

Finite Sample Bernstein – von Mises Theorem for Semiparametric Problems

Maxim Panov* and Vladimir Spokoiny†

Abstract. The classical parametric and semiparametric Bernstein – von Mises (BvM) results are reconsidered in a non-classical setup allowing finite samples and model misspecification. In the case of a finite dimensional nuisance parameter we obtain an upper bound on the error of Gaussian approximation of the posterior distribution for the target parameter which is explicit in the dimension of the nuisance and target parameters. This helps to identify the so called *critical dimension* p_n of the full parameter for which the BvM result is applicable. In the important i.i.d. case, we show that the condition “ p_n^3/n is small” is sufficient for the BvM result to be valid under general assumptions on the model. We also provide an example of a model with the phase transition effect: the statement of the BvM theorem fails when the dimension p_n approaches $n^{1/3}$. The results are extended to the case of infinite dimensional parameters with the nuisance parameter from a Sobolev class.

Keywords: prior, posterior, Bayesian inference, semiparametric, critical dimension.

1 Introduction

The prominent Bernstein – von Mises (BvM) theorem claims that the posterior measure is asymptotically normal with the mean close to the maximum likelihood estimator (MLE) and the posterior variance is nearly the inverse of the total Fisher information matrix. The BvM result provides a theoretical background for Bayesian computations of the MLE and its variance. Also it justifies usage of elliptic credible sets based on the first two moments of the posterior. The classical version of the BvM Theorem is stated for the standard parametric setup with a fixed parametric model and large samples; see Le Cam and Yang (1990); van der Vaart (2000) for a detailed overview. However, in modern statistics applications one often faces very complicated models involving a lot of parameters and with a limited sample size. This requires an extension of the classical results to this non-classical situation. We mention Cox (1993); Freedman (1999); Ghosal (1999); Johnstone (2010) and references therein for some special phenomena arising in the Bayesian analysis when the parameter dimension increases. Already consistency of the posterior distribution in nonparametric and semiparametric models is a nontrivial problem; cf. Schwartz (1965) and Barron et al. (1996). Asymptotic normality of the posterior measure for these classes of models is even more challenging; see e.g. Shen

*Moscow Institute of Physics and Technology, Institute for Information Transmission Problems of RAS, Datadvance Company, Pokrovsky blvd. 3 building 1B, 109028 Moscow, Russia, panov.maxim@gmail.com

†Weierstrass Institute and Humboldt University Berlin, Moscow Institute of Physics and Technology, Institute for Information Transmission Problems of RAS, Mohrenstr. 39, 10117 Berlin, Germany, spokoiny@wias-berlin.de

(2002). Some results for particular semi and nonparametric problems are available from Kim and Lee (2004); Kim (2006); Leahu (2011); Castillo and Nickl (2013). Cheng and Kosorok (2008) obtained a version of the BvM statement based on a high order expansion of the profile sampler. The recent paper by Bickel and Kleijn (2012) extends the BvM statement from the classical parametric case to a rather general i.i.d. framework. Castillo (2012) studies the semiparametric BvM result for Gaussian process functional priors. In Rivoirard and Rousseau (2012) a semiparametric BvM theorem is derived for linear functionals of density and in forthcoming work (Castillo and Rousseau, 2013) the result is generalized to a broad class of models and functionals. However, all these results are limited to the asymptotic setup and to some special classes of models like i.i.d. or Gaussian.

In this paper we reconsider the BvM result for the parametric component of a general semiparametric model. An important feature of the study is that the sample size is fixed, we proceed with just one sample. A finite sample theory is especially challenging because most notions, methods and tools in the classical theory are formulated in the asymptotic setup with growing sample size. Only a few finite sample general results are available; see e.g. the recent paper by Boucheron and Massart (2011). This paper focuses on the semiparametric problem when the full parameter is large or infinite dimensional but the target is low dimensional. In the Bayesian framework, the aim is the marginal of the posterior corresponding to the target parameter; cf. Castillo (2012). Typical examples are provided by functional estimation, estimation of a function at a point, or simply by estimating a given subvector of the parameter vector. An interesting feature of the semiparametric BvM result is that the nuisance parameter appears only via the effective score and the related efficient Fisher information; cf. Bickel and Kleijn (2012). The methods of study heavily rely on the notion of the hardest parametric submodel. In addition, one assumes that an estimate of the nuisance parameter is available which ensures a certain accuracy of estimation; see Cheng and Kosorok (2008) or Bickel and Kleijn (2012). This essentially simplifies the study but does not allow to derive a qualitative relation between the full dimension of the parameter space and the total available information in the data.

Some recent results study the impact of a growing parameter dimension p_n on the quality of Gaussian approximation of the posterior. We mention Ghosal (1999, 2000), Boucheron and Gassiat (2009), Johnstone (2010) and Bontemps (2011) for specific examples. See the discussion after Theorem 4 below for more details.

In this paper we show that the *bracketing* approach of Spokoiny (2012) can be used for obtaining a finite sample semiparametric version of the Bernstein – von Mises theorem even if the full parameter dimension grows with the sample size. The ultimate goal of this paper is to quantify the so called critical parameter dimension for which the BvM result can be applied. Our approach neither relies on a pilot estimate of the nuisance and target parameter nor involves the notion of the hardest parametric submodel. In the case of finite dimensional nuisance the obtained results only require some smoothness of the log-likelihood function, its finite exponential moments, and some identifiability conditions. Further we specify this result to the i.i.d. setup and show that the imposed conditions are satisfied if p_n^3/n is small. We present an example showing

that the dimension $p_n = O(n^{1/3})$ is indeed critical and the BvM result starts to fail if p_n grows over $n^{1/3}$. If the nuisance is infinite dimensional then additionally some smoothness of the nonparametric part is required. We state the general BvM results and show the validity of the BvM result for linear and generalized linear models with the nuisance parameter from the Sobolev class and a uniform sieve prior distribution on the parameter space.

Now we describe our setup. Let \mathbf{Y} denote the observed random data, and \mathbb{P} denote the data distribution. The parametric statistical model assumes that the unknown data distribution \mathbb{P} belongs to a given parametric family $(\mathbb{P}_{\mathbf{v}})$:

$$\mathbf{Y} \sim \mathbb{P} = \mathbb{P}_{\mathbf{v}^*} \in (\mathbb{P}_{\mathbf{v}}, \mathbf{v} \in \mathcal{T}),$$

where \mathcal{T} is some parameter space and $\mathbf{v}^* \in \mathcal{T}$ is the true value of parameter. In the semiparametric framework, one attempts to recover only a low dimensional component $\boldsymbol{\theta}$ of the whole parameter \mathbf{v} . This means that the target of estimation is $\boldsymbol{\theta}^* \stackrel{\text{def}}{=} \Pi_0 \mathbf{v}^*$ for some mapping $\Pi_0 : \mathcal{T} \rightarrow \mathbb{R}^q$, and $q \in \mathbb{N}$ stands for the dimension of the target. Usually in the classical semiparametric setup, the vector \mathbf{v} is represented as $\mathbf{v} = (\boldsymbol{\theta}, \boldsymbol{\eta})$, where $\boldsymbol{\theta}$ is the target of analysis while $\boldsymbol{\eta}$ is the *nuisance parameter*. We refer to this situation as the $(\boldsymbol{\theta}, \boldsymbol{\eta})$ -setup and our presentation follows this setting. An extension to the \mathbf{v} -setup with $\boldsymbol{\theta} = \Pi_0 \mathbf{v}$ is straightforward. Also for simplicity we first develop our results for the case when the total parameter space \mathcal{T} is a subset of the Euclidean space of dimensionality p .

Another issue addressed in this paper is the model misspecification. In the most of practical problems, it is unrealistic to expect that the model assumptions are exactly fulfilled, even if some rich nonparametric models are used. This means that the true data distribution \mathbb{P} does not belong to the considered family $(\mathbb{P}_{\mathbf{v}}, \mathbf{v} \in \mathcal{T})$. The “true” value \mathbf{v}^* of the parameter \mathbf{v} can be defined by

$$\mathbf{v}^* \stackrel{\text{def}}{=} \operatorname{argmax}_{\mathbf{v} \in \mathcal{T}} \mathbb{E} \mathcal{L}(\mathbf{v}), \tag{1}$$

where $\mathcal{L}(\mathbf{v}) = \log \frac{d\mathbb{P}_{\mathbf{v}}}{d\boldsymbol{\mu}_0}(\mathbf{Y})$ is the log-likelihood function of the family $(\mathbb{P}_{\mathbf{v}})$ for some dominating measure $\boldsymbol{\mu}_0$. Under model misspecification, \mathbf{v}^* defines the best parametric fit to \mathbb{P} by the considered family; cf. Chernozhukov and Hong (2003), Kleijn and van der Vaart (2006, 2012) and references therein. The target $\boldsymbol{\theta}^*$ is defined by the mapping Π_0 :

$$\boldsymbol{\theta}^* \stackrel{\text{def}}{=} \Pi_0 \mathbf{v}^*.$$

Now we switch to the Bayesian set-up. Let π be a prior measure on the parameter set \mathcal{T} . Below we study the properties of the posterior measure which is the random measure on \mathcal{T} describing the conditional distribution of \mathbf{v} given \mathbf{Y} and obtained by normalization of the product $\exp\{\mathcal{L}(\mathbf{v})\} \pi(d\mathbf{v})$. This relation is usually written as

$$\mathbf{v} \mid \mathbf{Y} \propto \exp\{\mathcal{L}(\mathbf{v})\} \pi(d\mathbf{v}). \tag{2}$$

An important feature of our analysis is that $\mathcal{L}(\boldsymbol{v})$ is not assumed to be the true log-likelihood. This means that a model misspecification is possible and the underlying data distribution can be beyond the considered parametric family. In this sense, the Bayes formula (2) describes a *quasi posterior*; Chernozhukov and Hong (2003). Below we show that smoothness of the log-likelihood function $\mathcal{L}(\boldsymbol{v})$ ensures a kind of a Gaussian approximation of the posterior measure. Our focus is to describe the accuracy of such approximation as a function of the parameter dimension p and the other important characteristics of the model.

We suppose that the prior measure π has a positive density $\pi(\boldsymbol{v})$ w.r.t. to the Lebesgue measure on \mathcal{Y} : $\pi(d\boldsymbol{v}) = \pi(\boldsymbol{v})d\boldsymbol{v}$. Then (2) can be written as

$$\boldsymbol{v} \mid \boldsymbol{Y} \propto \exp\{\mathcal{L}(\boldsymbol{v})\} \pi(\boldsymbol{v}). \quad (3)$$

The famous Bernstein – von Mises (BvM) theorem claims that the posterior centered by any efficient estimator $\tilde{\boldsymbol{v}}$ of the parameter \boldsymbol{v}^* (for example the MLE) and scaled by the total Fisher information matrix is nearly standard normal:

$$\mathcal{D}_0(\boldsymbol{v} - \tilde{\boldsymbol{v}}) \mid \boldsymbol{Y} \xrightarrow{w} \mathcal{N}(0, I_p),$$

where I_p is an identity matrix of dimension p .

An important feature of the posterior distribution is that it is entirely known and can be numerically assessed. If we know in addition that the posterior is nearly normal, it suffices to compute its mean and variance for building the concentration and credible sets. The BvM result does not require the prior distribution to be proper and the phenomenon can be observed in the case of improper priors as well (for examples, see Bochkina and Green (2014)).

In this work we investigate the properties of the posterior distribution for the target parameter $\boldsymbol{\vartheta} \mid \boldsymbol{Y} = \Pi_0 \boldsymbol{v} \mid \boldsymbol{Y}$. In this case (3) can be written as

$$\boldsymbol{\vartheta} \mid \boldsymbol{Y} \propto \int \exp\{\mathcal{L}(\boldsymbol{v})\} \pi(\boldsymbol{v}) d\boldsymbol{\eta}. \quad (4)$$

The BvM result in this case transforms into

$$\check{D}_0(\boldsymbol{\vartheta} - \tilde{\boldsymbol{\theta}}) \mid \boldsymbol{Y} \xrightarrow{w} \mathcal{N}(0, I_q),$$

where I_q is an identity matrix of dimension q , $\tilde{\boldsymbol{\theta}} = \Pi_0 \tilde{\boldsymbol{v}}$, and \check{D}_0^2 is given in (6).

We consider two important classes of priors, namely non-informative and flat Gaussian priors. Our goal is to show under mild conditions that the posterior distribution of the target parameter (4) is close to a prescribed Gaussian law even for finite samples. The other important issue is to specify the conditions on the sample size and the dimension of the parameter space for which the BvM result is still applicable.

2 BvM Theorem with a finite dimensional nuisance

This section presents our main results for the case of a finite dimensional parameter \mathbf{v} , i.e. $\dim(\mathcal{Y}) = p < \infty$. One of the main elements of our construction is a $p \times p$ matrix \mathcal{D}_0^2 which is defined similarly to the Fisher information matrix:

$$\mathcal{D}_0^2 \stackrel{\text{def}}{=} -\nabla^2 \mathbb{E}\mathcal{L}(\mathbf{v}^*). \tag{5}$$

Here and in what follows we work under conditions which are close to the classical conditions of a *regular parametric family* (see Ibragimov and Khas'minskij (1981)) and assume that the log-likelihood function $\mathcal{L}(\mathbf{v})$ is sufficiently smooth in \mathbf{v} , $\nabla\mathcal{L}(\mathbf{v})$ stands for its gradient and $\nabla^2\mathbb{E}\mathcal{L}(\mathbf{v})$ for the Hessian of the expectation $\mathbb{E}\mathcal{L}(\mathbf{v})$, and the true value \mathbf{v}^* is due to (1). Also define the score

$$\boldsymbol{\xi} \stackrel{\text{def}}{=} \mathcal{D}_0^{-1}\nabla\mathcal{L}(\mathbf{v}^*).$$

Under our conditions we can permute expectation and differentiation of the likelihood and thus the definition of \mathbf{v}^* implies $\nabla\mathbb{E}\mathcal{L}(\mathbf{v}^*) = 0$ and hence, $\mathbb{E}\boldsymbol{\xi} = 0$.

For the $(\boldsymbol{\theta}, \boldsymbol{\eta})$ -setup, we consider the block representation of the vector $\nabla\mathcal{L}(\mathbf{v}^*)$ and of the matrix and \mathcal{D}_0^2 from (5):

$$\nabla\mathcal{L}(\mathbf{v}^*) = \begin{pmatrix} \nabla_{\boldsymbol{\theta}} \\ \nabla_{\boldsymbol{\eta}} \end{pmatrix}, \quad \mathcal{D}_0^2 = \begin{pmatrix} D_0^2 & A_0 \\ A_0^\top & H_0^2 \end{pmatrix}.$$

Define also the $q \times q$ matrix \check{D}_0^2 and random vectors $\check{\nabla}_{\boldsymbol{\theta}}, \check{\boldsymbol{\xi}} \in \mathbb{R}^q$ as

$$\check{D}_0^2 \stackrel{\text{def}}{=} D_0^2 - A_0 H_0^{-2} A_0^\top, \tag{6}$$

$$\check{\nabla}_{\boldsymbol{\theta}} \stackrel{\text{def}}{=} \nabla_{\boldsymbol{\theta}} - A_0 H_0^{-2} \nabla_{\boldsymbol{\eta}},$$

$$\check{\boldsymbol{\xi}} \stackrel{\text{def}}{=} \check{D}_0^{-1} \check{\nabla}_{\boldsymbol{\theta}}. \tag{7}$$

The $q \times q$ matrix \check{D}_0^2 is usually called the efficient Fisher information matrix, while the random vector $\check{\boldsymbol{\xi}} \in \mathbb{R}^q$ is the efficient score. Everywhere in the text for a vector \mathbf{a} we denote by $\|\mathbf{a}\|$ its Euclidean norm and for a matrix A we denote by $\|A\|$ its operator norm.

2.1 Conditions

Our results assume a number of conditions to be satisfied. The list is essentially as in Spokoiny (2012), one can find there some discussion and examples showing that the conditions are not restrictive and are fulfilled in most classical models used in statistical studies like i.i.d., regression or generalized linear models. The conditions are split into

local and global. The local conditions only describe the properties of the process $\mathcal{L}(\mathbf{v})$ for $\mathbf{v} \in \mathcal{Y}_0(\mathbf{r}_0)$ with some fixed value \mathbf{r}_0 :

$$\mathcal{Y}_0(\mathbf{r}_0) \stackrel{\text{def}}{=} \{\mathbf{v} \in \mathcal{Y} : \|\mathcal{D}_0(\mathbf{v} - \mathbf{v}^*)\| \leq \mathbf{r}_0\}.$$

The global conditions have to be fulfilled on the whole \mathcal{Y} . Define the stochastic component $\zeta(\mathbf{v})$ of $\mathcal{L}(\mathbf{v})$:

$$\zeta(\mathbf{v}) \stackrel{\text{def}}{=} \mathcal{L}(\mathbf{v}) - \mathbb{E}\mathcal{L}(\mathbf{v})$$

and introduce the notation

$$L(\mathbf{v}, \mathbf{v}^*) \stackrel{\text{def}}{=} \mathcal{L}(\mathbf{v}) - \mathcal{L}(\mathbf{v}^*)$$

for the (quasi) log-likelihood ratio. We start with some exponential moments conditions.

(ED₀) There exists a constant $\nu_0 > 0$, a positive symmetric $p \times p$ matrix \mathcal{V}_0^2 satisfying $\text{Var}\{\nabla\zeta(\mathbf{v}^*)\} \leq \mathcal{V}_0^2$, and a constant $\mathbf{g} > 0$ such that

$$\sup_{\gamma \in \mathbb{R}^p} \log \mathbb{E} \exp \left\{ \mu \frac{\langle \nabla\zeta(\mathbf{v}^*), \gamma \rangle}{\|\mathcal{V}_0\gamma\|} \right\} \leq \frac{\nu_0^2 \mu^2}{2}, \quad |\mu| \leq \mathbf{g}.$$

(ED₂) There exists a constant $\nu_0 > 0$, a constant $\omega > 0$ and for each $\mathbf{r} > 0$ a constant $\mathbf{g}(\mathbf{r}) > 0$ such that for all $\mathbf{v} \in \mathcal{Y}_0(\mathbf{r})$:

$$\sup_{\gamma_1, \gamma_2 \in \mathbb{R}^p} \log \mathbb{E} \exp \left\{ \frac{\mu}{\omega} \frac{\gamma_1^\top \nabla^2 \zeta(\mathbf{v}) \gamma_2}{\|\mathcal{D}_0 \gamma_1\| \cdot \|\mathcal{D}_0 \gamma_2\|} \right\} \leq \frac{\nu_0^2 \mu^2}{2}, \quad |\mu| \leq \mathbf{g}(\mathbf{r}).$$

The next condition is needed to ensure some smoothness properties of expected log-likelihood $\mathbb{E}\mathcal{L}(\mathbf{v})$ in the local zone $\mathbf{v} \in \mathcal{Y}_0(\mathbf{r}_0)$. Define

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

Then $\mathcal{D}_0^2 = \mathcal{D}_0^2(\mathbf{v}^*)$.

(L₀) There exists a constant $\delta(\mathbf{r})$ such that it holds on the set $\mathcal{Y}_0(\mathbf{r})$ for all $\mathbf{r} \leq \mathbf{r}_0$

$$\|\mathcal{D}_0^{-1} \mathcal{D}_0^2(\mathbf{v}) \mathcal{D}_0^{-1} - \mathbf{I}_p\| \leq \delta(\mathbf{r}).$$

The global identification condition is:

(L_r) For any \mathbf{r} there exists a value $\mathbf{b}(\mathbf{r}) > 0$, such that $\mathbf{r}\mathbf{b}(\mathbf{r}) \rightarrow \infty$, $\mathbf{r} \rightarrow \infty$ and

$$-\mathbb{E}L(\mathbf{v}, \mathbf{v}^*) \geq \mathbf{r}^2 \mathbf{b}(\mathbf{r}) \quad \text{for all } \mathbf{v} \text{ with } \mathbf{r} = \|\mathcal{D}_0(\mathbf{v} - \mathbf{v}^*)\|.$$

Finally we specify the identifiability conditions. We begin by representing the information and the covariance matrices in block form:

$$\mathcal{D}_0^2 = \begin{pmatrix} D_0^2 & A_0 \\ A_0^\top & H_0^2 \end{pmatrix}, \quad \mathcal{V}_0^2 = \begin{pmatrix} V_0^2 & B_0 \\ B_0^\top & Q_0^2 \end{pmatrix}.$$

The *identifiability conditions* in Spokoiny (2012) ensure that the matrix \mathcal{D}_0^2 is positive and satisfied $\alpha^2 \mathcal{D}_0^2 \geq \mathcal{V}_0^2$ for some $\alpha > 0$. Here we restate these conditions in the special block form which is specific for the $(\boldsymbol{\theta}, \boldsymbol{\eta})$ -setup.

(\mathcal{I}) There are constants $\alpha > 0$ and $\nu < 1$ such that

$$\alpha^2 D_0^2 \geq V_0^2, \quad \alpha^2 H_0^2 \geq Q_0^2, \quad \alpha^2 \mathcal{D}_0^2 \geq \mathcal{V}_0^2. \tag{8}$$

and

$$\|D_0^{-1} A_0 H_0^{-2} A_0^\top D_0^{-1}\| \leq \nu. \tag{9}$$

The quantity ν bounds the angle between the target and nuisance subspaces in the tangent space. The regularity condition (\mathcal{I}) ensures that this angle is not too small and hence, the target and nuisance parameters are identifiable. In particular, the matrix \check{D}_0^2 from (6) is well posed under (\mathcal{I}). The bounds in (8) are given with the same constant α only for simplifying the notation. One can show that the last bound on \mathcal{D}_0^2 follows from the first two and (9) with another constant α' depending on α and ν only.

2.2 The main results

First we state the BvM result about the properties of the $\boldsymbol{\vartheta}$ -posterior given by (4) in the case of a uniform prior that is, $\pi(\boldsymbol{v}) \equiv 1$ on \mathcal{Y} . Define

$$\bar{\boldsymbol{\vartheta}} \stackrel{\text{def}}{=} \mathbb{E}(\boldsymbol{\vartheta} | \mathbf{Y}), \quad \mathfrak{S}^2 \stackrel{\text{def}}{=} \text{Cov}(\boldsymbol{\vartheta} | \mathbf{Y}) \stackrel{\text{def}}{=} \mathbb{E}\{(\boldsymbol{\vartheta} - \bar{\boldsymbol{\vartheta}})(\boldsymbol{\vartheta} - \bar{\boldsymbol{\vartheta}})^\top | \mathbf{Y}\}. \tag{10}$$

Also define

$$\boldsymbol{\theta}^\circ \stackrel{\text{def}}{=} \boldsymbol{\theta}^* + \check{D}_0^{-1} \check{\boldsymbol{\xi}},$$

where \check{D}_0 and $\check{\boldsymbol{\xi}}$ are defined by (6) and (7) respectively. The random point $\boldsymbol{\theta}^\circ$ can be viewed as a first order approximation of the profile MLE $\tilde{\boldsymbol{\theta}}$. Below we present a version of the BvM result in the considered nonasymptotic setup which claims that $\bar{\boldsymbol{\vartheta}}$ is close to $\boldsymbol{\theta}^\circ$, \mathfrak{S}^2 is nearly equal to \check{D}_0^{-2} , and $\check{D}_0(\boldsymbol{\vartheta} - \boldsymbol{\theta}^\circ)$ is nearly standard normal conditionally on \mathbf{Y} .

