



# From robust tests to Bayes-like posterior distributions

Yannick Baraud<sup>1</sup> 

Received: 3 July 2022 / Revised: 25 April 2023 / Accepted: 16 June 2023 /  
Published online: 13 July 2023

© The Author(s), under exclusive licence to Springer-Verlag GmbH Germany, part of Springer Nature 2023

## Abstract

In the Bayes paradigm, given a loss function and an  $n$ -sample, we present the construction of a new type of posterior distribution, that extends the classical Bayes one. The loss functions we have in mind are either those derived from the total variation and Hellinger distances or some  $\mathbb{L}_j$ -ones for  $j > 1$ . We prove that, with a probability close to one, this new posterior distribution concentrates its mass in a neighbourhood (for the chosen loss function) of the law of the data, provided that this law belongs to the support of the prior or, at least, lies close enough to it. We therefore establish that the new posterior distribution enjoys some robustness properties with respect to a possible misspecification of the prior, or more precisely, its support. We also show that the posterior distribution is stable with respect to the equidistribution assumption we started from. Besides, when the model is regular and well-specified and one uses the squared Hellinger loss, we show that our credible regions possess, at least for  $n$  sufficiently large, the same ellipsoidal shapes and approximately the same sizes as those we would derive from the classical Bayesian posterior distribution by using the Bernstein–von Mises theorem. Then we use our Bayesian-like approach to solve the following problems. We first consider the estimation of a location parameter or both the location and scale parameters of a density in a nonparametric framework. Then we tackle the problem of estimating a density, with the squared Hellinger loss, in a high-dimensional parametric model under some sparsity conditions on the parameter. Importantly, the results established in this paper are nonasymptotic and provide, as much as possible, bounds with explicit constants.

**Keywords** Bayes procedure · Gibbs estimator · Posterior distribution · Robustness · Hellinger distance · Total variation distance

---

This project has received funding from the European Union’s Horizon 2020 research and innovation programme under Grant Agreement No. 811017.

---

✉ Yannick Baraud  
yannick.baraud@uni.lu

<sup>1</sup> Department of Mathematics, University of Luxembourg, Maison du nombre, 6 Avenue de la Fonte, 4364 Esch-sur-Alzette, Grand Duchy of Luxembourg

**Mathematics Subject Classification** Primary 62G05 · 62G35 · 62F35 · 62F15

## 1 Introduction

Observe  $n$  i.i.d. random variables  $X_1, \dots, X_n$  with values in a measurable space  $(E, \mathcal{E})$  and assume that their common distribution  $P^*$  belongs to a family  $\mathcal{M}$  of candidate probabilities, or at least lies close enough to it in a suitable sense. We consider the problem of estimating  $P^*$  from the observation of  $X = (X_1, \dots, X_n)$  and we evaluate the performance of an estimator with values in  $\mathcal{M}$  by means of a given loss function  $\ell: \mathcal{P} \times \mathcal{M} \rightarrow \mathbb{R}_+$ , where  $\mathcal{P}$  denotes a set of probabilities containing  $P^*$ .

Our approach to solve this estimation problem has a Bayesian flavour. We endow  $\mathcal{M}$  with a  $\sigma$ -algebra  $\mathcal{A}$  and a probability measure  $\pi$  that plays the same role as the prior in the classical Bayes paradigm. Our aim is to design a posterior distribution  $\hat{\pi}_X$ , solely based on  $X$  and the choice of  $\ell$ , that concentrates its mass, with a probability close to one, on an  $\ell$ -ball, namely a set of the form

$$\mathcal{B}(P^*, r) = \{P \in \mathcal{M}, \ell(P^*, P) \leq r\} \quad \text{with } r > 0. \quad (1)$$

This means that with a probability close to 1, a point  $\hat{P}$  which is randomly drawn according to our (random) distribution  $\hat{\pi}_X$  is likely to estimate  $P^*$  with an accuracy (with respect to the chosen loss  $\ell$ ) not larger than  $r$ . Our objective is to design  $\hat{\pi}_X$  in such a way that this concentration property holds for small values of  $r$  and under mild assumptions on  $P^*$  and  $\mathcal{M}$ .

In the literature, many authors have studied the concentration properties of the classical Bayes posterior distribution on Hellinger balls. We refer to the pioneering papers by van der Vaart and his co-authors—see for example Ghosal, Ghosh and van der Vaart [19]. They show that the concentration property around  $P^*$  holds, as  $n$  tends to infinity, provided that the prior  $\pi$  puts enough mass on sets of the form  $\mathcal{K}(P^*, \varepsilon) = \{P \in \mathcal{M}, K(P^*, P) < \varepsilon\}$  where  $\varepsilon$  is a positive number and  $K(P^*, P)$  the Kullback–Leibler divergence between  $P^*$  and  $P$ . This assumption may, however, be quite restrictive even in the favorable situation where  $P^*$  belongs to the model  $\mathcal{M}$ . Such sets may indeed be empty, and the condition therefore unsatisfied, when the probabilities in  $\mathcal{M}$  are not equivalent. This is for example the case when  $\mathcal{M}$  is the set of all uniform distributions  $P_\theta$  on  $[\theta - 1/2, \theta + 1/2]$ , with  $\theta \in \mathbb{R}$ , although the problem of estimating  $P^* \in \mathcal{M}$  in this setting is quite easy, even in the Bayesian paradigm. The assumption appears even more restrictive when the probability  $P^*$  does not belong to  $\mathcal{M}$ , that is when the model is misspecified. For example, if the distributions in  $\mathcal{M}$  are all equivalent and  $R$  is singular with respect to  $\bar{P} \in \mathcal{M}$ ,  $\mathcal{K}(P^*, \varepsilon)$  is empty for  $P^* = (1 - 10^{-10})\bar{P} + 10^{-10}R$  although  $P^*$  and  $\bar{P} \in \mathcal{M}$  are statistically indistinguishable from any  $n$ -sample of realistic size.

Unfortunately, it is in general impossible to get rid of the restrictive conditions we have mentioned above. It is well known that the Bayes posterior distribution can be unstable in case of a misspecification of the model. Examples that illustrate this weakness have been given in Jiang and Tanner [21] and Baraud and Birgé [6] for instance. This instability is due to the fact that the Bayes posterior distribution is

based on the log-likelihood function and similar issues are known for the maximum likelihood estimator.

In order to obtain the concentration and stability properties we look for, we replace the log-likelihood function by a more stable one. Substituting another function to the log-likelihood one is not new in the literature and leads to what is called *quasi-posterior distributions*. The resulting estimators, called *quasi-Bayesian estimators* or *Laplace type estimators*, have been studied by various statisticians among which Chernozhukov and Hong [18] and Bissiri et al. [16] (we also refer to the references therein). These papers, however, do not address the problem of misspecification. In contrast, it is addressed in Jiang and Tanner [21] for performing variable selection in the logistic model. The authors show that the classical Bayesian approach is no longer reliable when the model is slightly misspecified while their Gibbs posterior distribution performs well and offers thus a much safer alternative. The problem of estimating a high-dimensional parameter  $\theta \in \mathbb{R}^d$  under a sparsity condition was considered in Atchadé [2]. His quasi-posterior distribution is obtained by replacing the joint density of the data by a more suitable one and by using some specific prior that forces sparsity. He proves that the so-defined posterior distribution contracts around the true parameter  $\theta^*$  at rate  $\sqrt{(s^* \log d)/n}$  (where  $s^*$  is the number of nonzero coordinates of  $\theta^*$ ) when both  $d$  and  $n$  tend to infinity. A common feature of the papers we have cited above lies in their asymptotic nature. This is not the case for Bhattacharya et al. [8] who replaced the likelihood function in the expression of the posterior distribution by the *fractional likelihood*, that is a suitable power of the likelihood function. The authors also consider the situation where the model is possibly misspecified but their result involves the  $\alpha$ -divergence which, as the Kullback one, can be infinite even when the true distribution of the data is close to the model for the total variation distance or the Hellinger one.

Baraud and Birgé [6] propose a surrogate to the Bayes posterior distribution that is called the  $\rho$ -posterior distribution in reference to the theory of  $\rho$ -estimation that was developed in Baraud et al. [7] and Baraud and Birgé [5]. In the frequentist paradigm, this theory aimed at solving the various problems connected to the instability of the maximum likelihood method. The  $\rho$ -posterior distribution preserves some of the nice features of the classical Bayes one but also possesses the robustness property we are interested in. The authors show that their posterior distribution concentrates on a Hellinger ball around  $P^*$  as soon as the prior puts enough mass around a point which is close enough to  $P^*$ . However their approach applies to specific dominated models  $\mathcal{M} = \{P = p \cdot \mu, p \in \mathcal{M}\}$  only. They assume that the family  $\mathcal{M}$  of densities that defines their model possesses some special combinatorial structure which is either met when  $\mathcal{M}$  is finite or when it satisfies some VC-type condition (see their Section 5). As a consequence, the concentration radius they obtain not only depends on the choice of the prior but also on a complexity term that is linked to this structure. Unlike theirs, our approach makes no such assumptions on  $\mathcal{M}$  and we are therefore able to get rid of this unpleasant complexity term while retaining a similar dependency with respect to the choice of the prior. Baraud and Birgé's posterior distribution has also the drawback to involve the supremum over the family  $\mathcal{M}$  of an empirical process. Their posterior distribution is therefore difficult to calculate in practice, unless  $\mathcal{M}$  is finite with a reasonable size. From a more theoretical point of view, it also raises some

unpleasant issues with regard to the measurability of this supremum in the situation where the family  $\mathcal{M}$  is uncountable, which is the typical case. Finally, Baraud and Birgé's approach restricts to the squared Hellinger loss while ours applies to many others.

Closer to our approach are the aggregation methods and PAC-Bayesian techniques that have been popularized by Olivier Catoni in statistical learning (see Catoni [17]). This approach has mainly been applied for the purpose of empirical risk minimization and statistical learning (see for example Alquier [1]). Our aim is to extend these techniques toward a versatile tool that can solve our Bayes-like estimation problem for various loss functions simultaneously.

The problem of designing a good estimator of  $P^*$  for a given loss function  $\ell$  was tackled in the frequentist paradigm in Baraud [4]. There, the author provides a general framework that enables one to deal with various loss functions of interest, among which the total variation, 1-Wasserstein, Hellinger, and  $\mathbb{L}_j$ -losses among others. His approach relies on the construction of a suitable family of robust tests and lies in the line of the former work of Le Cam [22], Birgé [9] and Birgé [11]. The aim of the present paper is to transpose this theory from the frequentist to the Bayesian paradigm. If  $\ell$  is the Kullback–Leibler divergence, our construction recovers the classical Bayes posterior distribution even though this is not the choice we would recommend for the reasons we have explained before.

Quite surprisingly, the concentration properties that we establish here require almost no assumption on  $\mathcal{M}$  and the distribution of the data (apart from independence). They mostly depend on the choices of the prior  $\pi$  and the loss function  $\ell$ . For a suitable element  $P$  which belongs to the model  $\mathcal{M}$  and lies close enough to  $P^*$ , these concentration properties depend on the minimal value of the radius  $r$  over which the log-ratio  $V(P, r) = \log[\pi(\mathcal{B}(P, 2r))/\pi(\mathcal{B}(P, r))]$  (with  $\mathcal{B}$  defined in (1)) increases at most linearly with  $r$ . This log-ratio was introduced in Birgé [12] for the purpose of analyzing the behaviour of the classical Bayes posterior distribution. In our Bayes-like paradigm, we show that the behaviour of the quantities  $V(P, r)$  for  $P \in \mathcal{M}$  and  $r > 0$  completely encapsulates the complexity of the model  $\mathcal{M}$ . We prove that our posterior distribution  $\hat{\pi}_X$  concentrates on an  $\ell$ -ball centered at  $P^*$  and the radius  $r = r(n)$  of which is usually of minimax order as  $n$  tends to infinity when the model is well-specified. From a nonasymptotic point of view, we show that  $\hat{\pi}_X$  retains its nice concentration properties as long as  $P^*$  remains close enough to an element  $P$  in  $\mathcal{M}$  around which the prior puts enough mass, that is, even in the situation where the model is slightly misspecified. Actually, we establish the stronger result that even when the data are only independent but not i.i.d., the above conclusion remains true for the average  $\bar{P}^*$  of their marginal distributions in place of  $P^*$ . We therefore show that the posterior distribution  $\hat{\pi}_X$  enjoys some robustness properties with respect to the equidistribution assumption we started from. The main theorems involve as much as possible explicit numerical constants. We illustrate our results with examples which are deliberately chosen to be as general and simple as possible. Our aim is to give a flavour of the results that can be established with our Bayes-like posterior, avoiding as much as possible the technicalities that would result from the choice of *ad-hoc* priors introduced to solve specific problems. Instead, we wish to discuss the optimality and robustness properties of our construction for solving general parametric and nonparametric esti-

mation problems in the density framework under assumptions that we wish to be as weak as possible. These posterior distributions will therefore provide a benchmark for comparison with other methods. Their practical implementation will be the subject of future work.

Of special interest is the choice of  $\ell$  given by the total variation distance or the Hellinger one. As we shall see, for such losses the stability of our posterior distribution automatically leads to estimators  $\widehat{P} \sim \widehat{\pi}_X$  that are naturally robust to the presence of outliers or contaminating data among the sample. These results contrast sharply with the instability of the classical Bayes posterior distribution we underlined earlier. Nevertheless, our posterior distribution also shares some similarities with the classical Bayes one. When the model is well-specified and one uses the squared Hellinger loss, we show that the credible regions of our posterior distribution asymptotically possess the same ellipsoidal shapes and approximately the same sizes as the ones we derive from the classical Bayes posterior by means of the Bernstein–von Mises theorem. Establishing an analogue of this theorem for our Bayes-like posterior distribution is, however, beyond the scope of the present paper.

Our paper is organized as follows. We present our statistical setting in Sect. 2. We consider there independent but not necessarily i.i.d. data in order to analyse later on the behaviour of our posterior distribution with respect to a possible departure from equidistribution. The construction of the posterior distribution is described in Sect. 3. In this section, we also show how more classical constructions based on the likelihood or the fractional likelihoods are particular cases of ours. We complete this section with some heuristics which, we hope, help understanding the main ideas of our approach. In particular, we bridge there the problem of designing robust posterior distributions to that of testing between two disjoint  $\ell$ -balls. Section 4 is devoted to the main theorems. We describe there the concentration properties of our posterior distribution. The applications of these results to classical loss functions are presented in Sect. 5. We put a special emphasis on the cases of the total variation distance and the squared Hellinger loss. In the remaining part of the paper, we only focus on these two losses. In Sect. 6 we highlight some similarities and differences between the classical Bayes posterior and ours for the squared Hellinger loss. In Sect. 7 we explain how our posterior distribution can be used to solve the problem of estimating a density, or a parameter associated with it, in several statistical frameworks of interest. We discuss there how the concentration properties of our posterior distribution deteriorate in the case of a misspecification of the model by the prior. We also consider the problems of estimating a density in a location-scale family and a high-dimensional parameter in a parametric model under a sparsity constraint. We also show how our estimation strategy leads to unusual rates of convergence for estimating a translation parameter in a non-regular statistical model. In Sect. 8, we provide an evaluation of the concentration radius of our posterior distributions in the parametric framework. Finally, Sect. 9 is devoted to the proofs of the main theorems and Sect. 10 to the other proofs.

## 2 The statistical setting

Let  $X = (X_1, \dots, X_n)$  be an  $n$ -tuple of independent random variables with values in a measurable space  $(E, \mathcal{E})$  and joint distribution  $\mathbf{P}^* = \otimes_{i=1}^n P_i^*$ . Even though this might not be true, we pretend that the  $X_i$  are i.i.d. and our aim is to estimate their (presumed) common distribution  $P^*$  from the observation of  $X$ . To do so, we introduce a family  $\mathcal{M}$  that consists of candidate probabilities (or merely finite signed measures in the case of the  $\mathbb{L}_j$ -loss). The reason for considering finite signed measures lies in the fact that statisticians sometimes estimate probability densities by integrable functions that are not necessarily densities but elements of a suitable linear space for instance (think of the case of projection estimators). We endow  $\mathcal{M}$  with a  $\sigma$ -algebra  $\mathcal{A}$  and a probability measure  $\pi$ , that we call a *prior* by analogy to the classical Bayesian framework, and we refer to the resulting pair  $(\mathcal{M}, \pi)$  as our *model*. The model  $(\mathcal{M}, \pi)$  plays here a similar role as in the classical Bayes paradigm. It encapsulates the *a priori* information that the statistician has on  $P^*$ . Nevertheless, we do not assume that  $P^*$ , if it ever exists, belongs to  $\mathcal{M}$  nor that the true marginals  $P_i^*$  do. We rather assume that the model  $(\mathcal{M}, \pi)$  is approximately correct in the sense that the average distribution

$$\bar{P}^* = \frac{1}{n} \sum_{i=1}^n P_i^*$$

is close enough to some point  $P$  in  $\mathcal{M}$  around which the prior  $\pi$  puts enough mass. We assume that  $\bar{P}^*$  belongs to a given set  $\mathcal{P}$  of probability measures on  $(E, \mathcal{E})$  and we measure the estimation accuracy by means of a loss function  $\ell : (\mathcal{M} \cup \mathcal{P}) \times \mathcal{M} \rightarrow \mathbb{R}_+$  which is not identical to 0 in order to avoid trivialities. Even though  $\ell$  may not be a genuine distance in general, we assume that it shares some similar features and we interpret it as if it were. For this reason, we call  $\ell$ -ball (or *ball* for short) centered at  $P \in \mathcal{P} \cup \mathcal{M}$  with radius  $r > 0$  the subset of  $\mathcal{M}$

$$\mathcal{B}(P, r) = \{Q \in \mathcal{M}, \ell(P, Q) \leq r\}.$$

Our aim is to built a *posterior distribution* (or posterior for short)  $\hat{\pi}_X$  on  $(\mathcal{M}, \mathcal{A})$ , depending on our observation  $X$ , which concentrates with a probability close to 1 on an  $\ell$ -ball of the form  $\mathcal{B}(\bar{P}^*, r_n)$  where we wish the value of  $r_n > 0$  to be small.

### 2.1 The special case of parametrized models

In many situations we consider statistical models  $\mathcal{M} = \{P_\theta, \theta \in \Theta\}$  which are parametrized via a one-to-one mapping  $\theta \mapsto P_\theta$ . When  $(\Theta, \mathfrak{B}, \nu)$  is a measurable space, we endow  $\mathcal{M}$  with the  $\sigma$ -algebra  $\mathcal{A} = \{A, \{\theta \in \Theta, P_\theta \in A\} \in \mathfrak{B}\}$ . This choice possesses several advantages. First, the mapping  $\theta \mapsto P_\theta$  is measurable from  $(\Theta, \mathfrak{B})$  onto  $(\mathcal{M}, \mathcal{A})$  and we may therefore define the prior  $\pi$  on  $(\mathcal{M}, \mathcal{A})$  as the image of  $\nu$  by this mapping. Besides, a function  $F$  is measurable on  $(\mathcal{M}, \mathcal{A})$  if and only if the mapping  $\theta \mapsto F \circ P_\theta$  is measurable on  $(\Theta, \mathfrak{B})$ . This property makes the measurability of  $F$  easier to check in general. In particular, the mapping  $F : P_\theta \mapsto \theta$  is measurable on  $(\mathcal{M}, \mathcal{A})$  because  $\theta \mapsto F \circ P_\theta = \theta$  is measurable on  $(\Theta, \mathfrak{B})$  and we may then define a posterior  $\hat{\nu}_X$  on  $(\Theta, \mathfrak{B})$  as the image by  $F$  of our posterior  $\hat{\pi}_X$  on

$(\mathcal{M}, \mathcal{A})$ . By definition of  $\widehat{\nu}_X$ , for all  $\theta \in \Theta$  and  $r > 0$

$$\widehat{\pi}_X(\mathcal{B}(P_\theta, r)) = \widehat{\nu}_X(\{\theta' \in \Theta, \ell(\theta, \theta') \leq r\}) \tag{2}$$

where  $\ell(\theta, \theta')$  denotes, slightly abusively,  $\ell(P_\theta, P_{\theta'})$  for  $\theta, \theta' \in \Theta$ . The concentration of  $\widehat{\pi}_X$  on an  $\ell$ -ball centered at  $P_\theta$  with radius  $r > 0$  is then equivalent to the concentration of  $\widehat{\nu}_X$  on the set  $\{\theta' \in \Theta, \ell(\theta, \theta') \leq r\}$ . Every time we consider a parametrized model, we assume that it is identifiable and implicitly use the construction that we presented above as well as its consequences.

