

Inexact Laplace Approximation and the Use of Posterior Mean in Bayesian Inference*

Vladimir Spokoiny[†]

Abstract. The prominent Bernstein – von Mises (BvM) Theorem claims a kind of approximation of the posterior distribution by a Gaussian one with the covariance close to the inverse of the total Fisher information matrix. A more general Laplace approximation result states a similar Gaussian approximation of the posterior with the parameters depending on the prior. These two results build a basis for Bayesian inference and uncertainty quantification in a rather general situation. Spokoiny and Panov (2021) offered a new look at this problem which allows to state rather strong results on the quality of Gaussian approximation in non-asymptotic and dimension free form assuming linearity and concavity of log-likelihood function which can be misspecified. The established results provide explicit non-asymptotic bounds on the quality of a Gaussian approximation of the posterior distribution in total variation distance in terms of the so called *effective dimension* p_G defined as interplay between information contained in the data and in the prior distribution. This paper substantially improves and further develops the results from Spokoiny and Panov (2021) using the recent progress on high dimensional Laplace approximation. We address the question of *effective and critical dimension* in Bayesian inference, the relations between Laplace approximation and Bernstein–von Mises Theorem, and, particularly, the use of *posterior mean* instead of *Maximum A Posteriori Probability* estimator in Bayesian inference. The results are illustrated for the case of log-density estimation.

MSC2020 subject classifications: Primary 62F15; secondary 60F25, 62C10.

Keywords: Gaussian approximation, effective dimension, posterior mean.

1 Introduction

The prominent Bernstein – von Mises (BvM) Theorem is one of the most fundamental and most mysterious results in Bayesian inference stating asymptotic normality of the posterior distribution. It is often considered as a Bayesian counterpart of the Fisher Theorem claiming asymptotic normality of the maximum likelihood estimator (MLE). The inverse of the Fisher information matrix describes the variance of the MLE in the Fisher Theorem and the variance of the posterior in the BvM Theorem. These two results build the basis for statistical inference and uncertainty quantification of majority of statistical procedures. Parametric BvM theory is well-understood Van der Vaart (2000). One of corollaries of the BvM result is the fact that the choice of the prior is unimportant, it does not show up in the limiting distribution and washes out of the

*Financial support by the German Research Foundation (DFG) through the Collaborative Research Center 1294 “Data assimilation” is gratefully acknowledged.

[†]Weierstrass Institute and Humboldt University Berlin, Mohrenstr. 39, 10117 Berlin, Germany; HSE and IITP RAS Moscow, spokoiny@wias-berlin.de

posterior distribution as the sample size increases. The situation changes dramatically when modern statistical problems with a high dimensional parameter space and limited sample size are considered. The prior becomes crucial, its choice is an important issue as for inference problems as for uncertainty quantification; see e.g. Knapik et al. (2016), Szabó et al. (2015), Rousseau and Szabo (2020) and references therein. Lu (2017) studied a parametric BvM theorem for nonlinear Bayesian inverse problems with an increasing number of parameters.

Consider a Bayesian inference problem for a pseudo log-likelihood function $L(\mathbf{v}) = L(\mathbf{Y}, \mathbf{v})$ with data \mathbf{Y} , a parameter $\mathbf{v} \in \mathbb{R}^p$, and a prior π on \mathbb{R}^p . This paper focuses on the case of a Gaussian prior $\pi_G \sim \mathcal{N}(\mathbf{v}_0, G^{-2})$ with a symmetric positive definite covariance matrix G^{-2} . An extension to non-Gaussian priors is commented in Section D in the Supplement (Spokoiny (2023)); later (S2023). The posterior density $\pi_G(\cdot)$ of \mathbf{v} given \mathbf{Y} can be written in the form

$$\pi_G | \mathbf{Y} \sim \pi_G(\mathbf{v}) \propto \exp\{L(\mathbf{v}) - \|G(\mathbf{v} - \mathbf{v}_0)\|^2/2\},$$

where the sign \propto means equality up to a normalizing multiplicative constant. Assume that the penalized maximum likelihood estimator (pMLE) $\tilde{\mathbf{v}}_G$ is well defined:

$$\tilde{\mathbf{v}}_G = \underset{\mathbf{v}}{\operatorname{argmax}}\{L(\mathbf{v}) - \|G(\mathbf{v} - \mathbf{v}_0)\|^2/2\}. \quad (1.1)$$

Clearly $\tilde{\mathbf{v}}_G$ is maximizer of $\pi_G(\mathbf{v})$ and it is often referred to as maximum a posteriori probability (MAP) estimator. Define also the penalized Fisher information matrix

$$\mathbb{F}_G(\mathbf{v}) = -\nabla^2 \mathbb{E}L(\mathbf{v}) + G^2,$$

where \mathbb{E} means the expectation w.r.t. the underlying data distribution \mathbb{P} . An important step in understanding the impact of the prior is made by the results on Laplace's approximation claiming that the posterior distribution π_G is close to the Gaussian distribution $\mathcal{N}(\tilde{\mathbf{v}}_G, \mathbb{F}_G^{-1}(\tilde{\mathbf{v}}_G))$. The Bernstein–von Mises phenomenon formally corresponds to the non-penalized case $G^2 = 0$. A number of papers discuss the BvM phenomenon for nonlinear inverse problems; see e.g. Nickl (2020); Monard et al. (2019); Giordano and Kekkonen (2020), where the convergence is quantified in a distance that metrizes the weak convergence. Schillings et al. (2020) showed that the Laplace approximation error in Hellinger distance converges to zero in the order of the noise level. The recent paper Helin and Kretschmann (2022) provides a finite sample error of Laplace approximation for the total variation (TV) distance with an explicit dependence on the dimension and on the nonlinearity of the forward mapping for Bayesian inverse problems. A common drawback of all these and similar results is that the error bounds depend implicitly or explicitly of the dimension of the parameter space. If the dimension grows this dependence may become crucial, thus questioning the range of applicability of Laplace approximation and of BvM Theorem. Spokoiny (2017) discussed general properties of the pMLE $\tilde{\mathbf{v}}_G$ in terms of the so called *effective dimension* which can be small or moderate even if the true parameter dimension is large. Spokoiny and Panov (2021) established similar and even stronger results under an additional assumption of

linearity of the stochastic term $L(\boldsymbol{v}) - \mathbb{E}L(\boldsymbol{v})$ in \boldsymbol{v} . Spokoiny (2019) explained how a non-linear inverse problem can be reduced to the stochastically linear case by extending the parameter space without significant increase of the effective dimension.

Another challenge of applying the BvM-type results is that the parameters of the approximating Gaussian distribution are defined through the pMLE $\tilde{\boldsymbol{v}}_G$ from (1.1). This is a high dimensional optimization problem for a random objective function. A closed form analytic solution is available only in very special and simple situations, otherwise it can only be obtained by high-tech optimization methods with some error. This leads to an open problem of justifying a Laplace approximation of the posterior with inexact parameters for Bayesian inference, Durmus and Moulines (2019). Finding $\tilde{\boldsymbol{v}}_G$ could be especially difficult if computing $L(\boldsymbol{v})$ and its gradient is costly. This leads to gradient free methods Nesterov and Spokoiny (2017) or Bayesian optimization approach Mockus (1989), Frazier (2018). Ma et al. (2019) argued that Monte Carlo Markov Chain (MCMC) sampling can be more efficient than high dimensional optimization. A particular issue for applying MCMC type methods for Bayesian inference is to justify the use of posterior mean instead of posterior mode.

This paper aims at addressing the mentioned challenges. Below the list of the most important achievements in the paper.

Effective dimension and dimension free guarantees for penalized MLE $\tilde{\boldsymbol{v}}_G$.

The *effective dimension* \mathfrak{p}_G is defined by an interplay between the information delivered by the data and information contained in the penalty; see Spokoiny and Panov (2021) or Section 2.4 for more details. Section 2 establishes explicit *non-asymptotic* and *dimension free* results on the accuracy of the pMLE $\tilde{\boldsymbol{v}}_G$ including concentration, Fisher and Wilks expansions, loss and risk bounds under assumptions of stochastic linearity and concavity of the log-likelihood function. The bounds apply well under model misspecification and are stated under the same *critical dimension* condition $\mathfrak{p}_G \ll n$, where n is the *effective sample size*. The use of self-concordance type conditions from Section A of the Supplement (S2023) on $f(\boldsymbol{v}) = \mathbb{E}L(\boldsymbol{v})$ allows to obtain sharper and more transparent results than in Spokoiny and Panov (2021). Section C of the Supplement (S2023) explains a rate optimal choice of penalty/prior and derives usual minimax rate results over Sobolev smoothness classes from the obtained bounds.

Laplace's approximation of the posterior. Section 3 of the paper presents new results on concentration and Gaussian approximation of the posterior. It is important to note that the analysis of the posterior distribution requires very different analytic tools than used for the pMLE study. We make use of the recent progress in high dimensional Laplace approximation from Section E of the Supplement (S2023). It appears that some rather sharp bounds on posterior concentration and contraction can be obtained under the critical dimension condition $\mathfrak{p}_G \ll n$ similar to the case of pMLE; see Proposition 3.1. This is a substantial improvement over Spokoiny and Panov (2021) where $\mathfrak{p}_G^3 \ll n$ was assumed. The main results of Theorem 3.4 provide some bounds on the accuracy of Gaussian approximation of the posterior. These results still require $\mathfrak{p}_G^3 \ll n$. For the *total variation* distance, the accuracy of approximation is of order $\sqrt{\mathfrak{p}_G^3/n}$. It can be improved to \mathfrak{p}_G^3/n if we limit ourselves to the class of centrally symmetric sets.

The use of posterior mean in place of MAP. The result on Gaussian approximation of the posterior justifies the use of elliptic credible set centered at the MAP $\tilde{\mathbf{v}}_G$. The possibility of using the posterior mean $\bar{\mathbf{v}}_G$ in place of the MAP (1.1) is an important challenging question answered by Theorem 3.10: this use is justified under the same condition $p_G^3 \ll n$, however, only after restricting to the class of elliptic credible sets. The proof involves some recent advances in Gaussian comparison Götze et al. (2019).

Log-density estimation. Section 4 specifies the general results to the case of log-density estimation. In particular, we provide finite sample explicit and sharp bounds on posterior concentration and contraction under the condition $s_0 > 0$ on the smoothness degree of the density while Rousseau and Szabo (2017) required $s_0 > 1/2$.

2 Properties of the pMLE $\tilde{\mathbf{v}}_G$

This section collects general results about concentration and expansion of the pMLE which substantially improve the bounds from Spokoiny and Panov (2021). We assume to be given a pseudo log-likelihood random function $L(\mathbf{v})$, $\mathbf{v} \in \mathcal{Y} \subseteq \mathbb{R}^p$, $p < \infty$. Given a quadratic penalty $\|G\mathbf{v}\|^2/2$, define

$$L_G(\mathbf{v}) = L(\mathbf{v}) - \|G\mathbf{v}\|^2/2.$$

Typical examples of choosing G^2 are given in Section B.1 of the Supplement (S2023). Consider the penalized MLE $\tilde{\mathbf{v}}_G$ and its population counterpart \mathbf{v}_G^*

$$\tilde{\mathbf{v}}_G = \underset{\mathbf{v}}{\operatorname{argmax}} L_G(\mathbf{v}), \quad \mathbf{v}_G^* = \underset{\mathbf{v}}{\operatorname{argmax}} \mathbb{E}L_G(\mathbf{v}).$$

The corresponding Fisher information matrix $\mathbb{F}_G(\mathbf{v})$ is given by

$$\mathbb{F}(\mathbf{v}) = -\nabla^2 \mathbb{E}L(\mathbf{v}), \quad \mathbb{F}_G(\mathbf{v}) = -\nabla^2 \mathbb{E}L_G(\mathbf{v}) = \mathbb{F}(\mathbf{v}) + G^2.$$

We assume $\mathbb{F}_G(\mathbf{v})$ to be positive definite for all considered \mathbf{v} . By $D_G(\mathbf{v})$ we denote a positive symmetric matrix with $D_G^2(\mathbf{v}) = \mathbb{F}_G(\mathbf{v})$, and $\mathbb{F}_G = \mathbb{F}_G(\mathbf{v}_G^*)$, $D_G = \mathbb{F}_G^{1/2}$.

2.1 Conditions

Now we present our conditions. The most important one is about linearity of the stochastic component $\zeta(\mathbf{v}) = L(\mathbf{v}) - \mathbb{E}L(\mathbf{v}) = L_G(\mathbf{v}) - \mathbb{E}L_G(\mathbf{v})$.

- (ζ) *The stochastic component $\zeta(\mathbf{v}) = L(\mathbf{v}) - \mathbb{E}L(\mathbf{v})$ of the process $L(\mathbf{v})$ is linear in \mathbf{v} . We denote by $\nabla\zeta \equiv \nabla\zeta(\mathbf{v}) \in \mathbb{R}^p$ its gradient.*

Below we assume some concentration properties of the stochastic vector $\nabla\zeta$; see (F.43) of Theorem F.15 of the Supplement (S2023).

($\nabla\zeta$) Let $V^2 = \text{Var}(\nabla\zeta)$ and $D_G^2 = D_G^2(\mathbf{v}_G^*)$. Then for any considered $\mathbf{x} > 0$

$$\mathbb{P}(\|D_G^{-1}\nabla\zeta\| \geq \mathbf{r}_G(\mathbf{x})) \leq 3e^{-\mathbf{x}}, \tag{2.1}$$

where for $\mathbf{p}_G = \text{tr}(D_G^{-2}V^2)$ and $\lambda_G = \|D_G^{-1}V^2D_G^{-1}\|$

$$\mathbf{r}_G(\mathbf{x}) \stackrel{\text{def}}{=} \sqrt{\mathbf{p}_G} + \sqrt{2\mathbf{x}\lambda_G}. \tag{2.2}$$

This condition can be effectively checked if the errors in the data exhibit sub-Gaussian or sub-exponential behaviour; see Section F.3 of the Supplement (S2023). The important value $\mathbf{p}_G = \text{tr}(D_G^{-2}V^2)$ can be called the *effective dimension*; see Spokoiny (2017).

We also assume that the deterministic part $\mathbb{E}L_G(\mathbf{v})$ of the penalized log-likelihood is a concave function. It can be relaxed using localization; see Spokoiny (2019).

(\mathcal{C}_G) \mathcal{Y} is an open and convex set in \mathbb{R}^p . The function $\mathbb{E}L_G(\mathbf{v})$ is concave on \mathcal{Y} .

In Section 3.1 we consider a stronger condition of semi-concavity of $\mathbb{E}L(\mathbf{v})$. Further we will also need some smoothness conditions on the function $f(\mathbf{v}) = \mathbb{E}L(\mathbf{v})$. The class of models satisfying the conditions (ζ), ($\nabla\zeta$) with a smooth function $f(\mathbf{v}) = \mathbb{E}L(\mathbf{v})$ will be referred to as *stochastically linear smooth* (SLS). This class includes linear regression, generalized linear models (GLM) and log-density models; see Spokoiny and Panov (2021). However, this class is much larger. For instance, nonlinear regression and nonlinear inverse problems can be adapted to the SLS framework by an extension of the parameter space; see Spokoiny (2019).

2.2 Concentration of the pMLE $\tilde{\mathbf{v}}_G$

This section discusses some concentration properties of the pMLE $\tilde{\mathbf{v}}_G = \text{argmax}_{\mathbf{v}} L_G(\mathbf{v})$.