We suppose that a large constant \mathbf{x} is fixed which specifies random events $\Omega(\mathbf{x})$ of *dominating probability*. We say that a generic random set $\Omega(\mathbf{x})$ is of dominating probability if

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

The notation C is for a generic absolute constant and x is for a positive value ensuring that e^{-x} is negligible. The exact values of C will be specified in each particular case. The formulation of the results also involve the radius r_0 and the spread $\Delta(r_0, x)$. The radius r_0 separates the local zone $\mathcal{Y}_o(r_0)$ which is a vicinity of the central point \boldsymbol{v}^* , and its complement $\mathcal{Y} \setminus \mathcal{Y}_o(r_0)$ for which we establish a large deviation result. The spread value $\Delta(r_0, x)$ measures the quality of local approximation of the log-likelihood $L(\boldsymbol{v}, \boldsymbol{v}^*)$ by a quadratic process $\mathbb{L}(\boldsymbol{v}, \boldsymbol{v}^*)$:

$$\Delta(r_0, x) \stackrel{\text{def}}{=} \{\delta(r_0) + 6\nu_0 z_{\mathbb{H}}(x)\omega\} r_0^2.$$

Here the term $\delta(r_0)r_0^2$ measures the error of a quadratic approximation of the expected log-likelihood $L(\boldsymbol{v})$ due to (\mathcal{L}_0) , while the second term $6\nu_0 z_{\mathbb{H}}(x)\omega r_0^2$ controls the stochastic term and involves the entropy of the parameter space which is involved in the definition of $z_{\mathbb{H}}(x)$. A precise formulation is given in Theorem 9 below.

Theorem 1. *Suppose the conditions of Section 2.1. Let the prior be uniform on \mathcal{Y} . Then there exists a random event $\Omega(x)$ of probability at least $1 - 4e^{-x}$ such that it holds on $\Omega(x)$*

$$\begin{aligned} \|\check{D}_0(\bar{\boldsymbol{\vartheta}} - \boldsymbol{\theta}^\circ)\|^2 &\leq 4\Delta(r_0, x) + 16e^{-x}, \\ \|\mathbf{I}_q - \check{D}_0\mathfrak{S}^2\check{D}_0\| &\leq 4\Delta(r_0, x) + 16e^{-x}, \end{aligned}$$

where $\bar{\boldsymbol{\vartheta}}$ and \mathfrak{S}^2 are from (10).

Moreover, on $\Omega(x)$ for any measurable set $A \subseteq \mathbb{R}^q$

$$\begin{aligned} \exp(-2\Delta(r_0, x) - 8e^{-x})\mathbb{P}(\boldsymbol{\gamma} \in A) - e^{-x} \\ \leq \mathbb{P}(\check{D}_0(\boldsymbol{\vartheta} - \boldsymbol{\theta}^\circ) \in A \mid \mathbf{Y}) \\ \leq \exp(2\Delta(r_0, x) + 5e^{-x})\mathbb{P}(\boldsymbol{\gamma} \in A), \end{aligned}$$

where $\boldsymbol{\gamma}$ is a standard Gaussian vector in \mathbb{R}^q .

The condition “ $\Delta(r_0, x)$ is small” yields the desirable BvM result, that is, the posterior measure after centering and standardization is close in total variation to the standard normal law. The classical asymptotic results immediately follow for many classical models (see discussion in Section 4). The next corollary extends the previous result by using empirically computable objects.

Corollary 1. *Under the conditions of Theorem 1 for any measurable set $A \subseteq \mathbb{R}^q$ a random event $\Omega(x)$ of a dominating probability at least $1 - 4e^{-x}$*

$$\begin{aligned} \exp(-2\Delta(r_0, x) - 8e^{-x})\{\mathbb{P}(\boldsymbol{\gamma} \in A) - \tau\} - e^{-x} \\ \leq \mathbb{P}(\mathfrak{S}^{-1}(\boldsymbol{\vartheta} - \bar{\boldsymbol{\vartheta}}) \in A \mid \mathbf{Y}) \\ \leq \exp(2\Delta(r_0, x) + 5e^{-x})\{\mathbb{P}(\boldsymbol{\gamma} \in A) + \tau\}, \end{aligned}$$

where γ is a standard Gaussian vector in \mathbb{R}^q and

$$\tau \stackrel{\text{def}}{=} \frac{1}{2} \left(q\Delta^2(\mathbf{r}_0, \mathbf{x}) + \{1 + \Delta(\mathbf{r}_0, \mathbf{x})\}^2 \Delta^2(\mathbf{r}_0, \mathbf{x}) \right).$$

This corollary is important as in practical applications we do not know matrix \check{D}_0 and vector θ° , but matrix \mathfrak{S}^{-1} and vector $\check{\boldsymbol{\theta}}$ can be computed by numerical computations. If dimension q is fixed the result becomes informative under the condition “ $\Delta(\mathbf{r}_0, \mathbf{x})$ is small”. Moreover, the statement can be extended to situations when the target dimension q grows but $\Delta(\mathbf{r}_0, \mathbf{x})q^{1/2}$ is still small.

2.3 Extension to a flat Gaussian prior

The previous results for a non-informative prior can be extended to the case of a flat prior $\Pi(d\mathbf{v})$. To be more specific we restrict ourselves to the case of a Gaussian prior. This is a prototypic situation because any smooth prior can be locally approximated by a Gaussian one. Without loss of generality the prior mean will be set to zero:

$$\Pi = \mathcal{N}(0, G^{-2})$$

with the density

$$\pi(\mathbf{v}) \propto \exp\{-\|G\mathbf{v}\|^2/2\}$$

for some positive symmetric matrix G^2 .

The non-informative prior can be viewed as a limiting case of a Gaussian prior as $G \rightarrow 0$. We are interested in quantifying this relation. How small should G be to ensure the BvM result? To explain the result, we first consider the Gaussian case when $P_{\mathbf{v}} = \mathcal{N}(\mathbf{v}, \mathcal{D}_0^{-2})$ and \mathbf{v}^* is the true point. It is well known that in this situation the non-informative prior leads to the Gaussian posterior $\mathcal{N}(\mathbf{v}^*, \mathcal{D}_0^{-2})$, while the Gaussian prior Π yields again the Gaussian posterior with the covariance \mathcal{D}_G^{-2} for $\mathcal{D}_G^2 = \mathcal{D}_0^2 + G^2$ and mean $\mathbf{v}_G^* = \mathcal{D}_G^{-2}\mathcal{D}_0^2\mathbf{v}^*$. Therefore, the prior does not significantly affect the posterior if two Gaussian measures $\mathcal{N}(\mathbf{v}^*, \mathcal{D}_0^{-2})$ and $\mathcal{N}(\mathbf{v}_G^*, \mathcal{D}_G^{-2})$ are nearly equivalent. The corresponding condition is represented in Lemma 8. It requires the values $\|\mathcal{D}_0^{-1}\mathcal{D}_G^2\mathcal{D}_0^{-1} - I_q\| = \|\mathcal{D}_0^{-1}G^2\mathcal{D}_0^{-1}\|$, $\text{tr}(\mathcal{D}_0^{-1}\mathcal{D}_G^2\mathcal{D}_0^{-1} - I_q)^2 = \text{tr}(\mathcal{D}_0^{-1}G^2\mathcal{D}_0^{-1})^2$ and $\|\mathcal{D}_G(\mathbf{v}^* - \mathbf{v}_G^*)\| \asymp \|\mathcal{D}_G^{-1}G^2\mathbf{v}^*\|$ to be small.

Theorem 2. *Suppose the conditions of Theorem 1. Let also $\Pi = \mathcal{N}(0, G^{-2})$ be a Gaussian prior measure on \mathbb{R}^p such that*

$$\|\mathcal{D}_0^{-1}G^2\mathcal{D}_0^{-1}\| \leq \epsilon \leq 1/2, \quad \text{tr}(\mathcal{D}_0^{-1}G^2\mathcal{D}_0^{-1})^2 \leq \delta^2, \quad \|\mathcal{D}_G^{-1}G^2\mathbf{v}^*\| \leq \beta,$$

where δ and β are given constants. Then it holds on a set $\Omega(\mathbf{x})$ of probability $1 - 5e^{-x}$

$$P(\check{D}_0(\boldsymbol{\theta} - \theta^\circ) \in A \mid \mathbf{Y}) \geq \exp(-2\Delta(\mathbf{r}_0, \mathbf{x}) - 8e^{-x})\{P(\gamma \in A) - \tau\} - e^{-x},$$

$$P(\check{D}_0(\boldsymbol{\theta} - \theta^\circ) \in A \mid \mathbf{Y}) \leq \exp(2\Delta(\mathbf{r}_0, \mathbf{x}) + 5e^{-x})\{P(\gamma \in A) + \tau\} + e^{-x},$$

where

$$\tau \stackrel{\text{def}}{=} \frac{1}{2} \sqrt{(1 + \epsilon)(3\beta + \epsilon z_B(\mathbf{x}))^2 + \delta^2}$$

and the quantile function $z_B(\mathbf{x})$ is defined in (33).

Similar conditions and results can be found in the literature for more specific models. In particular, Bontemps (2011) or Johnstone (2010) explored the Gaussian case; see Section 5.1 below for a more detailed comparison.

3 Infinite dimensional nuisance

This section describes how previous results can be extended to the case where the nuisance is infinite dimensional. More specifically, we consider the $(\boldsymbol{\theta}, \mathbf{f})$ -setup, where $\boldsymbol{\theta} \in \Theta \subseteq \mathbb{R}^q$ and $\mathbf{f} \in \mathcal{H}$ for some Hilbert space \mathcal{H} . Suppose that in \mathcal{H} exists a countable basis $\mathbf{e}_1, \mathbf{e}_2, \dots$. Then

$$\mathbf{f} = \mathbf{f}(\boldsymbol{\phi}) = \sum_{j=1}^{\infty} \phi_j \mathbf{e}_j \in \mathcal{H},$$

where a vector $\boldsymbol{\phi} = \{\phi_j\}_{j=1}^{\infty} \in \ell_2$ and $\phi_j = \langle \mathbf{f}, \mathbf{e}_j \rangle$.

Let the likelihood for the full model be $\mathcal{L}(\boldsymbol{\theta}, \mathbf{f})$. Denote with $\mathbf{v} = (\boldsymbol{\theta}, \boldsymbol{\phi})$

$$\mathcal{L}(\mathbf{v}) = \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}) = \mathcal{L}(\boldsymbol{\theta}, \mathbf{f}(\boldsymbol{\phi})).$$

The underlying full and target parameters can be defined by maximizing the expected log-likelihood:

$$\mathbf{v}^* \stackrel{\text{def}}{=} \underset{\mathbf{v}=(\boldsymbol{\theta}, \boldsymbol{\phi})}{\text{argmax}} \mathbb{E} \mathcal{L}(\mathbf{v}), \quad \boldsymbol{\theta}^* \stackrel{\text{def}}{=} \Pi_0 \mathbf{v}^*. \quad (11)$$

Also define the information matrix \mathcal{D}_0^2 for the full parameter \mathbf{v} and the efficient information matrix $\check{\mathcal{D}}^2$ for the target $\boldsymbol{\theta}$:

$$\begin{aligned} \mathcal{D}_0^2 &\stackrel{\text{def}}{=} \nabla^2 \mathbb{E}[\mathcal{L}(\mathbf{v}^*)] \in \text{Lin}(\ell_2, \ell_2), \\ \check{\mathcal{D}}_0^2 &\stackrel{\text{def}}{=} (\Pi_0 \mathcal{D}_0^{-2} \Pi_0^\top)^{-1} \in \mathbb{R}^{q \times q}, \end{aligned}$$

where $\text{Lin}(\ell_2, \ell_2)$ is the space of linear operators from ℓ_2 to ℓ_2 .

We apply the sieve approach and use an uninformative finite dimensional prior for the parameters $\boldsymbol{\theta}$ and $\boldsymbol{\phi}$. The question under consideration is how the sieve truncation affects the posterior properties.

Let $\boldsymbol{\eta} = \{\eta_j\}_{j=1}^m$ be the sieve projection of the full parameter $\boldsymbol{\phi}$. For notational convenience represent $\boldsymbol{\phi} = (\boldsymbol{\eta}, \boldsymbol{\varkappa})$ and the sieve approximation corresponds to $\boldsymbol{\varkappa} \equiv 0$. Write the “true” point \boldsymbol{v}^* in the form $\boldsymbol{v}^* = (\boldsymbol{\theta}^*, \boldsymbol{\eta}^*, \boldsymbol{\varkappa}^*)$. Approximation of the functional nuisance parameter $\boldsymbol{\phi}$ by the m -dimensional parameter $\boldsymbol{\eta}$ leads to two sources of bias due to projection of the functional parameter onto the finite dimensional space spanned by the first m basis functions. The first one is caused by ignoring the trimmed component $\boldsymbol{\varkappa}$. The “sieved” target $\boldsymbol{\theta}_s^*$ defined as

$$\boldsymbol{\theta}_s^* \stackrel{\text{def}}{=} \operatorname{argmax}_{\boldsymbol{\theta}} \max_{\boldsymbol{\eta}} \mathbb{E}L(\boldsymbol{\theta}, \boldsymbol{\eta}, 0),$$

can be different from the truth $\boldsymbol{\theta}^*$. The other one is caused by change of the efficient Fisher information \check{D}^2 by its sieve counterpart \check{D}_s^2 . The bias terms can be bounded under smoothness assumptions on the model and on the functional nuisance parameter \boldsymbol{f} using the standard methods of approximation theory. To avoid tedious calculus we simply assume a kind of consistency of the sieve approximation. The notation simplifies if we also assume that the basis \boldsymbol{e}_m in the space \mathcal{H} is selected to provide orthogonality of the corresponding block H^2 of the Fisher information matrix, that is, $H^2 = \boldsymbol{I}_f$. Of course, this implies the same structure of its blocks $H_{\boldsymbol{\eta}}^2$ and $H_{\boldsymbol{\varkappa}}^2$. Note that the general situation can be reduced to this orthogonal case by a simple linear transformation of the nuisance parameter $\boldsymbol{\phi}$. This allows to represent the total Fisher matrix in the form

$$\mathcal{D}_0^2 = \begin{pmatrix} D_{\boldsymbol{\theta}}^2 & A_{\boldsymbol{\theta}\boldsymbol{\eta}} & A_{\boldsymbol{\theta}\boldsymbol{\varkappa}} \\ A_{\boldsymbol{\theta}\boldsymbol{\eta}}^\top & \boldsymbol{I}_{\boldsymbol{\eta}} & 0 \\ A_{\boldsymbol{\theta}\boldsymbol{\varkappa}}^\top & 0 & \boldsymbol{I}_{\boldsymbol{\varkappa}} \end{pmatrix}. \tag{12}$$

First we specify the required conditions. The first one is a kind of semi parametric identifiability and it allows to separate the target and the nuisance parameters. Formally it requires that the angle between two tangent subspaces for these parameters is separated away from zero:

(\mathcal{I}^s) For $\nu < 1$, it holds

$$\|D_{\boldsymbol{\theta}}^{-1}A_{\boldsymbol{\theta}\boldsymbol{\eta}}A_{\boldsymbol{\theta}\boldsymbol{\eta}}^\top D_{\boldsymbol{\theta}}^{-1}\| \leq \nu.$$

The smoothness conditions on $\boldsymbol{\theta}$ and $\boldsymbol{\phi}$ are expressed via the component $\boldsymbol{\varkappa}^*$ of \boldsymbol{v}^* and the block $A_{\boldsymbol{\theta}\boldsymbol{\varkappa}}$ of \mathcal{D}_0^2 .

(\mathcal{B}) It holds with $\rho_s, b_s \leq 1/2$

$$\begin{aligned} \|D_{\boldsymbol{\theta}}^{-1}A_{\boldsymbol{\theta}\boldsymbol{\varkappa}}\boldsymbol{\varkappa}^*\| &\leq \rho_s, \\ \|D_{\boldsymbol{\theta}}^{-1}A_{\boldsymbol{\theta}\boldsymbol{\varkappa}}A_{\boldsymbol{\theta}\boldsymbol{\varkappa}}^\top D_{\boldsymbol{\theta}}^{-1}\| &\leq b_s \leq 1/2. \end{aligned}$$

For validity of our results we will need that the value of m is fixed in a proper way ensuring that the values of ρ_s and b_s are sufficiently small. These values can be upper bounded under the usual smoothness conditions on \mathbf{f} , e.g. if \mathbf{f} belongs to a Sobolev ball with a certain regularity; cf. Bontemps (2011); Bickel and Kleijn (2012); Castillo (2012). See also an example of computing the quantities ρ_s and b_s in Section 5.3 below.

Consider a non-informative sieve prior which is uniform on $(\boldsymbol{\theta}, \boldsymbol{\eta})$ sieve component of the full parameter and puts singular mass at point 0 for the remaining components of the nuisance parameter $\{\eta_j\}_{j=m+1}^\infty$. We focus on the posterior distribution of the target parameter. It is assumed that the conditions of Theorem 1 and Corollary 1 are fulfilled for the sieve prior. For the efficient Fisher information matrix \check{D}_s^2 and the vectors $\boldsymbol{\theta}_s^*$ and $\boldsymbol{\theta}_s^\circ$ defined as

$$\begin{aligned}\check{D}_s^2 &\stackrel{\text{def}}{=} D_\theta^2 - A_{\boldsymbol{\theta}\boldsymbol{\eta}}A_{\boldsymbol{\theta}\boldsymbol{\eta}}^\top, \\ \boldsymbol{\theta}_s^\circ &\stackrel{\text{def}}{=} \boldsymbol{\theta}_s^* + \check{D}_s^{-1}\check{\boldsymbol{\xi}}_s = \boldsymbol{\theta}_s^* + \check{D}_s^{-1}(\nabla_\theta - A_{\boldsymbol{\theta}\boldsymbol{\eta}}\nabla_\eta),\end{aligned}$$

Theorem 1 ensures the BvM result for the sieve non-informative prior on $\boldsymbol{\theta}$ and $\boldsymbol{\eta}$: the $\boldsymbol{\theta}$ -posterior is approximated by the Gaussian measure $\mathcal{N}(\boldsymbol{\theta}_s^\circ, \check{D}_s^{-2})$. The main question is whether the side truncation introduces a significant bias in the posterior distribution. For the full semiparametric model, define

$$\begin{aligned}\check{D}^2 &\stackrel{\text{def}}{=} D_\theta^2 - A_{\boldsymbol{\theta}\boldsymbol{\phi}}A_{\boldsymbol{\theta}\boldsymbol{\phi}}^\top, \\ \boldsymbol{\theta}^\circ &\stackrel{\text{def}}{=} \boldsymbol{\theta}^* + \check{D}^{-1}\check{\boldsymbol{\xi}} = \boldsymbol{\theta}^* + \check{D}^{-1}(\nabla_\theta - A_{\boldsymbol{\theta}\boldsymbol{\phi}}\nabla_\phi).\end{aligned}$$

The vector $\boldsymbol{\theta}^\circ$ is the efficient score and \check{D}^2 is the efficient Fisher matrix and they naturally appear in the infinite dimensional Gaussian case as the posterior mean and influence matrix for the improper non-informative prior. The next result accomplishes Theorem 1. Under our identifiability condition (\mathcal{I}^s) and the smoothness condition (B), it allows to measure the distance between the Gaussian measure $\mathcal{N}(\boldsymbol{\theta}_s^\circ, \check{D}_s^{-2})$ which approximates the sieve posterior, with the Gaussian measure $\mathcal{N}(\boldsymbol{\theta}^\circ, \check{D}^{-2})$ corresponding to the full-dimensional prior. Due to Lemma 8 these two measures are close to each other if the ratio of two matrices $\check{D}_s^{-1}\check{D}^2\check{D}_s^{-1}$ is close to identity and the normalized mean difference $\check{D}(\boldsymbol{\theta}^\circ - \boldsymbol{\theta}_s^\circ)$ is small.

Theorem 3. Consider a semiparametric model with a quasi log-likelihood $L(\boldsymbol{\theta}, \boldsymbol{\phi})$. The true value $(\boldsymbol{\theta}^*, \boldsymbol{\phi}^*)$ is given by (11). Let $\mathcal{D}_0^2 \stackrel{\text{def}}{=} -\nabla^2 \mathbb{E}L(\boldsymbol{\theta}^*, \boldsymbol{\phi}^*)$ be the corresponding Fisher operator. Suppose that the nuisance parameter $\boldsymbol{\phi}$ is rescaled to ensure that the corresponding $\boldsymbol{\phi}$ -block of \mathcal{D}_0^2 is identity. Let $(\boldsymbol{\eta}, 0)$ be a sieve approximation of the functional nuisance parameter $\boldsymbol{\phi} = (\boldsymbol{\eta}, \boldsymbol{\varkappa})$, and (12) be the related block representation of \mathcal{D}_0^2 . Suppose the identifiability condition (\mathcal{I}^s) and the smoothness condition (B). Then the efficient Fisher information matrices $\check{D}^2 = D_\theta^2 - A_{\boldsymbol{\theta}\boldsymbol{\phi}}A_{\boldsymbol{\theta}\boldsymbol{\phi}}^\top$ and $\check{D}_s^2 = D_\theta^2 - A_{\boldsymbol{\theta}\boldsymbol{\eta}}A_{\boldsymbol{\theta}\boldsymbol{\eta}}^\top$ in the full and sieve models are related by

$$\begin{aligned}\|\check{D}_s^{-1}\check{D}^2\check{D}_s^{-1} - I_q\| &\leq (1 - \nu)^{-1}\rho_s, \\ \text{tr}\{(\check{D}_s^{-1}\check{D}^2\check{D}_s^{-1} - I_q)^2\} &\leq (1 - \nu)^{-2}q\rho_s^2.\end{aligned}$$

The target parameter $\boldsymbol{\theta}^*$ and its sieve counterpart $\boldsymbol{\theta}_s^*$ are related by

$$\|\check{D}(\boldsymbol{\theta}^* - \boldsymbol{\theta}_s^*)\| \leq (1 - \nu)^{-1} b_s + \delta(\mathbf{r}_s) \mathbf{r}_s, \tag{13}$$

where $\mathbf{r}_s = \|\boldsymbol{\varkappa}^*\|$. Moreover, if $\boldsymbol{\theta}^\circ = \boldsymbol{\theta}^* + \check{D}^{-1} \check{\boldsymbol{\xi}}$ and $\boldsymbol{\theta}_s^\circ = \boldsymbol{\theta}_s^* + \check{D}_s^{-1} \check{\boldsymbol{\xi}}_s$, then it holds with $\boldsymbol{\xi}_\varkappa \stackrel{\text{def}}{=} \rho_s^{-1} D_{\boldsymbol{\theta}^*}^{-1} A_{\boldsymbol{\theta}^*} \nabla_{\boldsymbol{\varkappa}}$

$$\|\check{D}(\boldsymbol{\theta}^\circ - \boldsymbol{\theta}_s^\circ)\| \leq (1 - \nu)^{-1} \rho_s (\|\check{\boldsymbol{\xi}}_s\| + \|\boldsymbol{\xi}_\varkappa\|) + (1 - \nu)^{-1} b_s + \delta(\mathbf{r}_s) \mathbf{r}_s. \tag{14}$$

Finally, under (ED_0) and (\mathcal{I}) , it holds on a set $\Omega(\mathbf{x})$ with $\mathbb{P}(\Omega(\mathbf{x})) \geq 1 - 4e^{\mathbf{x}}$

$$\|\check{D}(\boldsymbol{\theta}^\circ - \boldsymbol{\theta}_s^\circ)\| \leq 2\alpha(1 - \nu)^{-1} \rho_s (q^{1/2} + 2\mathbf{x}) + (1 - \nu)^{-1} b_s + \delta(\mathbf{r}_s) \mathbf{r}_s. \tag{15}$$

We conclude that the sieve prior does a good job if the quantities $q^{1/2} \rho_s$ and b_s are small and $\|\boldsymbol{\varkappa}^*\|$ is not large.