### 2.2 Notation and conventions

Throughout this paper, we use the following notation and conventions. For  $a, b \in \mathbb{R}$ ,  $a \vee b$  and  $a \wedge b$  denote  $\min\{a, b\}$  and  $\max\{a, b\}$  respectively. For  $x \in \mathbb{R}$ ,  $(x)_+ = x \vee 0$  while  $(x)_- = (-x) \vee 0$ . The Euclidean spaces  $\mathbb{R}^k$  with  $k \geq 1$  are equipped with their Borel  $\sigma$ -algebras. The cardinality of a set  $A$  is denoted  $|A|$  and its complement  ${}^cA$ . In particular, for  $P \in \mathcal{P} \cup \mathcal{M}$  and  $r > 0$ ,  ${}^c\mathcal{B}(P, r) = \{Q \in \mathcal{M}, \ell(P, Q) > r\}$ . The elements of  $\mathbb{R}^k$  with  $k > 1$  are denoted with bold letters, e.g.  $\mathbf{x} = (x_1, \dots, x_k)$  and  $\mathbf{0} = (0, \dots, 0)$ . For  $\mathbf{x} \in \mathbb{R}^k$ ,  $|\mathbf{x}|_\infty = \max_{i \in \{1, \dots, k\}} |x_i|$  while  $|\mathbf{x}|$  denotes the Euclidean norm of  $\mathbf{x}$ . The inner product of  $\mathbb{R}^k$  is denoted by  $\langle \cdot, \cdot \rangle$  and the closed Euclidean ball centered at  $\mathbf{x}$  with radius  $r \geq 0$  by  $\mathcal{B}(\mathbf{x}, r)$ . By convention  $\inf_\emptyset = +\infty$  unless otherwise specified. We write  $f \equiv c$  when a function  $f$  is constant and equals  $c$  on its domain. For all suitable functions  $f$  on  $(E^n, \mathcal{E}^{\otimes n})$ ,  $\mathbb{E}[f(X)]$  means  $\int_{E^n} f d\mathbf{P}^*$  while for  $f$  on  $(E, \mathcal{E})$ ,  $\mathbb{E}_S[f(X)]$  denotes the integral  $\int_E f dS$  with respect to the measure  $S$  on  $(E, \mathcal{E})$ . For  $j \in [1, +\infty)$ , we denote by  $\mathcal{L}_j(E, \mathcal{E}, \mu)$ , the set of measurable functions  $f$  on  $(E, \mathcal{E})$  such that  $\|f\|_{j, \mu} = [\int_E |f|^j d\mu]^{1/j} < +\infty$  while  $\|f\|_\infty = \sup_{x \in E} |f(x)|$  is the supremum norm of a function  $f$  on  $E$ . If  $\pi'$  is a distribution on  $(\mathcal{M}, \mathcal{A})$ ,  $Q \sim \pi'$  means that  $Q$  is a random variable with distribution  $\pi'$ . Finally, all the measures that we consider are implicitly assumed to be  $\sigma$ -finite.

## 3 Construction of the posterior distribution

Throughout this section, the model  $(\mathcal{M}, \pi)$  is assumed to be fixed.

### 3.1 The properties of our loss functions

The construction of the posterior not only depends on the prior  $\pi$  but also on the choice of the loss function. We first assume that  $\ell$  satisfies some basic properties which are described below.

**Assumption 1** For all  $S \in \mathcal{P} \cup \mathcal{M}$ , the mapping

$$\ell(S, \cdot) : \begin{cases} (\mathcal{M}, \mathcal{A}) \longrightarrow \mathbb{R}_+ \\ P \longmapsto \ell(S, P) \end{cases}$$

is measurable.

Under such an assumption,  $\ell$ -balls are measurable and the quantities  $\pi(\mathcal{B}(P, r))$  for  $P \in \mathcal{P} \cup \mathcal{M}$  and  $r > 0$  are therefore well-defined.

**Assumption 2** There exists a positive number  $\tau$  such that, for all  $S \in \mathcal{P}$  and  $P, Q \in \mathcal{M}$ ,

$$\ell(S, Q) \leq \tau [\ell(S, P) + \ell(P, Q)] \tag{3}$$

$$\ell(S, Q) \geq \tau^{-1} \ell(P, Q) - \ell(S, P). \tag{4}$$

When  $\ell$  is a genuine distance, inequalities (3) and (4) are satisfied with  $\tau = 1$  since they correspond to the triangle inequality. When  $\ell$  is the square of a distance, these inequalities are satisfied with  $\tau = 2$ .

Importantly, we assume that  $\ell$  is associated with a family  $\mathcal{T}(\ell, \mathcal{M}) = \{t_{(P,Q)}, (P, Q) \in \mathcal{M}^2\}$  of test statistics on  $(E, \mathcal{E})$  which possesses the properties below. We shall see in Sect. 5 that many classical loss functions (among which the total variation distance, the squared Hellinger distance, etc.) can be associated with families  $\mathcal{T}(\ell, \mathcal{M})$  satisfying the following assumptions.

**Assumption 3** The elements  $t_{(P,Q)}$  of  $\mathcal{T}(\ell, \mathcal{M})$  satisfy:

(i) The mapping

$$t : \left( \begin{array}{l} (E \times \mathcal{M} \times \mathcal{M}, \mathcal{E} \otimes \mathcal{A} \otimes \mathcal{A}) \longrightarrow \mathbb{R} \\ (x, P, Q) \longmapsto t_{(P,Q)}(x) \end{array} \right)$$

is measurable.

(ii) For all  $P, Q \in \mathcal{M}$ ,  $t_{(P,Q)} = -t_{(Q,P)}$ .

(iii) there exist positive numbers  $a_0, a_1$  such that, for all  $S \in \mathcal{P}$  and  $P, Q \in \mathcal{M}$ ,

$$\mathbb{E}_S [t_{(P,Q)}(X)] \leq a_0 \ell(S, P) - a_1 \ell(S, Q). \tag{5}$$

(iv) For all  $P, Q \in \mathcal{M}$ ,

$$\sup_{x \in E} t_{(P,Q)}(x) - \inf_{x \in E} t_{(P,Q)}(x) \leq 1.$$

Under assumption (ii),  $t_{(P,P)} = 0$  and we deduce from (5) that  $(a_0 - a_1)\ell(S, P) \geq 0$ , hence that  $a_0 \geq a_1$  since  $\ell$  is not constantly equal to 0.

Some families  $\mathcal{T}(\ell, \mathcal{M})$  may satisfy the stronger

**Assumption 4** Additionally to Assumption 3, there exists  $a_2 > 0$  such that

(iv) for all  $S \in \mathcal{P}$  and  $P, Q \in \mathcal{M}$ ,

$$\text{Var}_S [t_{(P,Q)}(X)] \leq a_2 [\ell(S, P) + \ell(S, Q)].$$



### 3.2 Construction of the posterior

Let  $\mathcal{T}(\ell, \mathcal{M})$  be a family of test statistics that satisfies our Assumption 3 and let  $\beta$  and  $\lambda$  be two positive numbers such that

$$\lambda = (1 + c)\beta \quad \text{with } c > 0 \text{ satisfying } c_0 = (1 + c) - c(a_0/a_1) > 0. \quad (6)$$

We set

$$\mathbf{T}(X, P, Q) = \sum_{i=1}^n t_{(P,Q)}(X_i) \quad \text{for all } P, Q \in \mathcal{M}$$

and define  $\tilde{\pi}_X(\cdot|P)$  as the probability on  $(\mathcal{M}, \mathcal{A})$  with density

$$\frac{d\tilde{\pi}_X(\cdot|P)}{d\pi} : Q \mapsto \frac{\exp[\lambda \mathbf{T}(X, P, Q)]}{\int_{\mathcal{M}} \exp[\lambda \mathbf{T}(X, P, Q)] d\pi(Q)}.$$

Then, for  $P \in \mathcal{M}$  we set

$$\begin{aligned} \mathbf{T}(X, P) &= \int_{\mathcal{M}} \mathbf{T}(X, P, Q) d\tilde{\pi}_X(Q|P) \\ &= \int_{\mathcal{M}} \mathbf{T}(X, P, Q) \frac{\exp[\lambda \mathbf{T}(X, P, Q)]}{\int_{\mathcal{M}} \exp[\lambda \mathbf{T}(X, P, Q)] d\pi(Q)} d\pi(Q). \end{aligned}$$

Finally, we define  $\hat{\pi}_X$  as the posterior distribution on  $(\mathcal{M}, \mathcal{A})$  with density

$$\frac{d\hat{\pi}_X}{d\pi} : P \mapsto \frac{\exp[-\beta \mathbf{T}(X, P)]}{\int_{\mathcal{M}} \exp[-\beta \mathbf{T}(X, P)] d\pi(P)}. \quad (7)$$

Our Assumption 3-(i) ensures that  $d\tilde{\pi}_X(\cdot|P)/d\pi$  is a measurable function of  $(X, P, Q)$  and  $d\hat{\pi}_X/d\pi$  a measurable function of  $(X, P)$ .

The posterior distribution depends on our choice of  $\beta$  and  $\lambda$  (or equivalently  $c$ ) even though we drop this dependency with the notation  $\hat{\pi}_X$ .

### 3.3 Monte Carlo computation of functions of the posterior

Even though we focus on the concentration properties of the posterior  $\hat{\pi}_X$ , one may alternatively be interested in some estimators derived from it. For example, estimators of the form

$$I = \int_{\mathcal{M}} F(P) d\hat{\pi}_X(P)$$

where  $F$  is a real-valued  $\pi$ -integrable function on  $(\mathcal{M}, \mathcal{A})$ . For typical choices of  $F$ ,  $I$  gives the expected mean, mode or median of the posterior whenever these quantities

make sense. One may also choose  $F : P \mapsto \mathbb{1}_{P \in \mathcal{B}(P_0, \varepsilon)}$  with  $P_0 \in \mathcal{M}$  and  $\varepsilon > 0$  in order to compute the (posterior) probability that  $\ell(P_0, \widehat{P})$  is not larger than  $\varepsilon$  when  $\widehat{P} \sim \widehat{\pi}_X$ .

Interestingly, the integral  $I$  can be approximated by Monte Carlo as follows. Assume that the prior  $\pi$  admits a density of the form  $C^{-1}\Pi$  with respect to a given probability measure  $m$ , where  $\Pi$  is a nonnegative  $m$ -integrable function on  $(\mathcal{M}, \mathcal{A})$  and  $C = \int_{\mathcal{M}} \Pi(P)dm(P) > 0$  a positive normalizing constant (that will not be involved in our calculation). Let  $P_1, \dots, P_N$  be an  $N$ -sample with distribution  $m$  and for each  $i \in \{1, \dots, N\}$ ,  $Q_i^{(1)}, \dots, Q_i^{(N')}$  an independent  $N'$ -sample with the same distribution. We may approximate  $I$  by

$$\widehat{I}_{N, N'} = \sum_{i=1}^N F(P_i) \frac{\exp[-\beta W_{i, N'}(P_i)] \Pi(P_i)}{\sum_{i'=1}^{N'} \exp[-\beta W_{i', N'}(P_{i'})] \Pi(P_{i'})}$$

where for all  $i \in \{1, \dots, N\}$ ,

$$W_{i, N'}(P_i) = \sum_{j=1}^{N'} T(X, P_i, Q_i^{(j)}) \frac{\exp[\lambda T(X, P_i, Q_i^{(j)})] \Pi(Q_i^{(j)})}{\sum_{j'=1}^{N'} \exp[\lambda T(X, P_i, Q_i^{(j')})] \Pi(Q_i^{(j')})}.$$

It is then easy to check that, by the law of large numbers,

$$\lim_{N \rightarrow +\infty} \left[ \lim_{N' \rightarrow +\infty} \widehat{I}_{N, N'} \right] = I.$$

### 3.4 Connection with the classical Bayes posterior distribution

The classical Bayes posterior turns out to be a particular case of the posterior-type ones introduced in Sect. 3.2. As we shall see now, they are associated with the Kullback–Leibler divergence loss. We recall that the Kullback–Leibler divergence  $\ell(P, Q) = K(P, Q)$  between two probabilities  $P, Q$  on  $(E, \mathcal{E})$  is defined by

$$K(P, Q) = \begin{cases} \int_E \log \left( \frac{dP}{dQ} \right) dP & \text{when } P \ll Q; \\ +\infty & \text{otherwise.} \end{cases}$$

Let us consider now a family  $\mathcal{M}$  of probabilities that satisfy for some  $a > 0$  and suitable versions of their densities  $dQ/dP$  the following inequalities:

$$e^{-a} \leq \frac{dP}{dQ}(x) \leq e^a \quad \text{for all } x \in E \text{ and } P, Q \in \mathcal{M}. \tag{8}$$

It follows from Baraud [4, Proposition 12] that the families of functions

$$\mathcal{T}(\ell, \mathcal{M}) = \left\{ t_{(P,Q)} = \frac{1}{2a} \log \left( \frac{dQ}{dP} \right), P, Q \in \mathcal{M} \right\} \tag{9}$$

satisfies our Assumptions 3 and 4 with  $a_0 = a_1 = 1/(2a)$  and  $a_2 = 2a/[\tanh(a/2)]$ . Note that given  $P, Q \in \mathcal{M}, P \neq Q$ , the test based on the sign of  $t_{(P,Q)}$  is the classical likelihood ratio test between  $P$  and  $Q$ .

If we apply the construction described in Sect. 3.2 to the family  $\mathcal{T}(\ell, \mathcal{M})$  we obtain that for all  $P, Q, P_0 \in \mathcal{M}$ ,

$$\mathbf{T}(X, P, Q) = \mathbf{T}(X, P_0, Q) - \mathbf{T}(X, P_0, P).$$

For all  $\lambda > 0$ , the density of  $\tilde{\pi}_X(\cdot|P)$

$$Q \mapsto \frac{\exp[\lambda \mathbf{T}(X, P, Q)]}{\int_{\mathcal{M}} \exp[\lambda \mathbf{T}(X, P, Q)] d\pi(Q)} = \frac{\exp[\lambda \mathbf{T}(X, P_0, Q)]}{\int_{\mathcal{M}} \exp[\lambda \mathbf{T}(X, P_0, Q)] d\pi(Q)}$$

is independent of  $P$  and writing  $\tilde{\pi}_X(\cdot)$  in place of  $\tilde{\pi}_X(\cdot|P)$  we obtain that

$$\begin{aligned} \mathbf{T}(X, P) &= \int_{\mathcal{M}} \mathbf{T}(X, P, Q) d\tilde{\pi}_X(Q) \\ &= \int_{\mathcal{M}} \mathbf{T}(X, P_0, Q) d\tilde{\pi}_X(Q) - \mathbf{T}(X, P_0, P) \\ &= C - \frac{1}{2a} \sum_{i=1}^n \log \left( \frac{dP}{dP_0} \right) (X_i) \text{ with } C = \int_{\mathcal{M}} \mathbf{T}(X, P_0, Q) d\tilde{\pi}_X(Q). \end{aligned}$$

Finally, the density of our posterior  $\hat{\pi}_X$  at  $P \in \mathcal{M}$  is given by

$$\frac{d\hat{\pi}_X}{d\pi}(P) = \frac{\exp[-\beta \mathbf{T}(X, P)]}{\int_{\mathcal{M}} \exp[-\beta \mathbf{T}(X, P)] d\pi(P)} = \frac{[\prod_{i=1}^n (dP/dP_0)(X_i)]^{\beta/(2a)}}{\int_{\mathcal{M}} [\prod_{i=1}^n (dP/dP_0)(X_i)]^{\beta/(2a)} d\pi(P)}.$$

This is the density of the classical Bayes posterior when  $\beta = 2a$  while for other values of  $\beta$  it is that of fractional Bayes ones.

Nevertheless, in the present paper we restrict our study to loss functions that satisfy some triangle-type inequality – see Assumption 2. This excludes the Kullback–Leibler divergence unless one is ready to make strong assumptions on the unknown distribution of the data, which we do not want to do here.

### 3.5 Some heuristics

In this section, we present the basic ideas that underline our approach. In particular, we shall see how the estimation problem we want to solve is linked to the one of testing between two disjoint  $\ell$ -balls  $\mathcal{B}(P, r)$  and  $\mathcal{B}(Q, r)$  with  $P, Q \in \mathcal{M}$ .

In order to avoid unnecessary details, we assume here that we observe i.i.d. data  $X_1, \dots, X_n$  with distribution  $P^* \in \mathcal{P}$  and that we have at disposal a family  $\mathcal{T}(\ell, \mathcal{M})$  of functions that satisfies our Assumption 3. In particular it follows from Assumption 3-(iii) that

$$\mathbb{E} \left[ \frac{\mathbf{T}(X, P, Q)}{n} \right] = \frac{1}{n} \sum_{i=1}^n \mathbb{E} [t_{(P, Q)}(X_i)] \leq a_0 \ell(P^*, P) - a_1 \ell(P^*, Q).$$

The antisymmetric property required by Assumption 3-(ii) entails that

$$\mathbf{T}(X, P, Q) = -\mathbf{T}(X, Q, P)$$

and leads to the lower bound

$$\mathbb{E} \left[ \frac{\mathbf{T}(X, P, Q)}{n} \right] \geq a_1 \ell(P^*, P) - a_0 \ell(P^*, Q).$$

Assuming for the sake of simplicity that  $a_0 = a_1 = 1$ , these calculations show that  $n^{-1}\mathbf{T}(X, P, Q) = n^{-1} \sum_{i=1}^n t_{(P, Q)}(X_i)$  is an unbiased and consistent estimator of  $\ell(P^*, P) - \ell(P^*, Q)$ . In particular, if the two  $\ell$ -balls  $\mathcal{B}(P, r), \mathcal{B}(Q, r)$  are disjoint and  $P^*$  belongs to one of them, the sign of  $n^{-1}\mathbf{T}(X, P, Q) = n^{-1} \sum_{i=1}^n t_{(P, Q)}(X_i)$  provides a consistent test for deciding which one contains  $P^*$ . In fact, the test does not depend on the value of  $r$  and consequently chooses the element among  $\{P, Q\}$  which is the closest to  $P^*$  (with respect to  $\ell$ ), at least when  $n$  is large enough. As compared to the classical likelihood ratio test between  $P$  and  $Q$ , this test has the advantage not to assume that  $P^*$  is either  $P$  or  $Q$  but only that it lies in a small enough  $\ell$ -vicinity around one of these two probabilities. The test is said to be *robust* with respect to the model  $\{P, Q\}$ . Its nonasymptotic properties have been studied in Baraud [4].

