\} \leq \log \mathbb{E} \exp\left\{\frac{\lambda}{\|\nabla_{\mathbf{v}} \boldsymbol{\gamma}_1\|} (A \boldsymbol{\gamma}_1)^\top \mathcal{Y}(\mathbf{v})\right\} \leq \frac{v_0^2 \lambda^2}{2}.$$

Therefore, $(\mathcal{E}D_p)$ is fulfilled for $\mathbf{u}^\top A \mathcal{Y}(\mathbf{v})$. Now we apply the bound from Theorem B.13 and the entropy bound for the elliptic set $\|A\mathbf{u}\| \leq r$ from Theorem B.10.

Theorem B.15. Let a random vector process $\mathcal{Y}(\mathbf{v}) \in \mathbb{R}^q$ for $\mathbf{v} \in \mathcal{Y} \subseteq \mathbb{R}^p$ fulfill $\mathcal{Y}(\mathbf{v}^*) = 0$, $\mathbb{E}\mathcal{Y}(\mathbf{v}) \equiv 0$, and the condition (B.30) be satisfied. Let A fulfill $1/2 \leq \|AA^\top\|_{\text{op}} \leq 1$. Then for each r , it holds

$$\mathbb{P}\left\{\sup_{\mathbf{v} \in \mathcal{Y}_\circ(r)} \|A\mathcal{Y}(\mathbf{v})\| > \sqrt{8}v_0 r \mathfrak{z}_{\mathbb{H}}(x)\right\} \leq e^{-x},$$

where $\mathfrak{z}_{\mathbb{H}}(x)$ is given by (B.3) with $\mathbb{Q}_2 = C p_A + \mathbb{Q}_2(\mathcal{Y}_\circ(r))$ and $\mathbb{Q}_1 = C \sqrt{p_A} + \mathbb{Q}_1(\mathcal{Y}_\circ(r))$.

References

- [1] A. Andresen and V. Spokoiny. Critical dimension in profile semiparametric estimation. *Electron. J. Stat.* **8** (2) (2014) 3077–3125. MR3301302
- [2] A. Barron, L. Birgé and P. Massart. Risk bounds for model selection via penalization. *Probab. Theory Related Fields* **113** (3) (1999) 301–413. MR1679028