Given \mathbf{x} and $\mathbf{r}_G = \mathbf{r}_G(\mathbf{x})$ from (2.2), define for some $\nu < 1$ the set \mathcal{U}_G by

$$\mathcal{U}_G \stackrel{\text{def}}{=} \{\mathbf{u}: \|D_G\mathbf{u}\| \leq \nu^{-1}\mathbf{r}_G\}. \tag{2.3}$$

The result of this section states the concentration properties of the pMLE $\tilde{\mathbf{v}}_G$ in the local vicinity \mathcal{A}_G of \mathbf{v}_G^* of the form

$$\mathcal{A}_G \stackrel{\text{def}}{=} \mathbf{v}_G^* + \mathcal{U}_G = \{\mathbf{v} = \mathbf{v}_G^* + \mathbf{u}: \mathbf{u} \in \mathcal{U}_G\} \subseteq \mathcal{Y}^\circ.$$

Local Gateaux-regularity of $f(\mathbf{v}) = \mathbb{E}L(\mathbf{v})$ within \mathcal{A}_G will be measured by the error of the second order Taylor approximation

$$\begin{aligned} \delta_3(\mathbf{v}, \mathbf{u}) &= f(\mathbf{v} + \mathbf{u}) - f(\mathbf{v}) - \langle \nabla f(\mathbf{v}), \mathbf{u} \rangle - \frac{1}{2} \langle \nabla^2 f(\mathbf{v}), \mathbf{u}^{\otimes 2} \rangle, \\ \delta'_3(\mathbf{v}, \mathbf{u}) &= \langle \nabla f(\mathbf{v} + \mathbf{u}), \mathbf{u} \rangle - \langle \nabla f(\mathbf{v}), \mathbf{u} \rangle - \langle \nabla^2 f(\mathbf{v}), \mathbf{u}^{\otimes 2} \rangle. \end{aligned} \tag{2.4}$$

More precisely, define

$$\omega_G \stackrel{\text{def}}{=} \sup_{\mathbf{u} \in \mathcal{U}_G} \frac{2|\delta_3(\mathbf{v}_G^*, \mathbf{u})|}{\|D_G \mathbf{u}\|^2}, \quad \omega'_G \stackrel{\text{def}}{=} \sup_{\mathbf{u} \in \mathcal{U}_G} \frac{2|\delta'_3(\mathbf{v}_G^*, \mathbf{u})|}{\|D_G \mathbf{u}\|^2}. \quad (2.5)$$

The quantities ω_G and ω'_G can be effectively bounded under smoothness conditions (\mathcal{T}_3) or (\mathcal{S}_3) given in Section A of the Supplement (S2023). Under (\mathcal{T}_3) at $\mathbf{v} = \mathbf{v}_G^*$ with $D^2(\mathbf{v}_G^*) = D_G^2$ and $\mathbf{r} = \mathbf{r}_G$, by Lemma A.1 of the Supplement (S2023), it holds for a small constant τ_3

$$\omega'_G \leq \tau_3 \nu^{-1} \mathbf{r}_G, \quad \omega_G \leq \tau_3 \nu^{-1} \mathbf{r}_G / 3.$$

Furthermore, under (\mathcal{S}_3) , the same bounds apply with $\tau_3 = c_3 n^{-1/2}$; see Lemma A.2 of the Supplement (S2023).

Proposition 2.1. *Suppose (ζ) , $(\nabla\zeta)$, and (\mathcal{C}_G) . Let also*

$$1 - \nu - \omega'_G \geq 0; \quad (2.6)$$

see (2.5) and (2.3). Then $\tilde{\mathbf{v}}_G \in \mathcal{A}_G$ on a set $\Omega(\mathbf{x})$ with $\mathbb{P}(\Omega(\mathbf{x})) \geq 1 - 3e^{-x}$, i.e.

$$\|D_G(\tilde{\mathbf{v}}_G - \mathbf{v}_G^*)\| \leq \nu^{-1} \mathbf{r}_G. \quad (2.7)$$

Proof. By $(\nabla\zeta)$, on a the random set $\Omega(\mathbf{x})$ with $\mathbb{P}(\Omega(\mathbf{x})) \geq 1 - 3e^{-x}$, it holds $\|D_G^{-1} \nabla\zeta\| \leq \mathbf{r}_G$. Now the result follows from Proposition A.7 with $f(\mathbf{v}) = \mathbb{E}L_G(\mathbf{v})$, $g(\mathbf{v}) = L_G(\mathbf{v})$, $\mathbf{r} = \nu^{-1} \mathbf{r}_G$, and $\mathbf{A} = \nabla\zeta$. \square

Remark 2.1. The result (2.7) continues to apply with any matrix \mathbb{D} in place of D_G provided that $\mathbb{D} \leq D_G$ and $(\nabla\zeta)$ as well as (2.5), (2.6) hold after this change.

2.3 Fisher and Wilks expansions

This section presents some finite sample results about the behavior of the penalized MLE $\tilde{\mathbf{v}}_G$ and the excess $L_G(\tilde{\mathbf{v}}_G) - L_G(\mathbf{v}_G^*)$. Proposition 2.1 states the concentration properties of $\tilde{\mathbf{v}}_G$ around \mathbf{v}_G^* . Now we show that this concentration can be used to establish a version of the Fisher expansion for the estimation error $\tilde{\mathbf{v}}_G - \mathbf{v}_G^*$ and the Wilks expansion for the excess $L_G(\tilde{\mathbf{v}}_G) - L_G(\mathbf{v}_G^*)$.

Theorem 2.2. *Assume the conditions of Proposition 2.1 with $\nu = 2/3$. Then on $\Omega(\mathbf{x})$*

$$\begin{aligned} 2L_G(\tilde{\mathbf{v}}_G) - 2L_G(\mathbf{v}_G^*) - \|D_G^{-1} \nabla\zeta\|^2 &\leq \frac{\omega_G}{1 - \omega_G} \|D_G^{-1} \nabla\zeta\|^2, \\ 2L_G(\tilde{\mathbf{v}}_G) - 2L_G(\mathbf{v}_G^*) - \|D_G^{-1} \nabla\zeta\|^2 &\geq -\omega_G \|D_G^{-1} \nabla\zeta\|^2. \end{aligned}$$

Also

$$\begin{aligned} \|D_G(\tilde{\mathbf{v}}_G - \mathbf{v}_G^*) - D_G^{-1}\nabla\zeta\|^2 &\leq \frac{3\omega_G}{(1 - \omega_G)^2} \|D_G^{-1}\nabla\zeta\|^2, \\ \|D_G(\tilde{\mathbf{v}}_G - \mathbf{v}_G^*)\| &\leq \frac{1 + \sqrt{2\omega_G}}{1 - \omega_G} \|D_G^{-1}\nabla\zeta\|. \end{aligned} \tag{2.8}$$

Proof. The result follows from Proposition A.8 of the Supplement (S2023) similarly to Proposition 2.1. \square

2.4 Effective sample size and critical dimension in pMLE

This section discusses the important question of the critical parameter dimension still ensuring the validity of the presented results. A very important feature of our results is their dimension free and coordinate free form. The true parametric dimension p can be very large, it does not show up in the error terms. Neither do we use any spectral decomposition or sequence space structure, in particular, we do not require that the Fisher information matrix \mathbb{F} and the penalty matrix G^2 are diagonal or can be jointly diagonalized. The results are stated for the general data \mathbf{Y} and a quasi log-likelihood function. In particular, we do not assume independent or progressively dependent observations and additive structure of the log-likelihood. The *effective sample size* n can be defined via the smallest eigenvalue of the matrix $\mathbb{F}_G = D_G^2 = -\nabla^2 \mathbb{E}L_G(\mathbf{v}_G^*)$:

$$n^{-1} \stackrel{\text{def}}{=} \|\mathbb{F}_G^{-1}\|.$$

Our results apply as long as this value is sufficiently small. In typical examples like regression or density modeling such defined value is closely related to the sample size of the data.

For the concentration result of Proposition 2.1 we need the basic conditions (ζ) and (\mathcal{C}_G) . Further, $(\nabla\zeta)$ identifies the radius \mathbf{r}_G of the local vicinity \mathcal{A}_G . The final critical condition is given by (2.6). Essentially it says that the values ω_G and ω'_G are significantly smaller than 1. Under (\mathcal{S}_3) , $\omega'_G \leq c_3 \nu^{-1} \mathbf{r}_G n^{-1/2}$; see Lemma A.2 of the Supplement (S2023). So, (2.6) means $\mathbf{r}_G^2 \ll n$. Moreover, definition (2.1) of \mathbf{r}_G yields that $\mathbf{r}_G^2 \asymp \text{tr}(D_G^{-2}V^2) = \mathbf{p}_G$, where \mathbf{p}_G is the *effective dimension* of the problem. We conclude that the main properties of the pMLE $\tilde{\mathbf{v}}_G$ are valid under the condition $\mathbf{p}_G \ll n$ meaning sufficiently many observations per effective number of parameters.

2.5 The use of \tilde{D}_G^2 instead of D_G^2

The penalized information matrix $D_G^2 = D_G^2(\mathbf{v}_G^*) = -\nabla^2 \mathbb{E}L_G(\mathbf{v}_G^*)$ plays an important role in our results. In particular, D_G describes the shape of the concentration set $\mathcal{A}_G = \mathbf{v}_G^* + \mathcal{U}_G$. However, this matrix is not available as it involves the unknown point \mathbf{v}_G^* . If the matrix function $\mathbb{F}(\mathbf{v})$ is locally constant in \mathcal{A}_G , one can replace \mathbf{v}_G^* with its estimate $\tilde{\mathbf{v}}_G$. Variability of $\mathbb{F}(\mathbf{v})$, or, equivalently, $\mathbb{F}_G(\mathbf{v}) = \mathbb{F}(\mathbf{v}) + G^2$ can be

measured under the Fréchet smoothness of $f(\mathbf{v}) = \mathbb{E}L_G(\mathbf{v})$ by the value ω_G^+ from (A.4) of the Supplement (S2023) with $\mathbf{v} = \mathbf{v}_G^*$, $D(\mathbf{v}) = D_G$, and $\mathbf{r} = \nu^{-1}\mathbf{r}_G$.

Proposition 2.3. *Assume the conditions of Proposition 2.1 and let $\omega_G^+ \leq 1/2$; see (A.4). The random matrix $\tilde{D}_G^2 = \mathbb{F}_G(\tilde{\mathbf{v}}_G)$ fulfills on $\Omega(\mathbf{x})$ for any $\mathbf{u} \in \mathbb{R}^p$*

$$\begin{aligned} \|D_G^{-1}\tilde{D}_G^2 D_G^{-1} - I_p\| &\leq \omega_G^+, & \|D_G \tilde{D}_G^{-2} D_G - I_p\| &\leq \frac{\omega_G^+}{1 - \omega_G^+}, \\ (1 - \omega_G^+) \|D_G \mathbf{u}\|^2 &\leq \|\tilde{D}_G \mathbf{u}\|^2 \leq (1 + \omega_G^+) \|D_G \mathbf{u}\|^2. \end{aligned} \quad (2.9)$$

Proof. The value $\tilde{\mathbf{v}}_G - \mathbf{v}_G^*$ belongs to \mathcal{U}_G on $\Omega(\mathbf{x})$ and (2.9) follows from (A.5). \square

2.6 Smoothness and bias

Due to Proposition 2.1, the penalized MLE $\tilde{\mathbf{v}}_G$ is in fact an estimator of the vector \mathbf{v}_G^* . However, \mathbf{v}_G^* depends on penalization which introduces some bias. This section discusses whether one can use $\tilde{\mathbf{v}}_G$ for estimating the underlying truth \mathbf{v}^* defined as the maximizer of the expected log-likelihood: $\mathbf{v}^* = \operatorname{argmax}_{\mathbf{v}} \mathbb{E}L(\mathbf{v})$. First we describe the bias $\mathbf{b}_G = \mathbf{v}_G^* - \mathbf{v}^*$ induced by penalization. It is important to mention that the previous results about the properties of the pMLE $\tilde{\mathbf{v}}_G$ require strong concavity of the expected log-likelihood function $\mathbb{E}L_G(\mathbf{v})$ at least in a vicinity of the point \mathbf{v}_G^* . In some sense, this strong concavity is automatically forced by the penalizing term in the definition of \mathbf{v}_G^* . However, the underlying truth $\mathbf{v}^* = \operatorname{argmax}_{\mathbf{v}} \mathbb{E}L(\mathbf{v})$ is the maximizer of the non-penalized expected log-likelihood, and the corresponding Hessian $\mathbb{F}(\mathbf{v}^*) = -\nabla^2 \mathbb{E}L(\mathbf{v}^*)$ can degenerate. This makes evaluation of the bias more involved. To bypass this situation, we assume later in this section that the Hessian $\nabla^2 \mathbb{E}L_G(\mathbf{v})$ cannot change much in a reasonably large vicinity of \mathbf{v}^* . This allows to establish an accurate quadratic approximation of $f(\mathbf{v})$ and to evaluate the bias $\mathbf{b}_G = \mathbf{v}_G^* - \mathbf{v}^*$.

Define \mathbb{D}_G by $\mathbb{D}_G^2 = \mathbb{F}_G(\mathbf{v}^*)$; cf. $D_G^2 = \mathbb{F}_G(\mathbf{v}_G^*)$. Let also Q be a symmetric matrix satisfying $Q^2 \leq \mathbb{D}_G^2$. Typical examples include $Q = D_G$, $Q = \mathbb{D}_G$, and $Q^2 = nI_p$. Later we bound the norm $\|Q\mathbf{b}_G\|$. Denote with $\nu = 2/3$

$$\begin{aligned} \mathbf{b}_G &\stackrel{\text{def}}{=} \|Q \mathbb{D}_G^{-2} G^2 \mathbf{v}^*\|, \\ \omega_G^* &\stackrel{\text{def}}{=} \sup_{\mathbf{u}: \|Q\mathbf{u}\| \leq \nu^{-1} \mathbf{b}_G} \|\mathbb{D}_G^{-1} \mathbb{F}_G(\mathbf{v}^* + \mathbf{u}) \mathbb{D}_G^{-1} - I_p\|; \end{aligned} \quad (2.10)$$

cf. (A.4) and (A.5) of the Supplement (S2023) for $f(\mathbf{v}) = \mathbb{E}L_G(\mathbf{v})$. Note that the definition of ω_G^+ in Proposition 2.3 uses another $\mathbf{r} = \nu^{-1}\mathbf{r}_G$, therefore, different notation. Proposition A.11 of the Supplement (S2023) yields the following result.

Proposition 2.4. *Let $\mathbb{D}_G^2 = \mathbb{F}_G(\mathbf{v}^*)$, $\nu \leq 2/3$, and $\mathbf{b}_G = \|Q \mathbb{D}_G^{-2} G^2 \mathbf{v}^*\|$. Let also*

$\omega_G^* \leq 1/3$; see (2.10). Then the bias $\mathbf{b}_G = \mathbf{v}_G^* - \mathbf{v}^*$ fulfills

$$\begin{aligned} \|Q \mathbf{b}_G\| &\leq \frac{\mathbf{b}_G}{1 - \omega_G^*} = \frac{1}{1 - \omega_G^*} \|Q \mathcal{D}_G^{-2} G^2 \mathbf{v}^*\|, \quad (2.11) \\ \|Q(\mathbf{b}_G + \mathcal{D}_G^{-2} G^2 \mathbf{v}^*)\| &\leq \frac{\omega_G^*}{1 - \omega_G^*} \mathbf{b}_G = \frac{\omega_G^*}{1 - \omega_G^*} \|Q \mathcal{D}_G^{-2} G^2 \mathbf{v}^*\|. \end{aligned}$$

Corollary 2.5. *Assume the conditions of Proposition 2.4. Then*

$$\|\mathcal{D}_G \mathbf{b}_G\| \leq \frac{1}{1 - \omega_G^*} \|\mathcal{D}_G^{-1} G^2 \mathbf{v}^*\|, \quad \|\mathbf{b}_G\| \leq \frac{1}{1 - \omega_G^*} \|\mathcal{D}_G^{-2} G^2 \mathbf{v}^*\|. \quad (2.12)$$

The same bounds apply with $D_G^2 = \mathbb{F}_G(\mathbf{v}_G^*)$ in place of $\mathcal{D}_G^2 = \mathbb{F}_G(\mathbf{v}^*)$.