4 The i.i.d. case and critical dimension

This section comments on how the previously obtained general results can be linked to the classical asymptotic results in the statistical literature. The nice feature of the whole approach based on the local bracketing is that all the results are stated under the same list of conditions: once checked one can directly apply any of the mentioned results. Typical examples include i.i.d., generalized linear models (GLM), and median regression models. Here we briefly discuss how the BvM result can be applied to one typical case, namely, to an i.i.d. experiment.

Let $\mathbf{Y} = (Y_1, \dots, Y_n)^\top$ be an i.i.d. sample from a measure P . Here we suppose the conditions of Section 5.1 in Spokoiny (2012) on P and $(P_{\mathbf{v}})$ to be fulfilled. We admit that the parametric assumption $P \in (P_{\mathbf{v}}, \mathbf{v} \in \mathcal{Y})$ can be misspecified and consider the asymptotic setup with the full dimension $p = p_n$ which depends on n and grows to infinity as $p_n \rightarrow \infty$.

Theorem 4. *Suppose the conditions of Section 5.1 in Spokoiny (2012). Let also $p_n \rightarrow \infty$ and $p_n^3/n \rightarrow 0$. Then the result of Theorem 1 holds with $\Delta(\mathbf{r}_0, \mathbf{x}) = C\sqrt{p_n^3/n}$, $\mathcal{D}_0^2 = n\mathbb{F}_{\mathbf{v}^*}$, where $\mathbb{F}_{\mathbf{v}^*}$ is the Fisher information of $(P_{\mathbf{v}})$ at \mathbf{v}^* .*

A similar result about asymptotic normality of the posterior in a linear regression model can be found in Ghosal (1999). However, the convergence is proved under the condition $p_n^4 \log(p_n)/n \rightarrow 0$ which appears to be too strong. Ghosal (2000) showed that the dimensionality constraint can be relaxed to $p_n^3/n \rightarrow 0$ for exponential models with a product structure. Boucheron and Gassiat (2009) proved the BvM result in a specific class of i.i.d. model with discrete probability distribution under the condition $p_n^3/n \rightarrow 0$. Further examples and the related conditions for Gaussian models are presented in Johnstone (2010).

4.1 Critical dimension

This section discusses the issue of a *critical dimension*. Namely we show that the condition $p_n = o(n^{1/3})$ in Theorem 4 for the validity of the BvM result cannot be dropped or relaxed in a general situation. Namely, we present an example for which $p_n^3/n \geq \beta^2 > 0$ and the posterior distribution does not concentrate around MLE.

Let n and p_n be such that $M_n = n/p_n$ is an integer. We consider a simple Poissonian model with $Y_i \sim \text{Poisson}(v_j)$ for $i \in \mathcal{I}_j$, where $\mathcal{I}_j \stackrel{\text{def}}{=} \{i : \lceil i/M_n \rceil = j\}$ for $j = 1, \dots, p_n$ and $\lceil x \rceil$ is the nearest integer greater or equal to x . Let also $u_j = \log v_j$ be the canonical parameter. The log-likelihood $\mathcal{L}(\mathbf{u})$ with $\mathbf{u} = (u_1, \dots, u_{p_n})$ reads as

$$\mathcal{L}(\mathbf{u}) = \sum_{j=1}^{p_n} (Z_j u_j - M_n e^{u_j}),$$

where

$$Z_j \stackrel{\text{def}}{=} \sum_{i \in \mathcal{I}_j} Y_i.$$

We consider the problem of estimating the mean of the u_j 's:

$$\theta = \frac{1}{p_n} (u_1 + \dots + u_{p_n}).$$

Below we study this problem in the asymptotic setup with $p_n \rightarrow \infty$ as $n \rightarrow \infty$ when the underlying measure \mathbb{P} corresponds to $u_1^* = \dots = u_{p_n}^* = u^*$ for some u^* yielding $\theta^* = u^*$. The value u^* will be specified later. We consider an i.i.d. exponential prior on the parameters v_j of Poisson distribution:

$$v_j \sim \text{Exp}(\mu).$$

Below we allow that μ may depend on n . Our results are valid for $\mu \leq \mathfrak{C} \sqrt{\frac{n}{\log n}}$. The posterior is Gamma distributed:

$$v_j \mid \mathbf{Y} \sim \text{Gamma}(\alpha_j, \mu_j),$$

where $\alpha_j = 1 + \sum_{i \in \mathcal{I}_j} Y_i$, $\mu_j = \frac{\mu}{M_n \mu + 1}$.

First we describe the profile maximum likelihood estimator $\tilde{\theta}_n$ of the target parameter θ . The MLE for the full parameter \mathbf{v} reads as $\tilde{\mathbf{v}} = (\tilde{v}_1, \dots, \tilde{v}_{p_n})^\top$ with

$$\tilde{v}_j = Z_j / M_n.$$

Thus, the profile MLE $\tilde{\theta}_n$ reads as

$$\tilde{\theta}_n = \frac{1}{p_n} \sum_{j=1}^{p_n} \log(\tilde{v}_j).$$

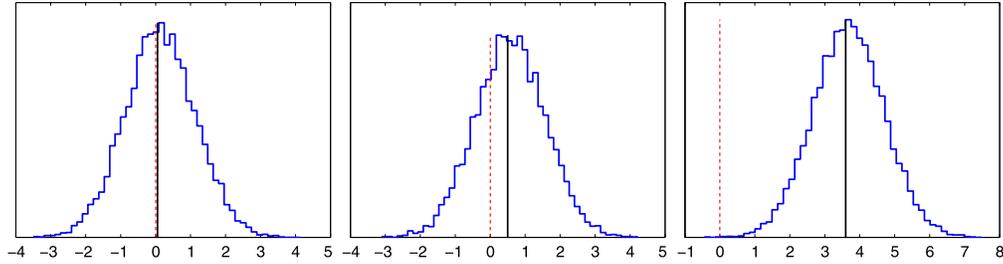


Figure 1: Posterior distribution of $\beta_n^{-1} p_n (\theta - \tilde{\theta}_n)$ for $\beta_n = 1/\log(p_n)$, $\beta_n = 1$, and $\beta_n = \log(p_n)$. Solid line is for posterior mean and dashed line is for true mean.

Furthermore, the efficient Fisher information \check{D}_0^2 is equal to $p_n^{-1} n$; see Lemma 11 below. As $\tilde{\theta}_n$ is the profile MLE it is efficient with the asymptotic variance equal to \check{D}_0^{-2} .

Theorem 5. *Let $Y_i \sim \text{Poisson}(v^*)$ for all $i = 1, \dots, n$, $v^* = 1/p_n$. Then*

1. *If $p_n^3/n \rightarrow 0$ as $p_n \rightarrow \infty$, then*

$$p_n^{1/2} n^{1/2} (\theta - \tilde{\theta}_n) \mid \mathbf{Y} \xrightarrow{w} \mathcal{N}(0, 1).$$

2. *Let $p_n^3/n \equiv \beta > 0$. Then*

$$p_n^{1/2} n^{1/2} (\theta - \tilde{\theta}_n) \mid \mathbf{Y} \xrightarrow{w} \mathcal{N}(\beta/2, 1).$$

3. *If $p_n^3/n \rightarrow \infty$, but $p_n^4/n^{3/2} \rightarrow 0$, then*

$$p_n^{1/2} n^{1/2} (\theta - \tilde{\theta}_n) \mid \mathbf{Y} \xrightarrow{w} \infty.$$

We carried out a series of experiments to numerically demonstrate the results of Theorem 5. The dimension of parameter space was fixed $p_n = 10000$. Three cases were considered:

1. $p_n^{3/2}/n^{1/2} = \frac{1}{\log p_n}$, which corresponds to $p_n^3/n \rightarrow 0, n \rightarrow \infty$.
2. $p_n^{3/2}/n^{1/2} \equiv 1$.
3. $p_n^{3/2}/n^{1/2} = \log p_n$, which corresponds to $p_n^3/n \rightarrow \infty, n \rightarrow \infty$.

For each sample 10000 realizations of \mathbf{Y} were generated from the exponential distribution $\text{Exp}(v_*)$ and so were corresponding posterior values $\theta \mid \mathbf{Y}$. The resulting posterior distribution for three cases is demonstrated in Figure 1. It can be easily seen that the results of Theorem 5 are numerically confirmed.

5 Examples

This section presents a number of examples illustrating the general results of Section 2.

5.1 Linear Gaussian regression and a flat Gaussian prior

Let $\mathbf{Y} = (Y_1, \dots, Y_n)^\top$ be a random vector in \mathbb{R}^n following the equation

$$\mathbf{Y} = \mathbf{f} + \boldsymbol{\varepsilon} = \boldsymbol{\Psi}^\top \mathbf{v}^* + \boldsymbol{\varepsilon}, \quad (16)$$

where the errors $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)^\top$ are independent zero mean, and $\boldsymbol{\Psi}$ is a given $p \times n$ design matrix. The mean vector $\mathbf{f} \in \mathbb{R}^n$ is unknown and the second equation in (16) means that it belongs to some given p -dimensional linear subspace in \mathbb{R}^n : $\mathbf{f} = \boldsymbol{\Psi}^\top \mathbf{v}^*$ for an unknown target vector $\mathbf{v}^* \in \mathbb{R}^p$. We write the matrix $\boldsymbol{\Psi}$ in the form $\boldsymbol{\Psi} = \{\boldsymbol{\Psi}_1, \dots, \boldsymbol{\Psi}_n\}$, so that $f_i = \boldsymbol{\Psi}_i^\top \mathbf{v}^*$. Below we suppose that $n > p$ and the rank of $\boldsymbol{\Psi}$ is p , or, equivalently, the rows of $\boldsymbol{\Psi}$ are linearly independent vectors in \mathbb{R}^n .

First, we consider the Gaussian case, i.e. $\varepsilon_i \in \mathcal{N}(0, \sigma_n^2 \mathbf{I}_n)$, $i = 1, \dots, n$, with \mathbf{I}_n being the $n \times n$ identity matrix. The variance of observations σ_n^2 is known but may depend on sample size n . For ease of comparison, we assume that the design matrix $\boldsymbol{\Psi}$ fulfills $\boldsymbol{\Psi}^\top \boldsymbol{\Psi} = \mathbf{I}_p$. In our notation, this implies $\mathcal{D}_0^2 = \sigma_n^{-2} \mathbf{I}_p$.

For a Gaussian prior, the posterior is exactly Gaussian and the Fisher and Wilks approximations are also exact. If a non-informative prior is used then the only necessary condition for validity of the BvM result is $p = p_n = o(n)$; see Bontemps (2011). Johnstone (2010) showed that the BvM result can be extended to the situation with a growing parameter dimension $p = p_n$ and a flat Gaussian prior with a covariance matrix $\tau_n^2 \mathbf{I}_n$ under the condition “ $(\sigma_n/\tau_n)^4 p_n$ is small” and “ $(\sigma_n/\tau_n^2) \|\mathbf{v}^*\|$ is small”. Our results from Section 2 cover the case of a flat Gaussian prior and the only conditions to check for the BvM Theorem are that “ $\text{tr}(\mathcal{D}_0^{-2} G^2) = (\sigma_n/\tau_n)^4 p_n$ is small” and “ $\|\mathcal{D}_G^{-1} G^2 \mathbf{v}^*\| = (\sigma_n/\tau_n^2) \|\mathbf{v}^*\|$ is small”. One can see that our general results well apply to the case of growing dimension and are as sharp as the existing results obtained for a very special Gaussian case.

5.2 Linear non-Gaussian regression

Now consider a more general situation, when errors ε_i in (16) follow a general distribution with a given density $f(\cdot)$: $\varepsilon_i \sim P_f$. Denote $h(x) = \log f(x)$. The log-likelihood function for this problem reads as

$$\mathcal{L}(\mathbf{v}) = \sum_{i=1}^n \log f(Y_i - \boldsymbol{\Psi}_i^\top \mathbf{v}) = \sum_{i=1}^n h(Y_i - \boldsymbol{\Psi}_i^\top \mathbf{v}). \quad (17)$$

Suppose that $h(z)$ is twice continuously differentiable and let

$$h^2 \stackrel{\text{def}}{=} \int h''(z) f(z) dz < \infty.$$

If the model is correctly specified, then

$$\mathcal{D}_0^2 = -\nabla^2 \mathbb{E} \sum_{i=1}^n h(Y_i - \Psi_i^\top \mathbf{v}^*) = \int h''(z) f(z) dz \cdot \sum_{i=1}^n \Psi_i \Psi_i^\top = h^2 \sum_{i=1}^n \Psi_i \Psi_i^\top. \quad (18)$$

Similarly,

$$\mathcal{D}^2(\mathbf{v}) = -\nabla^2 \mathbb{E} \sum_{i=1}^n h(Y_i - \Psi_i^\top \mathbf{v}) = \sum_{i=1}^n \Psi_i \Psi_i^\top \int h''(z - \Psi_i^\top (\mathbf{v} - \mathbf{v}^*)) f(z) dz.$$

Conditions from Section 2.1 need to be checked in order to apply general results from Section 2.2. We assume some exponential moment conditions on distribution P_f :

(e₀) *There exist some constants ν_0 and $\mathbf{g}_1 > 0$ such that for $\varepsilon \sim P_f$ it holds*

$$\log \mathbb{E} \exp(\mu h'(\varepsilon)/h) \leq \nu_0^2 \mu^2 / 2, \quad |\mu| \leq \mathbf{g}_1.$$

Condition (e₀) effectively means that the distribution of error ε has exponentially decreasing tail. Under e₀ condition (ED₀) is satisfied due to the following lemma.

Lemma 1. *Assume (e₀) and let $\mathcal{V}_0^2 = \mathcal{D}_0^2 = h^2 \sum_{i=1}^n \Psi_i \Psi_i^\top$. Then condition (ED₀) follows from (e₀) with this \mathcal{V}_0^2 and $\mathbf{g} = \mathbf{g}_1 N_1^{1/2}$, where*

$$N_1^{-1/2} \stackrel{\text{def}}{=} \max_i \sup_{\gamma \in \mathbb{R}^p} \frac{h |\Psi_i^\top \gamma|}{\|\mathcal{D}_0 \gamma\|}. \quad (19)$$

(\mathcal{L}_0) and (\mathcal{L}_r) are also valid under the Lipschitz continuity assumption on $h''(\cdot)$.

Lemma 2. *Let*

$$|h''(z) - h''(z_0)| \leq L |z - z_0|, \quad z, z_0 \in \mathbb{R}.$$

Then

$$\delta(\mathbf{r}) \leq \frac{L \mathbf{r}}{h N_1^{1/2}},$$

where N_1 is defined by (19).

The next interesting question is checking the condition (ED₂). The stochastic part of the likelihood reads as follows:

$$\zeta(\mathbf{v}) = \sum_{i=1}^n h(Y_i - \Psi_i^\top \mathbf{v}) - \mathbb{E} h(Y_i - \Psi_i^\top \mathbf{v}).$$

To check the condition (ED_2) we need to compute the Hessian of the $\zeta(\mathbf{v})$:

$$\nabla^2 \zeta(\mathbf{v}) = \sum_{i=1}^n \{h''(Y_i - \Psi_i^\top \mathbf{v}) - \mathbb{E}h''(Y_i - \Psi_i^\top \mathbf{v})\} \Psi_i \Psi_i^\top.$$

Finally, we need to impose the mild condition on marginal likelihood:

(e₂) *There exists some constant ν_0 and for every $\mathbf{r} > 0$ there exists $\mathbf{g}(\mathbf{r}) > 0$, such that for all numbers δ with $|\delta| \leq N_2^{-1/2} \mathbf{r}$*

$$\log \mathbb{E} \exp\left(\frac{\mu}{\mathfrak{s}_i} \{h''(Y_i + \delta) - \mathbb{E}h''(Y_i + \delta)\}\right) \leq \frac{\nu_0^2 \mu^2}{2}, \quad |\mu| \leq \mathbf{g}(\mathbf{r}),$$

where \mathfrak{s}_i are some known values and

$$N_2^{-1/2} \stackrel{\text{def}}{=} \max_i \sup_{\gamma \in \mathbb{R}^p} \frac{\mathfrak{s}_i |\Psi_i^\top \gamma|}{\|\mathcal{D}_0 \gamma\|}.$$

Then

$$\begin{aligned} & \sup_{\gamma_1, \gamma_2 \in \mathbb{R}^p} \log \mathbb{E} \exp \left\{ \frac{\mu}{\omega} \frac{\gamma_1^\top \nabla^2 \zeta(\mathbf{v}) \gamma_2}{\|\mathcal{D}_0 \gamma_1\| \cdot \|\mathcal{D}_0 \gamma_2\|} \right\} \\ &= \sup_{\gamma_1, \gamma_2 \in \mathbb{R}^p} \log \mathbb{E} \exp \left\{ \frac{\mu}{\omega} \frac{\sum_{i=1}^n (h''(Y_i - \Psi_i^\top \mathbf{v}) - \mathbb{E}h''(Y_i - \Psi_i^\top \mathbf{v})) \gamma_1^\top \Psi_i \Psi_i^\top \gamma_2}{\|\mathcal{D}_0 \gamma_1\| \cdot \|\mathcal{D}_0 \gamma_2\|} \right\} \\ &\leq \sup_{\gamma_1, \gamma_2 \in \mathbb{R}^p} \frac{\nu_0^2 \mu^2}{2\omega^2} \frac{\sum_{i=1}^n \mathfrak{s}_i^2 |\Psi_i^\top \gamma_1|^2 \cdot |\Psi_i^\top \gamma_2|^2}{\|\mathcal{D}_0 \gamma_1\|^2 \cdot \|\mathcal{D}_0 \gamma_2\|^2} \\ &\leq \frac{\nu_0^2 \mu^2}{2\omega^2} \frac{n}{N_2^2}. \end{aligned}$$

It easily follows that $\omega = \frac{\sqrt{n}}{N_2}$ does the job, and (ED_2) follows. In regular situations N_2 is of order of sample size n and then $\omega \sim n^{-1/2}$.

Finally, turn to the semiparametric setup with $\mathbf{v} = (\boldsymbol{\theta}, \boldsymbol{\eta}) \in \mathbb{R}^p$. Assume the identifiability condition (\mathcal{I}) from Section 2.1. Then our general results yield the following semiparametric BvM Theorem for the linear model (16).

Theorem 6. *Let (e_0) , (e_2) , conditions of Lemma 2 and condition (\mathcal{I}) for matrix \mathcal{D}_0^2 from (18) hold. Then the results of Theorem 1 hold for linear model (16) with $\mathbf{r}_0^2 \geq \mathbb{C}(p + \mathbf{x})$ and*

$$\Delta(\mathbf{r}_0, \mathbf{x}) = \{\delta(\mathbf{r}_0) + 6\nu_0 z_{\mathbb{H}}(\mathbf{x}) \omega\} \mathbf{r}_0^2 \leq \left\{ \frac{L}{h^2} \frac{\mathbf{r}_0}{N_1^{1/2}} + 6\nu_0 z_{\mathbb{H}}(\mathbf{x}) \frac{\sqrt{n}}{N_2} \right\} \mathbf{r}_0^2.$$

5.3 Semiparametric non-Gaussian linear regression

This section specifies the obtained results for the linear non-Gaussian model (16) with a semiparametric regression function

$$\mathbf{f}^* = \Psi^\top \boldsymbol{\theta}^* + \mathbf{g}^*, \tag{20}$$

where $\boldsymbol{\theta}^* \in \mathbb{R}^q$ is an unknown target vector and $\Psi = (\Psi_1, \dots, \Psi_n)$ is a $q \times n$ matrix with $\Psi_i = (\psi_1(X_i), \dots, \psi_q(X_i))^\top \in \mathbb{R}^q$ for given basis functions $\{\psi_j(\cdot), j = 1, \dots, q\}$ and the design points $X_i, i = 1, \dots, n$. Without loss of generality we can assume that these basis functions are design orthonormal:

$$\sum_i \psi_j(X_i)\psi_{j'}(X_i) = \delta_{j,j'}. \tag{21}$$

The general case can be reduced to this one by usual rotation and rescaling. Similarly, we suppose that the entries of the nuisance vector $\mathbf{g}^* = \{g^*(X_1), \dots, g^*(X_n)\}^\top$ are the values at the design points X_i of a function g^* which is an element of a functional space. This means that $g^* = g(x) = \sum_{k=1}^\infty \eta_k \varphi_k(x)$ for a given functional basis $\{\varphi_k\}_{k=1}^\infty$ (e.g. Fourier, wavelet, etc.) and an infinite dimensional nuisance parameter vector $\boldsymbol{\eta} = \{\eta_1, \eta_2, \dots, \eta_m, \dots\}^\top$.

Furthermore, we assume that g^* is smooth, that is, it can be well approximated by finite sums $g_m^*(\cdot) = \sum_{k=1}^m \eta_k \varphi_k(\cdot)$ in the sense that

$$\|g^* - g_m^*\| \leq \gamma_m. \tag{22}$$

For instance, if \mathcal{F}_s is a Sobolev ball, that is,

$$g^* \in \mathcal{F}_s(\mathbb{C}) \stackrel{\text{def}}{=} \left\{ g(x) = \sum_{k=1}^\infty \eta_k \varphi_k(x) : \sum_{k=1}^\infty \eta_k^2 k^{2s} \leq \mathbb{C} \right\},$$

then $\gamma_m \leq (m + 1)^{-s}$.

In addition we suppose the same smoothness condition for each basis function ψ_j for $j = 1, \dots, q$, that is,

$$\|\psi_j - \psi_{j,m}\| \leq \gamma_m, \tag{23}$$

where $\varphi_{j,m} = \Pi_m \varphi_j$ and Π_m is the projector on the spaced spanned by the first m basis functions $\varphi_1, \dots, \varphi_m$.