- [3] A. Belloni and V. Chernozhukov. On the computational complexity of MCMC-based estimators in large samples. *Ann. Statist.* **37** (4) (2009) 2011–2055. [MR2533478](#)
- [4] L. Birgé and P. Massart. Minimum contrast estimators on sieves: Exponential bounds and rates of convergence. *Bernoulli* **4** (3) (1998) 329–375. [MR1653272](#)
- [5] L. Birgé and P. Massart. Gaussian model selection. *J. Eur. Math. Soc. (JEMS)* **3** (3) (2001) 203–268. [MR1848946](#)
- [6] L. Birgé and P. Massart. Minimal penalties for Gaussian model selection. *Probab. Theory Related Fields* **138** (1–2) (2007) 33–73. [MR2288064](#)
- [7] S. Boucheron and P. Massart. A high-dimensional Wilks phenomenon. *Probab. Theory Related Fields* **150** (2011) 405–433. [MR2824862](#)
- [8] J. Fan, C. Zhang and J. Zhang. Generalized likelihood ratio statistics and Wilks phenomenon. *Ann. Statist.* **29** (1) (2001) 153–193. [MR1833962](#)
- [9] S. Ghosal. Asymptotic normality of posterior distributions in high-dimensional linear models. *Bernoulli* **5** (2) (1999) 315–331. [MR1681701](#)
- [10] S. Ghosal. Asymptotic normality of posterior distributions for exponential families when the number of parameters tends to infinity. *J. Multivariate Anal.* **74** (1) (2000) 49–68. [MR1790613](#)
- [11] Y. Golubev and V. Spokoiny. Exponential bounds for minimum contrast estimators. *Electron. J. Stat.* **3** (2009) 712–746. [MR2534199](#)
- [12] P. J. Green and B. W. Silverman. *Nonparametric Regression and Generalized Linear Models: A Roughness Penalty Approach*. Chapman & Hall, London, 1994. [MR1270012](#)
- [13] P. J. Huber. The behavior of maximum likelihood estimates under nonstandard conditions. In *Proc. 5th Berkeley Symp. Math. Stat. Probab., Univ. Calif. 1965/66*, 1 221–233. Univ. California Press, Berkeley, CA, 1967. [MR0216620](#)
- [14] I. A. Ibragimov and R. Z. Khas'minskij. *Statistical Estimation. Asymptotic Theory*. Springer, New York, 1981. Transl. from the Russian by Samuel Kotz. [MR0620321](#)
- [15] Y. Kim. The Bernstein–von Mises theorem for the proportional hazard model. *Ann. Statist.* **34** (4) (2006) 1678–1700. [MR2283713](#)
- [16] R. Koenker, P. Ng and S. Portnoy. Quantile smoothing splines. *Biometrika* **81** (4) (1994) 673–680. [MR1326417](#)
- [17] B. Laurent and P. Massart. Adaptive estimation of a quadratic functional by model selection. *Ann. Statist.* **28** (5) (2000) 1302–1338.
- [18] E. Mammen. Empirical process of residuals for high-dimensional linear models. *Ann. Statist.* **24** (1) (1996) 307–335. [MR1389892](#)
- [19] S. Portnoy. Asymptotic behavior of M -estimators of p regression parameters when p^2/n is large. I. Consistency. *Ann. Statist.* **12** (4) (1984) 1298–1309. [MR0760690](#)
- [20] S. Portnoy. Asymptotic behavior of M estimators of p regression parameters when p^2/n is large. II. Normal approximation. *Ann. Statist.* **13** (4) (1985) 1403–1417. [MR0811499](#)
- [21] S. Portnoy. Asymptotic behavior of the empiric distribution of M -estimated residuals from a regression model with many parameters. *Ann. Statist.* **14** (1986) 1152–1170. [MR0856812](#)
- [22] S. Portnoy. Asymptotic behavior of likelihood methods for exponential families when the number of parameters tends to infinity. *Ann. Statist.* **16** (1) (1988) 356–366. [MR0924876](#)
- [23] X. Shen. On methods of sieves and penalization. *Ann. Statist.* **25** (6) (1997) 2555–2591. [MR1604416](#)
- [24] X. Shen and W. H. Wong. Convergence rate of sieve estimates. *Ann. Statist.* **22** (2) (1994) 580–615. [MR1292531](#)
- [25] V. Spokoiny. Parametric estimation. Finite sample theory. *Ann. Statist.* **40** (6) (2012) 2877–2909. [MR3097963](#)
- [26] V. Spokoiny, W. Wang and W. Härdle. Local quantile regression (with rejoinder). *J. Statist. Plann. Inference* **143** (7) (2013) 1109–1129. [MR3049611](#)
- [27] V. Spokoiny and M. Zhilova. Bootstrap confidence sets under model misspecification. *Ann. Statist.* **43** (2015) 2653–2675. [MR3405607](#)
- [28] S. van de Geer. M -estimation using penalties or sieves. *J. Statist. Plann. Inference* **108** (1–2) (2002) 55–69. [MR1947391](#)
- [29] S. A. Van de Geer. *Applications of Empirical Process Theory. Cambridge Series in Statistical and Probabilistic Mathematics*. Cambridge University Press, Cambridge, 2000. [MR1739079](#)
- [30] A. van der Vaart and J. A. Wellner. *Weak Convergence and Empirical Processes. With Applications to Statistics. Springer Series in Statistics*. Springer, New York, 1996. [MR1385671](#)
- [31] A. Zaitsev, E. Burnaev and V. Spokoiny. Properties of the posterior distribution of a regression model based on Gaussian random fields. *Autom. Remote Control* **74** (10) (2013) 1645–1655. [MR3219856](#)