Let us now explain how such families  $\{\mathbf{T}(X, P, Q), (P, Q) \in \mathcal{M}^2\}$  of test statistics can be used to build robust estimators and not only tests. In the frequentist paradigm, the construction of  $\ell$ -estimators is based on the following heuristics. If, with a probability close to 1,  $n^{-1}\mathbf{T}(X, P, Q)$  is close to its expectation  $\ell(P^*, P) - \ell(P^*, Q)$  uniformly with respect to  $(P, Q) \in \mathcal{M}^2$  then  $n^{-1}\mathbf{T}'(X, P) = \sup_{Q \in \mathcal{M}} [n^{-1}\mathbf{T}(X, P, Q)]$  is close to

$$\sup_{Q \in \mathcal{M}} [\ell(P^*, P) - \ell(P^*, Q)] = \ell(P^*, P) - \inf_{Q \in \mathcal{M}} \ell(P^*, Q).$$

We therefore expect that a minimizer over  $\mathcal{M}$  of the function  $P \in \mathcal{M} \mapsto n^{-1}\mathbf{T}'(X, P)$  be close to a minimizer over  $\mathcal{M}$  of the function  $P \in \mathcal{M} \mapsto \ell(P^*, P) - \inf_{Q \in \mathcal{M}} \ell(P^*, Q)$ , that is an element that minimizes the loss  $\ell(P^*, P)$  among the probabilities  $P \in \mathcal{M}$ .

In the Bayesian paradigm, we may argue in a similar way as follows. Replacing  $n^{-1}\mathbf{T}(X, P, Q)$  by its expectation  $\ell(P^*, P) - \ell(P^*, Q)$ , as we did before, amounts

to replacing  $\mathbf{T}(X, \mathbf{P})$  by

$$\begin{aligned} &\bar{\mathbf{T}}(X, \mathbf{P}) \\ &= n \int_{\mathcal{M}} (\ell(P^*, P) - \ell(P^*, Q)) \frac{\exp[n\lambda (\ell(P^*, P) - \ell(P^*, Q))] d\pi(Q)}{\int_{\mathcal{M}} \exp[n\lambda (\ell(P^*, P) - \ell(P^*, Q))] d\pi(Q)} \\ &= n\ell(P^*, P) - n \int_{\mathcal{M}} \ell(P^*, Q) \frac{\exp[-n\lambda\ell(P^*, Q)]}{\int_{\mathcal{M}} \exp[-n\lambda\ell(P^*, Q)] d\pi(Q)} d\pi(Q). \end{aligned}$$

Note that the second term in the right-hand side does not depend on  $P$ . Consequently, replacing  $\mathbf{T}(X, \mathbf{P})$  by  $\bar{\mathbf{T}}(X, \mathbf{P})$  in the expression (7) of the density of  $\hat{\pi}_X$  leads to the density

$$P \mapsto \frac{\exp[-\beta\bar{\mathbf{T}}(X, P)]}{\int_{\mathcal{M}} \exp[-\beta\bar{\mathbf{T}}(X, P)] d\pi(P)} = \frac{\exp[-n\beta\ell(P^*, P)]}{\int_{\mathcal{M}} \exp[-n\beta\ell(P^*, P)] d\pi(P)}.$$

We recognize here the density of a Gibbs measure associated with the energy  $\ell(P^*, P)$  at point  $P \in \mathcal{M}$  and inverse temperature  $n\beta > 0$ . We know that when the temperature goes to 0 (or equivalently  $n\beta$  to infinity), Gibbs measures concentrate their masses in vicinities of low energy points in  $\mathcal{M}$ . In our case, these low energy points are those for which  $\ell(P^*, P)$  is minimal.

Similar ideas can be found in Catoni’s work and more specifically in his construction of Gibbs estimators—see Catoni [17, Chapter 4]. There, Catoni shows how to aggregate a continuous family of estimators in order to minimize a risk. In the present paper, we do not aim at aggregating estimators but we use similar ideas and tools that are due to Catoni and his co-authors for the construction of our robust posterior distribution.

## 4 The main results

### 4.1 Linking the prior to the complexity of the model

For  $P \in \mathcal{M}$  and  $r > 0$ , we recall that

$$V(P, r) = \log \left( \frac{\pi(\mathcal{B}(P, 2r))}{\pi(\mathcal{B}(P, r))} \right)$$

where we use the convention  $a/0 = +\infty$  for all  $a \geq 0$ . We said in the Introduction that such quantities encapsulate in some sense the complexity of the model  $(\mathcal{M}, \pi)$  and we shall now explain why. If  $\mathcal{M} = \{P_\theta, \theta \in \mathbb{R}^k\}$  is a parametric model endowed with a loss  $\ell$  such that  $\ell(\theta, \theta') = |\theta - \theta'|$ , so that  $(\mathcal{M}, \ell)$  is isometric to  $(\mathbb{R}^k, |\cdot|)$ , and if the prior  $\nu$  on  $\Theta = \mathbb{R}^k$  is improper and given by the Lebesgue measure, we obtain that for all  $P \in \mathcal{M}$  and  $r > 0$

$$V(P, r) = \log \left( \frac{\pi(\mathcal{B}(P, 2r))}{\pi(\mathcal{B}(P, r))} \right) = \log \left( \frac{(2r)^k}{r^k} \right) = k \log 2. \tag{10}$$

We observe that  $V(P, r)$  corresponds in this case to the usual dimension of  $\mathbb{R}^k$  (up the factor  $\log 2$ ). For more general models  $(\mathcal{M}, \pi)$  and loss functions  $\ell$ , we may interpret  $V(P, r)$  as some notion of dimension (or complexity) associated with the element  $P \in \mathcal{M}$  at the scale  $r > 0$ . As we do not consider improper priors but probability distributions,  $\lim_{r \rightarrow +\infty} \pi(\mathcal{B}(P, r)) = 1$  and consequently  $\lim_{r \rightarrow +\infty} V(P, r) = 0$ . This means that the connection with the notion of “dimension” is only relevant for values of  $r$  which are not too large.

Given  $\gamma \in (0, 1]$ , the set

$$\mathcal{R}(\beta, P) = \left\{ r \geq \frac{1}{n\beta a_1}, \text{ such that } \sup_{r' \geq r} \frac{V(P, r')}{r'} \leq \gamma n\beta a_1 \right\}$$

is the subinterval of  $\mathbb{R}_+$  on which the mapping  $r \mapsto V(P, r)$  is not larger than  $r \mapsto \gamma n\beta a_1 r$ . We denote by

$$r_n(\beta, P) = \inf \mathcal{R}(\beta, P) \tag{11}$$

the left endpoint of  $\mathcal{R}(\beta, P)$ . Since  $\mathcal{R}(\beta, P)$  is increasing with  $\beta$  with respect to set inclusion,  $r_n(\beta, P)$  is a nonincreasing function of  $\beta$ . For example, in the ideal situation given in (10) where  $V(P, r) \equiv k \log 2$  with  $k \log 2 \geq 1$ ,  $r_n(\beta, P) = (\gamma a_1)^{-1} [k \log 2 / (n\beta)]$ . When the model  $\mathcal{M} = \{P_\theta, \theta \in \Theta\}$  is parametric and the parameter space  $\Theta$  is an open subset of  $\mathbb{R}^k$  endowed with a prior  $\nu$ , we shall see in Sect. 8.2 that under suitable assumptions  $r_n(\beta, P_\theta)$  is indeed of order  $k/(n\beta)$ , at least for  $n$  sufficiently large.

The Bayesian paradigm offers the possibility to favour some elements of  $\mathcal{M}$  as compared to others. The order of magnitude of  $r_n(\beta, P)$  allows one to quantify how much the prior  $\pi$  advantages or disadvantages  $P \in \mathcal{M}$ . It follows from the definition of  $r_n(\beta, P)$  that

$$0 < \pi(\mathcal{B}(P, 2r)) \leq \exp(\gamma n\beta a_1 r) \pi(\mathcal{B}(P, r)) \text{ for all } r > r_n(\beta, P). \tag{12}$$

Letting  $r$  decrease to  $r_n(\beta, P)$ , we derive that (12) also holds for  $r = r_n(\beta, P)$ . In particular,  $\pi(\mathcal{B}(P, r)) > 0$  for  $r = r_n(\beta, P)$ . If the prior puts no mass on the  $\ell$ -ball  $\mathcal{B}(P, r)$ , which clearly corresponds to a situation where the prior disadvantages  $P$ ,  $r_n(\beta, P) > r$  and  $r_n(\beta, P)$  is therefore large if  $r$  is large. In the opposite case, if the prior puts enough mass on  $\mathcal{B}(P, r)$  in the sense that

$$\pi(\mathcal{B}(P, r)) \geq \exp(-\gamma n\beta a_1 r), \tag{13}$$

then for all  $r' \geq r$ ,

$$\begin{aligned} \pi(\mathcal{B}(P, r')) &\geq \exp(-\gamma n\beta a_1 r) \geq \exp(-\gamma n\beta a_1 r') \\ &\geq \exp(-\gamma n\beta a_1 r') \pi(\mathcal{B}(P, 2r')) \end{aligned}$$

hence,

$$\frac{\pi(\mathcal{B}(P, 2r'))}{\pi(\mathcal{B}(P, r'))} \leq \exp(\gamma n \beta a_1 r') \quad \text{and } r_n(\beta, P) \leq r.$$

The quantity  $r_n(\beta, P)$  is therefore small if  $r$  is small. Although (13) is not equivalent to (12) (it is actually stronger), the previous arguments provide a partial view on the relationship between  $\pi$  and  $r_n$  and conditions to decide whether  $P$  is favoured by  $\pi$  or not, according to the size of  $r_n(\beta, P)$ .

#### 4.2 A general result on the concentration property of the posterior distribution

According to the discussion of Sect. 4.1, we see that, when the set

$$\mathcal{M}(\beta) = \left\{ P \in \mathcal{M}, r_n(\beta, P) \leq a_1^{-1} \beta \right\} \tag{14}$$

is nonempty, it contains the most favoured elements of the model  $(\mathcal{M}, \pi)$  at level  $a_1^{-1} \beta$ . Since  $r_n(\beta, P)$  is nonincreasing with  $\beta$ , the set  $\mathcal{M}(\beta)$  is increasing with  $\beta$  with respect to set inclusion. If  $a_1^{-1} \beta \geq (n \beta a_1)^{-1}$  or equivalently  $\beta \geq 1/\sqrt{n}$ , the set  $\mathcal{M}(\beta)$  can alternatively be defined from  $V(P, r)$  as follows:

$$\mathcal{M}(\beta) = \left\{ P \in \mathcal{M}, V(P, r) \leq \gamma n \beta a_1 r \text{ for all } r \geq a_1^{-1} \beta \right\}. \tag{15}$$

This set plays a crucial role in our first result.

**Theorem 1** *Assume that the model  $(\mathcal{M}, \pi)$  and the loss  $\ell$  satisfy Assumptions 1 and 2 and the family  $\mathcal{T}(\ell, \mathcal{M})$  Assumption 3. Let  $\gamma < (c_0 \wedge c)/(2\tau)$  and  $\beta \geq 1/\sqrt{n}$  be chosen in such a way that the set  $\mathcal{M}(\beta)$  defined by (14) is not empty. Then, the posterior  $\widehat{\pi}_X$  defined by (7) possesses the following property. There exists  $\kappa_0 > 0$  only depending on  $c, \tau, \gamma$  and the ratio  $a_0/a_1$  such that, for all  $\xi > 0$  and any distribution  $\mathbf{P}^*$  with marginals in  $\mathcal{P}$ ,*

$$\mathbb{E} \left[ \widehat{\pi}_X \left( {}^c \mathcal{B}(\overline{P}^*, \kappa_0 r) \right) \right] \leq 2e^{-\xi} \tag{16}$$

with

$$r = \inf_{P \in \mathcal{M}(\beta)} \ell(\overline{P}^*, P) + \frac{1}{a_1} \left( \beta + \frac{2\xi}{n\beta} \right). \tag{17}$$

In particular,

$$\mathbb{P} \left[ \widehat{\pi}_X \left( {}^c \mathcal{B}(\overline{P}^*, \kappa_0 r) \right) \geq e^{-\xi/2} \right] \leq 2e^{-\xi/2}.$$

The value of  $\kappa_0$  is given by (119) in the proof. It only depends on the choice of the family  $\mathcal{T}(\ell, \mathcal{M})$  but not on the prior  $\pi$ . Hence, for a given family  $\mathcal{T}(\ell, \mathcal{M})$ ,  $\kappa_0$  is a numerical constant.

Let us now comment on Theorem 1. When  $X_1, \dots, X_n$  are truly i.i.d. with distribution  $P^*$  and the prior puts enough mass around  $P^*$ , in the sense that  $P^* \in \mathcal{M}(\beta)$ , then  $r = a_1^{-1}[\beta + 2\xi/(n\beta)]$  in (17). When this ideal situation is not met, either because the data are not identically distributed or because  $P^*$  does not belong to  $\mathcal{M}(\beta)$ ,  $r$  increases by at most an additive term of order  $\inf_{P \in \mathcal{M}(\beta)} \ell(\overline{P}^*, P)$ . When this approximation term remains small as compared to  $a_1^{-1}\beta$ , the value of  $r$  does not deteriorate too much as compared to the previous situation.

The value of  $r$  given by (17) depends on the choice of the parameter  $\beta$ . Since the set  $\mathcal{M}(\beta)$  is increasing (with respect to set inclusion) as  $\beta$  gets larger, the two terms  $\inf_{P \in \mathcal{M}(\beta)} \ell(\overline{P}^*, P)$  and  $a_1^{-1}\beta$  vary in opposite directions as  $\beta$  increases. The set  $\mathcal{M}(\beta)$  must be large enough to provide a suitable approximation of  $\overline{P}^*$  while  $\beta$  must not be too large in order to keep  $a_1^{-1}\beta$  to a reasonable size. Practically, we recommend to choose  $\beta = \beta(\alpha) \geq 1/\sqrt{n}$  such that

$$\pi(\mathcal{M}(\beta)) \geq 1 - \alpha \quad \text{for } \alpha \in (0, 1/10). \tag{18}$$

In Example 1 below and in Sect. 7.1, we give some examples of choices of  $\beta$ .

**Example 1** Let  $(\mathcal{M}, \pi)$  be a model where the prior  $\pi$  satisfies for some  $k \geq 1$  and constants  $0 < A \leq (2/e)B$ ,

$$(Ar)^k \wedge 1 \leq \pi(\mathcal{B}(P, r)) \leq (Br)^k \wedge 1 \quad \text{for all } P \in \mathcal{M} \text{ and } r > 0. \tag{19}$$

This means that the prior  $\pi$  behaves like the Lebesgue measure on an Euclidean space of dimension  $k$  for small enough values of  $r$ . Then,

$$V(P, r) = \log \frac{\pi(\mathcal{B}(P, 2r))}{\pi(\mathcal{B}(P, r))} \leq k \log \left( \frac{2B}{A} \right) \quad \text{for all } P \in \mathcal{M} \text{ and } r > 0 \tag{20}$$

which implies that for all  $P \in \mathcal{M}$

$$r_n(P, \beta) \leq \frac{k}{\gamma a_1 n \beta} \log \left( \frac{2B}{A} \right). \tag{21}$$

The right-hand side is not larger than  $a_1^{-1}\beta$  for

$$\beta = \sqrt{\frac{k \log(2B/A)}{\gamma n}} \tag{22}$$

which is larger than  $1/\sqrt{n}$  since  $(2B/A) \geq e$  and  $\gamma \in (0, 1]$ . For such a value of  $\beta$ , which does not depend on the distribution of the data, the element  $P$  belongs to  $\mathcal{M}(\beta)$  given by (15), and since  $P$  is arbitrary we derive that  $\mathcal{M}(\beta) = \mathcal{M}$ . Applying



Theorem 1 we conclude that the distribution  $\widehat{\pi}_X$  concentrates on an  $\ell$ -ball centered at  $\overline{P}^*$  with a radius  $r$  of order

$$\bar{r}_n = \inf_{P \in \mathcal{M}} \ell(\overline{P}^*, P) + \frac{1}{a_1} \left( \sqrt{\frac{k}{n}} + \frac{2\xi}{\sqrt{nk}} \right). \tag{23}$$

### 4.3 A refined result under Assumption 4

Let us assume now that the family  $\mathcal{T}(\ell, \mathcal{M})$  satisfies the stronger Assumption 4. We introduce the mapping

$$\phi : \begin{cases} (0, +\infty) \longrightarrow \mathbb{R}_+ \\ z \longmapsto \phi(z) = \frac{2(e^z - 1 - z)}{z^2}. \end{cases} \tag{24}$$

The function  $\phi$  is increasing on  $(0, +\infty)$  and tends to 1 when  $z$  tends to 0. Given  $\beta > 0$  and a family  $\mathcal{T}(\ell, \mathcal{M})$  that satisfies Assumption 4, we define

$$\bar{c}_1 = c_0 - \beta a_2 a_1^{-1} \tau^2 \phi[\beta(1 + 2c)](1 + 2c(1 + c)); \tag{25}$$

$$\bar{c}_2 = c - \beta a_2 a_1^{-1} \tau^2 \phi[\beta(1 + 2c)]c^2; \tag{26}$$

$$\bar{c}_3 = (2 + c) - \beta a_2 a_1^{-1} \tau^2 \phi[\beta(3 + 2c)](2 + c)^2. \tag{27}$$

Note that the value of  $\bar{c}_1 \wedge \bar{c}_2 \wedge \bar{c}_3$  is positive for  $\beta = 0$  and decreases continuously to  $-\infty$  when  $\beta$  grows to infinity. Consequently, there exists some  $\beta_0 > 0$  for which  $\bar{c}_1 \wedge \bar{c}_2 \wedge \bar{c}_3 = 0$  and  $\bar{c}_1 \wedge \bar{c}_2 \wedge \bar{c}_3$  is positive for all values  $\beta \in (0, \beta_0)$ .

Let us now present our second result on the concentration property of our posterior  $\widehat{\pi}_X$ .

**Theorem 2** *Assume that the model  $(\mathcal{M}, \pi)$  and the loss  $\ell$  satisfy Assumptions 1 and 2 and the family  $\mathcal{T}(\ell, \mathcal{M})$  Assumption 4. For  $\beta \in (0, \beta_0)$  and  $\gamma < (\bar{c}_1 \wedge \bar{c}_2 \wedge \bar{c}_3)/(2\tau)$ , the posterior  $\widehat{\pi}_X$  defined by (7) satisfies the following property. There exists  $\kappa_0 > 0$  only depending on  $a_0/a_1, a_2/a_1, c, \tau, \beta$  and  $\gamma$  such that, for all  $\xi > 0$  and any distribution  $\mathbf{P}^*$  with marginals in  $\mathcal{P}$ ,*

$$\mathbb{E} \left[ \widehat{\pi}_X \left( {}^c\mathcal{B}(\overline{P}^*, \kappa_0 r) \right) \right] \leq 2e^{-\xi} \tag{28}$$

with

$$r = \inf_{P \in \mathcal{M}} \left[ \ell(\overline{P}^*, P) + r_n(\beta, P) \right] + \frac{2\xi}{n\beta a_1}. \tag{29}$$

In particular,

$$\mathbb{P} \left[ \widehat{\pi}_X \left( {}^c\mathcal{B}(\overline{P}^*, \kappa_0 r) \right) \geq e^{-\xi/2} \right] \leq 2e^{-\xi/2}.$$

The value of  $\kappa_0$  is given by (132) in the proof. Note that the constraints on  $\beta$  and  $\gamma$ , that are required in our Theorem 2, and that on  $c$  given in (6) only depend on  $a_0, a_1$  and  $a_2$ , hence on the choice of the family  $\mathcal{T}(\ell, \mathcal{M})$ . When  $a_0, a_1$  and  $a_2$  do not depend on  $\mathcal{M}$ , the value of  $\beta$  can be chosen as a universal constant. In particular, it neither depends on the model  $(\mathcal{M}, \pi)$  nor on the sample size  $n$ .