To avoid the identifiability problem, we restrict the expansion of the function g to the first M coefficients for a large number M which may depend on the sample size n . For instance, one can take $M = n/\log(n)$. Alternatively, $M = n^a$ for some $a < 1$. Also, to simplify the presentation, we suppose that basis functions $\varphi_k(\cdot)$ are orthonormal in the sense that

$$\sum_i \varphi_{k'}(X_i)\varphi_k(X_i) = \delta_{k',k}. \tag{24}$$

Define also $\Phi_i = \{\varphi_1(X_i), \dots, \varphi_q(X_i)\}^\top \in \mathbb{R}^q$ for $i = 1, 2, \dots, n$. Then the full parameter of the model is $\mathbf{v} = (\boldsymbol{\theta}, \boldsymbol{\eta}) \in \mathbb{R}^p$ for $p = q + M$, and the decomposition (20) can be rewritten as

$$\mathbf{f} = \Xi \mathbf{v},$$

where $\Xi_i = (\Psi_i^\top, \Phi_i^\top)^\top$, $i = 1, \dots, n$, and $\Xi = (\Xi_1, \dots, \Xi_n)^\top$. The full Fisher information matrix reads as

$$\mathcal{D}_0^2 = \begin{pmatrix} D_0^2 & A_0 \\ A_0^\top & H_0^2 \end{pmatrix} = h^2 \begin{pmatrix} \sum_{i=1}^n \Psi_i \Psi_i^\top & \sum_{i=1}^n \Psi_i \Phi_i^\top \\ \sum_{i=1}^n \Phi_i \Psi_i^\top & \sum_{i=1}^n \Phi_i \Phi_i^\top \end{pmatrix}.$$

Due to the orthogonality conditions (21) and (24), the blocks D_0^2 and H_0^2 are proportional to identity: $D_0^2 = h^2 I_q$, $H_0^2 = h^2 I_M$. In what follows, we assume $h = 1$, the extension to the general case is trivial. The identifiability condition (9) can be written as

$$\|A_0 A_0^\top\| \leq \nu. \quad (25)$$

The condition (22) implies by orthogonality of the basis functions φ_m

$$\|A_{\boldsymbol{\theta} \times \boldsymbol{x}^*}\|^2 = \|g^* - g_m^*\|^2 \leq \gamma_m^2.$$

Similarly, the smoothness condition (23) implies that each row of $A_{\boldsymbol{\theta} \times \boldsymbol{x}}$ is bounded in norm by γ_m . Therefore,

$$\|A_{\boldsymbol{\theta} \times \boldsymbol{x}} A_{\boldsymbol{\theta} \times \boldsymbol{x}}^\top\| \leq \text{tr}(A_{\boldsymbol{\theta} \times \boldsymbol{x}} A_{\boldsymbol{\theta} \times \boldsymbol{x}}^\top) \leq q \gamma_m^2.$$

Theorem 7. Consider the model (16) with the log-likelihood (17) and the semi parametric regression function from (20). Suppose the indentifiability condition (25) and the smoothness conditions (22), (23). Then the result of Theorem 3 holds with $\rho_s = b_s = (1 - \nu)^{-1} q^{1/2} \gamma_m^2$.

Remark 1. In the case where all the functions $g^*, \psi_j, j = 1, \dots, q$ are from the Sobolev ball with smoothness s , we can conclude that $\gamma_m \leq (m + 1)^{-s}$ and $\rho_s = b_s \leq (1 - \nu)^{-1} q^{1/2} (m + 1)^{-2s}$.

5.4 Generalized linear modeling

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

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

$$Y_i \sim P_{\Psi_i^\top \mathbf{v}^*}. \tag{26}$$

The corresponding log-density of a GLM reads as

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

Under $\mathbb{P}_{\boldsymbol{\theta}^*}$ each observation Y_i follows (26), in particular, $\mathbb{E}Y_i = d'(\Psi_i^\top \mathbf{v}^*)$. However, similarly to the previous sections, it is accepted that the parametric model (26) is misspecified. Response misspecification means that the vector $\mathbf{f} \stackrel{\text{def}}{=} \mathbb{E}\mathbf{Y}$ cannot be represented in the form $d'(\Psi^\top \mathbf{v})$ whatever \mathbf{v} is. The other sort of misspecification concerns the data distribution. The model (26) assumes that the Y_i 's are independent and the marginal distribution belongs to the given parametric family P . In what follows, we only assume independent data having certain exponential moments. The target of estimation \mathbf{v}^* is defined by

$$\mathbf{v}^* \stackrel{\text{def}}{=} \underset{\mathbf{v}}{\operatorname{argmax}} \mathbb{E}\mathcal{L}(\mathbf{v}).$$

The quasi MLE $\tilde{\mathbf{v}}$ is defined by maximization of $\mathcal{L}(\mathbf{v})$:

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

Convexity of $d(\cdot)$ implies that $\mathcal{L}(\mathbf{v})$ is a concave function of \mathbf{v} , so that the optimization problem has a unique solution and can be effectively solved. However, a closed form solution is only available for the constant regression or for the linear Gaussian regression. The corresponding target \mathbf{v}^* is the maximizer of the expected log-likelihood:

$$\mathbf{v}^* = \underset{\mathbf{v}}{\operatorname{argmax}} \mathbb{E}\mathcal{L}(\mathbf{v}) = \underset{\mathbf{v}}{\operatorname{argmax}} \sum \{f_i \Psi_i^\top \mathbf{v} - d(\Psi_i^\top \mathbf{v})\}$$

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

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

(e0) *There exist some constants ν_0 and $\mathbf{g}_1 > 0$, and for every i a constant \mathfrak{s}_i such that $\mathbb{E}(\varepsilon_i/\mathfrak{s}_i)^2 \leq 1$ and*

$$\log \mathbb{E} \exp(\mu \varepsilon_i/\mathfrak{s}_i) \leq \nu_0^2 \mu^2/2, \quad |\mu| \leq \mathbf{g}_1. \tag{27}$$

A natural candidate for \mathfrak{s}_i is σ_i where $\sigma_i^2 = \mathbb{E}\varepsilon_i^2$ is the variance of ε_i . Under (27), introduce a $q \times q$ matrix \mathcal{V}_0 defined by

$$\mathcal{V}_0^2 \stackrel{\text{def}}{=} \sum \mathfrak{s}_i^2 \Psi_i \Psi_i^\top. \quad (28)$$

Condition (e_0) effectively means that each error term $\varepsilon_i = Y_i - \mathbb{E}Y_i$ has some bounded exponential moments: for $|\lambda| \leq \mathbf{g}_1$, it holds $f(\lambda) \stackrel{\text{def}}{=} \log \mathbb{E} \exp(\lambda \varepsilon_i / \mathfrak{s}_i) < \infty$. In words, condition (e_0) requires a light (exponentially decreasing) tail for the marginal distribution of each ε_i .

Define also

$$N_1^{-1/2} \stackrel{\text{def}}{=} \max_i \sup_{\gamma \in \mathbb{R}^p} \frac{\mathfrak{s}_i |\Psi_i^\top \gamma|}{\|\mathcal{V}_0 \gamma\|}. \quad (29)$$

Now conditions are satisfied due to following lemma, see Spokoiny (2012) for proof.

Lemma 3. Assume (e_0) and let \mathcal{V}_0^2 be defined by (28) and N_1 by (29). Then condition (ED_0) follows from (e_0) with this \mathcal{V}_0^2 and $\mathbf{g} = \mathbf{g}_1 N_1^{1/2}$. Moreover, the stochastic component $\zeta(\mathbf{v})$ is linear in \mathbf{v} and the condition (ED_2) is fulfilled with $\omega(\mathbf{r}) \equiv 0$.

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

Lemma 4. It holds

$$-\mathbb{E}L(\mathbf{v}, \mathbf{v}^*) = \sum \{d(\Psi_i^\top \mathbf{v}) - d(\Psi_i^\top \mathbf{v}^*) - d'(\Psi_i^\top \mathbf{v}^*) \Psi_i^\top (\mathbf{v} - \mathbf{v}^*)\} = \mathcal{K}(\mathbb{P}_{\mathbf{v}^*}, \mathbb{P}_{\mathbf{v}}),$$

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

$$-\mathbb{E}L(\mathbf{v}, \mathbf{v}^*) = \|\mathcal{D}(\mathbf{v}^\circ)(\mathbf{v} - \mathbf{v}^*)\|^2/2,$$

where $\mathbf{v}^\circ \in [\mathbf{v}^*, \mathbf{v}]$ and

$$\mathcal{D}^2(\mathbf{v}^\circ) = \sum d''(\Psi_i^\top \mathbf{v}^\circ) \Psi_i \Psi_i^\top.$$

The proof of this lemma can also be found in Spokoiny (2012). Define now the matrix \mathcal{D}_0^2 by

$$\mathcal{D}_0^2 \stackrel{\text{def}}{=} \mathcal{D}^2(\mathbf{v}^*) = \sum d''(\Psi_i^\top \mathbf{v}^*) \Psi_i \Psi_i^\top.$$

Let also \mathcal{V}_0^2 be defined by (28). Note that the matrices \mathcal{D}_0^2 and \mathcal{V}_0^2 coincide if the model $Y_i \sim P_{\Psi_i^\top \mathbf{v}^*}$ is correctly specified and $\mathfrak{s}_i^2 = d''(\Psi_i^\top \mathbf{v}^*)$. The matrix \mathcal{D}_0^2 describes a local

elliptic neighborhood of the central point \mathbf{v}^* in the form $\mathcal{Y}_0(\mathbf{r}) = \{\mathbf{v} : \|\mathcal{D}_0(\mathbf{v} - \mathbf{v}^*)\| \leq \mathbf{r}\}$. If the matrix function $\mathcal{D}^2(\mathbf{v})$ is continuous in this vicinity $\mathcal{Y}_0(\mathbf{r})$ then the value $\delta(\mathbf{r})$ measuring the approximation quality of $-\mathbb{E}L(\mathbf{v}, \mathbf{v}^*)$ by the quadratic function $\|\mathcal{D}_0(\mathbf{v} - \mathbf{v}^*)\|^2/2$ is small and the identifiability condition (\mathcal{L}_0) is fulfilled on $\mathcal{Y}_0(\mathbf{r})$. The following lemma gives bounds for $\delta(\mathbf{r})$.

Lemma 5. *Let $d''(z)$ be Lipschitz continuous:*

$$|d''(z) - d''(z_0)| \leq L|z - z_0|, \quad z, z_0 \in \mathbb{R}$$

Then

$$\delta(\mathbf{r}) \leq L \frac{\mathbf{r}}{N_2^{1/2}},$$

where

$$N_2^{-1/2} \stackrel{\text{def}}{=} \max_i \sup_{\boldsymbol{\gamma} \in \mathbb{R}^p} \frac{|\Psi_i^\top \boldsymbol{\gamma}|}{d''(\Psi_i^\top \mathbf{v}^*) \cdot \|\mathcal{D}_0 \boldsymbol{\gamma}\|}.$$

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

Theorem 8. *Let (e_0) and conditions of Lemma 5 hold. Then the results of Theorem 1 hold for GLM with*

$$\Delta(\mathbf{r}_0, \mathbf{x}) \leq L \frac{\mathbf{r}_0^3}{N_2^{1/2}}.$$

If the function $d(w)$ is quadratic then the approximation error δ vanishes as well and then quadratic approximation is valid globally, a localization step is not required. However, if $d(w)$ is not quadratic, the result applies only locally and it has to be accomplished with a large deviation bound. The GLM structure is helpful in the large deviation zone as well. Indeed, the identifiability condition $(\mathcal{L}\mathbf{r})$ easily follows from Lemma 4: it suffices to bound from below the matrix $\mathcal{D}(\mathbf{v})$ for $\mathbf{v} \in \mathcal{Y}_0(\mathbf{r})$:

$$\mathcal{D}(\mathbf{v}) \geq \mathbf{b}(\mathbf{r})\mathcal{D}_0, \quad \mathbf{v} \in \mathcal{Y}_0(\mathbf{r}).$$

An interesting question, similarly to the i.i.d. case, is the minimal radius \mathbf{r}_0 of the local vicinity $\mathcal{Y}_0(\mathbf{r}_0)$ ensuring the desirable concentration property. The required value conditions are fulfilled for $\mathbf{r}^2 \geq \mathbf{r}_0^2 = \mathbf{C}(\mathbf{x} + p)$, where \mathbf{C} only depends on ν_0, \mathbf{b} , and \mathbf{g} . Thus, the results are valid if

$$\delta(\mathbf{r}_0)\mathbf{r}_0^2 = \mathbf{C} \frac{\mathbf{r}_0^3}{N_1^{1/2}} = \mathbf{C} \frac{(\mathbf{x} + p)^{3/2}}{N_1^{1/2}}$$

is small.

The GLM model also allows a semiparametric extension, i.e.

$$w_i = \Psi_i^\top \boldsymbol{\theta}^* + g^*(X_i),$$

where $g^*(\cdot)$ is from a Sobolev class. This setup differs from Section 5.3 only by few technicalities and leads to similar theoretical results.

6 Supplementary

This section contains the imposed conditions and some supplementary statements which are of some interest by themselves.

6.1 Bracketing and upper function devices

This section briefly overviews the main constructions of Spokoiny (2012) including the bracketing bound and the upper function results. The bracketing bound describes the quality of quadratic approximation of the log-likelihood process $\mathcal{L}(\mathbf{v})$ in a local vicinity of the point \mathbf{v}^* , while the upper function method is used to show that the full MLE $\tilde{\mathbf{v}}$ belongs to this vicinity with a dominating probability. Given $\mathbf{r} > 0$, define the local set

$$\mathcal{Y}_0(\mathbf{r}) \stackrel{\text{def}}{=} \{\mathbf{v} : (\mathbf{v} - \mathbf{v}^*)^\top \mathcal{D}_0^2(\mathbf{v} - \mathbf{v}^*) \leq \mathbf{r}^2\}. \quad (30)$$

Define the quadratic processes $\mathbb{L}(\mathbf{v}, \mathbf{v}^*)$:

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

The next result states the local bracketing bound. The formulation assumes that some value \mathbf{x} is fixed such that $e^{-\mathbf{x}}$ is sufficiently small. If the dimension p is large, one can select $\mathbf{x} = C \log(p)$. We assume that a value $\mathbf{r} = \mathbf{r}_0$ is fixed which separates the local and global zones.

Theorem 9. *Suppose the conditions (ED_0) , (ED_2) , (\mathcal{L}_0) , and (\mathcal{I}) from Section 2.1 hold for some $\mathbf{r}_0 > 0$. Then on a random set $\Omega_{\mathbf{r}_0}(\mathbf{x})$ of dominating probability at least $1 - e^{-\mathbf{x}}$*

$$|L(\mathbf{v}, \mathbf{v}^*) - \mathbb{L}(\mathbf{v}, \mathbf{v}^*)| \leq \Delta(\mathbf{r}_0, \mathbf{x}), \quad \mathbf{v} \in \mathcal{Y}_0(\mathbf{r}_0), \quad (31)$$

where

$$\begin{aligned} \Delta(\mathbf{r}_0, \mathbf{x}) &\stackrel{\text{def}}{=} \{\delta(\mathbf{r}_0) + 6\nu_0 z_{\mathbb{H}}(\mathbf{x})\omega\} \mathbf{r}_0^2, \\ z_{\mathbb{H}}(\mathbf{x}) &\stackrel{\text{def}}{=} 2p^{1/2} + \sqrt{2\mathbf{x}} + \mathbf{g}^{-1}(\mathbf{g}^{-2}\mathbf{x} + 1)4p, \end{aligned} \quad (32)$$

and $\mathcal{Y}_0(\mathbf{r}_0)$ is defined in (30). Moreover, the random vector $\boldsymbol{\xi} = \mathcal{D}_0^{-1} \nabla \mathcal{L}(\mathbf{v}^*)$ fulfills on a random set $\Omega_B(\mathbf{x})$ of dominating probability at least $1 - 2e^{-\mathbf{x}}$

$$\|\boldsymbol{\xi}\|^2 \leq z_B^2(\mathbf{x}), \quad (33)$$

where $z_B^2(\mathbf{x}) \stackrel{\text{def}}{=} \mathbf{p}_B + 6\lambda_B \mathbf{x}$,

$$B \stackrel{\text{def}}{=} \mathcal{D}_0^{-1} \mathcal{V}_0^2 \mathcal{D}_0^{-1}, \quad \mathbf{p}_B \stackrel{\text{def}}{=} \text{tr}(B), \quad \lambda_B \stackrel{\text{def}}{=} \lambda_{\max}(B).$$

Furthermore, assume $(\mathcal{L}\mathbf{r})$ with $\mathbf{b}(\mathbf{r}) \equiv \mathbf{b}$ yielding

$$-EL(\mathbf{v}, \mathbf{v}^*) \geq \mathbf{b} \|\mathcal{D}_0(\mathbf{v} - \mathbf{v}^*)\|^2$$

for each $\mathbf{v} \in \mathcal{Y} \setminus \mathcal{Y}_0(\mathbf{r}_0)$. Let also

$$\mathbf{r} \geq \frac{2}{\mathbf{b}} \left\{ z_B(\mathbf{x}) + 6\nu_0 z_{\mathbb{H}}(\mathbf{x} + \log(2\mathbf{r}/\mathbf{r}_0)) \omega \right\}, \quad \mathbf{r} \geq \mathbf{r}_0$$

with $z_{\mathbb{H}}(\mathbf{x})$ from (32). Then,

$$L(\mathbf{v}, \mathbf{v}^*) \leq -\mathbf{b} \|\mathcal{D}_0(\mathbf{v} - \mathbf{v}^*)\|^2/2, \quad \mathbf{v} \in \mathcal{Y} \setminus \mathcal{Y}_0(\mathbf{r}_0). \tag{34}$$

holds on a random set $\Omega(\mathbf{x})$ of probability at least $1 - 4e^{-\mathbf{x}}$.

The result (31) is an improved version of the approximation bound obtained in Spokoiny (2012), Theorem 3.1. The result (33) can be found in the supplement to Spokoiny (2012). The result (34) is very similar to Theorem 4.2 from Spokoiny (2012).

6.2 Tail posterior probability for full parameter space

The next step in our analysis is to check that \mathbf{v} concentrates in a small vicinity $\mathcal{Y}_0(\mathbf{r}_0)$ of the central point \mathbf{v}^* with a properly selected \mathbf{r}_0 . The concentration properties of the posterior will be described by using the random quantity

$$\rho^*(\mathbf{r}_0) = \frac{\int_{\mathcal{Y} \setminus \mathcal{Y}_0(\mathbf{r}_0)} \exp\{L(\mathbf{v}, \mathbf{v}^*)\} d\mathbf{v}}{\int_{\mathcal{Y}_0(\mathbf{r}_0)} \exp\{L(\mathbf{v}, \mathbf{v}^*)\} d\mathbf{v}}.$$

Theorem 10. *Suppose the conditions of Theorem 9. Then it holds on $\Omega_{\mathbf{r}_0}(\mathbf{x})$*

$$\rho^*(\mathbf{r}_0) \leq \exp\{2\Delta(\mathbf{r}_0, \mathbf{x}) + \nu(\mathbf{r}_0)\} \mathbf{b}^{-p/2} \mathbb{P}(\|\boldsymbol{\gamma}\|^2 \geq \mathbf{b}\mathbf{r}_0^2), \tag{35}$$

with

$$\nu(\mathbf{r}_0) \stackrel{\text{def}}{=} -\log \mathbb{P}(\|\boldsymbol{\gamma} + \boldsymbol{\xi}\| \leq \mathbf{r}_0 \mid \mathbf{Y}).$$

If $\mathbf{r}_0 \geq z_B(\mathbf{x}) + z(p, \mathbf{x})$, then on $\Omega(\mathbf{x})$

$$\nu(\mathbf{r}_0) \leq 2e^{-\mathbf{x}}. \tag{36}$$

This result together with Theorem 10 and Lemma 7 yields simple sufficient conditions on the value \mathbf{r}_0 which ensures the concentration of the posterior on $\mathcal{Y}_0(\mathbf{r}_0)$.

Corollary 2. *Assume the conditions of Theorem 10. Then the additional inequality $\text{br}_0^2 \geq z^2(p, \mathbf{x} + \frac{p}{2} \log \frac{\varepsilon}{\delta})$ ensures on a random set $\Omega(\mathbf{x})$ of probability at least $1 - 4e^{-x}$*

$$\rho^*(\mathbf{r}_0) \leq \exp\{2\Delta(\mathbf{r}_0, \mathbf{x}) + 2e^{-x} - x\}.$$

6.3 Tail posterior probability for target parameter

The next major step in our analysis is to check that $\boldsymbol{\theta}$ concentrates in a small vicinity $\Theta_0(\mathbf{r}_0) = \{\boldsymbol{\theta}: \|\check{D}_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\| \leq \mathbf{r}_0\}$ of the central point $\boldsymbol{\theta}^* = \Pi_0 \mathbf{v}^*$ with a properly selected \mathbf{r}_0 . The concentration properties of the posterior will be described by using the random quantity

$$\rho(\mathbf{r}_0) \stackrel{\text{def}}{=} \frac{\int_{\mathcal{Y}} \exp\{L(\mathbf{v}, \mathbf{v}^*)\} \pi(\mathbf{v}) \mathbb{I}\{\boldsymbol{\theta} \notin \Theta_0(\mathbf{r}_0)\} d\mathbf{v}}{\int_{\mathcal{Y}} \exp\{L(\mathbf{v}, \mathbf{v}^*)\} \pi(\mathbf{v}) \mathbb{I}\{\boldsymbol{\theta} \in \Theta_0(\mathbf{r}_0)\} d\mathbf{v}}.$$

In what follows we suppose that prior is uniform, i.e. $\pi(\mathbf{v}) \equiv 1$, $\mathbf{v} \in \mathcal{Y}$. This results in the following representation for $\rho(\mathbf{r}_0)$:

$$\rho(\mathbf{r}_0) = \frac{\int_{\mathcal{Y}} \exp\{L(\mathbf{v}, \mathbf{v}^*)\} \mathbb{I}\{\boldsymbol{\theta} \notin \Theta_0(\mathbf{r}_0)\} d\mathbf{v}}{\int_{\mathcal{Y}} \exp\{L(\mathbf{v}, \mathbf{v}^*)\} \mathbb{I}\{\boldsymbol{\theta} \in \Theta_0(\mathbf{r}_0)\} d\mathbf{v}}. \quad (37)$$

Obviously $\mathbb{P}(\boldsymbol{\vartheta} \notin \Theta_0(\mathbf{r}_0) | \mathbf{Y}) \leq \rho(\mathbf{r}_0)$. Therefore, small values of $\rho(\mathbf{r}_0)$ indicate a small posterior probability of the large deviation set $\{\boldsymbol{\vartheta} \notin \Theta_0(\mathbf{r}_0)\}$.

Theorem 11. *Suppose (31). Then for $\text{br}_0^2 \geq z^2(p, \mathbf{x} + \frac{p}{2} \log \frac{\varepsilon}{\delta})$ on $\Omega(\mathbf{x})$ of probability at least $1 - 4e^{-x}$*

$$\rho(\mathbf{r}_0) \leq \rho^*(\mathbf{r}_0) \leq \exp\{2\Delta(\mathbf{r}_0, \mathbf{x}) + 2e^{-x} - x\}.$$

6.4 Local Gaussian approximation of the posterior: upper bound

It is convenient to introduce local conditional expectation: for a random variable η , define

$$\mathbb{E}^\circ \eta \stackrel{\text{def}}{=} \mathbb{E} \left[\eta \mathbb{I}\{\boldsymbol{\vartheta} \in \Theta_0(\mathbf{r}_0)\} \mid \mathbf{Y} \right].$$

The following theorem gives an exact statement about the upper bound of this posterior expectation. Let $\check{\boldsymbol{\xi}}$ be from (7) and $\boldsymbol{\theta}^\circ = \boldsymbol{\theta}^* + \check{D}_0^{-1} \check{\boldsymbol{\xi}}$.