This is a special cases of (2.11) with $Q = \mathcal{D}_G$ and $Q = I_p$. The last statement is due to Remark A.1 of the Supplement (S2023).

2.7 Loss and risk of the pMLE

Now we combine the previous results about the stochastic term $\tilde{\mathbf{v}}_G - \mathbf{v}_G^*$ and the bias term $\mathbf{b}_G = \mathbf{v}_G^* - \mathbf{v}^*$ to obtain the sharp bounds on the loss and risk of the pMLE $\tilde{\mathbf{v}}_G$.

Theorem 2.6. *Assume the conditions of Proposition 2.1 and 2.4. Then on $\Omega(\mathbf{x})$ with $\mathbb{P}(\Omega(\mathbf{x})) \geq 1 - 3e^{-x}$, it holds with $\boldsymbol{\xi}_G = D_G^{-1} \nabla \zeta$, \mathbf{r}_G from (2.2), and $n^{-1} = \|D_G^{-2}\|$*

$$\|D_G(\tilde{\mathbf{v}}_G - \mathbf{v}^*)\| \leq \frac{1 + \sqrt{2\omega_G}}{1 - \omega_G} \|\boldsymbol{\xi}_G\| + \frac{\|D_G^{-1} G^2 \mathbf{v}^*\|}{1 - \omega_G^*} \leq \frac{1 + \sqrt{2\omega_G}}{1 - \omega_G} \mathbf{r}_G + \frac{\|D_G^{-1} G^2 \mathbf{v}^*\|}{1 - \omega_G^*}, \quad (2.13)$$

$$\|\tilde{\mathbf{v}}_G - \mathbf{v}^*\| \leq \frac{1 + \sqrt{2\omega_G}}{\sqrt{n}(1 - \omega_G)} \|\boldsymbol{\xi}_G\| + \frac{\|D_G^{-2} G^2 \mathbf{v}^*\|}{1 - \omega_G^*} \leq \frac{3 \mathbf{r}_G}{\sqrt{n}} + 3 \|D_G^{-2} G^2 \mathbf{v}^*\|. \quad (2.14)$$

Proof. Let $\Omega(\mathbf{x})$ be the random set from $(\nabla \zeta)$ on which with $\|\boldsymbol{\xi}_G\| \leq \mathbf{r}_G$. It follows from (2.8) of Theorem 2.2 that on $\Omega(\mathbf{x})$ with $\mathbf{b}_G = \mathbf{v}_G^* - \mathbf{v}^*$

$$\|D_G(\tilde{\mathbf{v}}_G - \mathbf{v}^*) + D_G \mathbf{b}_G\| \leq \frac{1 + \sqrt{2\omega_G}}{1 - \omega_G} \|\boldsymbol{\xi}_G\|.$$

This and (2.12) imply (2.13). □

Now we state the results about the risk of the pMLE $\tilde{\mathbf{v}}_G$. To avoid technical burden, we fix a large \mathbf{x} , $\mathbf{r}_G = \mathbf{r}_G(\mathbf{x})$, and exclude an event $\{\|\boldsymbol{\xi}_G\| > \mathbf{r}_G\}$ having an exponentially small probability; see condition (2.1) of Proposition 2.1.

Theorem 2.7. Assume the conditions of Proposition 2.1 and Proposition 2.4. Then for a set $\Omega(\mathbf{x})$ with $\mathbb{P}(\Omega(\mathbf{x})) \geq 1 - 3e^{-x}$, it holds with $\mathbf{p}_G = \text{tr}(D_G^{-2} V^2)$

$$\begin{aligned} \|\mathbb{E}\{D_G(\tilde{\mathbf{v}}_G - \mathbf{v}^*) \mathbb{1}(\Omega(\mathbf{x}))\}\| &\leq \frac{1}{1 - \omega_G^*} \|D_G^{-1} G^2 \mathbf{v}^*\| + \frac{\sqrt{3\omega_G}}{1 - \omega_G} \mathbb{E}\|\boldsymbol{\xi}_G\| + \mathbf{c}_1 e^{-x} \\ &\leq \frac{1}{1 - \omega_G^*} \|D_G^{-1} G^2 \mathbf{v}^*\| + \frac{\sqrt{3\omega_G}}{1 - \omega_G} \sqrt{\mathbf{p}_G} + \mathbf{c}_1 e^{-x}, \end{aligned} \tag{2.15}$$

and

$$\begin{aligned} \mathbb{E}\{\|D_G(\tilde{\mathbf{v}}_G - \mathbf{v}^*)\|^2 \mathbb{1}(\Omega(\mathbf{x}))\} \\ \leq \left(1 + \frac{\sqrt{2\omega_G}}{1 - \omega_G}\right)^2 \mathbf{p}_G + \left(\frac{1}{1 - \omega_G^*} \|D_G^{-1} G^2 \mathbf{v}^*\| + \frac{\sqrt{3\omega_G}}{1 - \omega_G} \sqrt{\mathbf{p}_G} + \mathbf{c}_1 e^{-x}\right)^2. \end{aligned} \tag{2.16}$$

Remark 2.2. For ω_G^* , ω_G small, (2.16) yields classical bias-variance decomposition:

$$\mathbb{E}\{\|D_G(\tilde{\mathbf{v}}_G - \mathbf{v}^*)\|^2 \mathbb{1}(\Omega(\mathbf{x}))\} \leq (\mathbf{p}_G + \|D_G^{-1} G^2 \mathbf{v}^*\|^2) \{1 + o(1)\}. \tag{2.17}$$

With $n^{-1} = \|D_G^{-2}\|$, we also obtain

$$\mathbb{E}\{n\|\tilde{\mathbf{v}}_G - \mathbf{v}^*\|^2 \mathbb{1}(\Omega(\mathbf{x}))\} \leq (\mathbf{p}_G + n\|D_G^{-2} G^2 \mathbf{v}^*\|^2) \{1 + o(1)\}.$$

Moreover, under the small bias condition $\|D_G^{-2} G^2 \mathbf{v}^*\|^2 \ll \mathbf{p}_G/n$, the impact of the bias induced by penalization is negligible. The relation $\|D_G^{-2} G^2 \mathbf{v}^*\|^2 \asymp \mathbf{p}_G/n$ is usually referred to as “bias-variance trade-off”. Our bound is sharp in the sense that even for the special case of a linear models, (2.17) becomes equality.

Proof. Below we denote $\mathbb{E}_{\mathbf{x}}\eta = \mathbb{E}\{\eta \mathbb{1}(\Omega(\mathbf{x}))\}$ for any r.v. η . As $\mathbb{E}\boldsymbol{\xi}_G = 0$, we derive

$$\mathbb{E}_{\mathbf{x}}\{D_G(\tilde{\mathbf{v}}_G - \mathbf{v}^*)\} = \mathbb{E}_{\mathbf{x}}D_G(\tilde{\mathbf{v}}_G - \mathbf{v}^* - \boldsymbol{\xi}_G) - \mathbb{E}\boldsymbol{\xi}_G \mathbb{1}(\Omega^c(\mathbf{x})).$$

For the first term we apply (2.8) and (2.12) yielding

$$\|\mathbb{E}_{\mathbf{x}}D_G(\tilde{\mathbf{v}}_G - \mathbf{v}^* - \boldsymbol{\xi}_G)\| \leq \frac{\|D_G^{-1} G^2 \mathbf{v}^*\|}{1 - \omega_G^*} + \frac{\sqrt{3\omega_G \mathbf{p}_G}}{1 - \omega_G}.$$

To show (2.15), we also have to bound the tail moments of $\|\boldsymbol{\xi}_G\|$:

$$\|\mathbb{E}\boldsymbol{\xi}_G \mathbb{1}(\Omega^c(\mathbf{x}))\| \leq \mathbb{E}\|\boldsymbol{\xi}_G\| \mathbb{1}(\Omega^c(\mathbf{x})) \leq e^{-x/2}.$$

This can be easily done using deviation bounds for the quadratic form $\|\boldsymbol{\xi}_G\|^2$; see Theorem F.9 of the Supplement (S2023). Similarly one can bound the variance of $D_G \tilde{\mathbf{v}}_G$.

With $\mathbf{B}_G = D_G^{-1}V^2D_G^{-1}$

$$\begin{aligned} \text{Var}_{\mathbf{x}}(D_G \tilde{\mathbf{v}}_G) &\leq \mathbb{E}_{\mathbf{x}} \{D_G (\tilde{\mathbf{v}}_G - \mathbf{v}_G^*)\} \{D_G (\tilde{\mathbf{v}}_G - \mathbf{v}_G^*)\}^\top \\ &\leq \left(1 + \frac{\sqrt{2\omega_G}}{1 - \omega_G}\right)^2 \mathbb{E}(\boldsymbol{\xi}_G \boldsymbol{\xi}_G^\top) = \left(1 + \frac{\sqrt{2\omega_G}}{1 - \omega_G}\right)^2 \mathbf{B}_G. \end{aligned}$$

This yields for the quadratic risk $\mathbb{E}\|D_G(\tilde{\mathbf{v}}_G - \mathbf{v}^*)\|^2$

$$\mathbb{E}_{\mathbf{x}} \|D_G(\tilde{\mathbf{v}}_G - \mathbf{v}^*)\|^2 \leq \text{tr} \text{Var}_{\mathbf{x}}\{D_G(\tilde{\mathbf{v}}_G - \mathbf{v}^*)\} + \|\mathbb{E}_{\mathbf{x}} D_G(\tilde{\mathbf{v}}_G - \mathbf{v}^*)\|^2$$

and (2.16) follows. □

3 Laplace approximation of the posterior

This section studies the properties of the posterior $\mathbf{v}_G | \mathbf{Y}$. Our main result states Gaussian approximation of the posterior by $\mathcal{N}(\tilde{\mathbf{v}}_G, \tilde{D}_G^{-2})$. More specifically, our aim is, for any bounded measurable function g , to compare the conditional moments of $g(\mathbf{v}_G - \tilde{\mathbf{v}}_G)$ and of $g(\tilde{D}_G^{-1}\boldsymbol{\gamma})$, where $\boldsymbol{\gamma}$ is standard normal conditionally on \mathbf{Y} . The use of $\nabla L_G(\tilde{\mathbf{v}}_G) = 0$ yields

$$\begin{aligned} \mathbb{E}\{g(\mathbf{v}_G - \tilde{\mathbf{v}}_G) | \mathbf{Y}\} &= \frac{\int g(\mathbf{u} - \tilde{\mathbf{v}}_G) e^{L_G(\mathbf{u})} d\mathbf{u}}{\int e^{L_G(\mathbf{u})} d\mathbf{u}} = \frac{\int g(\mathbf{u}) e^{L_G(\tilde{\mathbf{v}}_G + \mathbf{u}) - L_G(\tilde{\mathbf{v}}_G)} d\mathbf{u}}{\int e^{L_G(\tilde{\mathbf{v}}_G + \mathbf{u}) - L_G(\tilde{\mathbf{v}}_G)} d\mathbf{u}} \\ &= \frac{\int g(\mathbf{u}) \exp\{L_G(\tilde{\mathbf{v}}_G + \mathbf{u}) - L_G(\tilde{\mathbf{v}}_G) - \langle \nabla L_G(\tilde{\mathbf{v}}_G), \mathbf{u} \rangle\} d\mathbf{u}}{\int \exp\{L_G(\tilde{\mathbf{v}}_G + \mathbf{u}) - L_G(\tilde{\mathbf{v}}_G) - \langle \nabla L_G(\tilde{\mathbf{v}}_G), \mathbf{u} \rangle\} d\mathbf{u}}. \end{aligned} \tag{3.1}$$

Now consider the Bregman divergence of the expected log-likelihood $f_G(\mathbf{v}) = \mathbb{E}L_G(\mathbf{v})$

$$f_G(\mathbf{v}; \mathbf{u}) = f_G(\mathbf{v} + \mathbf{u}) - f_G(\mathbf{v}) - \langle \nabla f_G(\mathbf{v}), \mathbf{u} \rangle, \quad \mathbf{u} \in \mathbb{R}^p.$$

As the stochastic term of $L(\mathbf{v})$ and thus, of $L_G(\mathbf{v})$ is linear in \mathbf{v} , it holds for any \mathbf{v}, \mathbf{u}

$$L_G(\mathbf{v} + \mathbf{u}) - L_G(\mathbf{v}) - \langle \nabla L_G(\mathbf{v}), \mathbf{u} \rangle = f_G(\mathbf{v} + \mathbf{u}) - f_G(\mathbf{u}) - \langle \nabla f_G(\mathbf{v}), \mathbf{u} \rangle = f_G(\mathbf{v}; \mathbf{u}).$$

Given $\tilde{\mathbf{v}}_G = \mathbf{v}$, we derive from (3.1)

$$\mathbb{E}\{g(\mathbf{v}_G - \tilde{\mathbf{v}}_G) | \mathbf{Y}\} = \mathbb{E}\{g(\mathbf{v}_G - \tilde{\mathbf{v}}_G) | \tilde{\mathbf{v}}_G = \mathbf{v}\} = \frac{\int g(\mathbf{u}) e^{f_G(\mathbf{v}; \mathbf{u})} d\mathbf{u}}{\int e^{f_G(\mathbf{v}; \mathbf{u})} d\mathbf{u}}. \tag{3.2}$$

This basic identity will be systematically used below. Laplace's approximation means nothing but the use of the second order Taylor approximation of the function $f_G(\cdot)$ at \mathbf{v} . Namely, $f_G(\mathbf{v}; \mathbf{u}) \approx -\|D_G(\mathbf{v}) \mathbf{u}\|^2/2$ and

$$\frac{\int g(\mathbf{u}) e^{f_G(\mathbf{v}; \mathbf{u})} d\mathbf{u}}{\int e^{f_G(\mathbf{v}; \mathbf{u})} d\mathbf{u}} \approx \frac{\int g(\mathbf{u}) e^{-\|D_G(\mathbf{v}) \mathbf{u}\|^2/2} d\mathbf{u}}{\int e^{-\|D_G(\mathbf{v}) \mathbf{u}\|^2/2} d\mathbf{u}}.$$

The analysis includes two major steps: posterior concentration and a Gaussian approximation of the posterior distribution.

3.1 Posterior concentration

We start with the important technical result describing the concentration sets of the posterior. In all our result, the value \mathbf{x} is fixed to ensure that $e^{-\mathbf{x}}$ is negligible.

Proposition 2.1 enables us to restrict the study to the case with $\tilde{\mathbf{v}}_G \in \mathcal{A}_G$. To describe the concentration properties of the posterior we need a slightly stronger concavity condition on $\mathbb{E}L(\mathbf{v})$, concavity of $\mathbb{E}L_G(\mathbf{v})$ is not sufficient.

(C) The function $\mathbb{E}L(\mathbf{v})$ is concave.

This condition can be relaxed to *weak concavity*.

(C_o) There exists $G_o^2 \leq G^2$ such that for any $\mathbf{v} \in \mathcal{A}_G$, the function $2\mathbb{E}L(\mathbf{v} + \mathbf{u}) - \|G_o \mathbf{u}\|^2$ is concave in \mathbf{u} .