**Example 2** (Example 1 continued) Let us go back to the framework of our Example 1 and assume that  $\mathcal{T}(\ell, \mathcal{M})$  satisfies the requirements of Theorem 2, hence Assumption 4. Applying our construction with some numerical value of  $\beta$  which satisfies the constraint of our Theorem 2, we deduce from (21) that  $\widehat{\pi}_X$  concentrates on an  $\ell$ -ball with radius of order

$$\bar{r} = \inf_{P \in \mathcal{M}} \ell(\bar{P}^*, P) + \frac{\log(2B/A)}{\gamma a_1 \beta} \frac{k}{n} + \frac{2}{a_1 \beta} \frac{\xi}{n}. \tag{30}$$

When the model is well-specified,  $\inf_{P \in \mathcal{M}} \ell(\bar{P}^*, P) = 0$  and the ball  $\mathcal{B}(P^*, \kappa_0 \bar{r})$  with radius  $\bar{r} = \bar{r}(n)$  contracts at the rate  $1/n$ . Applying our Theorem 1 under Assumption 3, ignoring the fact that the family  $\mathcal{T}(\ell, \mathcal{M})$  also satisfies Assumption 4, would lead to the weaker result that when the model is well-specified the posterior concentrates on an  $\ell$ -ball with radius of order  $\sqrt{k/n}$ , hence at a rate  $1/\sqrt{n}$ , as shown by (23).

### 4.4 Concentrated priors

Theorem 1 and 2 show that starting from a prior  $\pi$  that puts enough mass around most of the elements of  $\mathcal{M}$ , the posterior  $\widehat{\pi}_X$  concentrates on an  $\ell$ -ball with radius of order  $\inf_{P \in \mathcal{M}} \ell(\bar{P}^*, P) + r_n$  where  $r_n$  is small, at least under suitable assumptions and for  $n$  sufficiently large. The situation we want to investigate now is what happens when the prior is very concentrated on a small  $\ell$ -ball with radius  $\varepsilon > 0$  around an element  $\bar{Q} \in \mathcal{M}$  that might not be the true distribution of the data. More precisely, assume the following

**Assumption 5** For  $\bar{Q} \in \mathcal{M}$  and  $\varepsilon > 0$ ,

$$\pi \left( {}^c \mathcal{B}(\bar{Q}, \varepsilon) \right) \leq e^{-(2\xi+1)} \pi \left( \mathcal{B}(\bar{Q}, \varepsilon) \right) \quad \text{with } \xi > 0.$$

In this case, we establish the following result.

**Theorem 3** Assume that the model  $(\mathcal{M}, \pi)$  and the loss  $\ell$  satisfy Assumptions 1 and 2 and the family  $\mathcal{T}(\ell, \mathcal{M})$  Assumption 3. If Assumption 5 is satisfied, there exists  $\kappa_0 > 0$  only depending on  $c, \tau$  and the ratio  $a_0/a_1$  such that for any distribution  $P^*$  with marginals in  $\mathcal{P}$ ,

$$\mathbb{E} \left[ \widehat{\pi}_X \left( {}^c \mathcal{B}(\bar{P}^*, \kappa_0 r) \right) \right] \leq 2e^{-\xi} \quad \text{with } r = \ell(\bar{P}^*, \bar{Q}) \vee \frac{\beta}{a_1} \vee \varepsilon. \tag{31}$$

In particular, for the choice  $\beta = a_1 \varepsilon, r = \ell(\bar{P}^*, \bar{Q}) \vee \varepsilon$ .

If furthermore, Assumption 4 is satisfied and  $\beta \in (0, \beta_0)$  (where  $\beta_0$  is defined in Sect. 4.3), there exists  $\kappa'_0 > 0$  only depending on  $\tau, \beta, a_0/a_1$  and  $a_2/a_1$  such that for any distribution  $\mathbf{P}^*$  with marginals in  $\mathcal{P}$ ,

$$\mathbb{E} \left[ \widehat{\pi}_X \left( {}^c\mathcal{B}(\overline{\mathbf{P}}^*, \kappa'_0 r) \right) \right] \leq 2e^{-\xi} \quad \text{with } r = \ell(\overline{\mathbf{P}}^*, \overline{\mathbf{Q}}) \vee \varepsilon. \tag{32}$$

This result shows that for a suitable choice of  $\beta$ , the posterior  $\widehat{\pi}_X$  also concentrates on an  $\ell$ -ball centred at  $\overline{\mathbf{P}}^*$  with radius of order  $\varepsilon$  when the model is well-specified, that is, when the data are i.i.d. with distribution  $\overline{\mathbf{P}}^* = \overline{\mathbf{Q}}$ . When the model is misspecified, the radius of the ball is of order  $\ell(\overline{\mathbf{P}}^*, \overline{\mathbf{Q}}) \vee \varepsilon$  and therefore does not inflate more than the distance of  $\overline{\mathbf{P}}^*$  to the center  $\overline{\mathbf{Q}}$ . This result illustrates the stability of the posterior  $\widehat{\pi}_X$  with respect to misspecification.

### 5 Applications to classical loss functions

The aim of this section is to show how our general construction can be applied to loss functions  $\ell$  of interest. The propositions contained in this section about the corresponding families  $\mathcal{T}(\ell, \mathcal{M})$  have been established in Baraud [4] except for the squared Hellinger loss for which we refer to Baraud and Birgé [5, Proposition 3]. The list of loss functions we present here is not exhaustive. Our results also apply to all loss functions that derive from a variational formula of the form

$$\ell(P, Q) = \sup_{f \in \mathcal{F}} \left[ \int_E f dP - \int_E f dQ \right]$$

where  $\mathcal{F}$  is a suitable class of bounded functions. For such losses, we refer the reader to Baraud [4].

In this section, we consider models  $\mathcal{M} = \{P = p \cdot \mu, p \in \mathcal{M}\}$  which are dominated by a measure  $\mu$  on  $(E, \mathcal{E})$  and we denote by  $\mathcal{M} \subset \mathcal{L}_1(E, \mathcal{E}, \mu)$  the corresponding families of densities with respect to  $\mu$ . Elements  $P, Q, \dots$  in  $\mathcal{M}$  are associated with their densities in  $\mathcal{M}$  by using lower case letters  $p, q, \dots$ . In all the cases we consider,  $t_{(P, Q)}(x)$  is a measurable function of  $(p(x), q(x))$  for  $P, Q \in \mathcal{M}$  and  $x \in E$ . In order to satisfy our measurability Assumption 3-(i), it is therefore sufficient to assume that

$$\begin{aligned} (E \times \mathcal{M}, \mathcal{E} \otimes \mathcal{A}) &\longrightarrow \mathbb{R} \\ (x, P) &\longmapsto p(x) \end{aligned}$$

is measurable. In the case of a parametrized model  $\mathcal{M} = \{P_\theta = p_\theta \cdot \mu, \theta \in \Theta\}$ , as described in Sect. 2.1, this condition is satisfied as soon as the mapping

$$\begin{aligned} p : (E \times \Theta, \mathcal{E} \otimes \mathfrak{B}) &\longrightarrow \mathbb{R}_+ \\ (x, \theta) &\longmapsto p_\theta(x) \end{aligned}$$

is measurable. Throughout this section, we assume that such measurability assumptions are satisfied.

### 5.1 The case of the total variation distance

In this section,  $\mathcal{P}$  is the set of all probability measures on  $(E, \mathcal{E})$  and

$$\|P - Q\| = \frac{1}{2} \int_E |p - q| d\mu \tag{33}$$

denotes the total variation loss (TV-loss for short) between  $P, Q \in \mathcal{P}$ .

**Proposition 1** *The family  $\mathcal{T}(\ell, \mathcal{M})$  which consists of all the functions  $t_{(P, Q)}$  defined for  $P = p \cdot \mu$  and  $Q = q \cdot \mu$  in  $\mathcal{M}$  by*

$$t_{(P, Q)} = \frac{1}{2} [\mathbb{1}_{q > p} - Q(q > p)] - \frac{1}{2} [\mathbb{1}_{p > q} - P(p > q)] \tag{34}$$

satisfies Assumption 2 with  $\tau = 1$  and Assumption 3 with  $a_0 = 3/2$  and  $a_1 = 1/2$ .

It follows from Proposition 1 that we may apply our general construction to the so-defined family  $\mathcal{T}(\ell, \mathcal{M})$  with the values  $c = c_0 = 1/3$  (hence  $\lambda = 4/3$ ). The reader can check that the value  $\gamma = 1/100$  satisfies the requirement of our Theorem 1 and that (16) is satisfied with  $\kappa_0 = 220$ . Theorem 1 can therefore be rephrased as follows.

**Corollary 1** *Let  $\beta \geq 1/\sqrt{n}$ ,  $c = 1/3$  and  $\hat{\pi}_X^{TV}$  be the posterior defined by (7) and associated with the family  $\mathcal{T}(\ell, \mathcal{M})$  given in Proposition 1. For all  $\xi > 0$  and any distribution  $\mathbf{P}^*$ , with a probability at least  $1 - 2e^{-\xi/2}$ , the posterior  $\hat{\pi}_X^{TV}$  satisfies*

$$\begin{aligned} \hat{\pi}_X^{TV} \left( \left\{ P \in \mathcal{M}, \ell(\bar{P}^*, P) \leq 220 \left[ \inf_{P' \in \mathcal{M}(\beta)} \ell(\bar{P}^*, P') + 2 \left( \beta + \frac{2\xi}{n\beta} \right) \right] \right\} \right) \\ \geq 1 - e^{-\xi/2} \end{aligned} \tag{35}$$

where

$$\mathcal{M}(\beta) = \left\{ P \in \mathcal{M}, \sup_{r \geq 2\beta} \left[ \frac{200}{nr} \log \left( \frac{\pi(\mathcal{B}(P, 2r))}{\pi(\mathcal{B}(P, r))} \right) \right] \leq \beta \right\}.$$

By convexity, we may write that

$$\inf_{P \in \mathcal{M}(\beta)} \|P - \bar{P}^*\| \leq \inf_{P \in \mathcal{M}(\beta)} \left[ \frac{1}{n} \sum_{i=1}^n \|P - P_i^*\| \right]$$

and the left-hand side is therefore small when there exists  $P \in \mathcal{M}(\beta)$  that approximates well enough most of the marginals of  $\mathbf{P}^*$ . The concentration properties of  $\hat{\pi}_X^{TV}$  remain thus stable with respect to a possible misspecification of the model and a departure from the equidistribution assumption.

In fact, as we shall see in our Example 3 below, the average distribution  $\bar{P}^*$  may belong to  $\mathcal{M}(\beta)$  even when none of the marginals  $P_i^*$  does. This means that in good

situations, the posterior may concentrate around  $\bar{P}^*$ , as it would do in the i.i.d. case when the distribution of the data does belong to  $\mathcal{M}(\beta)$ , even when the data are non-i.i.d. and their marginals do not belong to  $\mathcal{M}(\beta)$ .

**Example 3** (Example 1 continued) Going back to Example 1 and taking for  $\ell$  the TV-loss (then  $a_1 = 1/2$ ), we deduce from (23) that

$$\bar{r}_n = \inf_{P \in \mathcal{M}} \left\| \bar{P}^* - P \right\| + 2 \left( \sqrt{\frac{k}{n}} + \frac{2\xi}{\sqrt{nk}} \right).$$

In particular, if for each  $i \in \{1, \dots, n\}$ ,  $P_i^*$  is the uniform distribution on  $[(i - 1)/n, i/n]$  and  $\mathcal{M}$  contains the uniform distribution  $\mathcal{U}([0, 1])$  on  $[0, 1]$ ,  $\mathcal{M}$  contains  $\bar{P}^* = \mathcal{U}([0, 1])$ , even if none of the marginals  $P_i^*$  belongs to  $\mathcal{M}$ . We then get that

$$\bar{r}_n = 2 \left( \sqrt{\frac{k}{n}} + \frac{2\xi}{\sqrt{nk}} \right)$$

and the posterior concentrates around  $\bar{P}^*$  at a parametric rate.

### 5.2 Case of the $\mathbb{L}_j$ -loss

Let  $j \in (1, +\infty)$ . We denote by  $\mathcal{P}_j$  the set of all finite and signed measures on  $(E, \mathcal{E}, \mu)$  which are of the form  $P = p \cdot \mu$  with  $p \in \mathcal{L}_j(E, \mu) \cap \mathcal{L}_1(E, \mu)$ . Let  $\ell_j$  be the loss defined by  $\ell_j(P, Q) = \|p - q\|_{\mu, j}$  for all  $P = p \cdot \mu$  and  $Q = q \cdot \mu$  in  $\mathcal{P}_j$ . In this section,  $\mathcal{P}$  is the subset that consists of all the probability measures in  $\mathcal{P}_j$ .

**Proposition 2** Let  $\mathcal{M} = \{P = p \cdot \mu, p \in \mathcal{M}\}$  be a subset of  $\mathcal{P}_j$  for which  $\mathcal{M}$  satisfies for some  $R > 0$

$$\|p - q\|_\infty \leq R \|p - q\|_{\mu, j} \quad \text{for all } p, q \in \mathcal{M}. \tag{36}$$

Define for  $P = p \cdot \mu$  and  $Q = q \cdot \mu$  in  $\mathcal{M}$ ,

$$f_{(P, Q)} = \frac{(p - q)_+^{j-1} - (p - q)_-^{j-1}}{\|p - q\|_{\mu, j}^{j-1}} \quad \text{when } P \neq Q \quad \text{and} \quad f_{(P, P)} = 0.$$

Then, the family  $\mathcal{T}(\ell_j, \mathcal{M})$  which contains the functions  $t_{(P, Q)}$  defined for  $P, Q \in \mathcal{M}$  by

$$t_{(P, Q)} = \frac{1}{2R^{j-1}} \left[ \int_E f_{(P, Q)} \frac{dP + dQ}{2} - f_{(P, Q)} \right] \tag{37}$$

satisfies Assumption 2 with  $\tau = 1$  and Assumption 3 with  $a_0 = 3/(4R^{j-1})$  and  $a_1 = 1/(4R^{j-1})$ .

When  $j = 2$ , (36) is typically satisfied when  $\mathcal{M}$  is a subset of a linear space enjoying good connections between the  $\mathbb{L}_2(\mu)$  and the supremum norms. Many finite dimensional linear spaces with good approximation properties do satisfy such connections (e.g. piecewise polynomials of a fixed degree on a regular partition of  $[0, 1]$ , trigonometric polynomials on  $[0, 1]$  etc.). We refer the reader to Birgé and Massart [14, Section 3] for additional examples. The property may also hold for infinite dimensional linear spaces as proven in Baraud [4].

It follows from Proposition 2 that one may choose  $c = c_0 = 1/3$  in (6) and  $\gamma = 1/100$  in Theorem 1. Besides, Theorem 1 applies with  $\kappa_0 = 220$ .

**Example 4** (Example 1 continued) Let us go back to our Example 1 with  $\ell = \ell_j$  and  $\mathcal{T}(\ell, \mathcal{M})$  given in Proposition 2. For the choice of  $\beta$  given in (22) and  $\gamma = 1/100$ , we deduce from (23) (with  $a_1 = 1/(4R^{j-1})$ ) that the resulting posterior  $\hat{\pi}_X$  concentrates on an  $\ell_j$ -ball around  $P^*$  with a radius of order

$$\bar{r}_n = \inf_{p \in \mathcal{M}} \left\| \frac{1}{n} \sum_{i=1}^n p_i^* - p \right\|_{\mu, j} + 4R^{j-1} \left( \sqrt{\frac{k}{n}} + \frac{2\xi}{\sqrt{nk}} \right).$$

### 5.3 The case of the squared Hellinger loss

Here,  $\mathcal{P}$  is the set of all probability measures on  $(E, \mathcal{E})$  and

$$\ell(P, Q) = h^2(P, Q) = \frac{1}{2} \int_E (\sqrt{p} - \sqrt{q})^2 d\mu, \tag{38}$$

is the squared Hellinger distance between two probabilities  $P, Q \in \mathcal{P}$ .

**Proposition 3** *Let  $\psi$  be the function defined by*

$$\psi : \begin{cases} [0, +\infty] \longrightarrow [-1, 1] \\ x \longmapsto \begin{cases} \frac{x-1}{x+1} & \text{if } x \in [0, +\infty) \\ 1 & \text{if } x = +\infty. \end{cases} \end{cases}$$

The family  $\mathcal{T}(\ell, \mathcal{M})$  containing the functions  $t_{(P, Q)}$  defined for  $P = p \cdot \mu$  and  $Q = q \cdot \mu$  in  $\mathcal{M}$  by

$$t_{(P, Q)} = \frac{1}{2} \psi \left( \sqrt{\frac{q}{p}} \right) \tag{39}$$

(with the conventions  $0/0 = 1$  and  $x/0 = +\infty$  for all  $x > 0$ ) satisfies Assumption 2 with  $\tau = 2$  and Assumption 4 with  $a_0 = 2, a_1 = 3/16, a_2 = 3\sqrt{2}/4$ .

With such a choice of family  $\mathcal{T}(\ell, \mathcal{M})$ , (6) is satisfied with  $c = 1/125$ , then  $c_0 \in [0.922, 0.923]$ , and the requirements of Theorem 2 are satisfied with  $\beta = 2\gamma =$

1/500. Then the value  $\kappa_0 = 1694$  suits. The definition (11) of  $r_n(\beta, P)$  for  $P \in \mathcal{M}$  becomes

$$r_n(\beta, P) = \inf \left\{ r \geq \frac{8000}{3n}, \frac{\pi(\mathcal{B}(P, 2r'))}{\pi(\mathcal{B}(P, r'))} \leq \exp\left(\frac{3nr'}{8.10^6}\right) \text{ for all } r' \geq r \right\}, \tag{40}$$

with the convention  $\sup \emptyset = 8000/(3n)$ . Theorem 2 can then be rephrased as follows.

**Corollary 2** *Let  $\pi_X^h$  be the posterior defined by (7) and associated with the family  $\mathcal{T}(\ell, \mathcal{M})$  given in Proposition 3 and the choices  $c = 1/125$  and  $\beta = 1/500$ . For all  $\xi > 0$  and any distribution  $\mathbf{P}^*$ , with a probability at least  $1 - 2e^{-\xi/2}$ ,*

$$\widehat{\pi}_X^h \left( \left\{ P \in \mathcal{M}, h^2(\overline{P}^*, P) \leq 1694r \right\} \right) \geq 1 - e^{-\xi/2}$$

where

$$r = \inf_{P \in \mathcal{M}} \left[ h^2(\overline{P}^*, P) + r_n(\beta, P) \right] + \frac{5334\xi}{n}$$

and  $r_n(\beta, P)$  is given by (40).

As for the total variation distance, we may write that

$$\inf_{P \in \mathcal{M}} h^2(\overline{P}^*, P) \leq \inf_{P \in \mathcal{M}} \left[ \frac{1}{n} \sum_{i=1}^n h^2(P_i^*, P) \right].$$

The left-hand side is small when there exists an element  $P \in \mathcal{M}$  that approximates well most of the marginal distribution  $P_i^*$ . If for such a  $P$ , the quantity  $r_n(\beta, P)$  is small enough, the posterior concentrates around  $\overline{P}^*$  just as it would do if the data were truly i.i.d. with distribution  $P \in \mathcal{M}$ .