Theorem 12. *Suppose (31). Then for any $f: \mathbb{R}^q \rightarrow \mathbb{R}_+$ it holds on $\Omega_{\mathbf{r}_0}(\mathbf{x})$*

$$\mathbb{E}^\circ f(\check{D}_0(\boldsymbol{\vartheta} - \boldsymbol{\theta}^\circ)) \leq \exp\{\Delta^+(\mathbf{r}_0, \mathbf{x})\} \mathbb{E} f(\boldsymbol{\gamma}), \quad (38)$$

where $\gamma \sim \mathcal{N}(0, I_q)$ and

$$\begin{aligned} \Delta^+(\mathbf{r}_0, \mathbf{x}) &\stackrel{\text{def}}{=} 2\Delta(\mathbf{r}_0, \mathbf{x}) + \nu(\mathbf{r}_0) + \rho_f(\mathbf{r}_0), \\ \rho_f(\mathbf{r}_0) &\stackrel{\text{def}}{=} \frac{\int_{\mathcal{R} \setminus \mathcal{r}_0(\mathbf{r}_0)} \exp\{L(\mathbf{v}, \mathbf{v}^*)\} f(\check{D}_0(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)) d\mathbf{v}}{\int_{\mathcal{r}_0(\mathbf{r}_0)} \exp\{L(\mathbf{v}, \mathbf{v}^*)\} f(\check{D}_0(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)) d\mathbf{v}}. \end{aligned}$$

Define for random event $\eta \in A \subseteq \mathbb{R}^q$:

$$\mathbb{P}^\circ(\eta \in A) = \mathbb{E}^\circ \mathbb{I}\{\eta \in A\}.$$

The next result considers a special case with $f(\mathbf{u}) = |\boldsymbol{\lambda}^\top \mathbf{u}|^2$ and $f(\mathbf{u}) = \mathbb{I}(\mathbf{u} \in A)$ for any measurable set A .

Corollary 3. For any $\boldsymbol{\lambda} \in \mathbb{R}^q$, it holds on $\Omega_{\mathbf{r}_0}(\mathbf{x})$

$$\mathbb{E}^\circ |\boldsymbol{\lambda}^\top \check{D}_0(\boldsymbol{\vartheta} - \boldsymbol{\theta}^\circ)|^2 \leq \exp\{\Delta^+(\mathbf{r}_0, \mathbf{x})\} \|\boldsymbol{\lambda}\|^2.$$

For any measurable set $A \subseteq \mathbb{R}^q$, it holds on $\Omega_{\mathbf{r}_0}(\mathbf{x})$

$$\mathbb{P}^\circ(\check{D}_0(\boldsymbol{\vartheta} - \boldsymbol{\theta}^\circ) \in A) \leq \exp\{\Delta^+(\mathbf{r}_0, \mathbf{x})\} \mathbb{P}(\gamma \in A). \tag{39}$$

On $\Omega(\mathbf{x})$ one obtains

$$\Delta^+(\mathbf{r}_0, \mathbf{x}) \leq 2\Delta(\mathbf{r}_0, \mathbf{x}) + 2e^{-\mathbf{x}} + 2 \exp\{\Delta(\mathbf{r}_0, \mathbf{x}) + 4e^{-\mathbf{x}} - \mathbf{x}\}.$$

The next corollary describes an upper bound for the posterior probability in case of changing of scaling.

Corollary 4. Let D_1 be a symmetric $q \times q$ matrix such that $\|I - D_1^{-1} \check{D}_0^2 D_1^{-1}\| \leq \alpha$. Let also $\hat{\boldsymbol{\theta}} \in \mathbb{R}^q$ be such that $\|\check{D}_0(\boldsymbol{\theta}^\circ - \hat{\boldsymbol{\theta}})\| \leq \beta$. Then for any measurable set $A \subseteq \mathbb{R}^q$, it holds on $\Omega(\mathbf{x})$ with $\boldsymbol{\delta}_0 \stackrel{\text{def}}{=} D_1(\boldsymbol{\theta}^\circ - \hat{\boldsymbol{\theta}})$

$$\begin{aligned} \mathbb{P}^\circ(D_1(\boldsymbol{\vartheta} - \hat{\boldsymbol{\theta}}) \in A) &\leq \exp\{\Delta^+(\mathbf{r}_0, \mathbf{x})\} \mathbb{P}(D_1 \check{D}_0^{-1} \gamma + \boldsymbol{\delta}_0 \in A) \\ &\leq \exp\{\Delta^+(\mathbf{r}_0, \mathbf{x})\} \left(\mathbb{P}(\gamma \in A) + \frac{1}{2} \sqrt{\alpha^2 q + (1 + \alpha)^2 \beta^2} \right). \end{aligned} \tag{40}$$

6.5 Local Gaussian approximation of the posterior: lower bound

Now we present a local lower bound for the posterior measure.

Theorem 13. Suppose (31). Then for any $f: \mathbb{R}^q \rightarrow \mathbb{R}_+$ it holds on $\Omega_{\mathbf{r}_0}(\mathbf{x})$

$$\mathbb{E}^\circ f(\check{D}_0(\boldsymbol{\vartheta} - \boldsymbol{\theta}^\circ)) \geq \exp\{-\Delta^-(\mathbf{r}_0, \mathbf{x})\} \mathbb{E}\{f(\gamma) \mathbb{I}(\|\gamma + \check{\boldsymbol{\xi}}\| \leq \mathbf{r}_0)\}, \tag{41}$$

where

$$\begin{aligned}\Delta^-(\mathbf{r}_0, \mathbf{x}) &\stackrel{\text{def}}{=} 2\Delta(\mathbf{r}_0, \mathbf{x}) + \nu(\mathbf{r}_0) + \rho^*(\mathbf{r}_0) + 2\tilde{\rho}_f(\mathbf{r}_0), \\ \tilde{\rho}_f(\mathbf{r}_0) &\stackrel{\text{def}}{=} \frac{\int_{\mathbb{R}^p \setminus \mathcal{r}_0(\mathbf{r}_0)} \exp\{\mathbb{L}(\mathbf{v}, \mathbf{v}^*)\} f(\check{D}_0(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)) d\mathbf{v}}{\int_{\mathcal{r}_0(\mathbf{r}_0)} \exp\{\mathbb{L}(\mathbf{v}, \mathbf{v}^*)\} f(\check{D}_0(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)) d\mathbf{v}}.\end{aligned}\quad (42)$$

This result means that posterior measure can be bounded from below by the standard normal law up to (small) multiplicative and additive constants. As a corollary, we state the result for quadratic and indicator functions $f(\mathbf{u})$. The proof is similar to Corollary 4 and Corollary 3.

Corollary 5. For any $\boldsymbol{\lambda} \in \mathbb{R}^q$, it holds on $\Omega_{\mathbf{r}_0}(\mathbf{x})$

$$\mathbb{E}^\circ |\boldsymbol{\lambda}^\top \check{D}_0(\boldsymbol{\vartheta} - \boldsymbol{\theta}^\circ)|^2 \geq \exp\{-\Delta^-(\mathbf{r}_0, \mathbf{x}) + e^{-\mathbf{x}}\} \|\boldsymbol{\lambda}\|^2.$$

For any measurable set $A \subseteq \mathbb{R}^q$, it holds on $\Omega_{\mathbf{r}_0}(\mathbf{x})$

$$\mathbb{P}^\circ(\check{D}_0(\boldsymbol{\vartheta} - \boldsymbol{\theta}^\circ) \in A) \geq \exp\{\Delta^-(\mathbf{r}_0, \mathbf{x})\} \mathbb{P}(\boldsymbol{\gamma} \in A) - e^{-\mathbf{x}}. \quad (43)$$

Let D_1^2 be a symmetric $q \times q$ matrix such that $\|I - D_1^{-1} \check{D}_0^2 D_1^{-1}\| \leq \alpha$ and let $\hat{\boldsymbol{\theta}} \in \mathbb{R}^q$ be such that $\|\check{D}_0(\boldsymbol{\theta}^\circ - \hat{\boldsymbol{\theta}})\| \leq \beta$. Define $\boldsymbol{\delta}_0 \stackrel{\text{def}}{=} D_1(\boldsymbol{\theta}^\circ - \hat{\boldsymbol{\theta}})$. Then for any measurable subset A in \mathbb{R}^q , it holds on $\Omega(\mathbf{x})$

$$\begin{aligned}\mathbb{P}^\circ(D_1(\boldsymbol{\vartheta} - \hat{\boldsymbol{\theta}}) \in A) &\geq \exp\{\Delta^-(\mathbf{r}_0, \mathbf{x})\} \mathbb{P}(D_1 \check{D}_0^{-1} \boldsymbol{\gamma} + \boldsymbol{\delta}_0 \in A) - e^{-\mathbf{x}} \\ &\geq \exp\{\Delta^-(\mathbf{r}_0, \mathbf{x})\} \left\{ \mathbb{P}(\boldsymbol{\gamma} \in A) - \frac{1}{2} \sqrt{\alpha^2 q + (1 + \alpha)^2 \beta^2} \right\} - e^{-\mathbf{x}}, \\ \Delta^-(\mathbf{r}_0, \mathbf{x}) &\leq 2\Delta(\mathbf{r}_0, \mathbf{x}) + 3e^{-\mathbf{x}} + 4 \exp\{\Delta(\mathbf{r}_0, \mathbf{x}) + 4e^{-\mathbf{x}} - \mathbf{x}\}.\end{aligned}$$

7 Proofs

This appendix collects the proofs of the results.

7.1 Some inequalities for the normal law

This section collects some simple but useful facts about the properties of the multivariate standard normal distribution. Many similar results can be found in the literature, we present the proofs to keep the presentation self-contained. Everywhere in this section $\boldsymbol{\gamma}$ means a standard normal vector in \mathbb{R}^q .

Lemma 6. For any $\mathbf{u} \in \mathbb{R}^q$, any unit vector $\mathbf{a} \in \mathbb{R}^q$, and any $z > 0$, it holds

$$\mathbb{P}(\|\boldsymbol{\gamma} - \mathbf{u}\| \geq z) \leq \exp\{-z^2/4 + q/2 + \|\mathbf{u}\|^2/2\}, \quad (44)$$

$$\mathbb{E}\{|\boldsymbol{\gamma}^\top \mathbf{a}|^2 \mathbb{I}(\|\boldsymbol{\gamma} - \mathbf{u}\| \geq z)\} \leq (2 + |\mathbf{u}^\top \mathbf{a}|^2) \exp\{-z^2/4 + q/2 + \|\mathbf{u}\|^2/2\}. \quad (45)$$

Proof. By the exponential Chebyshev inequality, for any $\lambda < 1$

$$\begin{aligned} \mathbb{P}(\|\boldsymbol{\gamma} - \mathbf{u}\| \geq z) &\leq \exp(-\lambda z^2/2) \mathbb{E} \exp(\lambda \|\boldsymbol{\gamma} - \mathbf{u}\|^2/2) \\ &= \exp\left\{-\frac{\lambda z^2}{2} - \frac{q}{2} \log(1 - \lambda) + \frac{\lambda}{2(1 - \lambda)} \|\mathbf{u}\|^2\right\}. \end{aligned}$$

In particular, with $\lambda = 1/2$, this implies (44). Further, for $\|\mathbf{a}\| = 1$

$$\begin{aligned} \mathbb{E}\{|\boldsymbol{\gamma}^\top \mathbf{a}|^2 \mathbb{I}(\|\boldsymbol{\gamma} - \mathbf{u}\| \geq z)\} &\leq \exp(-z^2/4) \mathbb{E}\{|\boldsymbol{\gamma}^\top \mathbf{a}|^2 \exp(\|\boldsymbol{\gamma} - \mathbf{u}\|^2/4)\} \\ &\leq (2 + |\mathbf{u}^\top \mathbf{a}|^2) \exp(-z^2/4 + q/2 + \|\mathbf{u}\|^2/2) \end{aligned}$$

and (45) follows. □

The next result explains the concentration effect for the norm $\|\boldsymbol{\xi}\|^2$ of a Gaussian vector. We use a version from Spokoiny (2012).

Lemma 7. *For each \mathbf{x} ,*

$$\mathbb{P}(\|\boldsymbol{\gamma}\| \geq z(q, \mathbf{x})) \leq \exp(-\mathbf{x}), \quad \mathbb{P}(\|\boldsymbol{\gamma}\| \leq z_1(q, \mathbf{x})) \leq \exp(-\mathbf{x}),$$

where

$$z^2(q, \mathbf{x}) \stackrel{\text{def}}{=} q + \sqrt{6.6q\mathbf{x}} \vee (6.6\mathbf{x}), \quad z_1^2(q, \mathbf{x}) \stackrel{\text{def}}{=} q - 2\sqrt{q\mathbf{x}}.$$

The next lemma bounds from above the Kullback-Leibler divergence between two normal distributions.

Lemma 8. *Let $\mathbb{P} = \mathcal{N}(\mathbf{b}, \Sigma)$ and $\mathbb{P}^\circ = \mathcal{N}(\mathbf{b}^\circ, \Sigma^\circ)$ for some non-degenerate matrices Σ and Σ° . If*

$$\|\Sigma^{-1/2} \Sigma^\circ \Sigma^{-1/2} - \mathbf{I}_q\| \leq \epsilon \leq 1/2, \quad \text{tr}(\Sigma^{-1/2} \Sigma^\circ \Sigma^{-1/2} - \mathbf{I}_q)^2 \leq \delta^2,$$

then

$$\begin{aligned} \mathcal{K}(\mathbb{P}, \mathbb{P}^\circ) &= -\mathbb{E}_0 \log \frac{d\mathbb{P}^\circ}{d\mathbb{P}} \leq \frac{\delta^2}{2} + \frac{1}{2} (\mathbf{b} - \mathbf{b}^\circ)^\top \Sigma^\circ (\mathbf{b} - \mathbf{b}^\circ) \\ &\leq \frac{\delta^2}{2} + \frac{1 + \epsilon}{2} (\mathbf{b} - \mathbf{b}^\circ)^\top \Sigma (\mathbf{b} - \mathbf{b}^\circ). \end{aligned}$$

For any measurable set $A \subset \mathbb{R}^q$, it holds

$$|\mathbb{P}(A) - \mathbb{P}^\circ(A)| \leq \sqrt{\mathcal{K}(\mathbb{P}, \mathbb{P}^\circ)/2} \leq \frac{1}{2} \sqrt{\delta^2 + (1 + \epsilon)(\mathbf{b} - \mathbf{b}^\circ)^\top \Sigma (\mathbf{b} - \mathbf{b}^\circ)}.$$

Proof. The change of variables $\mathbf{u} = \Sigma^{-1/2}(\mathbf{x} - \mathbf{b})$ reduces the general case to the situation when \mathbb{P} is standard normal in \mathbb{R}^q while $P_1 = \mathcal{N}(\boldsymbol{\beta}, B)$ with $\boldsymbol{\beta} = \Sigma^{1/2}(\mathbf{b}^\circ - \mathbf{b})$ and $B \stackrel{\text{def}}{=} \Sigma^{-1/2} \Sigma^\circ \Sigma^{-1/2}$

$$2 \log \frac{d\mathbb{P}^\circ}{d\mathbb{P}}(\boldsymbol{\gamma}) = \log \det(B) - (\boldsymbol{\gamma} - \boldsymbol{\beta})^\top B(\boldsymbol{\gamma} - \boldsymbol{\beta}) + \|\boldsymbol{\gamma}\|^2$$

with $\boldsymbol{\gamma}$ standard normal and

$$2\mathcal{K}(\mathbb{P}, \mathbb{P}^\circ) = -2\mathbb{E}_0 \log \frac{d\mathbb{P}^\circ}{d\mathbb{P}} = -\log \det(B) + \text{tr}(B - I_q) + \boldsymbol{\beta}^\top B\boldsymbol{\beta}.$$

Let a_j be the j th eigenvalue of $B - I_q$. From $\|B - I_q\| \leq 1/2$ it follows $|a_j| \leq 1/2$ and

$$\begin{aligned} 2\mathcal{K}(\mathbb{P}, \mathbb{P}^\circ) &= \boldsymbol{\beta}^\top B\boldsymbol{\beta} + \sum_{j=1}^q \{a_j - \log(1 + a_j)\} \leq \boldsymbol{\beta}^\top B\boldsymbol{\beta} + \sum_{j=1}^q a_j^2 \\ &\leq \boldsymbol{\beta}^\top B\boldsymbol{\beta} + \text{tr}(B - I_q)^2 \leq \boldsymbol{\beta}^\top B\boldsymbol{\beta} + \delta^2. \end{aligned}$$

This implies by Pinsker's inequality

$$\sup_A |\mathbb{P}(A) - \mathbb{P}^\circ(A)| \leq \sqrt{\frac{1}{2} \mathcal{K}(\mathbb{P}, \mathbb{P}^\circ)} \leq \frac{1}{2} \sqrt{\delta^2 + \boldsymbol{\beta}^\top B\boldsymbol{\beta}}$$

as required. \square

7.2 Proof of Theorem 10

Define $u(\mathbf{v}) = \mathbf{b} \|\mathcal{D}_0(\mathbf{v} - \mathbf{v}^*)\|^2/2$. Now, by a change of variables, one obtains

$$\begin{aligned} &\frac{\mathbf{b}^{p/2} \det(\mathcal{D}_0)}{(2\pi)^{p/2}} \int_{\mathcal{R} \setminus \mathcal{R}_0(\mathbf{r}_0)} \exp\{-u(\mathbf{v})\} d\mathbf{v} \\ &\leq \frac{\mathbf{b}^{p/2} \det(\mathcal{D}_0)}{(2\pi)^{p/2}} \int_{\mathcal{R} \setminus \mathcal{R}_0(\mathbf{r}_0)} \exp\{-\mathbf{b} \|\mathcal{D}_0(\mathbf{v} - \mathbf{v}^*)\|^2/2\} d\mathbf{v} = \mathbb{P}(\|\boldsymbol{\gamma}\|^2 \geq \mathbf{b}\mathbf{r}_0^2). \end{aligned}$$

For the integral in the numerator of (37), it holds on $\Omega(\mathbf{x})$ by (34)

$$\int_{\mathcal{R} \setminus \mathcal{R}_0(\mathbf{r}_0)} \exp\{L(\mathbf{v}, \mathbf{v}^*)\} d\mathbf{v} \leq \int_{\mathcal{R} \setminus \mathcal{R}_0(\mathbf{r}_0)} \exp\{-u(\mathbf{v})\} d\mathbf{v}.$$

For the integral in the denominator it holds

$$\begin{aligned} &\int_{\mathcal{R}_0(\mathbf{r}_0)} \exp\{L(\mathbf{v}, \mathbf{v}^*)\} d\mathbf{v} \\ &\geq \exp\{-\Delta(\mathbf{r}_0, \mathbf{x}) - m(\boldsymbol{\xi})\} \int_{\mathcal{R}_0(\mathbf{r}_0)} \exp\{\mathbb{L}(\mathbf{v}, \mathbf{v}^*) + m(\boldsymbol{\xi})\} d\mathbf{v}. \end{aligned} \quad (46)$$

Inequality (46) implies by definition of $\nu(\mathbf{r}_0)$:

$$\int_{\mathcal{Y}_0(\mathbf{r}_0)} \exp\{L(\mathbf{v}, \mathbf{v}^*)\} d\mathbf{v} \geq \exp\{-\Delta(\mathbf{r}_0, \mathbf{x}) - m(\boldsymbol{\xi}) - \nu(\mathbf{r}_0)\}. \tag{47}$$

The bound (47) for the local integral $\int_{\mathcal{Y}_0(\mathbf{r}_0)} \exp\{L(\mathbf{v}, \mathbf{v}^*)\} d\mathbf{v}$ implies that

$$\rho^*(\mathbf{r}_0) \leq \exp\{\Delta(\mathbf{r}_0, \mathbf{x}) + \nu(\mathbf{r}_0) + m(\boldsymbol{\xi})\} \int_{\mathcal{Y} \setminus \mathcal{Y}_0(\mathbf{r}_0)} \exp\{-u(\mathbf{v})\} d\mathbf{v}.$$

Finally

$$\exp\{m(\boldsymbol{\xi})\} = \exp\{-\|\boldsymbol{\xi}\|^2/2\} (2\pi)^{-p/2} \det(\mathcal{D}_0) \leq (2\pi)^{-p/2} \det(\mathcal{D}_0)$$

and the assertion (35) follows. The bound (36) is also straightforward:

$$\begin{aligned} \nu(\mathbf{r}_0) &= -\log \mathbb{P}(\|\boldsymbol{\gamma} + \boldsymbol{\xi}\| \leq \mathbf{r}_0 \mid \mathbf{Y}) \leq -\log \mathbb{P}(\|\boldsymbol{\gamma}\| + \|\boldsymbol{\xi}\| \leq \mathbf{r}_0 \mid \mathbf{Y}) \\ &\leq -\log \mathbb{P}(\|\boldsymbol{\gamma}\| \leq z(p, \mathbf{x}) \mid \mathbf{Y}) \leq 2e^{-\mathbf{x}}. \end{aligned}$$

7.3 Proof of Theorem 11

Obviously $\{\boldsymbol{\theta} \notin \Theta_0(\mathbf{r}_0), \mathbf{v} \in \mathcal{Y}\} \subset \{\mathcal{Y} \setminus \mathcal{Y}_0(\mathbf{r}_0)\}$. Therefore, it holds for the integral in the numerator of (37) in view of (34)

$$\int_{\mathcal{Y}} \exp\{L(\mathbf{v}, \mathbf{v}^*)\} \mathbb{I}\{\boldsymbol{\theta} \notin \Theta_0(\mathbf{r}_0)\} d\mathbf{v} \leq \int_{\mathcal{Y} \setminus \mathcal{Y}_0(\mathbf{r}_0)} \exp\{L(\mathbf{v}, \mathbf{v}^*)\} d\mathbf{v}.$$

For the denominator, the inclusion $\mathcal{Y}_0(\mathbf{r}_0) \subset \{\boldsymbol{\theta} \in \Theta_0(\mathbf{r}_0), \mathbf{v} \in \mathcal{Y}\}$ and (34) imply

$$\int_{\mathcal{Y}} \exp\{L(\mathbf{v}, \mathbf{v}^*)\} \mathbb{I}\{\boldsymbol{\theta} \in \Theta_0(\mathbf{r}_0)\} d\mathbf{v} \geq \int_{\mathcal{Y}_0(\mathbf{r}_0)} \exp\{L(\mathbf{v}, \mathbf{v}^*)\} d\mathbf{v}.$$

Finally

$$\rho(\mathbf{r}_0) = \frac{\int_{\mathcal{Y}} \exp\{L(\mathbf{v}, \mathbf{v}^*)\} \mathbb{I}\{\boldsymbol{\theta} \notin \Theta_0(\mathbf{r}_0)\} d\mathbf{v}}{\int_{\mathcal{Y}} \exp\{L(\mathbf{v}, \mathbf{v}^*)\} \mathbb{I}\{\boldsymbol{\theta} \in \Theta_0(\mathbf{r}_0)\} d\mathbf{v}} \leq \frac{\int_{\mathcal{Y} \setminus \mathcal{Y}_0(\mathbf{r}_0)} \exp\{L(\mathbf{v}, \mathbf{v}^*)\} d\mathbf{v}}{\int_{\mathcal{Y}_0(\mathbf{r}_0)} \exp\{L(\mathbf{v}, \mathbf{v}^*)\} d\mathbf{v}} = \rho^*(\mathbf{r}_0),$$

and the assertion follows from Theorem 10.