(C) is a special case of (C_o) with $G_o = 0$. In what follow we assume (C). However, all the results apply under (C_o) after replacing D^2 with $D_o^2 = D^2 + G_o^2$. Define

$$\mathbf{p}(\mathbf{v}) \stackrel{\text{def}}{=} \text{tr}\{D^2(\mathbf{v})D_G^{-2}(\mathbf{v})\}, \quad \mathbf{r}(\mathbf{v}) \stackrel{\text{def}}{=} 2\sqrt{\mathbf{p}(\mathbf{v})} + \sqrt{2\mathbf{x}}; \quad (3.3)$$

cf. (2.2) for \mathbf{p}_G and \mathbf{r}_G . This ensures with γ standard normal

$$\mathbb{P}(\|D(\mathbf{v})D_G^{-1}(\mathbf{v})\gamma\| > \mathbf{r}(\mathbf{v})) \leq e^{-\mathbf{x}},$$

see (F.11) of Corollary F.6 of the Supplement (S2023). With some fixed $\nu \leq 1$, e.g. $\nu = 2/3$, define for any $\mathbf{v} \in \mathcal{A}_G$

$$\mathcal{U}(\mathbf{v}) = \{\mathbf{u}: \|D(\mathbf{v})\mathbf{u}\| \leq \nu^{-1}\mathbf{r}(\mathbf{v})\}. \quad (3.4)$$

With $f(\mathbf{v}) = \mathbb{E}L(\mathbf{v})$ and $\delta_3(\mathbf{v}, \mathbf{u})$ from (2.4), local smoothness of $f(\cdot)$ at \mathbf{v} will be measured by the value $\omega(\mathbf{v})$:

$$\omega(\mathbf{v}) \stackrel{\text{def}}{=} \sup_{\mathbf{u} \in \mathcal{U}(\mathbf{v})} \frac{1}{\|D\mathbf{u}\|^2/2} |\delta_3(\mathbf{v}, \mathbf{u})|; \quad (3.5)$$

cf. (2.5). Under (S₃), it holds $\omega(\mathbf{v}) \leq \nu^{-1}c_3 \mathbf{r}(\mathbf{v}) n^{-1/2}/3$; see Lemma A.2 of the Supplement (S2023).

Proposition 3.1. *Suppose (C), (∇C), and (C). Let also $\mathbf{p}(\mathbf{v})$ and $\mathbf{r}(\mathbf{v})$ be defined by (3.3) and $\mathcal{U}(\mathbf{v})$ by (3.4). If $\omega(\mathbf{v})$ from (3.5) satisfies*

$$\omega(\mathbf{v}) \leq 1/3, \quad \mathbf{v} \in \mathcal{A}_G, \quad (3.6)$$

then on $\Omega(\mathbf{x})$, it holds with $\tilde{D} = D(\tilde{\mathbf{v}}_G)$ and $\tilde{\mathbf{r}} = \mathbf{r}(\tilde{\mathbf{v}}_G)$

$$\mathbb{P}\left(\mathbf{v}_G - \tilde{\mathbf{v}}_G \notin \tilde{\mathcal{U}} \mid \mathbf{Y}\right) = \mathbb{P}\left(\|\tilde{D}(\mathbf{v}_G - \tilde{\mathbf{v}}_G)\| > \tilde{\mathbf{r}} \mid \mathbf{Y}\right) \leq e^{-\mathbf{x}}. \tag{3.7}$$

Proof. Let us fix $\tilde{\mathbf{v}}_G = \mathbf{v}$ and apply (3.2) with $g(\mathbf{u}) = \mathbb{1}(\|D(\mathbf{v})\mathbf{u}\| \notin \mathcal{U}(\mathbf{v}))$. Then it suffices to bound uniformly in $\mathbf{v} \in \mathcal{A}_G$ the ratio

$$\rho(\mathbf{v}) \stackrel{\text{def}}{=} \frac{\int \mathbb{1}(D(\mathbf{v})\mathbf{u} \notin \mathcal{U}(\mathbf{v})) e^{f_G(\mathbf{v};\mathbf{u})} d\mathbf{u}}{\int e^{f_G(\mathbf{v};\mathbf{u})} d\mathbf{u}}.$$

Now (E.9) of Theorem E.1 of the Supplement (S2023) yields the result. □

3.2 Posterior contraction

Now we bring together all the previous results to bound the posterior deviations $\mathbf{v}_G - \mathbf{v}^*$. The difference $\mathbf{v}_G - \mathbf{v}^*$ can be decomposed as

$$\mathbf{v}_G - \mathbf{v}^* = (\mathbf{v}_G - \tilde{\mathbf{v}}_G) + (\tilde{\mathbf{v}}_G - \mathbf{v}^*). \tag{3.8}$$

Result (2.14) of Theorem 2.6 provides a deviation bound for $\|\tilde{\mathbf{v}}_G - \mathbf{v}^*\|$ while Proposition 3.1 claims concentration of the posterior on the set $\{\|\tilde{D}(\mathbf{v}_G - \tilde{\mathbf{v}}_G)\| \leq \tilde{\mathbf{r}}\}$. We conclude by the following result.

Proposition 3.2. *Assume the conditions of Theorem 2.6 and Proposition 3.1 and let $\|D^{-2}(\mathbf{v})\| \leq n^{-1}$ for $\mathbf{v} \in \mathcal{A}_G$. It holds on $\Omega(\mathbf{x})$*

$$\mathbb{P}\left(\|\mathbf{v}_G - \mathbf{v}^*\| \geq \|\tilde{\mathbf{v}}_G - \mathbf{v}^*\| + \nu^{-1}\tilde{\mathbf{r}}/\sqrt{n} \mid \mathbf{Y}\right) \leq 2e^{-\mathbf{x}}, \tag{3.9}$$

and $\|\tilde{\mathbf{v}}_G - \mathbf{v}^*\|$ satisfies (2.14), while $\tilde{\mathbf{r}} \leq (1 - \omega_G^+)^{-1/2} \mathbf{r}(\mathbf{v}_G^*)$.

Proof. Bound (3.9) follows from decomposition (3.8) and Proposition 3.1. Further, the use of (2.9) of Proposition 2.3 yields $\tilde{\mathbf{r}} \leq (1 - \omega_G^+)^{-1/2} \mathbf{r}(\mathbf{v}_G^*)$. □

The use of (2.14) of Theorem 2.6 implies that most of posterior mass is concentrated in the root- n vicinity of \mathbf{v}^* :

$$\mathbb{P}\left(\|\mathbf{v}_G - \mathbf{v}^*\| \geq 3\|D_G^{-2}G^2\mathbf{v}^*\| + \frac{3}{\sqrt{n}}(\mathbf{r}_G + \tilde{\mathbf{r}}) \mid \mathbf{Y}\right) \leq 2e^{-\mathbf{x}}.$$

A prior ensuring the bias-variance trade-off leads to the optimal contraction rate which corresponds to the optimal penalty choice in penalized maximum likelihood estimation.

3.3 Gaussian approximation of the posterior

This section presents our main results about the accuracy of Gaussian approximation of the posterior $\mathbf{v}_G \mid \mathbf{Y}$ in the total variation distance. The use of self-concordance type conditions from Section A helps to obtain very accurate and precise finite sample guarantees, which gradually improve the bounds from Spokoiny and Panov (2021).

Let $\mathcal{B}(\mathbb{R}^p)$ be the σ -field of all Borel sets in \mathbb{R}^p , while $\mathcal{B}_s(\mathbb{R}^p)$ stands for all centrally symmetric sets from $\mathcal{B}(\mathbb{R}^p)$.

Theorem 3.3. Assume (ζ) , $(\nabla\zeta)$, and (\mathcal{C}) . Furthermore, let

$$\omega(\mathbf{v}) \mathbf{p}(\mathbf{v}) \leq 2/3, \quad \mathbf{v} \in \mathcal{A}_G;$$

cf. (3.6). Then with

$$\diamond_2(\mathbf{v}) = \frac{0.75 \omega(\mathbf{v}) \mathbf{p}(\mathbf{v})}{1 - \omega(\mathbf{v})}$$

and $\tilde{\diamond} = \diamond_2(\tilde{\mathbf{v}}_G)$, it holds on $\Omega(\mathbf{x})$ with

$$\sup_{A \in \mathcal{B}(\mathbb{R}^p)} \left| \mathbb{P}(\mathbf{v}_G - \tilde{\mathbf{v}}_G \in A \mid \mathbf{Y}) - \mathbb{P}'(\tilde{D}_G^{-1} \gamma \in A) \right| \leq \frac{2(\tilde{\diamond} + e^{-\mathbf{x}})}{1 - \tilde{\diamond} - e^{-\mathbf{x}}} \leq 4(\tilde{\diamond} + e^{-\mathbf{x}}). \quad (3.10)$$

Here \mathbb{P}' means a standard Gaussian distribution of γ given \mathbf{Y} .

Now we present more advanced bounds on the error of Gaussian approximation under conditions (\mathcal{S}_3) and (\mathcal{S}_4) (resp. (\mathcal{T}_3) and (\mathcal{T}_4)) from Section A of the Supplement (S2023) for $\Upsilon^\circ = \mathcal{A}_G$.

Theorem 3.4. Assume (ζ) , $(\nabla\zeta)$, (\mathcal{C}) , (\mathcal{T}_3) , and let $\tau_3 \nu^{-1} \mathbf{r}(\mathbf{v}) \leq 3/4$ for $\mathbf{r}(\mathbf{v})$ from (3.3) and all $\mathbf{v} \in \mathcal{A}_G$. Then the concentration bound (3.7) holds. Moreover, let

$$\tau_3 \nu^{-1} \mathbf{r}(\mathbf{v}) \mathbf{p}(\mathbf{v}) \leq 2, \quad \mathbf{v} \in \mathcal{A}_G. \quad (3.11)$$

With $\omega(\mathbf{v}) \stackrel{\text{def}}{=} \tau_3 \mathbf{r}(\mathbf{v})/3 \leq 1/4$, define

$$\diamond_3(\mathbf{v}) \stackrel{\text{def}}{=} \frac{\tau_3}{4\{1 - \omega(\mathbf{v})\}^{3/2}} \{\mathbf{p}(\mathbf{v}) + 1\}^{3/2}.$$

Then the result (3.10) applies on $\Omega(\mathbf{x})$ with $\tilde{\diamond} = \diamond_3(\tilde{\mathbf{v}}_G)$. Moreover, under (\mathcal{T}_4)

$$\sup_{A \in \mathcal{B}_s(\mathbb{R}^p)} \left| \mathbb{P}(\mathbf{v}_G - \tilde{\mathbf{v}}_G \in A \mid \mathbf{Y}) - \mathbb{P}'(\tilde{D}_G^{-1} \gamma \in A) \right| \leq \frac{2(\tilde{\diamond}_4 + e^{-\mathbf{x}})}{1 - \tilde{\diamond}_4 - e^{-\mathbf{x}}} \leq 4(\tilde{\diamond}_4 + e^{-\mathbf{x}})$$

with $\tilde{\diamond}_4 = \diamond_4(\tilde{\mathbf{v}}_G)$ and

$$\diamond_4(\mathbf{v}) \stackrel{\text{def}}{=} \frac{1}{16\{1 - \omega(\mathbf{v})\}^2} \left[\tau_3^2 \{\mathbf{p}(\mathbf{v}) + 2\}^3 + 2\tau_4 \{\mathbf{p}(\mathbf{v}) + 1\}^2 \right].$$

The results continue to apply with (\mathcal{S}_3) (resp. (\mathcal{S}_4)) in place of (\mathcal{T}_3) (resp. (\mathcal{T}_4)) and $c_3 n^{-1/2}$ (resp. $c_4 n^{-1}$) in place of τ_3 (resp. τ_4).

Proof of Theorem 3.3 (resp. Theorem 3.4). Similarly to the proof of Proposition 3.1, we restrict ourselves to the event $\tilde{\mathbf{v}}_G \in \mathcal{A}_G$. Then we fix any possible value $\mathbf{v} \in \mathcal{A}_G$ of $\tilde{\mathbf{v}}_G$ and use (3.2) to represent the posterior probability of a set A in the form

$$\mathbb{P}(\mathbf{v}_G - \tilde{\mathbf{v}}_G \in A \mid \mathbf{Y}) = \frac{\int_A e^{f_G(\mathbf{v}; \mathbf{u})} d\mathbf{u}}{\int e^{f_G(\mathbf{v}; \mathbf{u})} d\mathbf{u}}.$$

Now the result follows by Theorem E.1 (resp. Theorem E.2). □

Under self-concordance conditions (\mathcal{S}_3) and (\mathcal{S}_4) , constraint (3.11) reads as

$$\sup_{\mathbf{v} \in \mathcal{A}_G} \frac{c_3 \nu^{-1} \mathbf{r}(\mathbf{v}) \mathbf{p}(\mathbf{v})}{n^{1/2}} \leq 2.$$

As $\omega(\mathbf{v}) \leq 1/4$, Theorem 3.4 yields on $\Omega(\mathbf{x})$ with $\tilde{\mathbf{p}} = \mathbf{p}(\tilde{\mathbf{v}}_G)$

$$\begin{aligned} \sup_{A \in \mathcal{B}(\mathbb{R}^p)} \left| \mathbb{P}(\mathbf{v}_G - \tilde{\mathbf{v}}_G \in A \mid \mathbf{Y}) - \mathbb{P}'(\tilde{D}_G^{-1} \boldsymbol{\gamma} \in A) \right| &\leq 2 c_3 \sqrt{\frac{(\tilde{\mathbf{p}} + 1)^3}{n}} + 4e^{-x}, \\ \sup_{A \in \mathcal{B}_s(\mathbb{R}^p)} \left| \mathbb{P}(\mathbf{v}_G - \tilde{\mathbf{v}}_G \in A \mid \mathbf{Y}) - \mathbb{P}'(\tilde{D}_G^{-1} \boldsymbol{\gamma} \in A) \right| &\leq \frac{c_3^2 (\tilde{\mathbf{p}} + 2)^3 + 2c_4 (\tilde{\mathbf{p}} + 1)^2}{2n} + 4e^{-x}. \end{aligned}$$

3.4 Critical dimension in Bayesian inference

Posterior concentration in Proposition 3.1 only requires $\omega(\mathbf{v}) \ll 1$ for all $\mathbf{v} \in \mathcal{A}_G$. Under (\mathcal{S}_3) , one can bound $\omega(\mathbf{v}) \asymp \sqrt{\mathbf{p}(\mathbf{v})/n}$ yielding the condition $\mathbf{p}(\mathbf{v}) \ll n$ on the critical dimension which is essentially the same as the condition $\mathbf{p}_G \ll n$ for the pMLE. This is an important finding and an essential improvement of Spokoiny and Panov (2021). The main result of Theorem 3.3 requires $\omega(\mathbf{v}) \mathbf{p}(\mathbf{v}) \ll 1$ which is much stronger because of the multiplicative factor $\mathbf{p}(\mathbf{v})$. Under (\mathcal{S}_3) , the remainder \diamond_3 is of order $\sqrt{\mathbf{p}^3(\mathbf{v})/n}$ while under (\mathcal{S}_4) , $\diamond_4 \asymp \mathbf{p}^3(\mathbf{v})/n$, still requiring $\mathbf{p}^3(\mathbf{v}) \ll n$. In some cases, e.g. for additive structure of the log-likelihood, it can be relaxed. However, it seems that the $\mathbf{p}^3(\mathbf{v}) \ll n$ condition is inherent in the problem and cannot be relaxed in general situation. We guess that in the region $n^{1/3} \ll \mathbf{p}(\mathbf{v}) \ll n$, another non-Gaussian type of limiting behavior of the posterior is well possible.