**Example 5** (Example 1 continued) Let us go back to Example 1, more precisely Example 2, with  $\ell = h^2$  and  $\mathcal{T}(\ell, \mathcal{M})$  given in Proposition 3. Inequality (21) is satisfied with  $\beta = 2\gamma = 1/500$  and  $a_1 = 3/16$ . It follows from (30) that  $\widehat{\pi}_X^h$  concentrates on an  $h^2$ -ball around  $\overline{P}^*$  with a radius of order

$$\bar{r} = \inf_{P \in \mathcal{M}} h^2(\overline{P}^*, P) + \frac{k + \xi}{n}.$$

## 6 Comparing the classical Bayesian approach to ours

In this section, our aim is to highlight some similarities and differences between the Bayesian posterior and ours. Throughout this section, we consider the squared

Hellinger loss  $\ell = h^2$  and denote by  $\widehat{\pi}_X^K$  the Bayes posterior associated with the model  $(\mathcal{M}, \pi)$ . The letter  $K$  in the notation  $\widehat{\pi}_X^K$  refers to the fact that the Bayesian posterior can be obtained from our general construction by using the Kullback–Leibler divergence as explained in Sect. 3.4. When  $\mathcal{M} = \{P_\theta, \theta \in \Theta\}$  is parametric with  $\Theta \subset \mathbb{R}^k$ , we denote by  $\widehat{\nu}_X^K$  the Bayesian posterior on the parameter space  $\Theta$  and  $\widehat{\nu}_X^h$  that associated to  $\widehat{\pi}_X^h$ .

### 6.1 Some classical concentration results for the Bayes posterior distribution

Most of the results that have been established about the concentration properties of the Bayesian posterior are asymptotic in nature. It seems difficult to establish a general nonasymptotic version of those as we do for our posterior. One of the only exceptions we are aware of is Birgé [13].

When the data are i.i.d. with a distribution  $P^* \in \mathcal{M}$ , a typical asymptotic form of these results is the following one (see Ghosal et al. [19] Theorems 2.1 and 2.4 for example). Let  $\varepsilon_n$  be a sequence of positive numbers that converges to zero when  $n$  goes to infinity. If  $P^*$  fulfils some suitable conditions, that we shall discuss later on and which depend on the prior  $\pi$  and  $\varepsilon_n$ , the following convergence in probability holds true

$$\widehat{\pi}_X^K(\{P \in \mathcal{M}, h^2(P^*, P) \geq M_n \varepsilon_n^2\}) \xrightarrow[n \rightarrow +\infty]{P} 0 \text{ under } P^*. \tag{41}$$

In (41),  $M_n = M$  denotes some large enough positive constant if  $n\varepsilon_n^2 \rightarrow +\infty$  as  $n \rightarrow +\infty$  while  $M_n$  is increasing to infinity as  $n \rightarrow +\infty$  if  $\liminf n\varepsilon_n^2 > 0$  as  $n \rightarrow +\infty$ . The first condition on  $\varepsilon_n$  is typically satisfied when  $\mathcal{M}$  is a nonparametric model while the second one generally applies to parametric ones.

In comparison, in this well-specified framework, our Corollary 2 leads to the following result. For all  $P^* \in \mathcal{M}$  and  $\xi > 0$

$$\begin{aligned} \mathbb{P} \left[ \widehat{\pi}_X^h \left( \left\{ P \in \mathcal{M}, h^2(P^*, P) \geq \kappa'_0 \left( r_n(\beta, P^*) + \frac{\xi}{n} \right) \right\} \right) \geq e^{-\xi/2} \right] \\ \leq 2e^{-\xi/2} \end{aligned} \tag{42}$$

for some numerical constant  $\kappa'_0 > 0$ . If  $P^*$  satisfies  $r_n(\beta, P^*) \leq \varepsilon_n^2$ , we recover (41) by setting  $\xi = \xi_n = (M_n/(\kappa'_0) - 1)n\varepsilon_n^2$ . However, our condition that  $r_n(\beta, P^*) \leq \varepsilon_n^2$  is not equivalent to that imposed on  $P^*$  by Ghosal, Ghosh and van der Vaart [19]. It is actually weaker. In their paper, this condition is fulfilled when the prior puts enough mass on Kullback–Leibler type balls around  $P^*$ . Our approach allows one to consider Hellinger balls only, which are larger and make our assumption weaker. In fact, as already underlined in the Introduction, these Kullback–Leibler type balls could be empty, and the condition unsatisfied, while our theorem would still apply.

The result established by Birgé [13] provides an improvement as compared to the one presented above and established by Ghosal, Ghosh and van der Vaart. Birgé shows that it is essentially possible to get rid of the Kullback–Leibler divergence (see



his Theorem 2) but only when the model is parametric and well-specified. Apart for the nonparametric framework, this result leaves little place for improvement since we know that the Bayesian posterior may fail to concentrate around the true parameter when the model becomes slightly ill-specified.

Another consequence of our Corollary 2, as compared to (41), is that it allows one to control

$$\widehat{\pi}_X^h \left( \left\{ P \in \mathcal{M}, h^2(P^*, P) \geq \kappa'_0 \left( \varepsilon_n^2 + \frac{\xi}{n} \right) \right\} \right)$$

uniformly over the set  $\{P^* \in \mathcal{M}, r_n(\beta, P^*) \leq \varepsilon_n^2\}$ . For example, in the framework of Example 2, for the choice  $\varepsilon_n^2 = ck/n$  with  $c = \log(2B/A)/(\gamma a_1 \beta)$ , we know that  $r_n(\beta, P^*) \leq \varepsilon_n^2$  for all  $P^* \in \mathcal{M}$  and we deduce from (42) that

$$\sup_{P^* \in \mathcal{M}} \mathbb{P} \left[ \widehat{\pi}_X^h \left( \left\{ P \in \mathcal{M}, h^2(P^*, P) \geq \kappa'_0 \left( \varepsilon_n^2 + \frac{\xi}{n} \right) \right\} \right) \geq e^{-\xi/2} \right] \leq 2e^{-\xi/2}.$$

The concentration properties of our posterior is therefore uniform over the statistical model  $\mathcal{M}$ .

### 6.2 About the shapes and sizes of the credible regions

A nice feature of the Bayesian approach lies in the fact that it allows one to build credible regions. In practice, they often play the same role as the confidence regions in the frequentist paradigm. When the data are i.i.d. with distribution  $P^* = P_{\theta^*}$  in a parametric model  $\mathcal{M} = \{P_\theta, \theta \in \Theta\}$ ,  $\Theta \subset \mathbb{R}^k$ , a credible set for the parameter  $\theta^*$  is a subset  $\widehat{\Theta}_{n,X} \subset \Theta$  (only depending on observable quantities) that satisfies  $\widehat{\nu}_X^K(\widehat{\Theta}_{n,X}) \geq 1 - e^{-\xi}$  for some choice of  $\xi > 0$ . When  $\mathcal{M}$  is a regular parametric model with a nonsingular Fisher information matrix  $\mathbf{J}$ , and provided that it satisfies additional assumptions—see van der Vaart [24]—the Bernstein–von Mises theorem applies and tells us that

$$\left\| \widehat{\nu}_X^K - \mathcal{N} \left( \widehat{\theta}_n, (n\mathbf{J}(\theta^*))^{-1} \right) \right\| \xrightarrow[n \rightarrow +\infty]{\text{P}} 0 \text{ under } P_{\theta^*}$$

where  $\widehat{\theta}_n$  denotes the Maximum Likelihood Estimator (MLE for short). Denoting by  $\overline{\chi}_k^{-1}(\xi)$  the  $(1 - e^{-\xi})$ -quantile of a chi-square random variable with  $k$  degrees of freedom and

$$\Theta_{n,X} = \left\{ \theta \in \Theta, n \left| \mathbf{J}^{1/2}(\theta^*) (\widehat{\theta}_n - \theta) \right|^2 \leq \overline{\chi}_k^{-1}(\xi) \right\}, \tag{43}$$

we deduce that

$$\left| \widehat{\nu}_X^K(\Theta_{n,X}) - (1 - e^{-\xi}) \right| \leq \left\| \widehat{\nu}_X^K - \mathcal{N} \left( \widehat{\theta}_n, (n\mathbf{J}(\theta^*))^{-1} \right) \right\| \xrightarrow[n \rightarrow +\infty]{\text{P}} 0$$

hence

$$\widehat{v}_X^K(\Theta_{n,X}) \xrightarrow[n \rightarrow +\infty]{P} 1 - e^{-\xi} \text{ under } P_{\theta^*}.$$

The asymptotic level of ‘‘credibility’’ of the set  $\Theta_{n,X}$  is therefore  $1 - e^{-\xi}$ . This set is not, however, a genuine credible region since it depends on the unknown parameter  $\theta^*$ . We would obtain a genuine credible region by replacing  $\theta^*$  by  $\widehat{\theta}_n$  in the expression of  $\Theta_{n,X}$ . This substitution would change the level of credibility but not the shape of the region, which is an ellipsoid centred at  $\widehat{\theta}_n$  and the axes of which are given by the eigenvectors of the Fisher information matrix.

The aim of this section is to show that our posterior concentrates its mass on regions that have the same shape and approximately the same size. The size of  $\Theta_{n,X}$  is determined by the value of the quantile  $\overline{\chi}_k^{-1}(\xi)$ . The aim of the following lemma is to specify the order of magnitude of this quantile as a function of  $k$  and  $\xi$ . In fact, we consider below the more general case of the quantiles of a gamma distribution  $\gamma(s, \sigma)$  with parameters  $s, \sigma > 0$ , that is, the distribution with density  $x \mapsto (x^{s-1}e^{-x/\sigma})/(\sigma^s \Gamma(s))$  with respect to the Lebesgue measure on  $\mathbb{R}_+$ . The proof is postponed to Sect. 10.1.

**Lemma 1** *For  $s, \sigma, \xi > 0$ , let  $\overline{\gamma}_{s,\sigma}^{-1}(\xi)$  be the  $(1 - e^{-\xi})$ -quantile of the gamma distribution  $\gamma(s, \sigma)$  and  $\overline{\Phi}^{-1}(\xi)$  that of a standard Gaussian random variable. Then,*

$$\overline{\gamma}_{s,\sigma}^{-1}(\xi) \leq \sigma \left( \sqrt{s} + \sqrt{\xi} \right)^2 \tag{44}$$

and for all  $s = t + 1 > 1$  and  $\xi \geq \log 2 + 1/(12t)$ ,

$$\overline{\gamma}_{s,\sigma}^{-1}(\xi) \geq \sigma \left[ t + \left[ \sqrt{t} \overline{\Phi}^{-1} \left( \xi - \frac{1}{12t} \right) \right] \vee \left[ \xi + \log \left( \frac{e^{-1/(12t)}}{\sqrt{2\pi t}} \right) \right] \right]. \tag{45}$$

Since  $\overline{\Phi}^{-1}(\xi)$  is equivalent to  $\sqrt{2\xi}$  for large values of  $\xi > 0$ , these two inequalities show that for  $s$  and  $\xi$  large enough,  $\overline{\gamma}_{s,\sigma}^{-1}(\xi)$  is of order  $\sigma [s + \xi]$ . In particular,  $\overline{\chi}_k^{-1}(\xi) = \overline{\gamma}_{k/2,2}^{-1}(\xi)$  is of order  $k + \xi$  for  $k$  and  $\xi$  large enough.

To compare ourselves with the classical Bayesian paradigm, we prove in Sect. 10.2 the result below for our posterior. This result is based on the assumption that the statistical model  $\mathcal{M}$  is regular in the sense that is defined in Ibragimov and Has’minskii [20]. In order to avoid too many technicalities here, we refer the reader to our Sect. 8.3, more precisely Corollary 4, for a complete description of the assumptions on the statistical model  $\mathcal{M}$ .

**Theorem 4** *Assume that the statistical model  $\mathcal{M}$  satisfies the assumptions of Corollary 4. If  $X_1, \dots, X_n$  are i.i.d. with distribution  $P_{\theta^*} \in \mathcal{M}$ , for all  $\xi > 0$  and  $n$  large enough, with a probability  $1 - 2e^{-\xi}$ ,*

$$\widehat{v}_X^h \left( \left\{ \theta \in \Theta, n \left| \mathbf{J}^{1/2}(\theta^*) (\theta - \theta^*) \right|^2 \leq \kappa^* (k + \xi) \right\} \right) \geq 1 - e^{-\xi} \tag{46}$$

where  $\kappa^*$  is a positive numerical constant.

The set

$$\left\{ \theta \in \Theta, n \left| \mathbf{J}^{1/2}(\theta^*) (\theta - \theta^*) \right|^2 \leq \kappa^* (k + \xi) \right\}$$

possesses the same shape and, by Lemma 1, approximately the same size as the set  $\Theta_{n,X}$  defined by (43). We deduce from Theorem 4 that the classical Bayes posterior and ours concentrate both on similar sets. If  $\widehat{\theta}_n$  is an asymptotically efficient estimator of  $\theta^*$ , it is therefore reasonable to look for a credible region of the form

$$\left\{ \theta \in \Theta, n \left| \mathbf{J}^{1/2}(\widehat{\theta}_n) (\theta - \widehat{\theta}_n) \right|^2 \leq t \right\}, t > 0$$

for  $\widehat{v}_X^h$  as we would do for the classical Bayes one.

### 6.3 Robustness

As already mentioned, our approach allows the statistician to design robust posteriors by choosing as a loss function the squared Hellinger loss or the total variation one. In this section, we illustrate this property on a concrete example. Consider the statistical model  $\mathcal{M} = \{P_\theta = \mathcal{N}(\theta, 1), \theta \in \mathbb{R}\}$  and the prior  $\pi$  associated with the distribution  $\nu = \mathcal{N}(0, 1)$  on  $\Theta = \mathbb{R}$ . Then, the Bayes posterior on  $\Theta$  is  $\widehat{v}_X^K = \mathcal{N}(\widehat{m}_n, \sigma_n^2)$  with  $\widehat{m}_n = (n + 1)^{-1} \sum_{i=1}^n X_i$  and  $\sigma_n^2 = 1/(n + 1)$ . It concentrates on intervals of the form  $[\widehat{m}_n - c/\sqrt{n + 1}, \widehat{m}_n + c/\sqrt{n + 1}]$  for  $c > 0$  large enough. If the distribution of the data is contaminated so that  $X_1, \dots, X_n$  are i.i.d. with distribution

$$P^* = \left(1 - \frac{1}{n}\right) P_0 + \frac{1}{n} \mathcal{N}(10^4(n + 1), 1/n),$$

then with a probability at least  $1 - (1 - 1/n)^n \geq 1 - 1/e > 63\%$ , the posterior concentrates around  $\widehat{m}_n \approx 10^4$ , hence far away from 0, even though  $P^*$  and  $P_0$  are close:  $\|P^* - P_0\| \leq 1/n$ .

In this specific framework, the model  $\mathcal{M}$  is regular, the Fisher information is constant and positive,  $\nu$  admits a positive density which is continuous at  $\theta^* = 0$  and for all  $\theta, \theta' \in \Theta, h^2(\theta, \theta') = 1 - e^{-|\theta - \theta'|^2/8}$ . We shall see in Sect. 8, more precisely in Corollary 4, that for such regular statistical models  $r_n(\beta, P_0) \leq \kappa^*/n$  for some numerical constant  $\kappa^* > 0$ , at least for  $n$  large enough. Since  $h^2(P^*, P_0) \leq \|P^* - P_0\| \leq 1/n$ , we deduce from Corollary 2 that the posterior  $\widehat{v}_X^h$  concentrates on a set of the form

$$\left\{ \theta \in \mathbb{R}, h^2(\theta, 0) \leq \frac{c}{n} \right\} = \left\{ \theta \in \mathbb{R}, |\theta| \leq \sqrt{8 \log \left( \frac{1}{1 - c/n} \right)} \right\}$$

with  $c > 0$ . This set is an interval around 0 of approximate length  $1/\sqrt{n}$ , at least for  $n$  sufficiently large. Despite the contamination of the data, the concentration property of  $\widehat{v}_X^h$  remains thus the same as in the well-specified case.

## 7 Applications

### 7.1 How to choose $\beta$ in Theorem 1 for a translation model?

In this section, we consider the translation model  $\mathcal{M} = \{P_\theta = p(\cdot - \theta) \cdot \mu, \theta \in \mathbb{R}\}$  where  $p$  is a density on  $\mathbb{R}$  with respect to the Lebesgue measure  $\mu$ . Our aim is to estimate the translation parameter  $\theta$  by using a prior  $\nu_\sigma$  on  $\Theta = \mathbb{R}$  with a density (with respect to  $\mu$ ) of the form  $q(\cdot/\sigma)/\sigma$  for some density  $q$  and positive number  $\sigma$ . We evaluate the estimation error by means of the total variation loss. In order to use our construction we need to tune the parameter  $\beta$ . In Sect. 4.2, we suggested to choose  $\beta \geq 1/\sqrt{n}$  satisfying (18). In order to find such a value of  $\beta = \beta(\alpha)$ , we may proceed as follows. Consider a symmetric bounded interval  $I = [-l/2, l/2] \subset \mathbb{R}$  of length  $l > 0$  satisfying  $\nu_\sigma(I) \geq 1 - \alpha$ , hence concentrating most of the mass of the prior  $\nu_\sigma$ . If the set  $\mathcal{M}(\beta)$  is large enough to contain  $\{P_\theta, \theta \in I\}$ ,

$$\pi(\mathcal{M}(\beta)) \geq \pi(\{P_\theta, \theta \in I\}) = \nu_\sigma(I) \geq 1 - \alpha \tag{47}$$

and  $\beta$  satisfies (18). We deduce from our Corollary 1 that the corresponding posterior  $\widehat{\pi}_X^{TV}$  concentrates with a probability at least  $1 - 2e^{-\xi/2}$  on a TV-ball with a radius of order

$$\inf_{P' \in \mathcal{M}(\beta)} \ell(\overline{P}^*, P') + 2 \left( \beta + \frac{2\xi}{n\beta} \right) \leq \inf_{\theta \in I} \ell(\overline{P}^*, P_\theta) + 2\beta + \frac{4\xi}{\sqrt{n}} = r(\beta). \tag{48}$$

The approximation term  $\inf_{\theta \in I} \ell(\overline{P}^*, P_\theta)$  is small as soon as  $\overline{P}^*$  is close enough to a distribution  $P_{\theta^*}$  whose parameter  $\theta^*$  belongs to  $I$ . If we want to prevent us from the situation where  $\operatorname{argmin}_{\theta \in \Theta} \ell(\overline{P}^*, P_\theta)$  is far from 0, we need to increase  $I$  (or equivalently diminish  $\alpha$ ). What would be the consequence on the value of  $\beta = \beta(\alpha)$ ? What if we increase  $\sigma$ , to make the prior distribution flatter, or diminish  $\sigma$  to make it more picky? Finally, what is the influence of the choice of the density  $q$  on the size of  $\beta$ ?

These are the questions we want to answer in this section. In order to simplify the presentation of our results and avoid technicalities, we make the change of variables  $l = 2\sigma t$ , or equivalently  $t = l/(2\sigma) > 0$ , and assume the following.