7.4 Proof of Theorem 12

We use that $L(\mathbf{v}, \mathbf{v}^*) = \boldsymbol{\xi}^\top \mathcal{D}_0(\mathbf{v} - \mathbf{v}^*) - \|\mathcal{D}_0(\mathbf{v} - \mathbf{v}^*)\|^2/2$ is proportional to the density of a Gaussian distribution. More precisely, define

$$m(\boldsymbol{\xi}) \stackrel{\text{def}}{=} -\|\boldsymbol{\xi}\|^2/2 + \log(\det \mathcal{D}_0) - p \log(\sqrt{2\pi}).$$

Then

$$m(\boldsymbol{\xi}) + \mathbb{L}(\mathbf{v}, \mathbf{v}^*) = -\|\mathcal{D}_0(\mathbf{v} - \mathbf{v}^*) - \boldsymbol{\xi}\|^2/2 + \log(\det \mathcal{D}_0) - p \log(\sqrt{2\pi})$$

is (conditionally on \mathbf{Y}) the log-density of the normal law with the mean $\mathbf{v}_0 = \mathbf{v}^* + \mathcal{D}_0^{-1}\boldsymbol{\xi}$ and the covariance matrix \mathcal{D}_0^{-2} . If we perform integration and leave only $\boldsymbol{\theta}$ part of \mathbf{v} then $m(\boldsymbol{\xi}) + \mathbb{L}(\mathbf{v}, \mathbf{v}^*)$ is (conditionally on \mathbf{Y}) the log-density of the normal law with the mean $\boldsymbol{\theta}^\circ = \check{D}_0^{-1}\check{\boldsymbol{\xi}} + \boldsymbol{\theta}^*$ and the covariance matrix \check{D}_0^{-2} . So, for any nonnegative function $f: \mathbb{R}^q \rightarrow \mathbb{R}_+$ we get

$$\begin{aligned} & \int_{\mathcal{I}} \exp\{L(\mathbf{v}, \mathbf{v}^*) + m(\boldsymbol{\xi})\} f(\check{D}_0(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)) d\mathbf{v} \\ &= \int_{\mathcal{I}_0(\mathbf{r}_0)} \exp\{L(\mathbf{v}, \mathbf{v}^*) + m(\boldsymbol{\xi})\} f(\check{D}_0(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)) d\mathbf{v} \\ & \quad + \int_{\mathcal{I} \setminus \mathcal{I}_0(\mathbf{r}_0)} \exp\{L(\mathbf{v}, \mathbf{v}^*) + m(\boldsymbol{\xi})\} f(\check{D}_0(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)) d\mathbf{v} \\ &= (1 + \rho_f(\mathbf{r}_0)) \int_{\mathcal{I}_0(\mathbf{r}_0)} \exp\{L(\mathbf{v}, \mathbf{v}^*) + m(\boldsymbol{\xi})\} f(\check{D}_0(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)) d\mathbf{v} \\ &\leq e^{\Delta(\mathbf{r}_0, \mathbf{x}) + \rho_f(\mathbf{r}_0)} \int_{\mathcal{I}_0(\mathbf{r}_0)} \exp\{\mathbb{L}(\mathbf{v}, \mathbf{v}^*) + m(\boldsymbol{\xi})\} f(\check{D}_0(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)) d\mathbf{v} \\ &\leq e^{\Delta(\mathbf{r}_0, \mathbf{x}) + \rho_f(\mathbf{r}_0)} \int_{\mathbb{R}^p} \exp\{\mathbb{L}(\mathbf{v}, \mathbf{v}^*) + m(\boldsymbol{\xi})\} f(\check{D}_0(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)) d\mathbf{v} \\ &= e^{\Delta(\mathbf{r}_0, \mathbf{x}) + \rho_f(\mathbf{r}_0)} \mathbb{E}f(\boldsymbol{\gamma}). \end{aligned}$$

Thus,

$$\int_{\mathcal{I}} \exp\{L(\mathbf{v}, \mathbf{v}^*)\} f(\check{D}_0(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)) d\mathbf{v} \leq \exp\{\Delta(\mathbf{r}_0, \mathbf{x}) - m(\boldsymbol{\xi}) + \rho_f(\mathbf{r}_0)\} \mathbb{E}f(\boldsymbol{\gamma}). \quad (48)$$

Now (48) and (47) imply

$$\frac{\int_{\mathcal{I}} \exp\{L(\mathbf{v}, \mathbf{v}^*)\} f(\check{D}_\epsilon(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)) d\mathbf{v}}{\int_{\mathcal{I}} \exp\{L(\mathbf{v}, \mathbf{v}^*)\} d\mathbf{v}} \leq \exp\{2\Delta(\mathbf{r}_0, \mathbf{x}) + \nu(\mathbf{r}_0) + \rho_f(\mathbf{r}_0)\} \mathbb{E}f(\boldsymbol{\gamma})$$

and (38) follows by definition of $\Delta^+(\mathbf{r}_0, \mathbf{x})$.

7.5 Proof of Corollary 3

As a direct implication of (38) one easily gets

$$\mathbb{E}^\circ |\boldsymbol{\lambda}^\top \check{D}_0(\boldsymbol{\vartheta} - \boldsymbol{\theta}^\circ)|^2 \leq \exp(\Delta^+(\mathbf{r}_0, \mathbf{x})) \|\boldsymbol{\lambda}\|^2.$$

The only important step is to show that $\rho_{x^2}(\mathbf{r}_0)$ is small. Denote

$$\boldsymbol{\lambda}_0 = \mathcal{D}_0^{-1} \begin{pmatrix} \check{D}_0 \boldsymbol{\lambda} \\ \mathbf{0} \end{pmatrix}$$

and $\mathbf{0}$ is a zero vector of dimension $(p - q)$. We proceed separately with the numerator and denominator. For the numerator by (44) and (45) on $\Omega(\mathbf{x})$

$$\begin{aligned} & \int_{\mathcal{R} \setminus \mathcal{R}_0(\mathbf{r}_0)} |\boldsymbol{\lambda}^\top \check{D}_0(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)|^2 \exp\{L(\mathbf{v}, \mathbf{v}^*)\} d\mathbf{v} \\ & \leq \int_{\mathcal{R} \setminus \mathcal{R}_0(\mathbf{r}_0)} |\boldsymbol{\lambda}^\top \check{D}_0(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)|^2 \exp\{-\mathbf{b}\|\mathcal{D}_0(\mathbf{v} - \mathbf{v}^*)\|^2/2\} d\mathbf{v} \\ & = \int_{\mathcal{R} \setminus \mathcal{R}_0(\mathbf{r}_0)} |\boldsymbol{\lambda}_0^\top \mathcal{D}_0(\mathbf{v} - \mathbf{v}_0)|^2 \exp\{-\mathbf{b}\|\mathcal{D}_0(\mathbf{v} - \mathbf{v}^*)\|^2/2\} d\mathbf{v} \\ & = \exp\{(-p/2 + 2) \log \mathbf{b} - \log(\det \mathcal{D}_0) + p \log(\sqrt{2\pi})\} \mathbb{E} |\boldsymbol{\lambda}_0^\top (\boldsymbol{\gamma} + \boldsymbol{\xi})|^2 \mathbb{I}(\|\boldsymbol{\gamma}\|^2 \geq \mathbf{b}\mathbf{r}_0^2) \\ & \leq (4 + 2\|\boldsymbol{\xi}\|^2) \exp\{-p/2 + 2 \log \mathbf{b} - \log(\det \mathcal{D}_0) + p \log(\sqrt{2\pi}) - \mathbf{b}\mathbf{r}_0^2/4 + p/2\} \|\boldsymbol{\lambda}_0\|^2 \\ & \leq \exp\{\|\boldsymbol{\xi}\|^2/2 + (p/2 + 2) \log(e/\mathbf{b}) - \log(\det \mathcal{D}_0) + p \log(\sqrt{2\pi}) - \mathbf{b}\mathbf{r}_0^2/4\} \|\boldsymbol{\lambda}_0\|^2 \\ & \leq \exp\{-\log(\det \mathcal{D}_0) + p \log(\sqrt{2\pi}) - \mathbf{x}\} \|\boldsymbol{\lambda}_0\|^2 \end{aligned}$$

for $\mathbf{b}\mathbf{r}_0^2 \geq (2p + 4) \log(e/\mathbf{b}) + 2z_B(\mathbf{x}) + 4\mathbf{x}$. For the denominator, it holds on $\Omega(\mathbf{x})$

$$\begin{aligned} & \int_{\mathcal{R}_0(\mathbf{r}_0)} |\boldsymbol{\lambda}^\top \check{D}_0(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)|^2 \exp\{L(\mathbf{v}, \mathbf{v}^*)\} d\mathbf{v} \\ & \geq e^{-\Delta(\mathbf{r}_0, \mathbf{x})} \int_{\mathcal{R}_0(\mathbf{r}_0)} |\boldsymbol{\lambda}^\top \check{D}_0(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)|^2 \exp\{\mathbb{L}(\mathbf{v}, \mathbf{v}^*)\} d\mathbf{v} \\ & = e^{-\Delta(\mathbf{r}_0, \mathbf{x})} \int_{\mathcal{R}_0(\mathbf{r}_0)} |\boldsymbol{\lambda}_0^\top \mathcal{D}_0(\mathbf{v} - \mathbf{v}_0)|^2 \exp\{\mathbb{L}(\mathbf{v}, \mathbf{v}^*)\} d\mathbf{v} \\ & = e^{-\Delta(\mathbf{r}_0, \mathbf{x}) - m(\boldsymbol{\xi})} \mathbb{E} |\boldsymbol{\lambda}_0^\top (\boldsymbol{\gamma} + \boldsymbol{\xi})|^2 \mathbb{I}(\|\boldsymbol{\gamma} + \boldsymbol{\xi}\|^2 \leq \mathbf{r}_0^2) \\ & = e^{-\Delta(\mathbf{r}_0, \mathbf{x}) - m(\boldsymbol{\xi})} \left\{ \mathbb{E} |\boldsymbol{\lambda}_0^\top (\boldsymbol{\gamma} + \boldsymbol{\xi})|^2 - \mathbb{E} |\boldsymbol{\lambda}_0^\top (\boldsymbol{\gamma} + \boldsymbol{\xi})|^2 \mathbb{I}(\|\boldsymbol{\gamma} + \boldsymbol{\xi}\|^2 \geq \mathbf{r}_0^2) \right\} \\ & \geq e^{-\Delta(\mathbf{r}_0, \mathbf{x}) - m(\boldsymbol{\xi})} \left\{ \|\boldsymbol{\lambda}_0\|^2 + |\boldsymbol{\lambda}_0^\top \boldsymbol{\xi}|^2 - 2\|\boldsymbol{\lambda}_0\|^2 \exp\{-\mathbf{r}_0^2/4 + p/2 + \|\boldsymbol{\xi}\|^2/2\} \right\} \\ & \geq e^{-\Delta(\mathbf{r}_0, \mathbf{x}) - m(\boldsymbol{\xi})} \|\boldsymbol{\lambda}_0\|^2 \{1 - 2e^{-\mathbf{x}}\} \geq \exp\{-\Delta(\mathbf{r}_0, \mathbf{x}) - m(\boldsymbol{\xi}) - 4e^{-\mathbf{x}}\} \|\boldsymbol{\lambda}_0\|^2 \end{aligned}$$

for $\mathbf{r}_0^2 \geq 2p + 2z_B(\mathbf{x}) + 4\mathbf{x}$ on $\Omega(\mathbf{x})$. This yields on $\Omega(\mathbf{x})$

$$\begin{aligned} \rho_{x^2}(\mathbf{r}_0) &= \frac{\int_{\mathcal{R} \setminus \mathcal{R}_0(\mathbf{r}_0)} |\boldsymbol{\lambda}^\top \check{D}_0(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)|^2 \exp\{L(\mathbf{v}, \mathbf{v}^*)\} d\mathbf{v}}{\int_{\mathcal{R}_0(\mathbf{r}_0)} |\boldsymbol{\lambda}^\top \check{D}_0(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)|^2 \exp\{L(\mathbf{v}, \mathbf{v}^*)\} d\mathbf{v}} \\ &\leq \frac{2 \exp\{-\log(\det \mathcal{D}_0) + p \log(\sqrt{2\pi})\} - \mathbf{x}}{\exp\{-\Delta(\mathbf{r}_0, \mathbf{x}) + m(\boldsymbol{\xi}) - 4e^{-\mathbf{x}}\}} \|\boldsymbol{\lambda}_0\|^2 \\ &= 2 \exp\{\Delta(\mathbf{r}_0, \mathbf{x}) - \|\boldsymbol{\xi}\|^2/2 + 4e^{-\mathbf{x}} - \mathbf{x}\} \leq 2 \exp\{\Delta(\mathbf{r}_0, \mathbf{x}) + 4e^{-\mathbf{x}} - \mathbf{x}\}. \end{aligned}$$

7.6 Proof of Corollary 4

The first statement follows from Theorem 12 with $f(\mathbf{u}) = \mathbb{I}(D_1 \check{D}_0^{-1} \mathbf{u} + \boldsymbol{\delta}_0 \in A)$. Further, it holds on $\Omega(\mathbf{x})$ for $\boldsymbol{\delta}_0 \stackrel{\text{def}}{=} D_1(\boldsymbol{\theta}^\circ - \widehat{\boldsymbol{\theta}})$

$$\|\boldsymbol{\delta}_0\|^2 = \|D_1(\boldsymbol{\theta}^\circ - \widehat{\boldsymbol{\theta}})\|^2 \leq (1 + \alpha) \|\check{D}_0(\boldsymbol{\theta}^\circ - \widehat{\boldsymbol{\theta}})\|^2 \leq (1 + \alpha) \beta^2.$$

For proving (40), we apply Pinsker's inequality to two normal distributions. Let γ be standard normal in \mathbb{R}^q . The random variable $D_1 \check{D}_0^{-1} \gamma + \boldsymbol{\delta}_0$ is normal with mean $\mathcal{N}(\boldsymbol{\delta}_0, B_1^{-1})$ with $B_1^{-1} \stackrel{\text{def}}{=} D_1 \check{D}_0^{-2} D_1$. Obviously $\|I_q - B_1\| = \|I_q - D_1^{-1} \check{D}_0^2 D_1^{-1}\| \leq \alpha$. Thus, by Lemma 8 for any measurable set A , it holds

$$\mathbb{P}(D_1 \check{D}_0^{-1} \gamma + \boldsymbol{\delta}_0 \in A \mid \mathbf{Y}) \leq \mathbb{P}(\gamma \in A) + \frac{1}{2} \sqrt{\alpha^2 q + (1 + \alpha)^2 \beta^2}.$$

7.7 Proof of Theorem 13

As in the proof of Theorem 12, for any nonnegative function $f: \mathbb{R}^q \rightarrow \mathbb{R}_+$, it holds

$$\begin{aligned} &\int_{\mathcal{R}} \exp\{L(\mathbf{v}, \mathbf{v}^*)\} f(\check{D}_0(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)) \mathbb{I}\{\boldsymbol{\theta} \in \Theta_0(\mathbf{r}_0)\} d\mathbf{v} \\ &\geq \int_{\mathcal{R}_0(\mathbf{r}_0)} \exp\{L(\mathbf{v}, \mathbf{v}^*)\} f(\check{D}_0(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)) d\mathbf{v} \\ &\geq \exp\{-\Delta(\mathbf{r}_0, \mathbf{x}) - m(\boldsymbol{\xi})\} \int_{\mathcal{R}_0(\mathbf{r}_0)} \exp\{\mathbb{L}(\mathbf{v}, \mathbf{v}^*)\} f(\check{D}_0(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)) d\mathbf{v} \\ &\geq \exp\{-\Delta(\mathbf{r}_0, \mathbf{x}) - m(\boldsymbol{\xi})\} \int_{\mathbb{R}^p} \exp\{\mathbb{L}(\mathbf{v}, \mathbf{v}^*)\} f(\check{D}_0(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)) d\mathbf{v} \\ &\quad - \exp\{-\Delta(\mathbf{r}_0, \mathbf{x}) - m(\boldsymbol{\xi})\} \int_{\mathbb{R}^p \setminus \mathcal{R}_0(\mathbf{r}_0)} \exp\{\mathbb{L}(\mathbf{v}, \mathbf{v}^*)\} f(\check{D}_0(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)) d\mathbf{v} \\ &= \exp\{-\Delta(\mathbf{r}_0, \mathbf{x}) - m(\boldsymbol{\xi})\} (1 - \tilde{\rho}_f(\mathbf{r}_0)) \int_{\mathbb{R}^p} \exp\{\mathbb{L}(\mathbf{v}, \mathbf{v}^*)\} f(\check{D}_0(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)) d\mathbf{v}. \end{aligned}$$

The definition (42) implies

$$\begin{aligned} & \int_{\mathcal{Y}} \exp\{L(\mathbf{v}, \mathbf{v}^*)\} f(\check{D}_0(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)) \mathbb{I}\{\boldsymbol{\theta} \in \Theta_0(\mathbf{r}_0)\} d\mathbf{v} \\ & \geq \exp\{-\Delta(\mathbf{r}_0, \mathbf{x}) - m(\boldsymbol{\xi})\} (1 - \tilde{\rho}_f(\mathbf{r}_0)) \int_{\Theta_0(\mathbf{r}_0) \times \mathbb{R}^{(p-q)}} \exp\{L(\mathbf{v}, \mathbf{v}^*)\} f(\check{D}_0(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)) d\mathbf{v} \\ & \geq \exp\{-\Delta(\mathbf{r}_0, \mathbf{x}) - m(\boldsymbol{\xi}) - 2\tilde{\rho}_f(\mathbf{r}_0)\} \mathbb{E}f(\boldsymbol{\gamma}) \mathbb{I}\{\|\boldsymbol{\gamma} + \check{\boldsymbol{\xi}}\| \leq \mathbf{r}_0\}. \end{aligned} \tag{49}$$

Here we used that $1 - \alpha \geq e^{-2\alpha}$ for $0 \leq \alpha \leq \frac{1}{2}$. Similarly,

$$\begin{aligned} \int_{\mathcal{Y}} \exp\{L(\mathbf{v}, \mathbf{v}^*)\} d\mathbf{v} &= \int_{\mathcal{Y}_0(\mathbf{r}_0)} \exp\{L(\mathbf{v}, \mathbf{v}^*)\} d\mathbf{v} + \int_{\mathcal{Y} \setminus \mathcal{Y}_0(\mathbf{r}_0)} \exp\{L(\mathbf{v}, \mathbf{v}^*)\} d\mathbf{v} \\ &= \{1 + \rho^*(\mathbf{r}_0)\} \int_{\mathcal{Y}_0(\mathbf{r}_0)} \exp\{L(\mathbf{v}, \mathbf{v}^*)\} d\mathbf{v} \\ &\leq \{1 + \rho^*(\mathbf{r}_0)\} \exp\{\Delta(\mathbf{r}_0, \mathbf{x}) - m(\boldsymbol{\xi})\} \mathbb{P}(\|\boldsymbol{\gamma} + \boldsymbol{\xi}\| \leq \mathbf{r}_0 \mid \mathbf{Y}), \end{aligned}$$

and finally

$$\int_{\mathcal{Y}} \exp\{L(\mathbf{v}, \mathbf{v}^*)\} d\mathbf{v} \leq \exp\{\Delta(\mathbf{r}_0, \mathbf{x}) - m(\boldsymbol{\xi}) + \nu(\mathbf{r}_0) + \rho^*(\mathbf{r}_0)\}. \tag{50}$$

The bounds (49) and (50) imply

$$\begin{aligned} & \frac{\int_{\mathcal{Y}} \exp\{L(\mathbf{v}, \mathbf{v}^*)\} f(\check{D}_0(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)) d\mathbf{v}}{\int_{\mathcal{Y}} \exp\{L(\mathbf{v}, \mathbf{v}^*)\} d\mathbf{v}} \\ & \geq \frac{\exp\{-\Delta(\mathbf{r}_0, \mathbf{x}) - m(\boldsymbol{\xi}) - 2\tilde{\rho}_f(\mathbf{r}_0)\} \mathbb{E}f(\boldsymbol{\gamma}) \mathbb{I}\{\|\boldsymbol{\gamma} + \check{\boldsymbol{\xi}}\| \leq \mathbf{r}_0\}}{\exp\{\Delta(\mathbf{r}_0, \mathbf{x}) - m(\boldsymbol{\xi}) + \nu(\mathbf{r}_0) + \rho^*(\mathbf{r}_0)\}} \\ & \geq \exp\{-2\Delta(\mathbf{r}_0, \mathbf{x}) - 2\tilde{\rho}_f(\mathbf{r}_0) - \nu(\mathbf{r}_0) - \rho^*(\mathbf{r}_0)\} \mathbb{E}f(\boldsymbol{\gamma}) \mathbb{I}\{\|\boldsymbol{\gamma} + \check{\boldsymbol{\xi}}\| \leq \mathbf{r}_0\}. \end{aligned}$$

This yields (41).

7.8 Proof of Theorem 1

Due to our previous results, it is convenient to decompose the r.v. $\boldsymbol{\vartheta}$ in the form

$$\boldsymbol{\vartheta} = \boldsymbol{\vartheta} \mathbb{I}\{\boldsymbol{\vartheta} \in \Theta_0(\mathbf{r}_0)\} + \boldsymbol{\vartheta} \mathbb{I}\{\boldsymbol{\vartheta} \notin \Theta_0(\mathbf{r}_0)\} = \boldsymbol{\vartheta}^\circ + \boldsymbol{\vartheta}^c.$$

The large deviation result yields that the posterior distribution of the part $\boldsymbol{\vartheta}^c$ is negligible provided a proper choice of \mathbf{r}_0 . Below we show that $\boldsymbol{\vartheta}^\circ$ is nearly normal which yields the BvM result. Define

$$\boldsymbol{\vartheta}^\circ \stackrel{\text{def}}{=} \mathbb{E}^\circ \boldsymbol{\vartheta}, \quad \mathfrak{S}_\circ^2 \stackrel{\text{def}}{=} \text{Cov}(\boldsymbol{\vartheta}^\circ) \stackrel{\text{def}}{=} \mathbb{E}^\circ \{(\boldsymbol{\vartheta} - \boldsymbol{\vartheta}^\circ)(\boldsymbol{\vartheta} - \boldsymbol{\vartheta}^\circ)^\top\}.$$

It suffices to show that holds on $\Omega(\mathbf{x})$

$$\begin{aligned}\|\check{D}_0(\boldsymbol{\vartheta}^\circ - \boldsymbol{\theta}^\circ)\|^2 &\leq 2\Delta^* \\ \|I_q - \check{D}_0\mathfrak{S}_\circ^2\check{D}_0\| &\leq 2\Delta^*,\end{aligned}$$

where $\Delta^* = \max\{\Delta^+, \Delta^-\}$.

Consider $\boldsymbol{\eta} \stackrel{\text{def}}{=} \check{D}_0(\boldsymbol{\vartheta} - \boldsymbol{\theta}^\circ)$. Corollaries 3 and 5 yield for any $\boldsymbol{\lambda} \in \mathbb{R}^q$ that

$$\|\boldsymbol{\lambda}\|^2 \exp(-\Delta^-) \leq \mathbb{E}^\circ |\boldsymbol{\lambda}^\top \boldsymbol{\eta}|^2 \leq \|\boldsymbol{\lambda}\|^2 \exp(\Delta^+). \quad (51)$$

Define the first two moments of $\boldsymbol{\eta}$:

$$\bar{\boldsymbol{\eta}} \stackrel{\text{def}}{=} \mathbb{E}^\circ \boldsymbol{\eta}, \quad S_\circ^2 \stackrel{\text{def}}{=} \mathbb{E}^\circ \{(\boldsymbol{\eta} - \bar{\boldsymbol{\eta}})(\boldsymbol{\eta} - \bar{\boldsymbol{\eta}})^\top\} = \check{D}_0\mathfrak{S}_\circ^2\check{D}_0.$$

Use the following technical statement.