3.5 Laplace approximation with inexact parameters

Our main result of Theorem 3.3 states an approximation of the posterior distribution by the Gaussian measure with parameters $\tilde{\mathbf{v}}_G$ and \tilde{D}_G^{-2} . However, the vector

$\tilde{\mathbf{v}}_G = \operatorname{argmax}_{\mathbf{v}} L_G(\mathbf{v})$ is typically hard to compute, because it solves a high dimensional optimization problem. If $\tilde{\mathbf{v}}_G$ and thus $\tilde{D}_G = D_G(\tilde{\mathbf{v}}_G)$ are not available, one would be interested to use something more simple in place of $\tilde{\mathbf{v}}_G$. Suppose to be given a vector $\hat{\mathbf{v}}_G$ close to $\tilde{\mathbf{v}}_G$ and a matrix \tilde{D}_G^2 close to \tilde{D}_G^2 . A typical example to keep in mind corresponds to $\hat{\mathbf{v}}$ being the numerically evaluated posterior mean and H^2 being the posterior covariance, also evaluated numerically. Below we aim at presenting some sufficient conditions that ensure a reasonable approximation of the posterior by $\mathcal{N}(\hat{\mathbf{v}}, H^{-2})$ using general results on Gaussian comparison; see Section G of the Supplement (S2023). For this result we need all the conditions of Theorem 3.3 corresponding to the special case with $\hat{\mathbf{v}} = \tilde{\mathbf{v}}_G$ and $H^2 = \tilde{D}_G^2$. We write $\hat{D} = D(\hat{\mathbf{v}})$. Also we restrict ourselves to the class $\mathcal{B}_{el}(\mathbb{R}^p)$ of elliptic sets A in \mathbb{R}^p of the form

$$A = \{\mathbf{v} \in \mathbb{R}^p : \|Q(\mathbf{v} - \hat{\mathbf{v}})\| \leq \mathbf{r}\}$$

for some linear mapping $Q: \mathbb{R}^p \rightarrow \mathbb{R}^q$ and $\mathbf{r} > 0$. Given two symmetric q -matrices Σ_1, Σ_2 and a vector $\mathbf{a} \in \mathbb{R}^q$, define

$$d(\Sigma_1, \Sigma_2, \mathbf{a}) \stackrel{\text{def}}{=} \left(\frac{1}{\|\Sigma_1\|_{\text{Fr}}} + \frac{1}{\|\Sigma_2\|_{\text{Fr}}} \right) \left(\|\Sigma_1 - \Sigma_2\|_1 + \|\mathbf{a}\|^2 \right); \quad (3.12)$$

see Section G. Obviously, if $\Sigma_1 \geq \Sigma_2$, then $\|\Sigma_1 - \Sigma_2\|_1 = \operatorname{tr}(\Sigma_1 - \Sigma_2)$.

Theorem 3.5. *Suppose the conditions of Theorem 3.3 to be fulfilled. Let also $\hat{\mathbf{v}} \in \mathcal{A}_G$. Given $Q: \mathbb{R}^p \rightarrow \mathbb{R}^q$, define $\Sigma_1 = Q\tilde{D}_G^{-2}Q^\top$, $\Sigma_2 = QH^{-2}Q^\top$, $\mathbf{a} = Q(\hat{\mathbf{v}} - \tilde{\mathbf{v}}_G)$ and suppose $\|\Sigma_j\|^2 \leq 3\|\Sigma_j\|_{\text{Fr}}^2$ for $j = 1, 2$. Then*

$$\sup_{\mathbf{r} > 0} \left| \mathbb{P}(\|Q(\mathbf{v}_G - \hat{\mathbf{v}})\| \leq \mathbf{r} \mid \mathbf{Y}) - \mathbb{P}'(\|QH^{-1}\boldsymbol{\gamma}\| \leq \mathbf{r}) \right| \leq \frac{2(\diamond + e^{-\mathbf{x}})}{1 - \diamond - e^{-\mathbf{x}}} + \mathbf{C}d(\Sigma_1, \Sigma_2, \mathbf{a}),$$

where \diamond is from Theorem 3.3, $d(\cdot)$ from (3.12), and \mathbf{C} is an absolute constant.

Proof. Use Theorem E.9 with $f(\mathbf{v}) = \mathbb{E}L_G(\mathbf{v})$, $\mathbf{x}^* = \tilde{\mathbf{v}}_G$, $\mathbf{x} = \hat{\mathbf{v}}$, and $D = \tilde{D}_G$. \square

The result is particularly transparent if $H = \tilde{D}_G$ or, if these two matrices are sufficiently close. Theorem E.10 yields the following bound.

Corollary 3.6. *Under the conditions of Theorem 3.5, it holds on $\Omega(\mathbf{x})$*

$$\sup_{\mathbf{r} > 0} \left| \mathbb{P}(\|Q(\mathbf{v}_G - \hat{\mathbf{v}})\| \leq \mathbf{r} \mid \mathbf{Y}) - \mathbb{P}'(\|Q\tilde{D}_G^{-1}\boldsymbol{\gamma}\| \leq \mathbf{r}) \right| \leq \frac{2(\diamond + e^{-\mathbf{x}})}{1 - \diamond - e^{-\mathbf{x}}} + \frac{\mathbf{C}\|Q(\hat{\mathbf{v}} - \tilde{\mathbf{v}}_G)\|^2}{\|Q\tilde{D}_G^{-2}Q^\top\|_{\text{Fr}}}.$$

3.6 Laplace approximation and Bernstein–von Mises Theorem

The prominent Bernstein–von Mises (BvM) Theorem claims asymptotic normality of the posterior distribution with the mean corresponding to the standard MLE $\tilde{\mathbf{v}} = \operatorname{argmax}_{\mathbf{v}} L(\mathbf{v})$ and the variance $\tilde{D}^{-2} = D^{-2}(\tilde{\mathbf{v}})$. In particular, the prior does not

show up in this result, its impact on the posterior distribution becomes negligible as the sample size n grows. In our setup, the situation is different. The main results of Theorems 3.3 through 3.4 state another Gaussian approximation of the posterior with the mean $\tilde{\mathbf{v}}_G$ and the variance \tilde{D}_G^{-2} both depending on the prior covariance G^{-2} . This dependence is important because the accuracy of approximation is given in terms of $\tilde{\mathbf{p}} = \mathbf{p}(\tilde{\mathbf{v}}_G)$ also depending on G^{-2} . It is of interest to describe a kind of phase transition from the classical BvM approximation by $\mathcal{N}(\tilde{\mathbf{v}}, \tilde{D}^{-2})$ to the prior-dependent Laplace approximation by $\mathcal{N}(\tilde{\mathbf{v}}_G, \tilde{D}_G^{-2})$. Intuitively it is clear that the prior impact can be measured by the relation between the model-based Fisher information matrix D^2 and the prior precision matrix G^2 . The main result below confirms this intuitive guess, however, the result is not trivial and requires a careful treatment based on Theorem 3.3 and the Gaussian comparison technique mentioned in Section 3.5. Implicitly we assume that the conditions ensuring concentration of the MLE $\tilde{\mathbf{v}}$ corresponding to $G^2 = 0$ to be fulfilled. In particular, we need that $D^2(\mathbf{v})$ is sufficiently large for all \mathbf{v} in the vicinity of \mathbf{v}^* . The radius of this vicinity is given by the value $\mathbf{r}_x = \sqrt{\text{tr}(D^{-2}V^2)} + \sqrt{2\mathbf{x}}$; see (2.2) of (∇ζ). Under correct model specification, it holds $V^2 \leq \mathbf{C}D^2$ and $\mathbf{r}_x^2 \leq \mathbf{C}p$.

Theorem 3.7. *Under the conditions of Theorem 3.3, it holds on $\Omega(\mathbf{x})$*

$$\begin{aligned} & \sup_{\mathbf{r}>0} \left| \mathbb{P}(\|Q(\mathbf{v}_G - \tilde{\mathbf{v}})\| \leq \mathbf{r} \mid \mathbf{Y}) - \mathbb{P}'(\|Q\tilde{D}^{-1}\boldsymbol{\gamma}\| \leq \mathbf{r}) \right| \\ & \leq \frac{2(\diamond + e^{-\mathbf{x}})}{1 - \diamond - e^{-\mathbf{x}}} + \frac{\mathbf{C}\|Q(\tilde{\mathbf{v}} - \tilde{\mathbf{v}}_G)\|^2}{\|Q\tilde{D}^{-2}Q^\top\|_{\text{Fr}}} + \mathbf{C}\|\tilde{D}_G^{-1}G^2\tilde{D}_G^{-1}\| \frac{\text{tr}(Q\tilde{D}^{-2}Q^\top)}{\|Q\tilde{D}^{-2}Q^\top\|_{\text{Fr}}}. \end{aligned}$$

Moreover, with $Q = \tilde{D}$, $\mathbf{r}_x = \sqrt{\text{tr}(D^{-2}V^2)} + \sqrt{2\mathbf{x}} \leq \mathbf{C}\sqrt{p}$,

$$\begin{aligned} & \sup_{\mathbf{r}>0} \left| \mathbb{P}(\|\tilde{D}(\mathbf{v}_G - \tilde{\mathbf{v}})\| \leq \mathbf{r} \mid \mathbf{Y}) - \mathbb{P}'(\|\boldsymbol{\gamma}\| \leq \mathbf{r}) \right| \\ & \lesssim \diamond + e^{-\mathbf{x}} + \|G\mathbf{v}^*\|^2/\sqrt{p} + \|D_G^{-1}G^2D_G^{-1}\|^2\sqrt{p}. \end{aligned}$$

Proof. Theorem 3.5 yields in view of $\tilde{D}^{-2} \geq \tilde{D}_G^{-2}$

$$\begin{aligned} & \sup_{\mathbf{r}>0} \left| \mathbb{P}(\|Q(\mathbf{v}_G - \tilde{\mathbf{v}})\| \leq \mathbf{r} \mid \mathbf{Y}) - \mathbb{P}'(\|Q\tilde{D}^{-1}\boldsymbol{\gamma}\| \leq \mathbf{r}) \right| \\ & \leq \frac{2(\diamond + e^{-\mathbf{x}})}{1 - \diamond - e^{-\mathbf{x}}} + \frac{\mathbf{C}\|Q(\tilde{\mathbf{v}} - \tilde{\mathbf{v}}_G)\|^2}{\|Q\tilde{D}^{-2}Q^\top\|_{\text{Fr}}} + \frac{\mathbf{C}\text{tr}\{Q(\tilde{D}^{-2} - \tilde{D}_G^{-2})Q^\top\}}{\|Q\tilde{D}^{-2}Q^\top\|_{\text{Fr}}}. \end{aligned}$$

The last term here can easily be bounded:

$$\begin{aligned} \text{tr}\{Q(\tilde{D}^{-2} - \tilde{D}_G^{-2})Q^\top\} &= \text{tr}\{Q\tilde{D}^{-1}(\mathbf{I}_p - \tilde{D}\tilde{D}_G^{-2}\tilde{D})\tilde{D}^{-1}Q^\top\} \\ &\leq \|\mathbf{I}_p - \tilde{D}\tilde{D}_G^{-2}\tilde{D}\| \text{tr}(Q\tilde{D}^{-2}Q^\top) = \|\tilde{D}_G^{-1}G^2\tilde{D}_G^{-1}\| \text{tr}(Q\tilde{D}^{-2}Q^\top). \end{aligned}$$

Here we used that

$$\|\tilde{D}_G^{-1}G^2\tilde{D}_G^{-1}\| = \|\tilde{D}_G^{-1}(\tilde{D}_G^2 - \tilde{D}^2)\tilde{D}_G^{-1}\| = \|\mathbf{I}_p - \tilde{D}_G^{-1}\tilde{D}^2\tilde{D}_G^{-1}\| = \|\mathbf{I}_p - \tilde{D}\tilde{D}_G^{-2}\tilde{D}\|.$$

For $Q = \tilde{D}$, we use that $\|Q\tilde{D}^{-2}Q^\top\|_{\text{Fr}} = \sqrt{p}$, $\text{tr}(Q\tilde{D}^{-2}Q^\top) = p$, and apply the Fisher expansion (2.8) of Theorem 2.2 to $\tilde{\mathbf{v}}$ and $\tilde{\mathbf{v}}_G$. On $\Omega(\mathbf{x})$, it holds $\|D^{-1}\nabla\| \leq \mathbf{r}_x$ and with ω from (2.5)

$$\begin{aligned} \|D(\tilde{\mathbf{v}} - \tilde{\mathbf{v}}_G)\| &\leq \|D(\mathbf{v}^* - \mathbf{v}_G^*)\| + \|(\mathbf{I}_p - DD_G^{-2}D)D^{-1}\nabla\| \\ &\quad + \|D(\tilde{\mathbf{v}} - \mathbf{v}^*) - D^{-1}\nabla\| + \|D(\tilde{\mathbf{v}}_G - \mathbf{v}_G^* - D_G^{-2}\nabla)\| \\ &\leq \|D(\mathbf{v}^* - \mathbf{v}_G^*)\| + \mathbf{c}\|D_G^{-1}G^2D_G^{-1}\| \mathbf{r}_x + \mathbf{C}\omega \mathbf{r}_x. \end{aligned}$$

By (2.11) of Proposition 2.4

$$\|D(\mathbf{v}^* - \mathbf{v}_G^*)\|^2 \leq \|DD_G^{-2}D^\top\| \|G\mathbf{v}^*\|^2 \leq \|G\mathbf{v}^*\|^2.$$

As $\mathbf{r}_x^2 \leq \mathbf{C}p$, $\omega\sqrt{p} \leq \mathbf{C}\diamond_3$, $\tilde{D}^2 \leq 2D^2$, and $\tilde{D}_G^{-2} \leq 2D_G^{-2}$, the assertion follows. \square

Remark 3.1. The use of BvM requires rather strong bounds on the penalizing matrix G^2 and the related bias $\|G\mathbf{v}^*\|$. We need the condition of “light penalization” $\|D_G^{-1}G^2D_G^{-1}\| \ll p^{-1/2}$ which is much stronger than $G^2 \ll D^2$. Similarly, the “light bias” condition $\|G\mathbf{v}^*\|^2 \ll p^{1/2}$ is more restrictive than the “small bias” or “under-smoothing” condition $\|G\mathbf{v}^*\|^2 \ll p$.

3.7 Posterior mean

This section addresses an important question of using the posterior mean in place of the MAP $\tilde{\mathbf{v}}_G$ for Bayesian inference. Our main result justifies the use of the posterior mean in place of the MAP under the same critical dimension condition $\mathfrak{p}(\mathbf{v}) \ll n^{1/3}$ which is required for the Gaussian approximation result. First we quantify the deviation of the posterior mean $\bar{\mathbf{v}}_G$ from $\tilde{\mathbf{v}}_G$. Then we apply Corollary 3.6 to measure the impact of using $\bar{\mathbf{v}}_G$ in place of $\tilde{\mathbf{v}}_G$. By definition

$$\bar{\mathbf{v}}_G - \tilde{\mathbf{v}}_G \stackrel{\text{def}}{=} \mathbf{E}(\mathbf{v}_G | \mathbf{Y}) - \tilde{\mathbf{v}}_G = \frac{\int (\mathbf{v} - \tilde{\mathbf{v}}_G) e^{L_G(\mathbf{v})} d\mathbf{v}}{\int e^{L_G(\mathbf{v})} d\mathbf{v}}.$$

More precisely, we consider a linear mapping $Q: \mathbb{R}^p \rightarrow \mathbb{R}^q$ and evaluate the value $\|Q(\bar{\mathbf{v}}_G - \tilde{\mathbf{v}}_G)\|$. The choice of Q is important. In particular, we cannot take $Q = D_G$ because this choice makes the bound dependent and linearly growing with p .