**Assumption 6** The density  $q$  is positive, symmetric and decreasing on  $\mathbb{R}_+$ . There exists some nonnegative and nondecreasing function  $\varphi : [0, 1) \rightarrow \mathbb{R}_+$  such that

$$\|P_0 - P_\theta\| \leq r \iff |\theta| \leq \varphi(r) \quad \text{for all } r \in [0, 1).$$

When  $p$  is symmetric and nonincreasing on  $\mathbb{R}_+$ , the total variation distance between  $P_0$  and  $P_\theta$  is given by

$$\|P_0 - P_\theta\| = 2P_0([0, |\theta|/2]) \quad \text{for all } \theta \in \mathbb{R}.$$

Our Assumption 6 is then satisfied with  $\varphi(r) = F_0^{-1}[(r + 1)/2]$  for all  $r \in [0, 1)$ , where  $F_0^{-1}$  denotes the quantile function of the distribution  $P_0$ . We set

$$\bar{\Gamma} = \max \left\{ \left[ \sup_{0 < r \leq 1/4} \frac{\varphi(2r)}{\varphi(r)} \right] q(0), \frac{1}{2\varphi(1/4)} \right\} \tag{49}$$

and assume that this quantity is finite. Note that it only depends on  $q(0)$  and  $p$ . For example, if  $p$  is the density  $x \mapsto (1/2)e^{-|x|}$ ,

$$\|P_0 - P_\theta\| = 1 - \exp[-|\theta|/2] \quad \text{and} \quad \varphi : r \mapsto -2 \log(1 - r).$$

Since the mapping  $r \mapsto [\varphi(2r)/\varphi(r)]$  is increasing, we obtain in this case

$$\bar{\Gamma} = \frac{1}{\log(4/3)} \max \left\{ q(0) \log 2, \frac{1}{4} \right\}.$$

If now  $p : x \mapsto (s/2)(1 - |x|)^{s-1} \mathbb{1}_{|x| < 1}$  with  $s > 0$ ,

$$\|P_0 - P_\theta\| = 1 - (1 - |\theta|/2)^s \quad \text{and} \quad \varphi : r \mapsto 2[1 - (1 - r)^{1/s}].$$

The mapping  $r \mapsto \varphi(2r)/\varphi(r)$  has a continuous extension on  $[0, 1/4]$  and is therefore bounded. Given  $q(0)$ ,  $\bar{\Gamma}$  is therefore a finite number.

The following result is proven in Sect. 10.3.

**Proposition 4** *Assume that Assumption 6 is satisfied and  $\bar{\Gamma}$  is finite. Let  $t$  be a  $(1 - \alpha/2)$ -quantile of  $q$  with  $\alpha \leq 1/2$ . The set  $\mathcal{M}(\beta)$  contains the subset  $\{P_\theta, \theta \in [-\sigma t, \sigma t]\}$  and therefore satisfies (47) if*

$$\beta \geq \bar{\beta} = \sqrt{\frac{1}{n\gamma} \max \left\{ \log \left( \frac{\bar{\Gamma}(\sigma \vee 1)}{q(2t)} \right), \log 4 \right\}}. \tag{50}$$

Let us now comment on this result. The quantity  $\bar{\beta}$  may be written as  $C/\sqrt{n}$  with

$$C = \sqrt{\frac{1}{\gamma} \max \left\{ \log \left( \frac{\bar{\Gamma}(\sigma \vee 1)}{q(2t)} \right), \log 4 \right\}}.$$

Increasing the value of  $\sigma$  or that of  $t$  enlarges the interval  $I = [-\sigma t, \sigma t]$ . It also makes the value of  $C = C(\sigma, t)$  larger. Increasing  $\sigma$  makes the prior  $\nu_\sigma$  flatter and for

a fixed value of  $t > 0$ ,  $C = C(\sigma)$  increases as  $\sqrt{\log \sigma}$  when  $\sigma$  is larger than 1. In the other case, for a fixed value of  $\sigma$ ,  $C = C(t)$  increases as  $\sqrt{\log(1/q(2t))}$ . For example, when  $q$  is the density of a standard Gaussian random variable,  $\sqrt{\log(1/q(2t))}$  is of order  $t$ , while for the Laplace and the Cauchy distributions it is of order  $\sqrt{t}$  and  $\sqrt{\log t}$  respectively. This result illustrates the fact that it is safer to use priors with heavy tails when the size of the location parameter is uncertain. In case of a light-tailed prior, it may be wise to introduce a scaling parameter  $\sigma > 1$ . By taking  $\sigma = 10$ , the concentration radius only increases by a factor less than 1.6, while the interval  $I$  is ten times longer.

### 7.2 Fast rates

We go back to the statistical framework described in Sect. 7.1 and consider the special case of the density  $p : x \mapsto sx^{s-1}\mathbb{1}_{(0,1]}$  with  $s \in (0, 1]$ . As before, we choose the TV-loss. In this specific situation,

$$\|P_\theta - P_{\theta'}\| = |\theta - \theta'|^s \wedge 1 \quad \text{for all } \theta, \theta' \in \mathbb{R} \tag{51}$$

and consequently,  $\varphi(r) = r^{1/s}$  for all  $r \in [0, 1)$ . Besides, the family  $\mathcal{T}(\ell, \mathcal{M})$  given by (34) satisfies not only Assumption 3 but also Assumption 4 with  $a_2 = 1$ . These two facts are proven in Baraud [4, Examples 5 and 6]. As a consequence, Theorem 2 applies. The reader can check that the constants  $c = \beta = 0.1$  and  $\gamma = 0.01$  satisfy the requirements of Theorem 2 and that its conclusion holds true with  $\kappa_0 = 144$ .

In order to be more specific about the concentration radius of our posterior  $\hat{\pi}_X^{\text{TV}}$ , the following proposition provides an upper bound for the quantity  $r_n(\beta, P_\theta)$ . The proof is postponed to Sect. 10.4.

**Proposition 5** *Let  $t_0$  be the third quartile of  $v_1$ . If the density  $q$  is positive, symmetric and decreasing on  $[0, +\infty)$ , for all  $\theta \in \mathbb{R}$  the quantity  $r_n(\beta, P_\theta)$  is not larger than*

$$\bar{r}_n(\beta, P_\theta) = \frac{2000}{n} \max \left\{ \log \left( \frac{\bar{\Gamma}(\sigma \vee 1)}{q \left[ 2 \left( \frac{|\theta|}{\sigma} \vee t_0 \right) \right]} \right), \log 4 \right\}. \tag{52}$$

Then, our Theorem 2 tells us that for all  $\xi > 0$ , with a probability at least  $1 - 2e^{-\xi/2}$ , the posterior satisfies

$$\hat{\pi}_X^{\text{TV}} \left( \mathcal{B}(\bar{P}^*, 144r) \right) \geq 1 - e^{-\xi/2}$$

with

$$r \leq \inf_{\theta \in \mathbb{R}} \left[ \|\bar{P}^* - P_\theta\| + \bar{r}_n(\beta, P_\theta) \right] + \frac{40\xi}{n}. \tag{53}$$

When the data are i.i.d. with distribution  $P_{\theta^*}$ , with probability close to 1, a randomized estimator  $P_{\hat{\theta}}$  with distribution  $\hat{\pi}_X^{TV}$  satisfies with high probability

$$|\theta^* - \hat{\theta}|^s \wedge 1 = \|P_{\theta^*} - P_{\hat{\theta}}\| \leq \frac{C(\xi, s, q, \theta^*, \sigma)}{n}.$$

This inequality implies, at least for  $n$  large enough, that

$$|\theta^* - \hat{\theta}| \leq \frac{C^{1/s}(\xi, s, q, \theta^*, \sigma)}{n^{1/s}},$$

which means that the parameter  $\theta^*$  is estimated at rate  $n^{-1/s}$ . This rate is much faster than the usual  $(1/\sqrt{n})$ -parametric one that is reached by an estimator based on a moment method for instance. For example, when  $s = 1/3$  and  $n = 100$ , a moment estimator provides an accuracy of order  $10^{-1}$  while that of  $\hat{\theta}$  is of order  $10^{-6}$ . Since  $p$  is unbounded, note that the maximum likelihood estimator for  $\theta^*$  does not exist and is therefore useless.

It follows from the work of Le Cam that in a translation model  $\mathcal{M}$  of the form  $\{P_\theta = p(\cdot - \theta) \cdot \mu, \theta \in \mathbb{R}\}$ , where  $p$  is a density with respect to the Lebesgue measure  $\mu$ , it is impossible to estimate a distribution  $P^* \in \mathcal{M}$  from an  $n$ -sample at a rate faster than  $1/n$  for the TV-loss. Because of (51), the rate we get is not only optimal for estimating the distribution  $P_{\theta^*}$  but also for estimating the parameter  $\theta^*$  with respect to the Euclidean distance.

An alternative rate-optimal estimator for estimating  $\theta^*$  is that given by the minimum of the observations. This estimator is unfortunately obviously non-robust to the presence of an outlier among the sample. Our construction provides an estimator which possesses the property of being both rate-optimal and robust.

It also interesting to see how the quantity  $\bar{r}_n(\beta, P_\theta)$  given in (52) deteriorates under a misspecification of the prior  $\nu_\sigma$ , that is, when the size of the parameter  $\theta^*$  is large compared to  $\sigma$ . When  $q$  is Gaussian,  $\bar{r}_n(\beta, P_{\theta^*})$  increases by a factor of order  $(\theta^*/\sigma)^2$  while for the Laplace and Cauchy distributions it is of order  $|\theta^*|/\sigma$  and  $\log(|\theta^*|/\sigma)$  respectively. From these results, we conclude as before that the Cauchy distribution possesses some advantages over the other two distributions when little information is available on the location of the parameter  $\theta^*$ .

### 7.3 A general result under entropy

In this section, we equip  $E = \mathbb{R}^k$  with the Lebesgue measure  $\mu$  and the norm  $|\cdot|_\infty$ . We consider the TV-loss and the location-scale family

$$\mathcal{M} = \left\{ P_{(p, \mathbf{m}, \sigma)} = \frac{1}{\sigma^k} p\left(\frac{\cdot - \mathbf{m}}{\sigma}\right) \cdot \mu, p \in \mathcal{M}_0, \mathbf{m} \in \mathbb{R}^k, \sigma > 0 \right\}, \tag{54}$$

where  $\mathcal{M}_0$  is a set of densities on  $\mathbb{R}^k$ . Given independent observations  $X_1, \dots, X_n$  with presumed distribution  $P^* = P_{(p^*, \mathbf{m}^*, \sigma^*)} \in \mathcal{M}$ , our aim is to estimate the density

$p^* \in \mathcal{M}_0$ , the location parameter  $\mathbf{m}^* \in \mathbb{R}^k$  and the scale parameter  $\sigma^* > 0$ , hence the parameter  $\theta^* = (p^*, \mathbf{m}^*, \sigma^*) \in \Theta = \mathcal{M}_0 \times \mathbb{R}^k \times (0, +\infty)$ . We assume that the set of densities  $\mathcal{M}_0$  satisfies the following conditions.

**Assumption 7** Let  $\tilde{D}$  be a continuous nonincreasing mapping from  $(0, +\infty)$  to  $[1, +\infty)$  such that  $\lim_{\eta \rightarrow +\infty} \eta^{-2} \tilde{D}(\eta) = 0$ . For all  $\eta > 0$ , there exists a finite subset  $\mathcal{M}_0[\eta] \subset \mathcal{M}_0$  satisfying

$$|\mathcal{M}_0[\eta]| \leq \exp[\tilde{D}(\eta)] \tag{55}$$

such that for all  $p \in \mathcal{M}_0$ , there exists  $\bar{p} \in \mathcal{M}_0[\eta]$  that satisfies

$$\|P_{(p, \mathbf{0}, 1)} - P_{(\bar{p}, \mathbf{0}, 1)}\| = \frac{1}{2} \int_{\mathbb{R}^k} |p - \bar{p}| d\mu \leq \eta. \tag{56}$$

Besides, we assume that there exist  $A, s > 0$  such that for all  $p \in \mathcal{M}_0, \mathbf{m} \in \mathbb{R}^k$  and  $\sigma \geq 1$ ,

$$\|P_{(p, \mathbf{0}, 1)} - P_{(p, \mathbf{m}, \sigma)}\| \leq \left[ A \left( \left( \frac{\mathbf{m}}{\sigma} \Big|_{\infty} \right)^s + \left( 1 - \frac{1}{\sigma} \right)^s \right) \right] \wedge 1. \tag{57}$$

The first part of Assumption 7, which corresponds to inequalities (55) and (56), aims at measuring the size of the set  $\mathcal{M}_0$  by means of its entropy. The entropy of a set controls its metric dimension and usually determines the minimax rate of convergence over it as shown in Birgé [9]. With the second part of Assumption 7, namely inequality (57), we require some regularity properties of the TV-loss with respect to the location and scale parameters. It will be commented on later. We shall see that this condition may be satisfied even when the densities in  $\mathcal{M}_0$  are not smooth.

Let us now turn to the choice of our prior. We first consider a countable subset of the parameter space  $\Theta$  that will be proven to possess good approximation properties. Namely, we define for  $\eta, \delta > 0$

$$\Theta[\eta, \delta] = \left\{ \left( \bar{p}, (1 + \delta)^{j_0} \delta \mathbf{j}, (1 + \delta)^{j_0} \right), (\bar{p}, j_0, \mathbf{j}) \in \mathcal{M}_0[\eta] \times \mathbb{Z} \times \mathbb{Z}^k \right\}$$

and we associate a positive weight  $L_\theta$  with any element  $\theta = \theta(\bar{p}, j_0, \mathbf{j}) \in \Theta[\eta, \delta]$  as follows

$$L_\theta = (k + 1)L + \log |\mathcal{M}_0[\eta]| + 2 \sum_{i=0}^k \log(1 + |j_i|) \tag{58}$$

with  $L = \log[(\pi^2/3) - 1]$ . It is not difficult to check that  $\sum_{\theta \in \Theta[\eta, \delta]} e^{-L_\theta} = 1$ , and we may therefore endow  $\mathcal{M}$  with the (discrete) prior  $\pi$  defined as

$$\pi(\{P_\theta\}) = e^{-L_\theta} \quad \text{for all } \theta \in \Theta[\eta, \delta]. \tag{59}$$



With such a prior, our posterior  $\widehat{\pi}_X^{TV}$  given in Corollary 1 possesses the following properties.

**Corollary 3** *Let  $\xi > 0$   $K > 1$ . Assume that  $\mathcal{M}_0$  satisfies Assumption 7 and define*

$$\eta = \eta_n = \inf \mathcal{D}_n \quad \text{with} \quad \mathcal{D}_n = \left\{ \eta > 0, \tilde{D}(\eta) \leq \frac{n\eta^2}{24} \right\} \tag{60}$$

$$\delta = \delta_n = \left( \frac{\eta_n}{2A} \right)^{1/s}, \tag{61}$$

$$\beta = \beta_n = \frac{1}{2} \left[ K\eta_n + 2\sqrt{\frac{18.6(k+1)}{n}} \right] \tag{62}$$

and the subset  $\mathcal{M}_n(K)$  of  $\mathcal{M}$  that consists of the elements  $P_{(p, \mathbf{m}, \sigma)}$  for which

$$|\log \sigma| \vee \left| \frac{\mathbf{m}}{\sigma} \right|_\infty \leq \Lambda_n = \exp \left[ \frac{(K^2 - 1)n\eta_n^2}{48(k+1)} + \log \log(1 + \delta_n) \right]. \tag{63}$$

Then, the posterior  $\widehat{\pi}_X^{TV}$  satisfies the following property: there exists a numerical constant  $\kappa'_0 > 0$  such that for all  $\xi > 0$ ,

$$\mathbb{E} \left[ \widehat{\pi}_X \left( {}^c\mathcal{B}(\overline{P}^*, \kappa'_0 r_n) \right) \right] \leq 2e^{-\xi} \tag{64}$$

with

$$r_n = \inf_{P \in \mathcal{M}_n(K)} \ell(\overline{P}^*, P) + K\eta_n + \sqrt{\frac{k+1}{n}} + \frac{\xi}{\sqrt{n(k+1)}} \wedge \frac{\xi}{Kn\eta_n}. \tag{65}$$

Let us now comment on this result. The radius  $r_n$  is the sum of three main terms, omitting the dependency with respect to  $\xi$ . The first one,  $\inf_{P \in \mathcal{M}_n(K)} \ell(\overline{P}^*, P)$ , corresponds to the approximation of  $\overline{P}^*$  by an element of  $\mathcal{M}$  whose location and scale parameters satisfy the constraints given in (63). The quantity  $\eta_n$ , involved in the second term, usually corresponds to the minimax rate for solely estimating a density  $p \in \mathcal{M}_0$  from an  $n$ -sample. Finally, the third term  $\sqrt{(k+1)/n}$  corresponds to the rate we would get for solely estimating the location and translation parameters  $(\mathbf{m}, \sigma) \in \mathbb{R}^{k+1}$  when the density  $p$  is known.

Let us now provide some examples for which our condition (57) is satisfied. We start with an example where the densities in  $\mathcal{M}_0$  are smooth.

**Lemma 2** *Assume that the set  $\mathcal{M}_0$  consists of densities  $p$  that are supported on  $[0, 1]^k$ , satisfy  $\sup_{p \in \mathcal{M}_0} \|p\|_\infty \leq L_0$  and*

$$\sup_{p \in \mathcal{M}_0} |p(\mathbf{x}) - p(\mathbf{x}')| \leq L_1 |\mathbf{x} - \mathbf{x}'|^s \quad \text{for all } \mathbf{x}, \mathbf{x}' \in \mathbb{R}^k, \tag{66}$$

with constants  $L_0, L_1 > 0$  and  $s \in (0, 1]$ . Then (57) is satisfied with  $A = L_1 \vee [(1 + L_1 k^{s/2} + L_0)/2]$ .

Nevertheless, condition (57) may also be satisfied for families  $\mathcal{M}_0$  of densities which are not smooth, as shown in Lemma 3 below. It makes it possible to consider the following example.

**Example 6** We consider here the situation where  $k = 1$  and  $\mathcal{M}_0$  is the set of all nonincreasing densities on  $[0, 1]$  that are bounded by  $B > 1$ . Then,  $\mathcal{M}$  consists of all the probabilities whose densities are supported on intervals  $I$  with positive lengths, nonincreasing on  $I$  and which are bounded by  $B/\mu(I)$ . Birman and Solomjak [15] proved that  $\mathcal{M}_0$  satisfies Assumption 7 with  $\tilde{D}(\eta)$  of order  $(1/\eta) \vee 1$  (up to some constant that depends on  $B$ ). We deduce from (60) that  $\eta_n$  is therefore of order  $n^{-1/3}$ . Besides, it follows from Lemma 3 below that (57) is satisfied with  $A = B$  and  $s = 1$ . We may therefore apply Corollary 3. For a value of  $K$  large enough compared to 1,  $\Lambda_n$  defined by (63) is larger than  $\exp [CK^2n^{1/3}]$  for some constant  $C > 0$  (depending on  $A$ ). In particular, if  $X_1, \dots, X_n$  are i.i.d. with a density of the form

$$x \mapsto p^*(x) = \frac{1}{\sigma^*} p\left(\frac{x - m^*}{\sigma^*}\right)$$

where  $p \in \mathcal{M}_0, |m^*/\sigma^*| \leq \exp [CK^2n^{1/3}]$  and

$$\exp \left[ -\exp [CK^2n^{1/3}] \right] \leq \sigma^* \leq \exp \left[ \exp [CK^2n^{1/3}] \right],$$

(64) is satisfied with  $r_n$  of order  $C'n^{-1/3}$  where the constant  $C' > 0$  only depends on  $\xi, K, B$  but not on  $m^*$  and  $\sigma^*$ . This means that the concentration properties of  $\hat{\pi}_X$  hold true uniformly over a huge range of translation and scale parameters  $\mathbf{m}$  and  $\sigma$  when  $n$  is large enough.