Lemma 9. Assume (51). Then with $\Delta^* = \max\{\Delta^+, \Delta^-\} \leq 1/2$

$$\|\bar{\boldsymbol{\eta}}\|^2 \leq 2\Delta^*, \quad \|S_\circ^2 - I_q\| \leq 2\Delta^*. \quad (52)$$

Proof. Let \mathbf{u} be any unit vector in \mathbb{R}^q . We obtain from (51)

$$\exp(-\Delta^-) \leq \mathbb{E}^\circ |\mathbf{u}^\top \boldsymbol{\eta}|^2 \leq \exp(\Delta^+).$$

Note now that

$$\mathbb{E}^\circ |\mathbf{u}^\top \boldsymbol{\eta}|^2 = \mathbf{u}^\top S_\circ^2 \mathbf{u} + |\mathbf{u}^\top \bar{\boldsymbol{\eta}}|^2.$$

Hence

$$\exp(-\Delta^-) \leq \mathbf{u}^\top S_\circ^2 \mathbf{u} + |\mathbf{u}^\top \bar{\boldsymbol{\eta}}|^2 \leq \exp(\Delta^+). \quad (53)$$

In a similar way with $\mathbf{u} = \bar{\boldsymbol{\eta}}/\|\bar{\boldsymbol{\eta}}\|$ and $\boldsymbol{\gamma} \sim \mathcal{N}(0, I_q)$

$$\mathbb{E}^\circ |\mathbf{u}^\top (\boldsymbol{\eta} - \bar{\boldsymbol{\eta}})|^2 \geq e^{-\Delta^-} \mathbb{E} |\mathbf{u}^\top (\boldsymbol{\gamma} - \bar{\boldsymbol{\eta}})|^2 = e^{-\Delta^-} (1 + \|\bar{\boldsymbol{\eta}}\|^2)$$

yielding

$$\mathbf{u}^\top S_\circ^2 \mathbf{u} \geq (1 + \|\bar{\boldsymbol{\eta}}\|^2) \exp(-\Delta^-).$$

This inequality contradicts (53) if $\|\bar{\boldsymbol{\eta}}\|^2 > 2\Delta^* > 1$, and (52) follows. \square

The bound for the first moment implies with $\boldsymbol{\vartheta}^\circ = \mathbb{E}^\circ \boldsymbol{\vartheta}$

$$\|\check{D}_0(\boldsymbol{\vartheta}^\circ - \boldsymbol{\theta}^\circ)\|^2 \leq 2\Delta^*$$

while the second bound yields with

$$\|\check{D}_0 \check{\mathfrak{S}}_0^2 \check{D}_0 - I_q\| \leq 2\Delta^*.$$

The last result follows from (39) and (43) with an additional assumption that \mathbf{x} is large enough to ensure $\Delta^+(\mathbf{r}_0, \mathbf{x}) \leq 2\Delta(\mathbf{r}_0, \mathbf{x}) + 5e^{-\mathbf{x}}$ and $\Delta^-(\mathbf{r}_0, \mathbf{x}) \geq 2\Delta(\mathbf{r}_0, \mathbf{x}) - 8e^{-\mathbf{x}}$.

7.9 Proof of Theorem 2

Define $L_G(\mathbf{v}) = L(\mathbf{v}) - \|G\mathbf{v}\|^2/2$. The stochastic component of $L_G(\mathbf{v})$ coincides with one of $L(\mathbf{v})$. Also the quadratic term $\|G\mathbf{v}\|^2/2$ does not deteriorate the smoothness properties of the expected process $\mathbb{E}L_G(\mathbf{v})$. In particular, one can locally approximate $\mathbb{E}L_G(\mathbf{v}_G^\circ, \mathbf{v})$ by a quadratic function $\|\mathcal{D}_G(\mathbf{v} - \mathbf{v}_G^\circ)\|^2/2$ with

$$\begin{aligned} \mathbf{v}_G^\circ &\stackrel{\text{def}}{=} \mathbf{v}_G^* + \mathcal{D}_G^{-1}\boldsymbol{\xi}_G, & \boldsymbol{\xi}_G &\stackrel{\text{def}}{=} \mathcal{D}_G^{-1}\nabla\mathbb{E}L_G(\mathbf{v}^*) = \mathcal{D}_G^{-1}\nabla\mathbb{E}L(\mathbf{v}^*) + \mathcal{D}_G^{-1}G^2\mathbf{v}^*, \\ \mathbf{v}_G^* &\stackrel{\text{def}}{=} \underset{\mathbf{v}}{\operatorname{argmax}} \mathbb{E}L_G(\mathbf{v}), & \mathcal{D}_G^2 &\stackrel{\text{def}}{=} -\nabla^2\mathbb{E}L(\mathbf{v}^*) + G^2 = \mathcal{D}_0^2 + G^2. \end{aligned}$$

Now one can easily see that all the conditions of Theorem 1 are fulfilled for the process $L_G(\mathbf{v})$ when \mathbf{v}° is replaced by \mathbf{v}_G° and \mathcal{D}_0 by \mathcal{D}_G . The result approximates the posterior $\mathbf{v} | \mathbf{Y}$ for the Gaussian prior Π by the normal law $\mathcal{N}(\mathbf{v}_G^\circ, \mathcal{D}_G^{-2})$. Now the final result follows by Lemma 8 if we can bound $\|\mathcal{D}_0^{-1}\mathcal{D}_G^2\mathcal{D}_0^{-1} - I_p\|$ and $\|\mathcal{D}_G(\mathbf{v}^\circ - \mathbf{v}_G^\circ)\|$. By definition

$$\mathcal{D}_0^{-1}\mathcal{D}_G^2\mathcal{D}_0^{-1} - I_p = \mathcal{D}_0^{-1}G^2\mathcal{D}_0^{-1}.$$

Further, the definition ensures $\nabla\mathbb{E}L_G(\mathbf{v}) = \nabla\mathbb{E}L(\mathbf{v}) - G^2\mathbf{v}$ for any \mathbf{v} and also $\nabla\mathbb{E}L(\mathbf{v}^*) = \nabla\mathbb{E}L_G(\mathbf{v}_G^*) = 0$. This implies

$$\nabla\mathbb{E}L_G(\mathbf{v}_G^*) - \nabla\mathbb{E}L_G(\mathbf{v}^*) = \mathcal{D}_G^2(\check{\mathbf{v}})(\mathbf{v}_G^* - \mathbf{v}^*),$$

where $\check{\mathbf{v}}$ is a point from the interval connecting \mathbf{v}^* and \mathbf{v}_G^* and $\mathcal{D}_G^2(\check{\mathbf{v}}) = -\nabla^2\mathbb{E}L(\check{\mathbf{v}}) + G^2$. Therefore,

$$\mathcal{D}_G^2(\check{\mathbf{v}})(\mathbf{v}_G^* - \mathbf{v}^*) = G^2\mathbf{v}^*. \quad (54)$$

Let $\mathbf{r}_G = \|\mathcal{D}_G(\mathbf{v}^\circ - \mathbf{v}_G^\circ)\|$. It holds by (\mathcal{L}_0) and (54) with probability $\geq 1 - 2e^{-\mathbf{x}}$

$$\begin{aligned} \mathbf{r}_G &= \|\mathcal{D}_G(\mathbf{v}^\circ - \mathbf{v}_G^\circ)\| \leq \|\mathcal{D}_G(\mathbf{v}^* - \mathbf{v}_G^*)\| + \|\boldsymbol{\xi}_G - \mathcal{D}_G\mathcal{D}_0^{-1}\boldsymbol{\xi}\| \\ &\leq \|\mathcal{D}_G\mathcal{D}_G^{-2}(\check{\mathbf{v}})\mathcal{D}_G\mathcal{D}_G^{-1}G^2\mathbf{v}^*\| + \|\mathcal{D}_G^{-1}\mathcal{D}_0\boldsymbol{\xi} + \mathcal{D}_G^{-1}G^2\mathbf{v}^* - \mathcal{D}_G\mathcal{D}_0^{-1}\boldsymbol{\xi}\| \\ &\leq \{1 - \delta(\mathbf{r}_G)\}^{-1}\beta + \beta + \|\mathcal{D}_G^{-1}(\mathcal{D}_0^2 - \mathcal{D}_G^2)\mathcal{D}_0^{-1}\| \cdot \|\boldsymbol{\xi}\| \\ &\leq \{1 - \delta(\mathbf{r}_G)\}^{-1}\beta + \beta + \epsilon z_B(\mathbf{x}), \end{aligned}$$

where $z_B(\mathbf{x})$ is defined by (33). In particular, $\delta(\mathbf{r}_G) \leq 1/2$ and $\epsilon \leq 1/2$ imply $\mathbf{r}_G \leq 3\beta + \epsilon z_B(\mathbf{x})$. This yields by Lemma 8 for any $A^* \subset \mathbb{R}^p$:

$$|\mathbb{P}_0(A^*) - \mathbb{P}_G(A^*)| \leq \frac{1}{2} \sqrt{\delta^2 + (1 + \epsilon)(3\beta + \epsilon z_B(\mathbf{x}))}, \quad (55)$$

where measure \mathbb{P}_0 stands for Gaussian measure $\mathcal{N}(\mathbf{v}^\circ, \mathcal{D}_0^{-2})$ and measure \mathbb{P}_G stands for Gaussian measure $\mathcal{N}(\mathbf{v}_G^\circ, \mathcal{D}_G^{-2})$. Now we can take $A^* = A \times \mathbb{R}^{(p-q)}$. This yields by (55):

$$|\mathbb{P}_0(A^*) - \mathbb{P}_G(A^*)| = |\check{\mathbb{P}}_0(A) - \check{\mathbb{P}}_G(A)| \leq \frac{1}{2} \sqrt{\delta^2 + (1 + \epsilon)(3\beta + \epsilon z_B(\mathbf{x}))},$$

where measure $\check{\mathbb{P}}_0$ stands for Gaussian measure $\mathcal{N}(\boldsymbol{\theta}^\circ, \check{\mathcal{D}}_0^{-2})$ and measure $\check{\mathbb{P}}_G$ stands for Gaussian measure $\mathcal{N}(\boldsymbol{\theta}_G^\circ, \check{\mathcal{D}}_G^{-2})$. This yields the result of the theorem.

7.10 Proof of Theorem 3

First we compare the deterministic quantities \check{D}_s^2 and $\boldsymbol{\theta}^*$ with their sieve counterparts. The identifiability condition (\mathcal{I}^s) guarantees for $\check{D}_s^2 = D_\theta^2 - A_{\theta\eta} A_{\theta\eta}^\top$ that

$$\|D_\theta^{-1} \check{D}_s^2 D_\theta^{-1}\| \geq 1 - \nu. \quad (56)$$

Further $\check{D}^2 = D_\theta^2 - A_{\theta\phi} A_{\theta\phi}^\top$ and $\check{D}_s^2 = D_\theta^2 - A_{\theta\eta} A_{\theta\eta}^\top$

$$\check{D}_s^2 - \check{D}^2 = D_\theta^2 - A_{\theta\eta} A_{\theta\eta}^\top - (D_\theta^2 - A_{\theta\phi} A_{\theta\phi}^\top) = A_{\theta\kappa} A_{\theta\kappa}^\top.$$

The use of the smoothness condition (B) implies

$$\begin{aligned} \|D_\theta^{-1} (\check{D}_s^2 - \check{D}^2) D_\theta^{-1}\| &= \|D_\theta^{-1} A_{\theta\kappa} A_{\theta\kappa}^\top D_\theta^{-1}\| \leq \rho_s, \\ \|\check{D}^{-1} (\check{D}_s^2 - \check{D}^2) \check{D}^{-1}\| &\leq (1 - \nu)^{-1} \|D_\theta^{-1} A_{\theta\kappa} A_{\theta\kappa}^\top D_\theta^{-1}\| \leq (1 - \nu)^{-1} \rho_s. \end{aligned} \quad (57)$$

Next we bound the bias $\boldsymbol{\theta}^* - \boldsymbol{\theta}_s^*$ introduced in $\boldsymbol{\theta}^*$ by truncation. Consider first the Gaussian likelihood with

$$\mathbb{E}L(\mathbf{v}^*, \boldsymbol{\varkappa}^*) - \mathbb{E}L(\mathbf{v}, \boldsymbol{\varkappa}) = \|\mathcal{D}_0\{(\mathbf{v}, \boldsymbol{\varkappa}) - (\mathbf{v}^*, \boldsymbol{\varkappa}^*)\}\|^2/2.$$

Define $\mathbf{v}_s^* = (\boldsymbol{\theta}_s^*, \boldsymbol{\eta}_s^*)$ as the minimizer of the quadratic function $\|\mathcal{D}_0((\mathbf{v}, 0) - (\mathbf{v}^*, \boldsymbol{\varkappa}^*))\|^2$ over the set of parameters $(\mathbf{v}, 0) = (\boldsymbol{\theta}, \boldsymbol{\eta}, 0)$:

$$\begin{aligned} \mathbf{v}_s^* &= \operatorname{argmin}_{\mathbf{v}=(\boldsymbol{\theta}, \boldsymbol{\eta})} \|\mathcal{D}_0((\mathbf{v}, 0) - (\mathbf{v}^*, \boldsymbol{\varkappa}^*))\|^2 \\ &= \operatorname{argmin}_{(\boldsymbol{\theta}, \boldsymbol{\eta})} \{(\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top D_\theta^2 (\boldsymbol{\theta} - \boldsymbol{\theta}^*) + \|\boldsymbol{\eta} - \boldsymbol{\eta}^*\|^2 \\ &\quad + 2(\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top A_{\theta\eta} (\boldsymbol{\eta} - \boldsymbol{\eta}^*) - 2(\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top A_{\theta\kappa} \boldsymbol{\varkappa}^*\}. \end{aligned}$$

This implies by direct calculus that $\mathbf{v}_s^* = (\boldsymbol{\theta}_s^*, \boldsymbol{\eta}_s^*)$ fulfills

$$\begin{aligned} \boldsymbol{\eta}_s^* - \boldsymbol{\eta}^* &= -A_{\boldsymbol{\theta}\boldsymbol{\eta}}^\top (\boldsymbol{\theta}_s^* - \boldsymbol{\theta}^*), \\ \boldsymbol{\theta}_s^* - \boldsymbol{\theta}^* &= (D_{\boldsymbol{\theta}}^2 - A_{\boldsymbol{\theta}\boldsymbol{\eta}}A_{\boldsymbol{\theta}\boldsymbol{\eta}}^\top)^{-1}A_{\boldsymbol{\theta}\boldsymbol{\varkappa}}\boldsymbol{\varkappa}^* = \check{D}_s^{-2}A_{\boldsymbol{\theta}\boldsymbol{\varkappa}}\boldsymbol{\varkappa}^*. \end{aligned} \tag{58}$$

The identifiability condition (56) and the smoothness conditions (22), (23) imply

$$\begin{aligned} \|\check{D}_s(\boldsymbol{\theta}_s^* - \boldsymbol{\theta}^*)\| &= \|\check{D}_s^{-1}A_{\boldsymbol{\theta}\boldsymbol{\varkappa}}\boldsymbol{\varkappa}^*\| \\ &\leq (1 - \nu)^{-1}\|D_{\boldsymbol{\theta}}^{-1}A_{\boldsymbol{\theta}\boldsymbol{\varkappa}}\boldsymbol{\varkappa}^*\| \leq (1 - \nu)^{-1}b_s. \end{aligned} \tag{59}$$

In the general non-Gaussian case, we use that

$$\|\mathcal{D}_0^{-1}\{\nabla EL(\mathbf{v}^*, \boldsymbol{\varkappa}^*) - \nabla EL(\mathbf{v}_s^*, 0)\} - \mathcal{D}_0\{(\mathbf{v}^*, \boldsymbol{\varkappa}^*) - (\mathbf{v}_s^*, 0)\}\| \leq \delta(\mathbf{r}_s)\mathbf{r}_s$$

with

$$\mathbf{r}_s = \|\mathcal{D}_0\{(\mathbf{v}^*, \boldsymbol{\varkappa}^*) - (\mathbf{v}_s^*, 0)\}\| \leq \|\mathcal{D}_0\{(\mathbf{v}^*, \boldsymbol{\varkappa}^*) - (\mathbf{v}^*, 0)\}\| = \|\boldsymbol{\varkappa}^*\|.$$

The definition of $(\boldsymbol{\theta}^*, \boldsymbol{\eta}^*, \boldsymbol{\varkappa}^*)$ and \mathbf{v}_s^* imply that $\nabla_{\mathbf{v}}EL(\boldsymbol{\theta}^*, \boldsymbol{\eta}^*, \boldsymbol{\varkappa}^*) = 0$ and $\nabla_{\mathbf{v}}EL(\mathbf{v}_s^*, 0) = 0$. Now, projecting on the $\boldsymbol{\theta}$ -subspace implies similarly to (58)

$$\|\check{D}_s(\boldsymbol{\theta}^* - \boldsymbol{\theta}_s^*) - \check{D}_s^{-1}A_{\boldsymbol{\theta}\boldsymbol{\varkappa}}\boldsymbol{\varkappa}^*\| \leq \delta(\mathbf{r}_s)\mathbf{r}_s.$$

This implies by (59)

$$\|\check{D}_s(\boldsymbol{\theta}^* - \boldsymbol{\theta}_s^*)\| \leq (1 - \nu)^{-1}b_s + \delta(\mathbf{r}_s)\mathbf{r}_s,$$

which completes the proof of (13).

Now we consider the stochastic part $\check{D}_s^{-1}\check{\boldsymbol{\xi}}_s - \check{D}^{-1}\check{\boldsymbol{\xi}}$ of $\boldsymbol{\theta}^\circ - \boldsymbol{\theta}_s^\circ$. The bounds (57) and (56) imply for $\check{\boldsymbol{\xi}} = \nabla_{\boldsymbol{\theta}} - A_{\boldsymbol{\theta}\boldsymbol{\eta}}\nabla_{\boldsymbol{\eta}} - A_{\boldsymbol{\theta}\boldsymbol{\varkappa}}\nabla_{\boldsymbol{\varkappa}}$ and $\check{\boldsymbol{\xi}}_s = \nabla_{\boldsymbol{\theta}} - A_{\boldsymbol{\theta}\boldsymbol{\eta}}\nabla_{\boldsymbol{\eta}}$

$$\begin{aligned} \|\check{D}\{\check{D}_s^{-1}\check{\boldsymbol{\xi}}_s - \check{D}^{-1}\check{\boldsymbol{\xi}}\}\| &= \|\check{D}\check{D}_s^{-2}(\nabla_{\boldsymbol{\theta}} - A_{\boldsymbol{\theta}\boldsymbol{\eta}}\nabla_{\boldsymbol{\eta}}) - \check{D}^{-1}(\nabla_{\boldsymbol{\theta}} - A_{\boldsymbol{\theta}\boldsymbol{\eta}}\nabla_{\boldsymbol{\eta}} - A_{\boldsymbol{\theta}\boldsymbol{\varkappa}}\nabla_{\boldsymbol{\varkappa}})\| \\ &= \|\{\check{D}\check{D}_s^{-2} - \check{D}^{-1}\}(\nabla_{\boldsymbol{\theta}} - A_{\boldsymbol{\theta}\boldsymbol{\eta}}\nabla_{\boldsymbol{\eta}}) + \check{D}^{-1}A_{\boldsymbol{\theta}\boldsymbol{\varkappa}}\nabla_{\boldsymbol{\varkappa}}\| \\ &\leq \|\{I_q - \check{D}^{-1}\check{D}_s^2\check{D}^{-1}\}\check{D}\check{D}_s^{-2}(\nabla_{\boldsymbol{\theta}} - A_{\boldsymbol{\theta}\boldsymbol{\eta}}\nabla_{\boldsymbol{\eta}})\| + \|\check{D}^{-1}A_{\boldsymbol{\theta}\boldsymbol{\varkappa}}\nabla_{\boldsymbol{\varkappa}}\| \\ &\leq (1 - \nu)^{-1}\rho_s\|\check{\boldsymbol{\xi}}_s\| + (1 - \nu)^{-1}\|D_{\boldsymbol{\theta}}^{-1}A_{\boldsymbol{\theta}\boldsymbol{\varkappa}}\nabla_{\boldsymbol{\varkappa}}\| \\ &= (1 - \nu)^{-1}\rho_s(\|\check{\boldsymbol{\xi}}_s\| + \|\boldsymbol{\xi}_{\boldsymbol{\varkappa}}\|). \end{aligned}$$

Here $\boldsymbol{\xi}_{\boldsymbol{\varkappa}} = \rho_s^{-1}D_{\boldsymbol{\theta}}^{-1}A_{\boldsymbol{\theta}\boldsymbol{\varkappa}}\nabla_{\boldsymbol{\varkappa}}$ and we have also used that $\check{D}_s^2 \geq \check{D}^2$. This proves (14). It remains to check (15). If the full model is true then $\text{Var}(\check{\boldsymbol{\xi}}_s) = I_q$ and under (ED_0) , it holds on a set of probability at least $1 - 2e^{-x}$

$$\|\check{\boldsymbol{\xi}}_s\| \leq \sqrt{q} + 2x; \tag{60}$$

see Spokoiny (2012). Similarly, under the correct model specification, it holds $\text{Var}(\nabla_{\boldsymbol{x}}) = \boldsymbol{I}_{\boldsymbol{x}}$. For $\boldsymbol{\xi}_{\boldsymbol{x}} = \rho_s^{-1} D_{\boldsymbol{\theta}}^{-1} A_{\boldsymbol{\theta}_{\boldsymbol{x}}} \nabla_{\boldsymbol{x}}$, it holds by (B)

$$\text{Var}(\boldsymbol{\xi}_{\boldsymbol{x}}) \leq \rho_s^{-2} \|D_{\boldsymbol{\theta}}^{-1} A_{\boldsymbol{\theta}_{\boldsymbol{x}}} A_{\boldsymbol{\theta}_{\boldsymbol{x}}}^{\top} D_{\boldsymbol{\theta}}^{-1}\| \leq 1,$$

and hence, the use of (ED_0) implies with a dominating probability $1 - 2e^{-\boldsymbol{x}}$

$$\|\boldsymbol{\xi}_{\boldsymbol{x}}\| \leq \sqrt{q} + 2\boldsymbol{x}$$

similarly to (60). In the general situation, if the parametric assumption is not exactly true, we still can use (ED_0) . It ensures $\text{Var}(\check{\boldsymbol{\xi}}_s) \leq \boldsymbol{\alpha}^2 \boldsymbol{I}_q$, $\text{Var}(\boldsymbol{\xi}_{\boldsymbol{x}}) \leq \boldsymbol{\alpha}^2 \boldsymbol{I}_q$, and

$$\|\check{\boldsymbol{\xi}}_s\| \leq \boldsymbol{\alpha}(\sqrt{q} + 2\boldsymbol{x}), \quad \|\boldsymbol{\xi}_{\boldsymbol{x}}\| \leq \boldsymbol{\alpha}(\sqrt{q} + 2\boldsymbol{x}).$$

This yields the last claim of the theorem.