Theorem 3.8. *Assume the conditions of Theorem 3.4 and let $Q^\top Q \leq D^2(\mathbf{v})$ for all $\mathbf{v} \in \mathcal{A}_G$. Then it holds with some absolute constant \mathbf{C}*

$$\|Q(\bar{\mathbf{v}}_G - \tilde{\mathbf{v}}_G)\| \leq 2.4\mathbf{c}_3 \|Q\tilde{D}^{-2}Q^\top\|^{1/2} (\tilde{\mathfrak{p}} + 1)^{3/2} n^{-1/2} + \mathbf{C}e^{-x}.$$

Proof. One can apply the same trick as before: by $\nabla L_G(\tilde{\mathbf{v}}_G) = 0$

$$Q(\bar{\mathbf{v}}_G - \tilde{\mathbf{v}}_G) = \frac{\int Q\mathbf{u} \exp\{L_G(\tilde{\mathbf{v}}_G + \mathbf{u}) - L_G(\tilde{\mathbf{v}}_G) - \langle \nabla L_G(\tilde{\mathbf{v}}_G), \mathbf{u} \rangle\} d\mathbf{u}}{\int \exp\{L_G(\tilde{\mathbf{v}}_G + \mathbf{u}) - L_G(\tilde{\mathbf{v}}_G) - \langle \nabla L_G(\tilde{\mathbf{v}}_G), \mathbf{u} \rangle\} d\mathbf{u}}.$$

For any particular value $\tilde{\mathbf{v}}_G = \mathbf{v}$, stochastic linearity allows to replace the Bregman divergence of the log-likelihood $L_G(\mathbf{v})$ by the similar one for the expected log-likelihood $f_G(\mathbf{v}) = \mathbb{E}L_G(\mathbf{v})$. This yields

$$\|Q(\bar{\mathbf{v}}_G - \tilde{\mathbf{v}}_G)\| \leq \left\| \frac{\int Q\mathbf{u} e^{f_G(\mathbf{v};\mathbf{u})} d\mathbf{u}}{\int e^{f_G(\mathbf{v};\mathbf{u})} d\mathbf{u}} \right\|.$$

Now we may apply (E.16) of Theorem E.6 of the Supplement (S2023) with $D = D_G(\mathbf{v})$. □

Corollary 3.9. *Assume the conditions of Theorem 3.4. Then*

$$\|\tilde{D}(\bar{\mathbf{v}}_G - \tilde{\mathbf{v}}_G)\| \leq 2.4c_3 (\tilde{\mathfrak{p}} + 1)^{3/2} n^{-1/2} + Ce^{-x}.$$

Now we put together the result of Theorem 3.8 and the accuracy bound from Theorem 3.5. To make the result more transparent, assume $Q = \tilde{D}$ and $H = \tilde{D}_G$.

Theorem 3.10. *Assume the conditions of Theorem 3.4. Then on $\Omega(\mathbf{x})$*

$$\sup_{\mathbf{r} > 0} \left| \mathbb{P}(\|\tilde{D}(\mathbf{v}_G - \bar{\mathbf{v}}_G)\| \leq \mathbf{r} \mid \mathbf{Y}) - \mathbb{P}'(\|\tilde{D}\tilde{D}_G^{-1}\gamma\| \leq \mathbf{r}) \right| \leq c \left(\frac{(\tilde{\mathfrak{p}} + 1)^{3/2}}{n^{1/2}} + e^{-x} \right).$$

We conclude that the use of posterior mean in place of posterior mode is possible under the same condition $\tilde{\mathfrak{p}}^3 \ll n$. This is a non-trivial result based on recent progress in Gaussian probability from Götze et al. (2019).

4 Log-density estimation

Suppose we are given a random sample X_1, \dots, X_n in \mathbb{R}^d . The density model assumes that all these random variables are independent identically distributed from some measure P with a density $f(\mathbf{x})$ with respect to a σ -finite measure μ_0 in \mathbb{R}^d . This density function is the target of estimation. By definition, the function f is non-negative, measurable, and integrates to one: $\int f(\mathbf{x}) d\mu_0(\mathbf{x}) = 1$. Here and below, the integral \int without limits means the integral over the whole space \mathbb{R}^d . If $f(\cdot)$ has a smaller support \mathcal{X} , one can restrict integration to this set. Below we parametrize the model by a linear decomposition of the log-density function. Let $\{\psi_j(\mathbf{x}), j = 1, \dots, p\}$ with $p \leq \infty$ be a collection of functions in \mathbb{R}^d (a dictionary). For each $\mathbf{v} = (v_j) \in \mathbb{R}^p$, define

$$\ell(\mathbf{x}, \mathbf{v}) \stackrel{\text{def}}{=} v_1\psi_1(\mathbf{x}) + \dots + v_p\psi_p(\mathbf{x}) - \phi(\mathbf{v}) = \langle \Psi(\mathbf{x}), \mathbf{v} \rangle - \phi(\mathbf{v}),$$

where $\Psi(\mathbf{x})$ is a vector with components $\psi_j(\mathbf{x})$ and $\phi(\mathbf{v})$ is given by

$$\phi(\mathbf{v}) \stackrel{\text{def}}{=} \log \int e^{\langle \Psi(\mathbf{x}), \mathbf{v} \rangle} d\mu_0(\mathbf{x}). \quad (4.1)$$

It is worth stressing that the data point \mathbf{x} only enters in the linear term $\langle \Psi(\mathbf{x}), \mathbf{v} \rangle$ of the log-likelihood $\ell(\mathbf{x}, \mathbf{v})$. The function $\phi(\mathbf{v})$ is entirely model-driven. Below we restrict \mathbf{v} to a subset \mathcal{Y} in \mathbb{R}^p such that $\phi(\mathbf{v})$ is well defined and the integral of $e^{\langle \Psi(\mathbf{x}), \mathbf{v} \rangle}$ is finite. Linear log-density modeling assumes

$$\log f(\mathbf{x}) = \ell(\mathbf{x}, \mathbf{v}^*) = \langle \Psi(\mathbf{x}), \mathbf{v}^* \rangle - \phi(\mathbf{v}^*) \quad (4.2)$$

for some $\mathbf{v}^* \in \mathcal{Y} \subseteq \mathbb{R}^p$. A nice feature of such representation is that the function $\log f(\mathbf{x})$ in the contrary to the density itself does not need to be non-negative. One more important benefit of using the log-density is that the stochastic part of the corresponding log-likelihood is *linear* w.r.t. the parameter \mathbf{v} . With $S = \sum_{i=1}^n \Psi(X_i)$, for a given penalty operator G^2 , the penalized log-likelihood $L_G(\mathbf{v})$ reads as

$$L_G(\mathbf{v}) = \sum_{i=1}^n \langle \Psi(X_i), \mathbf{v} \rangle - n\phi(\mathbf{v}) - \frac{1}{2} \|G\mathbf{v}\|^2 = \langle S, \mathbf{v} \rangle - n\phi(\mathbf{v}) - \frac{1}{2} \|G\mathbf{v}\|^2.$$

The penalized MLE $\tilde{\mathbf{v}}_G$ and its population counterpart \mathbf{v}_G^* are defined as

$$\tilde{\mathbf{v}}_G = \operatorname{argmax}_{\mathbf{v} \in \mathcal{Y}} L_G(\mathbf{v}), \quad \mathbf{v}_G^* = \operatorname{argmax}_{\mathbf{v} \in \mathcal{Y}} \mathbb{E}L_G(\mathbf{v}).$$

4.1 Conditions

For applying the general results of Section 2 and Section 3, it suffices to check the general conditions of Section 2 for the log-density model. First note that the generalized linear structure of the model automatically yields conditions (C) and (ζ). Indeed, convexity of $\phi(\cdot)$ implies that $\mathbb{E}L(\mathbf{v}) = \langle \mathbb{E}S, \mathbf{v} \rangle - n\phi(\mathbf{v})$ is concave. Further, for the stochastic component $\zeta(\mathbf{v}) = L(\mathbf{v}) - \mathbb{E}L(\mathbf{v})$, it holds

$$\nabla \zeta(\mathbf{v}) = \nabla \zeta = S - \mathbb{E}S = \sum_{i=1}^n [\Psi(X_i) - \mathbb{E}\Psi(X_i)],$$

and (ζ) follows. Further, the representation $\mathbb{E}L(\mathbf{v}) = \langle \mathbb{E}S, \mathbf{v} \rangle - n\phi(\mathbf{v})$ implies

$$\mathbb{F}(\mathbf{v}) = -\nabla^2 \mathbb{E}L(\mathbf{v}) = -\nabla^2 L(\mathbf{v}) = n\nabla^2 \phi(\mathbf{v}).$$

To simplify our presentation, we assume that X_1, \dots, X_n are indeed i.i.d. and the density $f(\mathbf{x})$ can be represented in the form (4.2) for some parameter vector \mathbf{v}^* . This can be easily extended to non i.i.d. case at cost of more complicated notations. Then

$$\mathbf{v}^* = \operatorname{argmax}_{\mathbf{v} \in \mathcal{Y}} \mathbb{E}L(\mathbf{v}) = \operatorname{argmax}_{\mathbf{v} \in \mathcal{Y}} \{ \langle \mathbb{E}S, \mathbf{v} \rangle - n\phi(\mathbf{v}) \} = \operatorname{argmax}_{\mathbf{v} \in \mathcal{Y}} \{ \langle \overline{\Psi}, \mathbf{v} \rangle - \phi(\mathbf{v}) \}, \quad (4.3)$$

where $\bar{\Psi} = \mathbb{E} \Psi(X_1)$ and $\mathbb{E} S = n \bar{\Psi}$. This yields the identity

$$\nabla \phi(\mathbf{v}^*) = \bar{\Psi}.$$

Moreover, by (4.1), $\nabla^2 \phi(\mathbf{v}^*) = \text{Var}\{\Psi(X_1)\}$ and

$$V^2 = \text{Var}(\nabla \zeta) = n \nabla^2 \phi(\mathbf{v}^*) = \mathbb{F}(\mathbf{v}^*). \tag{4.4}$$

Here we present our conditions. For any $\mathbf{v} \in \mathcal{Y}$ and $\varrho > 0$, define $\mathfrak{m}(\mathbf{v})$ by $\mathfrak{m}^2(\mathbf{v}) = \nabla^2 \phi(\mathbf{v})$ and consider the corresponding balls in \mathbb{R}^p

$$\mathcal{B}_\varrho(\mathbf{v}) \stackrel{\text{def}}{=} \{\mathbf{u} \in \mathbb{R}^p : \|\mathfrak{m}(\mathbf{v})\mathbf{u}\| \leq \varrho\} = \{\mathbf{u} \in \mathbb{R}^p : \langle \nabla^2 \phi(\mathbf{v}), \mathbf{u}^{\otimes 2} \rangle \leq \varrho^2\}.$$

- (f) X_1, \dots, X_n are i.i.d. from a density f satisfying $\log f(\mathbf{x}) = \Psi(\mathbf{x})^\top \mathbf{v}^* - \phi(\mathbf{v}^*)$.
- (Y) The set \mathcal{Y} is open and convex, the value $\phi(\mathbf{v})$ from (4.1) is finite for all $\mathbf{v} \in \mathcal{Y}$, \mathbf{v}^* from (4.3) is an internal point in \mathcal{Y} such that $\mathcal{B}_{2\varrho}(\mathbf{v}^*) \subset \mathcal{Y}$ for a fixed $\varrho > 0$.
- (ϕ) For the Bregman divergence $\phi(\mathbf{v}; \mathbf{u}) \stackrel{\text{def}}{=} \phi(\mathbf{v} + \mathbf{u}) - \phi(\mathbf{v}) - \langle \nabla \phi(\mathbf{v}), \mathbf{u} \rangle$, it holds

$$\sup_{\mathbf{v} \in \mathcal{B}_\varrho(\mathbf{v}^*)} \sup_{\mathbf{u} \in \mathcal{B}_{2\varrho}(\mathbf{v})} \exp \phi(\mathbf{v}; \mathbf{u}) \leq \mathbf{C}_\varrho. \tag{4.5}$$

Introduce a measure $P_{\mathbf{v}}$ by the relation:

$$\frac{dP_{\mathbf{v}}}{d\mu_0}(\mathbf{x}) = \exp\{\langle \Psi(\mathbf{x}), \mathbf{v} \rangle - \phi(\mathbf{v})\}. \tag{4.6}$$

Identity (4.1) ensures that $P_{\mathbf{v}}$ is a probabilistic measure. Moreover, under (4.2), the data generating measure \mathbb{P} coincides with $P_{\mathbf{v}^*}^{\otimes n}$.

(Ψ₄) There are $\mathbf{C}_{\Psi,3} \geq 0$ and $\mathbf{C}_{\Psi,4} \geq 3$ such that for all $\mathbf{v} \in \mathcal{B}_\varrho(\mathbf{v}^*)$ and $\gamma \in \mathbb{R}^p$

$$\begin{aligned} |E_{\mathbf{v}} \langle \Psi(X_1) - E_{\mathbf{v}} \Psi(X_1), \gamma \rangle^3| &\leq \mathbf{C}_{\Psi,3} E_{\mathbf{v}}^{3/2} \langle \Psi(X_1) - E_{\mathbf{v}} \Psi(X_1), \gamma \rangle^2, \\ E_{\mathbf{v}} \langle \Psi(X_1) - E_{\mathbf{v}} \Psi(X_1), \gamma \rangle^4 &\leq \mathbf{C}_{\Psi,4} E_{\mathbf{v}}^2 \langle \Psi(X_1) - E_{\mathbf{v}} \Psi(X_1), \gamma \rangle^2. \end{aligned}$$

In fact, conditions (ϕ) and (Ψ₄) follow from (Y) and can be considered as a kind of definition of important quantities \mathbf{C}_ϱ , $\mathbf{C}_{\Psi,3}$, and $\mathbf{C}_{\Psi,4}$ which will be used for describing the smoothness properties of $\phi(\mathbf{v})$. The matrix $\nabla^2 \phi(\mathbf{v})$ is supposed well conditioned for $\mathbf{v} \in \mathcal{B}_\varrho(\mathbf{v}^*)$.

(∇²ϕ) For the information matrix $\nabla^2 \phi(\mathbf{v})$, it holds with some $\mathbf{C}_{\mathbb{F}} \geq 1$

$$\mathbf{C}_{\mathbb{F}}^{-1} \mathbf{I}_p \leq \nabla^2 \phi(\mathbf{v}) \leq \mathbf{C}_{\mathbb{F}} \mathbf{I}_p, \quad \mathbf{v} \in \mathcal{B}_\varrho(\mathbf{v}^*). \tag{4.7}$$

Later we show that it suffices to check (4.7) at \mathbf{v}^* , then it will be fulfilled in $\mathcal{B}_\varrho(\mathbf{v}^*)$ with a slightly larger constant $\mathbf{C}_{\mathbb{F}}$.