**Lemma 3** *Let  $p$  be a nonincreasing density on  $(0, +\infty)$ . For all  $\sigma \geq 1$*

$$\frac{1}{2} \int_{\mathbb{R}} \left| \frac{1}{\sigma} p\left(\frac{x}{\sigma}\right) - p(x) \right| dx \leq \left(1 - \frac{1}{\sigma}\right). \tag{67}$$

*If, furthermore,  $p$  is bounded by  $B \geq 1$ , for all  $m \in \mathbb{R}$ ,*

$$\frac{1}{2} \int_{\mathbb{R}} |p(x) - p(x - m)| dx \leq (|m|B) \wedge 1. \tag{68}$$

*In particular, for all  $m \in \mathbb{R}$  and  $\sigma \geq 1$ ,*

$$\frac{1}{2} \int_{\mathbb{R}} \left| \frac{1}{\sigma} p\left(\frac{x - m}{\sigma}\right) - p(x) \right| dx \leq \left[ B \left| \frac{m}{\sigma} \right| + \left(1 - \frac{1}{\sigma}\right) \right] \wedge 1. \tag{69}$$

### 7.4 Estimating a parameter under sparsity

Let us consider a parametric dominated model  $\mathcal{M} = \{P_\theta = p_\theta \cdot \mu, \theta \in \mathbb{R}^k\}$  where the dimension  $k$  of the parameters is large. We presume, even though this might not

be true, that the data are i.i.d. with distribution  $P_{\theta^*} \in \mathcal{M}$  and that the coordinates of the true parameter  $\theta^* = (\theta_1^*, \dots, \theta_k^*)$  are all zero except for a small number of them. Our aim is to estimate  $P_{\theta^*}$  from the observation of  $X_1, \dots, X_n$  by using the squared Hellinger loss.

To tackle this problem, we partition the model  $\mathcal{M}$  into the sub-models  $\{\mathcal{M}_m, m \subset \{1, \dots, k\}\}$  where  $\mathcal{M}_m$  consists of those distributions  $P_\theta \in \mathcal{M}$  for which the coordinates of  $\theta = (\theta_1, \dots, \theta_k)$  are all zero except those with an index  $i \in m$ . We denote by  $\Theta_m$  the set of such parameters, so that  $\mathcal{M}_m = \{P_\theta, \theta \in \Theta_m\}$ , and we use the conventions  $\Theta_\emptyset = \{\mathbf{0}\}$  and  $\mathcal{M}_\emptyset = \{P_\mathbf{0}\}$ . Given some positive number  $R > 0$ , we equip each parameter space  $\Theta_m, m \subset \{1, \dots, k\}$ , with the uniform distribution  $\nu_m$  on  $\Theta_m(R) = [-R, R]^k \cap \Theta_m$  when  $m \neq \emptyset$  and the Dirac mass  $\nu_\emptyset = \delta_\mathbf{0}$  at  $\mathbf{0} \in \mathbb{R}^k$  when  $m = \emptyset$ . We may then define on  $\mathbb{R}^k = \bigcup_{m \subset \{1, \dots, k\}} \Theta_m$ , the hierarchical prior

$$v = \sum_{m \subset \{1, \dots, k\}} e^{-L_m} \nu_m \quad \text{with} \quad L_m = |m| \log k + k \log \left(1 + \frac{1}{k}\right). \tag{70}$$

We endow  $\mathcal{M}$  with the  $\sigma$ -algebra and the prior  $\pi$  as described in Sect. 2.1. Besides, we assume that there exists  $s \in (0, 1]$  and a positive number  $B_k = B_k(R)$ , possibly depending on  $k$  and  $R$  (although we drop the dependency with respect to  $R$ ), such that

$$h^2(P_\theta, P_{\theta'}) \leq B_k |\theta - \theta'|_\infty^s \quad \text{for all } \theta, \theta' \in [-R, R]^k. \tag{71}$$

The following result is proven in Sect. 10.8.

**Proposition 6** *Assume that*

$$p : \begin{cases} E \times \mathbb{R}^k \longrightarrow \mathbb{R}_+ \\ (x, \theta) \longmapsto p_\theta(x) \end{cases}$$

*is measurable. If  $RB_k^{1/s} \geq 1$  there exists a numerical constant  $\kappa'_0 > 0$  such that for any distribution  $\mathbf{P}^*$  and  $\xi > 0$*

$$\mathbb{E} \left[ \widehat{\pi}_X^h \left( c_{\mathcal{B}}(\overline{P}^*, \kappa'_0 r) \right) \right] \leq 2e^{-\xi}$$

where

$$r = \inf_{m \subset \{1, \dots, k\}} \left[ \inf_{\theta \in \Theta_m(R)} \ell(\overline{P}^*, P_\theta) + \frac{|m| \log(2kR(nB_k)^{1/s}) + \xi}{n} \right]. \tag{72}$$

Let us now comment on this result. First of all, the mapping

$$R \mapsto \sup \left\{ \frac{h^2(P_\theta, P_{\theta'})}{|\theta - \theta'|_\infty^s}, \theta \neq \theta', \theta, \theta' \in [-R, R]^k \right\}$$

being nondecreasing, our condition  $RB_k^{1/s} = R[B_k(R)]^{1/s} \geq 1$  is always satisfied for a value of  $R$  sufficiently large.

When  $B_k$  does not increase faster than a power of  $k$ , the radius  $r$  given in (72) only depends logarithmically on the dimension  $k$  of the parameter space, as expected.

Let us now illustrate Proposition 6 by choosing some specific models  $\mathcal{M} = \{P_\theta, \theta \in \mathbb{R}^k\}$ . If  $P_\theta$  is the Gaussian distribution with mean  $\theta \in \mathbb{R}^k$  and covariance matrix  $\sigma^2 I_k$ , where  $I_k$  denotes the  $k \times k$  identity matrix,

$$h^2(P_\theta, P_{\theta'}) = 1 - \exp\left[-\frac{|\theta - \theta'|^2}{8\sigma^2}\right] \leq \frac{|\theta - \theta'|^2}{8\sigma^2} \leq \frac{k|\theta - \theta'|_\infty^2}{8\sigma^2}.$$

Then, inequality (71) is satisfied with  $B_k = k/(8\sigma^2)$  and  $s = 2$ . In particular, our condition  $RB_k^{1/s} \geq 1$  is equivalent to  $R \geq 2\sigma\sqrt{(2/k)}$ . In this case, the value of  $r$  given by (72) is of order

$$\inf_{m \subset \{1, \dots, k\}} \left[ \inf_{\theta \in \Theta_m(R)} \ell(\bar{P}^*, P_\theta) + \frac{|m| \log(knR/\sigma) + \xi}{n} \right].$$

More generally, if  $\mathcal{M} = \{P_\theta, \theta \in \mathbb{R}^k\}$  is a regular statistical model with a nonsingular Fisher information matrix  $\mathbf{J}(\theta)$  for all  $\theta \in \mathbb{R}^k$ , we know from the book of Ibragimov and Has'minskiĭ [20, Theorem 7.1, p. 81] that for all  $\theta, \theta' \in \mathbb{R}^k$  such that  $\theta, \theta' \in [-R, R]^k$

$$h^2(P_\theta, P_{\theta'}) \leq \frac{|\theta - \theta'|^2}{8} \sup_{\theta'' \in \mathbb{R}^k, |\theta''|_\infty \leq R} \text{tr}(\mathbf{J}(\theta'')).$$

Then, Assumption (71) holds with  $s = 2$  and we may take

$$B_k = \frac{k^2}{8} \sup_{\theta'' \in \mathbb{R}^k, |\theta''|_\infty \leq R} \varrho(\mathbf{J}(\theta''))$$

where  $\varrho(\mathbf{J}(\theta''))$  denotes the largest eigenvalue of the matrix  $\mathbf{J}(\theta'')$ . This value is independent of  $\theta''$  when  $\mathcal{M}$  is a translation model.

Finally note that the second term in (72) only increases logarithmically with respect to  $R$ , at least when  $B_k = B_k(R)$  does not increase faster than a power of  $R$ . By taking larger values of  $R$  one may therefore considerably enlarge the sizes of the cubes  $\Theta_m(R)$ , and therefore diminish the approximation term in (72), while only slightly increasing the second term  $[|m| \log(2kR(nB_k)^{1/s}) + \xi]/n$ .

### 8 Some tools for evaluating $r_n(\beta, P)$

The aim of this section is to provide some mathematical results that allow one to bound the quantity  $r_n(\beta, P)$  from above, or at least evaluate its order of magnitude, when  $n$  is sufficiently large. Throughout this section, we consider a parametric statistical model

$\mathcal{M} = \{P_\theta, \theta \in \Theta\}$  where the parameter space  $\Theta \subset \mathbb{R}^k$  is endowed with a prior  $\nu$  which admits a density  $q$  with respect to the Lebesgue measure on  $\mathbb{R}^k$ . In order to use the definition (11) of the quantity  $r_n(\beta, P)$ , we assume that we have at disposal a family  $\mathcal{T}(\ell, \mathcal{M})$  that satisfies our Assumption 3, which provides us with a value of  $a_1 > 0$ , as well as a value  $\gamma$  that satisfy the requirements of our main theorems. Our aim is to bound  $r_n(\beta, P)$  as a function of  $a_1, \gamma, \beta, k$  and  $n$  under suitable assumptions on the density  $q$  and the behaviour of the loss  $\ell$ . Once  $\ell$  and  $\mathcal{T}(\ell, \mathcal{M})$  are given,  $a_1$  and  $\gamma$  can be considered as fixed numerical constants. The value of  $\beta$  can also be considered as a numerical constant when Theorem 2 applies. Otherwise, it can be chosen of order  $\sqrt{k/n}$  as in our Example 1.

### 8.1 Bounding $r_n(\beta, P_\theta)$ in parametric models

In what follows,  $|\cdot|_*$  denotes some arbitrary norm on  $\mathbb{R}^k$  and  $\mathcal{B}_*(x, z)$  the corresponding closed ball centered at  $x \in \mathbb{R}^k$  with radius  $z \geq 0$ .

**Assumption 8** Let  $\theta^*$  be an element of  $\Theta \subset \mathbb{R}^k$ .

(i) There exist positive numbers  $\underline{a}, \bar{a}$  and  $s$  such that

$$\underline{a} |\theta - \theta^*|_*^s \leq \ell(\theta, \theta^*) \leq \bar{a} |\theta - \theta^*|_*^s \quad \text{for all } \theta \in \Theta. \tag{73}$$

(ii) There exists a positive nonincreasing function  $\nu_\theta$  on  $\mathbb{R}_+$  such that

$$\nu(\mathcal{B}_*(\theta^*, 2x)) \leq \nu_\theta(x) \nu(\mathcal{B}_*(\theta^*, x)) \quad \text{for all } x > 0. \tag{74}$$

Under Assumption 8-(i), the loss function behaves like a power of a norm between the parameters.

The following result is an extension of Proposition 10 in Baraud and Birgé [6]. It was established there for the special case of the squared Hellinger loss and we provide here an extension to an arbitrary one. Since the proof follows the same lines, we omit it.

**Proposition 7** Under Assumption 8,

$$r_n(\beta, P_{\theta^*}) \leq \inf \left\{ r \geq \frac{1}{n\beta a_1}, r \geq \frac{\varrho_0 \log [\nu_{\theta^*} ([r/\bar{a}]^{1/s})]}{\gamma n \beta a_1} \right\} \tag{75}$$

with  $\varrho_0 = 1 + \log(2\bar{a}/\underline{a})/[s \log 2]$ . If  $\nu_{\theta^*} \equiv \nu > 0$ , then

$$r_n(\beta, P_{\theta^*}) \leq \frac{(\varrho_0 \log \nu) \vee 1}{a_1 n \gamma \beta}. \tag{76}$$

If Assumption 8-(i) is satisfied and if the parameter space  $\Theta$  is convex and  $q$  satisfies

$$\underline{b} \leq q(\theta) \leq \bar{b} \quad \text{for all } \theta \in \Theta \quad \text{with } 0 < \underline{b} \leq \bar{b}, \tag{77}$$

then Assumption 8-(ii) holds with  $v_{\theta^*} \equiv 2^k(\bar{b}/b)$ . Consequently,

$$r_n(\beta, P_{\theta^*}) \leq \frac{\varrho_1}{a_1\gamma} \frac{k}{n\beta} \quad \text{with} \quad \varrho_1 = \left[ \varrho_0 \log \left( 2 \left[ \bar{b}/b \right]^{1/k} \right) \right] \vee 1. \tag{78}$$

When Assumption 8-(i) is satisfied and  $\nu$  admits a density which is bounded away from 0 and infinity on a convex parameter space  $\Theta \subset \mathbb{R}^k$ ,  $r_n(\beta, P_\theta)$  is of order  $k/(n\beta)$  for all  $\theta \in \Theta$ . This result may also hold true when the density is not bounded away from infinity as shown in the following example. If  $k = 1$ ,  $\Theta = [-1, 1]$  and  $q : \theta \mapsto (t/2)|\theta|^{t-1} \mathbb{1}_{[-1,1]}(\theta)$  with  $t \in (0, 1)$ , Assumption 8-(ii) holds with  $v_\theta \equiv 2^{1+t} (2^t - 1)^{-1}$  for all  $\theta \in [-1, 1]$ —see Baraud and Birgé [6, Proposition 11]. Then (76) still applies. In the other direction, when the density  $q$  takes very small values in the neighbourhood of the parameter  $\theta$ , the function  $v_\theta$  may take large values around 0. This is for example the case when  $q$  is proportional to  $\theta \mapsto \exp[-1/(2|\theta|^t)] \mathbb{1}_{[-1,1]}(\theta)$ ,  $t > 0$ , and  $\theta = 0$ . It follows from Baraud and Birgé [6, Proposition 12] (and its proof) that Assumption 8-(ii) is satisfied with  $v_\theta : x \mapsto \exp(c(t)/x^t)$  for some quantity  $c(t) > 0$ . Applying (75) leads to an upper bound on  $r_n(\beta, P_\theta)$  of order  $(n\beta)^{-s/(s+t)}$ .

### 8.2 Some asymptotic order of magnitude

In Sect. 8.2, we have given some general tools for controlling the quantity  $r_n(\beta, P_\theta)$  for a given value of  $n$ . In this section, we present some sufficient conditions under which  $r_n(\beta, P_\theta)$  is of order  $k/(n\beta)$  at least when  $n$  is large enough. These conditions are not the weakest possible ones but they have the advantage to be relatively easy to check on many examples.

**Assumption 9** The density  $q$  is continuous and positive at  $\theta^* \in \Theta$ . The loss function  $\ell$  satisfies the following properties for some positive number  $s > 0$  and a norm  $|\cdot|_*$  on  $\mathbb{R}^k$ .

(i) For all  $\varepsilon > 0$ , there exists  $z = z(\varepsilon) > 0$  such that

$$(1 - \varepsilon) |\theta - \theta^*|_*^s \leq \ell(\theta, \theta^*) \leq (1 + \varepsilon) |\theta - \theta^*|_*^s \quad \text{for all } \theta \in \mathcal{B}_*(\theta^*, z).$$

(ii) There exists a subset  $\mathcal{K} \subset \Theta$ , the interior of which contains  $\theta^*$ , that satisfies for some positive numbers  $\underline{a}_\mathcal{K}$  and  $\eta$ :

$$\underline{a}_\mathcal{K} |\theta - \theta^*|_*^s \leq \ell(\theta, \theta^*) \quad \text{for } \theta \in \mathcal{K} \text{ and for } \theta \notin \mathcal{K} \quad \ell(\theta, \theta^*) \geq \eta > 0. \tag{79}$$

Under these assumptions, we establish the following proposition, the proof of which is postponed to Sect. 10.10.

**Proposition 8** Under Assumption 9, at least for  $n$  sufficiently large,

$$r_n(\beta, P_{\theta^*}) \leq \frac{(1 + 1/s) k}{a_1\gamma} \frac{1}{n\beta}. \tag{80}$$

### 8.3 The case of the squared Hellinger loss on a regular statistical model

Of particular interest is the situation where the statistical model  $\mathcal{M} = \{P_\theta, \theta \in \Theta\}$ ,  $\Theta \subset \mathbb{R}^k$ , is regular. There exist several ways of defining a regular model in statistics and we adopt here the definition of Ibragimov and Has'minskiĭ [20].

**Definition 1** Let  $\mu$  be a measure on  $(E, \mathcal{E})$  and  $\Theta$  an open subset of  $\mathbb{R}^k$ . The statistical model  $\mathcal{M} = \{P_\theta = p_\theta \cdot \mu, \theta \in \Theta\}$  is said to be *regular* if the family of functions  $\{\zeta_\theta = \sqrt{p_\theta}, \theta \in \Theta\} \subset \mathcal{L}_2(E, \mathcal{E}, \mu)$  satisfies the following properties.

- (i) For  $\mu$ -almost all  $x \in E$ ,  $\theta \mapsto \zeta_\theta(x)$  is continuous.
- (ii) For all  $\theta \in \Theta$ , there exists  $\dot{\zeta}_\theta = (\dot{\zeta}_{\theta,1}, \dots, \dot{\zeta}_{\theta,k}) : E \rightarrow \mathbb{R}^k$  such that

$$\int_E |\dot{\zeta}_\theta(x)|^2 d\mu(x) < +\infty$$

and

$$\int_E |\zeta_{\theta+\epsilon}(x) - \zeta_\theta(x) - \langle \dot{\zeta}_\theta(x), \epsilon \rangle|^2 d\mu(x) = o(|\epsilon|^2) \text{ when } |\epsilon| \rightarrow 0.$$

- (iii) For all  $i \in \{1, \dots, k\}$ , the mapping  $\theta \mapsto \dot{\zeta}_{\theta,i}$  is continuous in  $\mathcal{L}_2(E, \mathcal{E}, \mu)$ .

When the model is regular, the matrix

$$\mathbf{J}(\theta) = \left( 4 \int_E \dot{\zeta}_{\theta,i}(x) \dot{\zeta}_{\theta,j}(x) d\mu(x) \right)_{\substack{1 \leq i \leq k \\ 1 \leq j \leq k}},$$

is called the *Fisher information matrix*.

The matrix  $\mathbf{J}(\theta)$  is symmetric and nonnegative and we may therefore consider its square root  $\mathbf{J}^{1/2}(\theta)$ , that is, the symmetric  $(k \times k)$ -nonnegative matrix that satisfies  $\mathbf{J}^{1/2}(\theta)\mathbf{J}^{1/2}(\theta) = \mathbf{J}(\theta)$ .

Regular statistical models enjoy nice metric properties that are described in Proposition 9 below. For a proof we refer the reader to Ibragimov and Has'minskiĭ [20]—Lemma 7.1 page 65, Theorem 7.6 page 81 and its proof.

**Proposition 9** Let  $\Theta$  be an open subset of  $\mathbb{R}^k$  and  $\theta^* \in \Theta$ . If  $\mathcal{M} = \{P_\theta = p_\theta \cdot \mu, \theta \in \Theta\}$  is regular and the Fisher information matrix  $\mathbf{J}(\theta^*)$  nonsingular at  $\theta^* \in \Theta$ , Assumption 9-(i) is satisfied with  $\ell = h^2$ ,  $s = 2$  and for the norm  $|\cdot|_*$  defined by

$$|\mathbf{x}|_* = \frac{1}{\sqrt{8}} \left| \mathbf{J}^{1/2}(\theta^*)\mathbf{x} \right| \text{ for all } \mathbf{x} \in \mathbb{R}^k. \tag{81}$$

Besides, for any compact subset  $\mathcal{K} \subset \Theta$  there exist positive numbers  $\bar{a}_\mathcal{K}, \underline{a}_\mathcal{K}$  such that

$$\underline{a}_\mathcal{K} |\theta - \theta^*|_*^2 \leq h^2(\theta, \theta^*) \leq \bar{a}_\mathcal{K} |\theta - \theta^*|_*^2 \text{ for all } \theta \in \mathcal{K}. \tag{82}$$

Using Proposition 8, we immediately infer the following result.