7.11 Proof of Theorem 5

First we check that the required conditions of Section 2.1 are fulfilled in the considered example. This can be easily done if we slightly change the definition of the local set $\mathcal{Y}_0(\mathbf{r}_0)$. Namely, for $\mathbf{u}^* = (u_1^*, \dots, u_{p_n}^*)^{\top}$, define $\mathcal{Y}_0(\sqrt{\mathfrak{z}})$ as a rectangle

$$\mathcal{Y}_0(\sqrt{\mathfrak{z}}) \stackrel{\text{def}}{=} \{\mathbf{u} : M_n \mathcal{K}(u_j, u_j^*) \leq \mathfrak{z}, j = 1, \dots, p_n\}.$$

Here $\mathcal{K}(u, u^*)$ is the Kullback-Leibler divergence for the Poisson family:

$$\mathcal{K}(u, u^*) = e^u(u - u^*) - e^u + e^{u^*}.$$

Lemma 10. *Let \mathfrak{z}_n be such that $2p_n e^{-\mathfrak{z}_n} \leq 1/2$. Then it holds*

$$\mathbb{P}(\tilde{\mathbf{u}} \in \mathcal{Y}_0(\sqrt{\mathfrak{z}_n})) \geq 1 - 4p_n e^{-\mathfrak{z}_n}. \quad (61)$$

In particular, the choice $\mathfrak{z}_n = \mathbf{x}_n + \log(p_n)$ with $\mathbf{x}_n = \mathbf{C} \log n$ provides

$$\mathbb{P}(\tilde{\mathbf{u}} \in \mathcal{Y}_0(\sqrt{\mathfrak{z}_n})) \geq 1 - 4e^{-\mathbf{x}_n}. \quad (62)$$

Proof. We use the bound from Polzehl and Spokoiny (2006)

$$\mathbb{P}(M_n \mathcal{K}(\tilde{u}_j, u_j^*) > \mathfrak{z}_n) \leq 2e^{-\mathfrak{z}_n}.$$

This yields

$$\mathbb{P}(\tilde{\mathbf{u}} \in \mathcal{Y}_0(\sqrt{\mathfrak{z}_n})) \geq (1 - 2e^{-\mathfrak{z}_n})^{p_n}.$$

Now the elementary inequalities $\log(1 - \alpha) \geq -2\alpha$ for $0 \leq \alpha \leq 1/2$ and $e^{-\delta} \geq 1 - \delta$ for $\delta \geq 0$ applied with $\alpha_n = 2e^{-\mathfrak{z}_n}$ and $\delta_n = 2\alpha_n p_n$ imply

$$(1 - \alpha_n)^{p_n} = e^{\log(1 - \alpha_n)p_n} \geq e^{-2\alpha_n p_n} \geq 1 - 2\alpha_n p_n$$

and (61) follows. □

In the special case $u_1^* = \dots = u_{p_n}^* = u^*$, the set $\mathcal{Y}_0(\sqrt{\mathfrak{z}})$ is a cube which can be also viewed as a ball in the sup-norm. Moreover, if $\mathfrak{z}_n/(M_n e^{u^*}) \leq 1/2$, this cube is contained in the cube $\{\mathbf{u} : \|\mathbf{u} - \mathbf{u}^*\| \leq \sqrt{\mathfrak{z}_n/(M_n e^{u^*})}\}$ in view of $e^x - 1 - x \leq a^2 \leq 1/2$ for $|x| \leq a \leq 1$. The concentration bound (62) enables us to check the local conditions only on the cube $\mathcal{Y}_0(\sqrt{\mathfrak{z}_n})$. Especially the condition $(\mathcal{E}\mathcal{D}_1)$ is trivially fulfilled because $\zeta(\mathbf{u}) = \mathcal{L}(\mathbf{u}) - \mathbb{E}\mathcal{L}(\mathbf{u})$ is linear in \mathbf{u} and θ is a linear functional of \mathbf{u} . Condition (\mathcal{L}_0) can be checked on $\mathcal{Y}_0(\sqrt{\mathfrak{z}_n})$ with $\delta(\mathfrak{z}_n) = \sqrt{\mathfrak{z}_n/(M_n e^{u^*})}$.

It remains to compute the value \check{D}_0^2 . Define $\beta_n = p_n/M_n^{1/2} = p_n^{3/2}/n^{1/2}$. If $n = p_n^3$, then $\beta_n = 1$.

Lemma 11. *Let $v^* = 1/p_n$. Then it holds*

$$\check{D}_0^2 = p_n^2 \beta_n^{-2}.$$

Now we are ready to finalize the proof Theorem 5.

Proof. Let β_n be bounded. The definition implies

$$p_n(\theta - \tilde{\theta}_n) = \sum_{j=1}^{p_n} \log\left(\frac{v_j}{Z_j/M_n}\right).$$

The posterior distribution $v_j | \mathbf{Y}$ is $\text{Gamma}(\alpha_j, \mu_j)$ with $\alpha_j = 1 + Z_j$ and $\mu_j = \frac{\mu}{M_n \mu + 1}$. We use following decomposition

$$\frac{v_j}{Z_j/M_n} = \frac{M_n \mu_j \alpha_j}{\alpha_j - 1} (1 + \alpha_j^{-1/2} \gamma_j),$$

where

$$\gamma_j \stackrel{\text{def}}{=} (\alpha_j \mu_j^2)^{-1/2} (v_j - \alpha_j \mu_j)$$

has zero mean and unit variance. We can use the Taylor expansion

$$p_n(\theta - \tilde{\theta}_n) = \sum_{j=1}^{p_n} \log\left(1 - \frac{1}{M_n \mu + 1}\right) + \sum_{j=1}^{p_n} \log\left(1 + \frac{1}{\alpha_j - 1}\right) + \sum_{j=1}^{p_n} \log\left(1 + \alpha_j^{-1/2} \gamma_j\right).$$

Now take into account properties of the real data distribution.

$$\alpha_j = \frac{M_n}{p_n} \left(1 + \sqrt{\frac{p_n}{M_n}} \delta_j \right),$$

where δ_j is asymptotically standard normal.

Suppose now that $\beta_n^3/\sqrt{p_n} \rightarrow 0$ as $p_n \rightarrow \infty$. Then $M_n/p_n = (\sqrt{p_n}/\beta_n^3)^{2/3} p_n^{2/3} \rightarrow \infty$ as $p_n \rightarrow \infty$. Thus for p_n sufficient large, $\alpha_j \approx M_n/p_n$. Moreover, it holds for p_n sufficiently large that $\max_{j=1, \dots, p_n} \alpha_j^{-1/2} |\gamma_j| \leq 1/2$ with a high probability. Below we can restrict ourselves to the case when $\alpha_j^{-1/2} |\gamma_j| \leq 1/2$. This allows to use the Taylor expansion

$$\begin{aligned} p_n(\theta - \tilde{\theta}_n) &= \sum_{j=1}^{p_n} \log\left(1 - \frac{1}{M_n \mu + 1}\right) + \sum_{j=1}^{p_n} \log\left(1 + \frac{1}{\alpha_j - 1}\right) + \sum_{j=1}^{p_n} \log\left(1 + \frac{\gamma_j}{\sqrt{\alpha_j}}\right) \\ &= \sum_{j=1}^{p_n} \frac{1}{\alpha_j - 1} + \sum_{j=1}^{p_n} \frac{1}{\sqrt{\alpha_j}} \gamma_j - \sum_{j=1}^{p_n} \frac{1}{2\alpha_j} \gamma_j^2 + R. \end{aligned}$$

One can easily check that the remainder R is of order $\beta_n^3/\sqrt{p_n} \rightarrow 0$. Moreover, $p_n^{-1/2} \sum_{j=1}^{p_n} \gamma_j$ is asymptotically standard normal, while $p_n^{-1} \sum_{j=1}^{p_n} \gamma_j^2 \xrightarrow{\mathbb{P}} 1$. The central limit theorem here can be easily checked because of the Lyapunov condition being valid. Also $\sum_{j=1}^{p_n} (\alpha_j - 1)^{-1} = \frac{p_n^2}{M_n} + o_n(\beta_n^2)$. Now check what happens if $\beta_n \rightarrow 0$:

$$\beta_n^{-1} p_n(\theta - \tilde{\theta}_n) = \beta_n + \frac{1}{\sqrt{p_n}} \sum_{j=1}^{p_n} \gamma_j - \frac{\beta_n}{2p_n} \sum_{j=1}^{p_n} \gamma_j^2 + o_n(1) \xrightarrow{w} \mathcal{N}(0, 1).$$

Similarly, with $\beta_n \equiv \beta$,

$$\beta^{-1} p_n(\theta - \tilde{\theta}_n) = \beta + \frac{1}{\sqrt{p_n}} \sum_{j=1}^{p_n} \gamma_j - \frac{\beta}{2p_n} \sum_{j=1}^{p_n} \gamma_j^2 + o_n(1) \xrightarrow{w} \mathcal{N}(\beta/2, 1).$$

This proves the result for $\beta_n \equiv \beta$. Finally in the case when β_n grows to infinity, but $\beta_n^3/\sqrt{p_n} \rightarrow 0$, then $\beta_n^{-1}(\theta - \tilde{\theta}_n) \xrightarrow{\mathbb{P}} \infty$. \square

7.12 Proof of Lemma 11

Let $\bar{u}_j = u_j - u_j^*$. Then

$$L(\mathbf{u}, \mathbf{u}^*) = \mathcal{L}(\mathbf{u}) - \mathcal{L}(\mathbf{u}^*) = \sum_{j=1}^{p_n} \{Z_j \bar{u}_j - M_n p_n^{-1} (e^{\bar{u}_j} - 1)\}.$$

The expected value of Z_j is M_n/p_n which leads to the following expectation of likelihood:

$$\mathbb{E}L(\mathbf{u}, \mathbf{u}^*) = \frac{M_n}{p_n} \sum_{j=1}^{p_n} (\bar{u}_j - (e^{\bar{u}_j} - 1)) = -\frac{M_n}{p_n} \sum_{j=1}^{p_n} \frac{\bar{u}_j^2}{2} + O(\|\bar{\mathbf{u}}\|^3).$$

Then we substitute $\bar{u}_1 = p_n \bar{\theta} - \sum_{j=2}^{p_n} \bar{u}_j$, where $\bar{\theta} = \theta - \theta^*$. Thus we get

$$\mathbb{E}L(\mathbf{u}, \mathbf{u}^*) = -\frac{M_n}{p_n} \frac{1}{2} (p_n \bar{\theta} - \sum_{j=2}^{p_n} \bar{u}_j)^2 - \frac{M_n}{p_n} \sum_{j=2}^{p_n} \frac{\bar{u}_j^2}{2} + O(\|\bar{\mathbf{u}}\|^3).$$

This Taylor expansion allows us to compute components of the Fisher information matrix:

$$\mathcal{D}_0^2 = -\nabla^2 \mathbb{E}L(\mathbf{u}^*) = \frac{M_n}{p_n} \begin{pmatrix} p_n^2 & -p_n & \dots & \dots & -p_n \\ -p_n & 2 & 1 & \dots & 1 \\ \vdots & 1 & \ddots & \ddots & \vdots \\ \vdots & \vdots & \ddots & \ddots & 1 \\ -p_n & 1 & \dots & 1 & 2 \end{pmatrix}.$$

The Fisher information for the target parameter θ can be computed as follows:

$$\check{D}_0^2 = M_n p_n (1 - \mathbf{e}^\top \mathbb{H}^{-1} \mathbf{e}),$$

where $\mathbf{e} = (1, \dots, 1)^\top$ and $\mathbb{H} = I + E$ with $E = \mathbf{e}\mathbf{e}^\top$ being the matrix of ones of size $(p_n - 1) \times (p_n - 1)$. It follows

$$\mathbf{e}^\top \mathbb{H}^{-1} \mathbf{e} = \text{tr}(\mathbf{e}^\top \mathbb{H}^{-1} \mathbf{e}) = \text{tr}(\mathbb{H}^{-1} \mathbf{e}\mathbf{e}^\top) = \text{tr}((E + I)^{-1} E).$$

Further, $(E + I)^{-1} E = I - (E + I)^{-1}$ yielding

$$\mathbf{e}^\top \mathbb{H}^{-1} \mathbf{e} = \text{tr}\{I - (E + I)^{-1}\} = (p_n - 1) - \text{tr}\{(E + I)^{-1}\} = (p_n - 1) - \sum_{j=1}^{p_n} \lambda_j,$$

where λ_j are eigenvalues of matrix $(E + I)^{-1}$. It is easy to see that $\lambda_1 = p_n^{-1}$ while $\lambda_2 = \dots = \lambda_{p_n-1} = 1$. Thus

$$\begin{aligned} \mathbf{e}^\top \mathbb{H}^{-1} \mathbf{e} &= (p_n - 1) - \{p_n^{-1} + (p_n - 2)\} = 1 - p_n^{-1}, \\ \check{D}_0^2 &= M_n p_n (1 - \mathbf{e}^\top \mathbb{H}^{-1} \mathbf{e}) = M_n p_n \{1 - (1 - p_n^{-1})\} = M_n = p_n^2 \beta_n^{-2}, \end{aligned}$$

which completes the proof.

7.13 Proof of Lemma 1

It holds

$$\begin{aligned}
& \sup_{\gamma \in \mathbb{R}^p} \log \mathbb{E} \exp \left\{ \mu \frac{\gamma^\top \nabla \zeta(\mathbf{v}^*)}{\|\mathcal{D}_0 \gamma\|} \right\} \\
&= \sup_{\gamma \in \mathbb{R}^p} \log \mathbb{E} \exp \left\{ \mu \frac{\sum_{i=1}^n (h'(Y_i - \Psi_i^\top \mathbf{v}^*) - \mathbb{E} h'(Y_i - \Psi_i^\top \mathbf{v}^*)) \gamma^\top \Psi_i}{\|\mathcal{D}_0 \gamma\|} \right\} \\
&= \sup_{\gamma \in \mathbb{R}^p} \log \mathbb{E} \exp \left\{ \frac{\mu h \sum_{i=1}^n \gamma^\top \Psi_i}{\|\mathcal{D}_0 \gamma\|} h'(\varepsilon)/h \right\}.
\end{aligned}$$

By definition $h|\gamma^\top \Psi_i|/\|\mathcal{D}_0 \gamma\| \leq N_1^{-1/2}$ and hence $\mu h|\gamma^\top \Psi_i|/\|\mathcal{D}_0 \gamma\| \leq \mathbf{g}_1$. Thus,

$$\sup_{\gamma \in \mathbb{R}^p} \log \mathbb{E} \exp \left\{ \mu \frac{\gamma^\top \nabla \zeta(\mathbf{v}^*)}{\|\mathcal{D}_0 \gamma\|} \right\} \leq \sup_{\gamma \in \mathbb{R}^p} \frac{\nu_0^2 \mu^2 h^2 \sum_{i=1}^n |\Psi_i^\top \gamma|^2}{2 \|\mathcal{D}_0 \gamma\|^2} = \frac{\nu_0^2 \mu^2}{2}.$$

7.14 Proof of Lemma 5

It holds

$$\begin{aligned}
\|\mathcal{D}_0^{-1}(\mathcal{D}^2(\mathbf{v}) - \mathcal{D}_0^2)\mathcal{D}_0^{-1}\| &= \sup_{\gamma \in \mathbb{R}^p: \|\gamma\|=1} |\gamma^\top \mathcal{D}_0^{-1}(\mathcal{D}^2(\mathbf{v}) - \mathcal{D}_0^2)\mathcal{D}_0^{-1}\gamma| \\
&\leq \sup_{\gamma \in \mathbb{R}^p: \|\gamma\|=1} \left| \sum_{i=1}^n (d''(\Psi_i^\top \mathbf{v}) - d''(\Psi_i^\top \mathbf{v}^*)) \gamma^\top \mathcal{D}_0^{-1} \Psi_i \Psi_i^\top \mathcal{D}_0^{-1} \gamma \right| \\
&\leq \sup_{\gamma \in \mathbb{R}^p: \|\gamma\|=1} \sum_{i=1}^n |d''(\Psi_i^\top \mathbf{v}) - d''(\Psi_i^\top \mathbf{v}^*)| \gamma^\top \mathcal{D}_0^{-1} \Psi_i \Psi_i^\top \mathcal{D}_0^{-1} \gamma \\
&\leq \sup_{\gamma \in \mathbb{R}^p: \|\gamma\|=1} \sum_{i=1}^n L |\Psi_i^\top (\mathbf{v} - \mathbf{v}^*)| \gamma^\top \mathcal{D}_0^{-1} \Psi_i \Psi_i^\top \mathcal{D}_0^{-1} \gamma \\
&\leq LN^{-1/2} \|\mathcal{D}_0(\mathbf{v} - \mathbf{v}^*)\| \sup_{\gamma \in \mathbb{R}^p: \|\gamma\|=1} \gamma^\top \mathcal{D}_0^{-1} \left(\sum_{i=1}^n d''(\Psi_i^\top \mathbf{v}^*) \Psi_i \Psi_i^\top \right) \mathcal{D}_0^{-1} \gamma \\
&\leq L \frac{\mathbf{r}}{N_2^{1/2}}.
\end{aligned}$$

References

- Barron, A., Schervish, M. J., and Wasserman, L. (1996). "The Consistency of Posterior Distributions in Nonparametric Problems." *The Annals of Statistics*, 27: 536–561. [665](#)

- Bickel, P. J. and Kleijn, B. J. K. (2012). “The semiparametric Bernstein-von Mises theorem.” *The Annals of Statistics*, 40(1): 206–237. [666](#), [676](#)
- Bochkina, N. and Green, P. J. (2014). “The Bernstein-von Mises theorem and non-regular models.” *The Annals of Statistics*, 42(5): 1850–1878. [668](#)
- Bontemps, D. (2011). “Bernstein–von Mises theorem for Gaussian regression with increasing number of regressors.” *The Annals of Statistics*, 39(5): 2557–2584. [666](#), [674](#), [676](#), [680](#)
- Boucheron, S. and Gassiat, E. (2009). “A Bernstein-von Mises theorem for discrete probability distributions.” *Electronic Journal of Statistics*, 3: 114–148. [666](#), [677](#)
- Boucheron, S. and Massart, P. (2011). “A high-dimensional Wilks phenomenon.” *Probability Theory and Related Fields*, 150: 405–433. [666](#)
- Castillo, I. (2012). “A semiparametric Bernstein–von Mises theorem for Gaussian process priors.” *Probability Theory and Related Fields*, 152: 53–99. [666](#), [676](#)
- Castillo, I. and Nickl, R. (2013). “Nonparametric Bernstein–von Mises theorems in Gaussian white noise.” *The Annals of Statistics*, 41(4): 1999–2028. [666](#)
- Castillo, I. and Rousseau, J. (2013). “A General Bernstein–von Mises Theorem in semiparametric models.” Available at arXiv:[1305.4482](#) [math.ST]. [666](#)
- Cheng, G. and Kosorok, M. R. (2008). “General frequentist properties of the posterior profile distribution.” *The Annals of Statistics*, 36(4): 1819–1853. [666](#)
- Chernozhukov, V. and Hong, H. (2003). “An MCMC approach to classical estimation.” *Journal of Econometrics*, 115(2): 293–346. [667](#), [668](#)
- Cox, D. D. (1993). “An analysis of Bayesian inference for nonparametric regression.” *The Annals of Statistics*, 21(2): 903–923. [665](#)
- Freedman, D. (1999). “On the Bernstein-von Mises theorem with infinite-dimensional parameters.” *The Annals of Statistics*, 27(4): 1119–1140. [665](#)
- Ghosal, S. (1999). “Asymptotic normality of posterior distributions in high-dimensional linear models.” *Bernoulli*, 5(2): 315–331. [665](#), [666](#), [677](#)
- (2000). “Asymptotic normality of posterior distributions for exponential families when the number of parameters tends to infinity.” *Journal of Multivariate Analysis*, 74(1): 49–68. [666](#), [677](#)
- Ibragimov, I. and Khas'minskij, R. (1981). *Statistical estimation. Asymptotic theory. Translated from the Russian by Samuel Kotz*. New York–Heidelberg–Berlin: Springer-Verlag. [669](#)
- Johnstone, I. M. (2010). “High dimensional Bernstein–von Mises: simple examples.” In *Borrowing strength: theory powering applications—a Festschrift for Lawrence D. Brown*, volume 6 of *Institute of Mathematical Statistics Collections*, 87–98. Beachwood, OH: Institute of Mathematical Statistics. [665](#), [666](#), [674](#), [677](#), [680](#)
- Kim, Y. (2006). “The Bernstein–von Mises theorem for the proportional hazard model.” *The Annals of Statistics*, 34(4): 1678–1700. [666](#)

- Kim, Y. and Lee, J. (2004). “A Bernstein–von Mises theorem in the nonparametric right-censoring model.” *The Annals of Statistics*, 32(4): 1492–1512. [666](#)
- Kleijn, B. J. K. and van der Vaart, A. W. (2006). “Misspecification in infinite-dimensional Bayesian statistics.” *The Annals of Statistics*, 34(2): 837–877. [667](#)
- (2012). “The Bernstein-von-Mises theorem under misspecification.” *Electronic Journal of Statistics*, 6: 354–381. [667](#)
- Le Cam, L. and Yang, G. L. (1990). *Asymptotics in Statistics: Some Basic Concepts*. Springer in Statistics. [665](#)
- Leahu, H. (2011). “On the Bernstein-von Mises phenomenon in the Gaussian white noise model.” *Electronic Journal of Statistics*, 5: 373–404. [666](#)
- McCullagh, P. and Nelder, J. (1989). *Generalized linear models. 2nd ed.* Monographs on Statistics and Applied Probability. 37. London etc.: Chapman and Hall. [684](#)
- Polzehl, J. and Spokoiny, V. (2006). “Propagation-separation approach for local likelihood estimation.” *Probability Theory and Related Fields*, 135(3): 335–362. [704](#)
- Rivoirard, V. and Rousseau, J. (2012). “Bernstein–von Mises theorem for linear functionals of the density.” *The Annals of Statistics*, 40(3): 1489–1523. [666](#)
- Schwartz, L. (1965). “On Bayes Procedures.” *Probability Theory and Related Fields*, 4(1): 10–26. [665](#)
- Shen, X. (2002). “Asymptotic normality of semiparametric and nonparametric posterior distributions.” *Journal of the American Statistical Association*, 97(457): 222–235. [665](#)
- Spokoiny, V. (2012). “Parametric estimation. Finite sample theory.” *The Annals of Statistics*, 40(6): 2877–2909. [666](#), [669](#), [671](#), [677](#), [686](#), [688](#), [689](#), [693](#), [704](#)
- van der Vaart, A. W. (2000). *Asymptotic Statistics (Cambridge Series in Statistical and Probabilistic Mathematics)*. Cambridge University Press. [665](#)

Acknowledgments

The authors would like to thank the editors and the two anonymous reviewers for their constructive comments and suggestions on the original version of this paper. The research results of sections 3 and 5 are developed in IITP RAS with the support of the grant by the Russian Science Foundation (project 14-50-00150). The research results of sections 2 and 4 are partially supported by Laboratory for Structural Methods of Data Analysis in Predictive Modeling, MIPT, RF government grant, ag. 11.G34.31.0073, and German Research Foundation (DFG) through the Collaborative Research Center 649 “Economic Risk”.