Check of conditions (\mathcal{S}_3) and (\mathcal{S}_4)

Let $P_{\mathbf{v}}$ be defined by (4.6). It is straightforward to check that $E_{\mathbf{v}}\Psi(X_1) = \nabla\phi(\mathbf{v})$ and $\text{Var}_{\mathbf{v}}(\Psi(X_1)) = \nabla^2\phi(\mathbf{v})$. Further, if $\mathbf{u} \in \mathcal{B}_{\varrho}(\mathbf{v})$ and $\mathbf{v} + \mathbf{u} \in \mathcal{T}$, then

$$\phi(\mathbf{v} + \mathbf{u}) = \log E_0 \exp\{\langle \Psi(X_1), \mathbf{v} + \mathbf{u} \rangle\} = \log E_{\mathbf{v}} \exp\{\langle \Psi(X_1), \mathbf{u} \rangle + \phi(\mathbf{v})\}.$$

This yields in view of $E_{\mathbf{v}}\Psi(X_1) = \nabla\phi(\mathbf{v})$ that $\boldsymbol{\varepsilon} = \Psi(X_1) - E_{\mathbf{v}}\Psi(X_1)$ fulfills

$$\begin{aligned} \log E_{\mathbf{v}} \exp\langle \boldsymbol{\varepsilon}, \mathbf{u} \rangle &= \phi(\mathbf{v} + \mathbf{u}) - \phi(\mathbf{v}) - \langle E_{\mathbf{v}}\Psi(X_1), \mathbf{u} \rangle \\ &= \phi(\mathbf{v} + \mathbf{u}) - \phi(\mathbf{v}) - \langle \nabla\phi(\mathbf{v}), \mathbf{u} \rangle. \end{aligned} \quad (4.8)$$

Lemma 4.1. *The function $\phi(\mathbf{v})$ satisfies for any $\mathbf{v} \in \mathcal{B}_{\varrho}(\mathbf{v}^*)$ and $\boldsymbol{\gamma} \in \mathbb{R}^p$*

$$|\langle \nabla^3\phi(\mathbf{v}), \boldsymbol{\gamma}^{\otimes 3} \rangle| \leq \mathbf{C}_{\Psi,3} \langle \nabla^2\phi(\mathbf{v}), \boldsymbol{\gamma}^{\otimes 2} \rangle^{3/2}, \quad (4.9)$$

$$|\langle \nabla^4\phi(\mathbf{v}), \boldsymbol{\gamma}^{\otimes 4} \rangle| \leq (\mathbf{C}_{\Psi,4} - 3) \langle \nabla^2\phi(\mathbf{v}), \boldsymbol{\gamma}^{\otimes 2} \rangle^2. \quad (4.10)$$

Moreover, for any $\boldsymbol{\gamma}_1, \boldsymbol{\gamma} \in \mathbb{R}^p$

$$|\langle \nabla^3\phi(\mathbf{v}), \boldsymbol{\gamma}_1 \otimes \boldsymbol{\gamma}^{\otimes 2} \rangle| \leq \sqrt{\mathbf{C}_{\Psi,4} \langle \nabla^2\phi(\mathbf{v}), \boldsymbol{\gamma}_1^{\otimes 2} \rangle} \langle \nabla^2\phi(\mathbf{v}), \boldsymbol{\gamma}^{\otimes 2} \rangle.$$

Proof. Denote $\boldsymbol{\varepsilon} = X_1 - E_{\mathbf{v}}X_1$. By (4.8) with $\mathbf{u} = t\boldsymbol{\gamma}$ for t sufficiently small

$$\chi(t) \stackrel{\text{def}}{=} \log E_{\mathbf{v}} \exp(t\langle \boldsymbol{\varepsilon}, \boldsymbol{\gamma} \rangle) = \phi(\mathbf{v} + t\boldsymbol{\gamma}) - \phi(\mathbf{v}) - \langle \nabla\phi(\mathbf{v}), t\boldsymbol{\gamma} \rangle,$$

and by (Ψ_4) with $\mathbf{C}_{\Psi,4} \geq 3$

$$\begin{aligned} |\chi^{(3)}(0)| &= |E_{\mathbf{v}}\langle \boldsymbol{\varepsilon}, \boldsymbol{\gamma} \rangle^3| \leq \mathbf{C}_{\Psi,3} E_{\mathbf{v}}^{3/2}\langle \boldsymbol{\varepsilon}, \boldsymbol{\gamma} \rangle^2, \\ |\chi^{(4)}(0)| &= |E_{\mathbf{v}}\langle \boldsymbol{\varepsilon}, \boldsymbol{\gamma} \rangle^4 - 3E_{\mathbf{v}}^2\langle \boldsymbol{\varepsilon}, \boldsymbol{\gamma} \rangle^2| \leq (\mathbf{C}_{\Psi,4} - 3)E_{\mathbf{v}}^2\langle \boldsymbol{\varepsilon}, \boldsymbol{\gamma} \rangle^2. \end{aligned}$$

If $\boldsymbol{\gamma}_1 \neq \boldsymbol{\gamma}$ then we may proceed in a similar way with the bivariate function $\chi(t_1, t) = \log E_{\mathbf{v}} \exp\{t_1\langle \boldsymbol{\varepsilon}, \boldsymbol{\gamma}_1 \rangle + t\langle \boldsymbol{\varepsilon}, \boldsymbol{\gamma} \rangle\}$. Its mixed derivative at zero satisfies

$$\left| \frac{\partial^3}{\partial t_1 \partial t^2} \chi(0, 0) \right| = |E_{\mathbf{v}} \langle \boldsymbol{\varepsilon}, \boldsymbol{\gamma}_1 \rangle \langle \boldsymbol{\varepsilon}, \boldsymbol{\gamma} \rangle^2| \leq \{E_{\mathbf{v}} \langle \boldsymbol{\varepsilon}, \boldsymbol{\gamma}_1 \rangle^2 E_{\mathbf{v}} \langle \boldsymbol{\varepsilon}, \boldsymbol{\gamma} \rangle^4\}^{1/2}$$

and the result follows as well. \square

Lemma 4.2. *If $\mathbf{v} \in \mathcal{B}_{\varrho}(\mathbf{v}^*)$ then with $c_{\phi} = \sqrt{\mathbf{C}_{\Psi,4} \mathbf{C}_{\varrho}}$*

$$\sup_{\mathbf{u} \in \mathcal{B}_{\varrho}(\mathbf{v})} \sup_{\boldsymbol{\gamma} \in \mathbb{R}^p} \frac{\langle \nabla^2\phi(\mathbf{v} + \mathbf{u}), \boldsymbol{\gamma}^{\otimes 2} \rangle}{\langle \nabla^2\phi(\mathbf{v}), \boldsymbol{\gamma}^{\otimes 2} \rangle} \leq c_{\phi}. \quad (4.11)$$

Proof. Let $\langle \nabla^2 \phi(\mathbf{v}), \mathbf{u}^{\otimes 2} \rangle \leq \varrho^2$. By (4.8) with $\boldsymbol{\varepsilon} = X_1 - E_{\mathbf{v}} X_1$

$$\nabla^2 \phi(\mathbf{v} + \mathbf{u}) = \nabla^2 \log E_{\mathbf{v}} e^{\langle \boldsymbol{\varepsilon}, \mathbf{u} \rangle} = \frac{E_{\mathbf{v}} \{ \boldsymbol{\varepsilon} \boldsymbol{\varepsilon}^\top e^{\langle \boldsymbol{\varepsilon}, \mathbf{u} \rangle} \}}{(E_{\mathbf{v}} e^{\langle \boldsymbol{\varepsilon}, \mathbf{u} \rangle})^2} - \frac{E_{\mathbf{v}} \{ \boldsymbol{\varepsilon} e^{\langle \boldsymbol{\varepsilon}, \mathbf{u} \rangle} \} E_{\mathbf{v}} \{ \boldsymbol{\varepsilon} e^{\langle \boldsymbol{\varepsilon}, \mathbf{u} \rangle} \}^\top}{(E_{\mathbf{v}} e^{\langle \boldsymbol{\varepsilon}, \mathbf{u} \rangle})^2}$$

and by (4.10) and (4.5) in view of $E_{\mathbf{v}} e^{\langle \boldsymbol{\varepsilon}, \mathbf{u} \rangle} \geq 1$

$$\begin{aligned} \langle \nabla^2 \phi(\mathbf{v} + \mathbf{u}), \boldsymbol{\gamma}^{\otimes 2} \rangle &\leq E_{\mathbf{v}} \{ \langle \boldsymbol{\varepsilon}, \boldsymbol{\gamma} \rangle^2 e^{\langle \boldsymbol{\varepsilon}, \mathbf{u} \rangle} \} \\ &\leq E_{\mathbf{v}}^{1/2} \langle \boldsymbol{\varepsilon}, \boldsymbol{\gamma} \rangle^4 E_{\mathbf{v}}^{1/2} e^{2\langle \boldsymbol{\varepsilon}, \mathbf{u} \rangle} \leq \sqrt{C_{\Psi,4} C_{\varrho}} \langle \nabla^2 \phi(\mathbf{v}), \boldsymbol{\gamma}^{\otimes 2} \rangle \end{aligned}$$

and the assertion follows. □

Lemma 4.3. *Let $\mathbf{v} \in \mathcal{B}_{\varrho}(\mathbf{v}^*)$ and $\mathbf{r} \leq \varrho \sqrt{n}$. Then $f(\mathbf{v}) = \mathbb{E}_{\mathbf{v}^*} L(\mathbf{v})$ satisfies (\mathcal{S}_3) and (\mathcal{S}_4) with $h(\mathbf{v}) = \langle \nabla \phi(\mathbf{v}^*), \mathbf{v} \rangle - \phi(\mathbf{v})$, $\mathfrak{m}^2(\mathbf{v}) = \nabla^2 \phi(\mathbf{v})$, and constants c_3 and c_4 depending on C_{ϱ} , $C_{\Psi,3}$, and $C_{\Psi,4}$ only.*

Proof. Let $\mathbf{v} \in \mathcal{B}_{\varrho}(\mathbf{v}^*)$. For any \mathbf{u} with $\|\mathfrak{m}(\mathbf{v})\mathbf{u}\| \leq \mathbf{r}/\sqrt{n} \leq \varrho$, by (4.9) and (4.11)

$$\frac{|\langle \nabla^3 \phi(\mathbf{v} + t\mathbf{u}), \mathbf{u}^{\otimes 3} \rangle|}{\|\mathfrak{m}(\mathbf{v})\mathbf{u}\|^3} \leq \frac{C_{\Psi,3} \|\mathfrak{m}(\mathbf{v} + t\mathbf{u})\mathbf{u}\|^3}{\|\mathfrak{m}(\mathbf{v})\mathbf{u}\|^3} \leq C_{\Psi,3} c_{\phi}^{3/2},$$

and (\mathcal{S}_3) follows with $c_3 = C_{\Psi,3} c_{\phi}^{3/2}$. The proof of (\mathcal{S}_4) is similar. □

Check of $(\nabla \zeta)$

Now we check the deviation bound for $\nabla \zeta = \mathbf{S} - \mathbb{E} \mathbf{S}$ under (f) and (\mathcal{R}) . I.i.d. structure of $\mathbf{S} = \sum_i X_i$ and (4.4) yield $\text{Var}(\mathbf{S}) = V^2 = n \nabla^2 \phi(\mathbf{v}^*)$. Further, for any $\mathbf{u} \in \mathcal{B}_{\varrho}(\mathbf{v}^*)$, again by the i.i.d. assumption and by (4.8)

$$n^{-1} \log \mathbb{E}_{\mathbf{v}^*} \exp \{ \langle \nabla \zeta, \mathbf{u} \rangle \} = \log \mathbb{E}_{\mathbf{v}^*} e^{\langle \boldsymbol{\varepsilon}, \mathbf{u} \rangle} = \phi(\mathbf{v}^* + \mathbf{u}) - \phi(\mathbf{v}^*) - \langle \nabla \phi(\mathbf{v}^*), \mathbf{u} \rangle.$$

Fix $\mathbf{r} \leq \varrho n^{1/2}$ and consider all \mathbf{u} with $n \langle \nabla^2 \phi(\mathbf{v}^*), \mathbf{u}^{\otimes 2} \rangle \leq \mathbf{r}^2$. Then by (\mathcal{S}_3) and (A.3) of Lemma A.2 of the Supplement (S2023)

$$\phi(\mathbf{v}^* + \mathbf{u}) - \phi(\mathbf{v}^*) - \langle \nabla \phi(\mathbf{v}^*), \mathbf{u} \rangle \leq \frac{1 + c_3 \mathbf{r} n^{-1/2} / 3}{2} \langle \nabla^2 \phi(\mathbf{v}^*), \mathbf{u}^{\otimes 2} \rangle \leq \langle \nabla^2 \phi(\mathbf{v}^*), \mathbf{u}^{\otimes 2} \rangle$$

provided that $c_3 \mathbf{r} \leq 3n^{1/2}$. This implies (F.37) with $\mathbf{g} = \varrho \sqrt{n}$ and thus, the deviation bound (F.43) of Theorem F.15 implies $(\nabla \zeta)$ for $\varrho \sqrt{n}$ sufficiently large.

4.2 Smoothness constraints, bias-variance trade-off

To handle the bias term, we impose some smoothness conditions on the underlying density parameter \mathbf{v}^* ; see Section C of the Supplement (S2023). We also limit ourselves to the penalty matrices G^2 which ensure a kind of bias-variance trade-off.

$$(\mathbf{G}_0\mathbf{v}^*) \quad \|\mathbf{G}_0\mathbf{v}^*\|^2 \leq 1 \text{ for some fixed } \mathbf{G}_0^2.$$

Later we follow the suggestion of Section C and apply the penalizing matrix $G^2 = w\mathbf{G}_0^2$. Set $\mathbf{D}^2 = n\nabla^2\phi(\mathbf{v}^*)$, $\mathbf{D}_G^2 = \mathbf{D}^2 + G^2$. The particular value $w = w_*$ can be selected by the bias-variance relation (C.3). First we evaluate the bias term. This is important to ensure that the point \mathbf{v}_G^* is still in the local vicinity $\mathcal{B}_\varrho(\mathbf{v}^*)$.

Proposition 4.4. *Assume (\mathbf{f}) , (ϕ) , (Ψ_4) , (Υ) , $(\mathbf{G}_0\mathbf{v}^*)$. Let $G^2 = w\mathbf{G}_0^2$, $\mathbf{D}_G^2 = n\nabla^2\phi(\mathbf{v}^*) + G^2$, and $\nu^{-1}w^{1/2}n^{-1/2} \leq \varrho$. Then $\mathbf{v}_G^* \in \mathcal{B}_\varrho(\mathbf{v}^*)$ and*

$$\|\mathbf{D}_G(\mathbf{v}_G^* - \mathbf{v}^*)\| \leq \nu^{-1}w^{1/2}. \quad (4.12)$$

Proof. We intend to apply Proposition 2.4 with $Q = \mathbf{D}$ and $\mathbf{r}_o = \nu^{-1}\|\mathbf{D}_G^{-1}G^2\mathbf{v}^*\|$ for $\nu = 2/3$. It holds by $(\mathbf{G}_0\mathbf{v}^*)$ in view of $G^2 \leq \mathbf{D}_G^2$ and $G^2 = w\mathbf{G}_0^2$

$$\mathbf{r}_o = \nu^{-1}\|\mathbf{D}_G^{-1}G^2\mathbf{v}^*\| \leq \nu^{-1}w^{1/2}\|\mathbf{G}_0\mathbf{v}^*\| = \nu^{-1}w^{1/2}.$$

Further, $\nu^{-1}w^{1/2}n^{-1/2} \leq \varrho$ ensures that the set $\{\mathbf{v}: \|\mathbf{D}(\mathbf{v} - \mathbf{v}^*)\| \leq \mathbf{r}_o\}$ belongs to the ball $\mathcal{B}_\varrho(\mathbf{v}^*)$. Now Lemma 4.3 yields (\mathbf{S}_3) for all $\mathbf{v} \in \mathcal{B}_\varrho(\mathbf{v}^*)$. By Lemma A.3, it holds $\omega_G^* \leq \mathbf{c}_3\mathbf{r}_o \leq \mathbf{c}_3\nu^{-1}w^{1/2}n^{-1/2} \leq 1/3$ for $n \geq n_0$. Now Proposition 2.4 yields (4.12) and $\mathbf{v}_G^* \in \mathcal{B}_\varrho(\mathbf{v}^*)$. Also by Lemma A.3, the matrix $D_G^2 = D_G^2(\mathbf{v}_G^*)$ satisfies $(1 - \omega_G^*)\mathbf{D}_G^2 \leq D_G^2 \leq (1 + \omega_G^*)\mathbf{D}_G^2$. \square

For $G^2 = w\mathbf{G}_0^2$, the effective dimension \mathbf{p}_G is given by (C.2):

$$\mathbf{p}_G = \mathbf{p}(w) = \text{tr}(\mathbf{D}^2 \mathbf{D}_G^{-2}) = \text{tr}\{\mathbf{D}^2(\mathbf{D}^2 + w\mathbf{G}_0^2)^{-1}\}.$$

A particular value w_* is defined by the bias-variance relation reads $\mathbf{p}(w) \asymp w$; see (C.3). Fix some \mathbf{x} and consider

$$\mathbf{r}_G = \mathbf{r}(w) = \sqrt{\mathbf{p}(w)} + \sqrt{2\mathbf{x}}.$$

Now expansion (2.13) of Theorem 2.6 applies provided that the concentration set $\|D_G(\mathbf{v} - \mathbf{v}_G^*)\| \leq \nu^{-1}\mathbf{r}_G$ is contained in $\mathcal{B}_\varrho(\mathbf{v}_G^*)$. Proposition 4.4 allows to use \mathbf{D}_G in place of D_G in this condition. In our results, \mathbf{C} stands for a fixed constant depending on the other constants in our conditions like s_0 , \mathbf{C}_F , \mathbf{c}_ϕ , ϱ , and $\mathbf{C}_{\Psi,4}$.

Theorem 4.5. *Assume (\mathbf{f}) , (ϕ) , (Ψ_4) , (Υ) , $(\mathbf{G}_0\mathbf{v}^*)$, and let $G^2 = w\mathbf{G}_0^2$. If $\nu^{-1}w^{1/2}n^{-1/2} \leq \varrho$ with $\nu = 2/3$ and $\mathbf{r}(w)n^{-1/2} \leq \varrho$, then on $\Omega(\mathbf{x})$*

$$\|\mathbf{D}_G(\tilde{\mathbf{v}}_G - \mathbf{v}^*)\| \leq 2\nu^{-1}\mathbf{r}(w) + \nu^{-1}w^{1/2}.$$

If $\mathfrak{r}(w_*) = \mathfrak{C}_0 w_*^{1/2}$ then it holds on $\Omega(\mathbf{x})$ for $G^2 = w_* G_0^2$ and $\mathbb{D}_{w_*}^2 = \mathbb{D}^2 + w_* G_0^2$

$$\|\mathbb{D}_{w_*}(\tilde{\mathbf{v}}_{w_*} - \mathbf{v}^*)\| \leq \mathfrak{C} \mathfrak{C}_0 w_*^{1/2}.$$

Under the same conditions one can specify the results of Section 2 including the Fisher–Wilks expansions of Theorem 2.2. Moreover, for a Gaussian prior $\mathcal{N}(0, G^{-2})$ with $G^2 = w G_0^2$, one can derive the results about concentration and contraction of the posterior $\mathbf{v}_G \mid \mathbf{X}$. Only Laplace’s approximation of the posterior in Theorem 3.4 requires a stronger condition on critical dimension: $\mathfrak{c}_3 \nu^{-1} \mathfrak{r}(w) \mathfrak{p}(w) \leq 2n^{1/2}$.

4.3 Rate optimality under Sobolev smoothness

To state standard rate-optimal results and to compare our conclusions with the existing results in the literature, we consider the univariate case $d = 1$ and introduce the condition on Sobolev smoothness of the log-density $\log f(\mathbf{x})$.

$$(\mathbf{s}_0, \mathbf{w}_0) \quad \mathbf{v}^* \in \mathcal{B}(s_0, w_0) = \left\{ \mathbf{v} = (v_j) : \sum_j v_j^2 j^{2s_0} \leq w_0 \right\} \text{ for } s_0 > 0.$$

The assumption $\mathbf{v}^* \in \mathcal{B}(s_0, w_0)$ is standard in log-density estimation; cf. Castillo and Nickl (2014) or Spokoiny and Panov (2021) with $s_0 > 1$. Rousseau and Szabo (2017) requires $s_0 > 1/2$ while our results below are valid for $s_0 > 0$. The only exception is the final Gaussian approximation result requiring $s_0 > 1$. We also assume that $w_0 \asymp 1$.

Now we state the results under the smoothness condition $(\mathbf{s}_0, \mathbf{w}_0)$. The precision/penalization matrix G^2 is taken of the (s, w) -form: $G^2 = \text{diag}(g_1^2, \dots, g_p^2)$ with $g_j^2 = j^{2s}/w$. One can take any degree $s > 1/2$, $s \geq s_0$, we recommend a large value like $s = 4$ or $s = 5$. Only the factor w should be fixed carefully to get the optimal accuracy of estimation from the relation $(wn)^{1/(2s)} \approx (w_0 n)^{1/(2s_0+1)}$; see (C.12) or (C.13). Under $(\nabla^2 \phi)$, the corresponding effective dimension \mathfrak{p}_G and the Laplace effective dimension $\mathfrak{p}(\mathbf{v})$ are determined by the index m for which $g_m^2 \approx n$. Alternatively one can use a m_0 -truncation prior with $m_0 \approx (w_0 n)^{1/(2s_0+1)}$. Such a choice of the prior parameters is frequently used for nonparametric *rate optimal* results about concentration of the pMLE and posterior contraction; cf. Castillo and Nickl (2014), Castillo and Rousseau (2015), Rousseau and Szabo (2017). Note that the mentioned results require at least $s_0 > 1/2$, while our concentration and contraction results apply under $s_0 > 0$. Theorem C.4 and C.5 of the Supplement (S2023) yield the following results.

Theorem 4.6. Assume (f) , (\mathcal{Y}) , (ϕ) , (Ψ_4) , $(\nabla^2 \phi)$ and $(\mathbf{s}_0, \mathbf{w}_0)$. Fix $s > 1/2$, $s \geq s_0$ and define $g_j^2 = w^{-1} j^{2s}$ with w satisfying $(wn)^{1/(2s)} \approx (w_0 n)^{1/(2s_0+1)}$; see (C.12) or (C.13). Then on $\Omega(\mathbf{x})$ for $n \geq n_0$,

$$\|\tilde{\mathbf{v}}_G - \mathbf{v}^*\| \leq \mathfrak{C} n^{-s_0/(2s_0+1)},$$

and the posterior measure $\mathbf{v}_G \mid \mathbf{X}$ satisfies

$$\mathbb{P}\left(\|\mathbf{v}_G - \tilde{\mathbf{v}}_G\| > \mathfrak{C} n^{-s_0/(2s_0+1)} \mid \mathbf{X}\right) \leq e^{-x}.$$

Our main result about Gaussian approximation requires more smoothness of the log-density $\log f(\cdot)$. Namely, we require $(\mathbf{s}_0, \mathbf{w}_0)$ with $s_0 > 1$.

Theorem 4.7. *Assume the conditions of Theorem 4.6 and $(\mathbf{s}_0, \mathbf{w}_0)$ with $s_0 > 1$. Define $m_0 = (w_0 n)^{1/(2s_0+1)}$. For the m -truncation prior with $m = m_0$ or for the (s, w) -smooth prior with $(wn)^{1/(2s)} = (w_0 n)^{1/(2s_0+1)} = m_0$, on $\Omega(\mathbf{x})$ for $n \geq n_0$,*

$$\begin{aligned} \sup_{A \in \mathcal{B}(\mathbb{R}^p)} \left| \mathbb{P}(\mathbf{v}_G - \tilde{\mathbf{v}}_G \in A \mid \mathbf{X}) - \mathbb{P}'(\tilde{D}_G^{-1} \boldsymbol{\gamma} \in A) \right| &\leq \mathfrak{C} n^{(1-s_0)/(2s_0+1)}, \\ \sup_{A \in \mathcal{B}_s(\mathbb{R}^p)} \left| \mathbb{P}(\mathbf{v}_G - \tilde{\mathbf{v}}_G \in A \mid \mathbf{X}) - \mathbb{P}'(\tilde{D}_G^{-1} \boldsymbol{\gamma} \in A) \right| &\leq \mathfrak{C} n^{(2-2s_0)/(2s_0+1)}. \end{aligned} \quad (4.13)$$

The error terms in (4.13) tend to zero as $n \rightarrow \infty$ because $s_0 > 1$.

Supplementary Material

Supplement to “Laplace approximation and the use of posterior mean in Bayesian inference” (DOI: [10.1214/23-BA1391SUPP](https://doi.org/10.1214/23-BA1391SUPP); .pdf).

- A. Local smoothness conditions. Introduces smoothness conditions used in the text and their implications.
- B. Examples of priors. Typical examples of priors such as truncation and smooth priors.
- C. Smooth priors and rate over Sobolev classes. Classical rate results for smooth priors.
- D. Non-Gaussian priors. Explains how the results extend to non-Gaussian priors.
- E. Dimension free bounds for Laplace approximation. Presents further results and the proofs of the results from Section 3 about Laplace approximation of the posterior.
- F. Deviation bounds for quadratic forms. This section collects technical results about Gaussian and non-Gaussian quadratic forms.
- G. Gaussian comparison. This section presents an important result about Gaussian comparison from Götze et al. (2019).

References

- Castillo, I. and Nickl, R. (2014). On the Bernstein–von Mises phenomenon for non-parametric Bayes procedures. *Annals of Statistics*, 42(5):1941–1969. [MR3262473](https://doi.org/10.1214/14-AOS1246). doi: <https://doi.org/10.1214/14-AOS1246>. 1327
- Castillo, I. and Rousseau, J. (2015). A Bernstein–von Mises theorem for smooth functionals in semiparametric models. *Annals of Statistics*, 43(6):2353–2383. [MR3405597](https://doi.org/10.1214/15-AOS1336). doi: <https://doi.org/10.1214/15-AOS1336>. 1327
- Durmus, A. and Moulines, É. (2019). High-dimensional Bayesian inference via the unadjusted Langevin algorithm. *Bernoulli*, 25(4A):2854–2882. [MR4003567](https://doi.org/10.3150/18-BEJ1073). doi: <https://doi.org/10.3150/18-BEJ1073>. 1305
- Frazier, P. I. (2018). A tutorial on Bayesian optimization. <https://arxiv.org/1807.02811>. 1305

- Giordano, M. and Kekkonen, H. (2020). Bernstein–von Mises theorems and uncertainty quantification for linear inverse problems. *SIAM/ASA Journal on Uncertainty Quantification*, 8(1):342–373. MR4069334. doi: <https://doi.org/10.1137/18M1226269>. 1304
- Götze, F., Naumov, A., Spokoiny, V., and Ulyanov, V. (2019). Large ball probabilities, Gaussian comparison and anti-concentration. *Bernoulli*, 25(4A):2538–2563. arXiv:1708.08663. MR4003557. doi: <https://doi.org/10.3150/18-BEJ1062>. 1306, 1321, 1328
- Helin, T. and Kretschmann, R. (2022). Non-asymptotic error estimates for the Laplace approximation in Bayesian inverse problems. *Numerische Mathematik*, 150(2). MR4382587. doi: <https://doi.org/10.1007/s00211-021-01266-9>. 1304
- Knapik, B. T., Szabó, B. T., van der Vaart, A. W., and van Zanten, J. H. (2016). Bayes procedures for adaptive inference in inverse problems for the white noise model. *Probab. Theory Related Fields*, 164(3-4):771–813. MR3477780. doi: <https://doi.org/10.1007/s00440-015-0619-7>. 1304
- Lu, Y. (2017). On the Bernstein-von Mises theorem for high dimensional nonlinear Bayesian inverse problems. <https://arxiv.org/1706.00289>. 1304
- Ma, Y., Chen, Y., Jin, C., Flammarion, N., and Jordan, M. I. (2019). Sampling can be faster than optimization. *Proceedings of the National Academy of Sciences*, 42(116). <http://arxiv.org/1811.08413>. MR4025861. doi: <https://doi.org/10.1073/pnas.1820003116>. 1305
- Mockus, J. (1989). *Bayesian approach to global optimization. Theory and applications*, volume 37 of *Mathematics and Its Applications. Soviet Series* Dordrecht etc.: Kluwer Academic Publishers. MR1111483. doi: <https://doi.org/10.1007/978-94-009-0909-0>. 1305
- Monard, F., Nickl, R., Paternain, G. P., et al. (2019). Efficient nonparametric Bayesian inference for X-ray transforms. *The Annals of Statistics*, 47(2):1113–1147. MR3909962. doi: <https://doi.org/10.1214/18-AOS1708>. 1304
- Nesterov, Y. and Spokoiny, V. (2017). Random gradient-free minimization of convex functions. *Foundations of Computational Mathematics*, 17(2):527–566. MR3627456. doi: <https://doi.org/10.1007/s10208-015-9296-2>. 1305
- Nickl, R. (2020). Bernstein–von Mises theorems for statistical inverse problems I: Schrödinger equation. *Journal of the European Mathematical Society*, 22:2697–2750. MR4118619. doi: <https://doi.org/10.4171/JEMS/975>. 1304
- Rousseau, J. and Szabo, B. (2017). Asymptotic behaviour of the empirical Bayes posteriors associated to maximum marginal likelihood estimator. *The Annals of Statistics*, 45(2):833–865. MR3650402. doi: <https://doi.org/10.1214/16-AOS1469>. 1306, 1327
- Rousseau, J. and Szabo, B. (2020). Asymptotic frequentist coverage properties of Bayesian credible sets for sieve priors. *The Annals of Statistics*, 48(4):2155–2179. MR4134790. doi: <https://doi.org/10.1214/19-AOS1881>. 1304

- Schillings, C., Sprungk, B., and Wacker, P. (2020). On the convergence of the Laplace approximation and noise-level-robustness of Laplace-based Monte Carlo methods for Bayesian inverse problems. *Numerische Mathematik*, 145:915–971. MR4125981. doi: <https://doi.org/10.1007/s00211-020-01131-1>. 1304
- Spokoiny, V. (2017). Penalized maximum likelihood estimation and effective dimension. *AIHP*, 53(1):389–429. arXiv:1205.0498. MR3606746. doi: <https://doi.org/10.1214/15-AIHP720>. 1304, 1307
- Spokoiny, V. (2019). Bayesian inference for nonlinear inverse problems. <https://arxiv.org/1912.12694>. 1305, 1307
- Spokoiny, V. (2023). “Supplementary Material for “Inexact Laplace Approximation and the Use of Posterior Mean in Bayesian Inference““ *Bayesian Analysis*. doi: <https://doi.org/10.1214/23-BA1391SUPP>. 1304
- Spokoiny, V. and Panov, M. (2021). Accuracy of Gaussian approximation for high-dimensional posterior distributions. *Bernoulli*. in print. arXiv:1910.06028. 1303, 1304, 1305, 1306, 1307, 1316, 1317, 1327
- Szabó, B., van der Vaart, A. W., and van Zanten, J. H. (2015). Frequentist coverage of adaptive nonparametric Bayesian credible sets. *The Annals of Statistics*, 43(4):1391–1428. MR3357861. doi: <https://doi.org/10.1214/14-AOS1270>. 1304
- Van der Vaart, A. W. (2000). *Asymptotic Statistics*, volume 3. Cambridge university press. MR1652247. doi: <https://doi.org/10.1017/CB09780511802256>. 1303