**Corollary 4** *Let  $\Theta$  be an open subset of  $\mathbb{R}^k$ . Assume that  $\mathcal{M} = \{P_\theta = p_\theta \cdot \mu, \theta \in \Theta\}$  is regular and the Fisher information matrix  $\mathbf{J}(\theta^*)$  nonsingular at  $\theta^* \in \Theta \subset \mathbb{R}^k$ . Assume that there exists a compact set  $\mathcal{K} \subset \Theta$ , containing  $\theta^*$  in its interior, such that  $h(\theta, \theta^*) \geq \eta > 0$  for all  $\theta \notin \mathcal{K}$ . Assume furthermore that the density  $q$  is continuous and positive at  $\theta^*$ . Then,  $r_n(\beta, P_{\theta^*}) \leq [3/(2a_1\gamma\beta)](k/n)$ , at least for  $n$  sufficiently large.*

### 9 Proofs of Theorems 1, 2 and 3

Throughout this proof we fix some  $\bar{Q} \in \mathcal{M}, r, \beta > 0$  and use the following notation:  $c_1 = 1 + c, c_2 = 2 + c,$

$$\mathcal{V}(\pi, \bar{Q}) = \{r > 0, \pi(\mathcal{B}(\bar{Q}, r)) > 0\}$$

and for  $r \in \mathcal{V}(\pi, \bar{Q}), \mathcal{B} = \mathcal{B}(\bar{Q}, r)$  and  $\pi_{\mathcal{B}} = [\pi(\mathcal{B})]^{-1} \mathbb{1}_{\mathcal{B}} \cdot \pi.$

#### 9.1 Main parts of the proofs of Theorems 1 and 2

Throughout the proofs of these two theorems we fix some positive number  $z,$  that will be chosen later on,  $r \geq r_n(\beta, \bar{Q})$  and set

$$A = \left\{ \int_{\mathcal{M}} \exp[-\beta \mathbf{T}(X, P)] d\pi(P) > z \right\}.$$

It follows from the definition (7) of  $\hat{\pi}_X$  that for all  $J \in \mathbb{N}$

$$\begin{aligned} \mathbb{E} \left[ \hat{\pi}_X \left( {}^c\mathcal{B}(\bar{Q}, 2^J r) \right) \right] &= \mathbb{E} \left[ \hat{\pi}_X \left( {}^c\mathcal{B}(\bar{Q}, 2^J r) \right) \mathbb{1}_{\mathcal{C}_A} \right] + \mathbb{E} \left[ \hat{\pi}_X \left( {}^c\mathcal{B}(\bar{Q}, 2^J r) \right) \mathbb{1}_A \right] \\ &\leq \mathbb{P}(\mathcal{C}_A) + \frac{1}{z} \mathbb{E} \left[ \int_{{}^c\mathcal{B}(\bar{Q}, 2^J r)} \exp[-\beta \mathbf{T}(X, P)] d\pi(P) \right] \\ &= \mathbb{P}(\mathcal{C}_A) + \frac{1}{z} \int_{{}^c\mathcal{B}(\bar{Q}, 2^J r)} \mathbb{E} \left[ \exp[-\beta \mathbf{T}(X, P)] \right] d\pi(P). \end{aligned} \tag{83}$$

In a first step, we prove that for some well chosen values of  $\beta, z, r$  and for  $J$  large enough, each of the two terms in the right-hand side of (83) is not larger than  $e^{-\xi}.$  To achieve this goal, we bound the first term of the right-hand side of (83) by applying Markov’s inequality

$$\begin{aligned} \mathbb{P}(\mathcal{C}_A) &= \mathbb{P} \left[ \int_{\mathcal{M}} \exp[-\beta \mathbf{T}(X, P)] d\pi(P) \leq z \right] \\ &= \mathbb{P} \left[ \left[ \int_{\mathcal{M}} \exp[-\beta \mathbf{T}(X, P)] d\pi(P) \right]^{-1} \geq z^{-1} \right] \end{aligned}$$



$$\leq z \mathbb{E} \left[ \frac{1}{\int_{\mathcal{M}} \exp[-\beta \mathbf{T}(X, P)] d\pi(P)} \right]$$

and then by using Lemma 6, we obtain that

$$\mathbb{P}({}^cA) \leq \frac{z}{\pi^2(\mathcal{B})} \left[ \int_{\mathcal{B}^2} \exp[-\mathbf{L}(P, Q)] d\pi_{\mathcal{B}}(P) d\pi_{\mathcal{B}}(Q) \right]^{-1}. \tag{84}$$

We therefore have a control of  $\mathbb{P}({}^cA)$  by choosing  $z$  small enough. We bound the second term of (83) by using Lemma 5.

We then finish the proofs of Theorems 1 and 2 as follows. In the context of Theorem 1, we finally establish that for a suitable value of  $J$  and all  $\bar{Q} \in \mathcal{M}(\beta)$ ,

$$\mathbb{E} \left[ \hat{\pi}_X \left( {}^c\mathcal{B}(\bar{Q}, 2^J r) \right) \right] \leq 2e^{-\xi} \quad \text{with} \quad r = r(\bar{Q}) = \ell(\bar{P}^*, \bar{Q}) + a_1^{-1} \left( \beta + \frac{2\xi}{n\beta} \right).$$

By (3),  $\mathcal{B}(\bar{Q}, 2^J r) \subset \mathcal{B}(\bar{P}^*, \tau \ell(\bar{P}^*, \bar{Q}) + \tau 2^J r)$  for all  $\bar{Q} \in \mathcal{M}(\beta)$ , and consequently  $\mathbb{E} \left[ \hat{\pi}_X \left( {}^c\mathcal{B}(\bar{P}^*, \bar{r}) \right) \right] \leq 2e^{-\xi}$  with

$$\bar{r} = \bar{r}(\bar{Q}) = \tau \left[ \ell(\bar{P}^*, \bar{Q}) + 2^J r \right] = \tau \left[ (1 + 2^J) \ell(\bar{P}^*, \bar{Q}) + 2^J a_1^{-1} \left( \beta + \frac{2\xi}{n\beta} \right) \right].$$

We obtain (16) by monotone convergence, taking a sequence  $(\bar{Q}_N)_{N \geq 0} \subset \mathcal{M}(\beta)$  such that  $\ell(\bar{P}^*, \bar{Q}_N)$  is nonincreasing to  $\inf_{P \in \mathcal{M}(\beta)} \ell(\bar{P}^*, P)$ , so that

$$\begin{aligned} \lim_{N \rightarrow +\infty} \bar{r}(\bar{Q}_N) &= \tau \left[ (1 + 2^J) \inf_{\bar{Q} \in \mathcal{M}(\beta)} \ell(\bar{P}^*, \bar{Q}) + 2^J a_1^{-1} \left( \beta + \frac{2\xi}{n\beta} \right) \right] \\ &\leq \tau(1 + 2^J) \left[ \inf_{\bar{Q} \in \mathcal{M}(\beta)} \ell(\bar{P}^*, \bar{Q}) + a_1^{-1} \left( \beta + \frac{2\xi}{n\beta} \right) \right] \end{aligned}$$

and (16) holds provided that  $\kappa_0 \geq \tau(2^J + 1)$ .

In the context of Theorem 2, we show that for some suitable value of  $J$  and all  $\bar{Q} \in \mathcal{M}$ ,

$$\mathbb{E} \left[ \hat{\pi}_X \left( {}^c\mathcal{B}(\bar{Q}, 2^J r) \right) \right] \leq 2e^{-\xi} \quad \text{with} \quad r = \ell(\bar{P}^*, \bar{Q}) + r_n(\bar{Q}, \beta) + \frac{2\xi}{n\beta a_1},$$

and we get (28) by arguing similarly.

### 9.2 Preliminary results

In the proofs of Theorems 1 and 2, we use the following consequence of our Assumption 3. We may write

$$\frac{1}{n} \sum_{i=1}^n \mathbb{E} [t_{(P,Q)}(X_i)] = \mathbb{E}_S [t_{(P,Q)}(X)] \quad \text{with } S = \bar{P}^* = \frac{1}{n} \sum_{i=1}^n P_i^* \in \mathcal{P}$$

and we deduce from (5) that for all  $P, Q \in \mathcal{M}$ ,

$$\frac{1}{n} \sum_{i=1}^n \mathbb{E} [t_{(P,Q)}(X_i)] \leq a_0 \ell(\bar{P}^*, P) - a_1 \ell(\bar{P}^*, Q). \tag{85}$$

Besides, using the antisymmetry property (ii) we also obtain that

$$\frac{1}{n} \sum_{i=1}^n \mathbb{E} [t_{(P,Q)}(X_i)] \geq a_1 \ell(\bar{P}^*, P) - a_0 \ell(\bar{P}^*, Q). \tag{86}$$

For the proof of Theorems 2, we additionally use the following consequence of our Assumption 4. By taking  $S = \bar{P}^*$  and using the convexity of the mapping  $u \mapsto u^2$ , we deduce that for all  $P, Q \in \mathcal{M}$

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \text{Var} [t_{(P,Q)}(X_i)] &= \mathbb{E}_S [t_{(P,Q)}^2(X)] - \frac{1}{n} \sum_{i=1}^n (\mathbb{E} [t_{(P,Q)}(X_i)])^2 \\ &\leq \mathbb{E}_S [t_{(P,Q)}^2(X)] - (\mathbb{E}_S [t_{(P,Q)}(X)])^2 \\ &= \text{Var}_S [t_{(P,Q)}(X)] \end{aligned}$$

and it follows then from Assumption 4 (iv) that for all  $P, Q \in \mathcal{M}$

$$\frac{1}{n} \sum_{i=1}^n \text{Var} [t_{(P,Q)}(X_i)] \leq a_2 [\ell(\bar{P}^*, P) + \ell(\bar{P}^*, Q)]. \tag{87}$$

The proofs of our main results rely on the following lemmas.

**Lemma 4** *Let  $(U, V)$  be a pair of random variables with values in a product space  $(E \times F, \mathcal{E} \otimes \mathcal{F})$  and marginal distributions  $P_U$  and  $P_V$  respectively. For all measurable function  $h$  on  $(E \times F, \mathcal{E} \otimes \mathcal{F})$ ,*

$$\mathbb{E}_U \left[ \frac{1}{\mathbb{E}_V [\exp [-h(U, V)]]} \right] \leq \left[ \mathbb{E}_V \left[ \frac{1}{\mathbb{E}_U [\exp [h(U, V)]]} \right] \right]^{-1}.$$

This lemma is proven in Audibert and Catoni [3, Lemma 4.2, p. 28].

**Lemma 5** For  $P, Q \in \mathcal{M}$ , we set

$$\mathbf{M}(P, Q) = \log \left[ \int_{\mathcal{M}} \mathbb{E} \left[ \exp \left[ \beta \left( c\mathbf{T}(X, P, Q') - c_1\mathbf{T}(X, P, Q) \right) \right] d\pi(Q') \right] \right].$$

For all  $r \in \mathcal{V}(\pi, \overline{Q})$  and  $P \in \mathcal{M}$ ,

$$\mathbb{E} \left[ \exp \left[ -\beta\mathbf{T}(X, P) \right] \right] \leq \frac{1}{\pi(\mathcal{B})} \left[ \int_{\mathcal{B}} \exp \left[ -\mathbf{M}(P, Q) \right] d\pi_{\mathcal{B}}(Q) \right]^{-1}. \tag{88}$$

**Proof** Let  $r \in \mathcal{V}(\pi, \overline{Q})$ . For  $P, Q \in \mathcal{M}$ , we set

$$I(X, P, Q) = c_1\beta\mathbf{T}(X, P, Q) - \log \int_{\mathcal{M}} \exp \left[ c\beta\mathbf{T}(X, P, Q') \right] d\pi(Q').$$

Then,

$$\begin{aligned} & \mathbb{E} \left[ \exp \left[ -I(X, P, Q) \right] \right] \\ &= \mathbb{E} \left[ \exp \left[ -c_1\beta\mathbf{T}(X, P, Q) + \log \int_{\mathcal{M}} \exp \left[ c\beta\mathbf{T}(X, P, Q') \right] d\pi(Q') \right] \right] \\ &= \mathbb{E} \left[ \int_{\mathcal{M}} \exp \left[ c\beta\mathbf{T}(X, P, Q') - c_1\beta\mathbf{T}(X, P, Q) \right] d\pi(Q') \right] \\ &= \exp \left[ \mathbf{M}(P, Q) \right]. \end{aligned} \tag{89}$$

Since  $\lambda = c_1\beta = (1 + c)\beta$ , it follows from the convexity of the exponential that

$$\begin{aligned} \mathbb{E} \left[ \exp \left[ -\beta\mathbf{T}(X, P) \right] \right] &= \mathbb{E} \left[ \exp \left[ \int_{\mathcal{M}} [-\beta\mathbf{T}(X, P, Q)] d\tilde{\pi}_X(Q|P) \right] \right] \\ &\leq \mathbb{E} \left[ \int_{\mathcal{M}} \exp \left[ -\beta\mathbf{T}(X, P, Q) \right] d\tilde{\pi}_X(Q|P) \right] \\ &= \mathbb{E} \left[ \frac{\int_{\mathcal{M}} \exp \left[ c\beta\mathbf{T}(X, P, Q) \right] d\pi(Q)}{\int_{\mathcal{M}} \exp \left[ c_1\beta\mathbf{T}(X, P, Q) \right] d\pi(Q)} \right] \\ &\leq \mathbb{E} \left[ \frac{\int_{\mathcal{M}} \exp \left[ c\beta\mathbf{T}(X, P, Q) \right] d\pi(Q)}{\int_{\mathcal{B}} \exp \left[ c_1\beta\mathbf{T}(X, P, Q) \right] d\pi(Q)} \right]. \end{aligned}$$

Hence,

$$\begin{aligned} \mathbb{E} \left[ \exp \left[ -\beta\mathbf{T}(X, P) \right] \right] &\leq \mathbb{E} \left[ \frac{1}{\int_{\mathcal{B}} \exp \left[ I(X, P, Q) \right] d\pi(Q)} \right] \\ &= \frac{1}{\pi(\mathcal{B})} \mathbb{E} \left[ \frac{1}{\int_{\mathcal{B}} \exp \left[ I(X, P, Q) \right] d\pi_{\mathcal{B}}(Q)} \right]. \end{aligned}$$

Applying Lemma 4 with  $U = X, V = Q$  with distribution  $\pi_{\mathcal{B}}$ , and  $h(U, V) = -I(X, P, Q)$ , we obtain that

$$\mathbb{E}[\exp[-\beta \mathbf{T}(X, P)]] \leq \frac{1}{\pi(\mathcal{B})} \left[ \int_{\mathcal{B}} \frac{1}{\mathbb{E}[\exp[-I(X, P, Q)]]} d\pi_{\mathcal{B}}(Q) \right]^{-1}$$

and (88) follows from (89). □

**Lemma 6** For  $P, Q \in \mathcal{M}$ , we set

$$\mathbf{L}(P, Q) = \log \int_{\mathcal{M}} \mathbb{E}[\exp[\beta (c_2 \mathbf{T}(X, P, Q') - c_1 \mathbf{T}(X, P, Q))] ] d\pi(Q').$$

For all  $r \in \mathcal{V}(\pi, \overline{Q})$ ,

$$\begin{aligned} & \mathbb{E} \left[ \frac{1}{\int_{\mathcal{M}} \exp[-\beta \mathbf{T}(X, P)] d\pi(P)} \right] \\ & \leq \frac{1}{\pi^2(\mathcal{B})} \left[ \int_{\mathcal{B}^2} \exp[-\mathbf{L}(P, Q)] d\pi_{\mathcal{B}}(P) d\pi_{\mathcal{B}}(Q) \right]^{-1}. \end{aligned}$$

**Proof** For  $P, Q \in \mathcal{M}$ , we set

$$H(X, P, Q) = \beta c_1 \mathbf{T}(X, P, Q) - \log \left[ \int_{\mathcal{M}} \exp [c_2 \beta \mathbf{T}(X, P, Q')] d\pi(Q') \right].$$

Then,

$$\begin{aligned} & \mathbb{E}[\exp[-H(X, P, Q)]] \\ & = \mathbb{E} \left[ \exp[-\beta c_1 \mathbf{T}(X, P, Q)] \int_{\mathcal{M}} \exp [c_2 \beta \mathbf{T}(X, P, Q')] d\pi(Q') \right] \\ & = \mathbb{E} \left[ \int_{\mathcal{M}} \exp [\beta (c_2 \mathbf{T}(X, P, Q') - c_1 \mathbf{T}(X, P, Q))] d\pi(Q') \right] \\ & = \exp [\mathbf{L}(P, Q)]. \end{aligned} \tag{90}$$

It follows from the convexity of the exponential and the fact that  $\lambda = c_1 \beta$  that for all  $P \in \mathcal{M}$ ,

$$\begin{aligned} \mathbb{E}[\exp[\beta \mathbf{T}(X, P)]] & = \mathbb{E} \left[ \exp \left[ \int_{\mathcal{M}} [\beta \mathbf{T}(X, P, Q)] d\tilde{\pi}_X(Q|P) \right] \right] \\ & \leq \mathbb{E} \left[ \int_{\mathcal{M}} \exp [\beta \mathbf{T}(X, P, Q)] d\tilde{\pi}_X(Q|P) \right] \\ & = \mathbb{E} \left[ \frac{\int_{\mathcal{M}} \exp [c_2 \beta \mathbf{T}(X, P, Q)] d\pi(Q)}{\int_{\mathcal{M}} \exp [c_1 \beta \mathbf{T}(X, P, Q)] d\pi(Q)} \right] \end{aligned}$$

$$= \mathbb{E} \left[ \frac{1}{\int_{\mathcal{M}} \exp [H(\mathbf{X}, P, Q)] d\pi(Q)} \right].$$

Applying Lemma 4 with  $U = \mathbf{X}$ ,  $V = Q$  with distribution  $\pi$ , and  $h(U, V) = -H(\mathbf{X}, P, Q)$  we obtain that

$$\mathbb{E} [\exp [\beta \mathbf{T}(\mathbf{X}, P)]] \leq \left[ \int_{\mathcal{M}} \frac{1}{\mathbb{E} [\exp [-H(\mathbf{X}, P, Q)]]} d\pi(Q) \right]^{-1}.$$

We deduce from (90) that for all  $P \in \mathcal{M}$

$$\begin{aligned} \mathbb{E} [\exp [\beta \mathbf{T}(\mathbf{X}, P)]] &\leq \left[ \int_{\mathcal{M}} \exp [-\mathbf{L}(P, Q)] d\pi(Q) \right]^{-1} \\ &\leq \frac{1}{\pi(\mathcal{B})} \left[ \int_{\mathcal{B}} \exp [-\mathbf{L}(P, Q)] d\pi_{\mathcal{B}}(Q) \right]^{-1}. \end{aligned} \tag{91}$$

Applying Lemma 4 with  $U = \mathbf{X}$ ,  $V = P$  with distribution  $\pi$  and  $h(U, V) = \beta \mathbf{T}(\mathbf{X}, P)$ , gives

$$\begin{aligned} \mathbb{E} \left[ \frac{1}{\int_{\mathcal{M}} \exp [-\beta \mathbf{T}(\mathbf{X}, P)] d\pi(P)} \right] &\leq \left[ \int_{\mathcal{M}} \frac{1}{\mathbb{E} [\exp [\beta \mathbf{T}(\mathbf{X}, P)]]} d\pi(P) \right]^{-1} \\ &\leq \frac{1}{\pi(\mathcal{B})} \left[ \int_{\mathcal{B}} \frac{1}{\mathbb{E} [\exp [\beta \mathbf{T}(\mathbf{X}, P)]]} d\pi_{\mathcal{B}}(P) \right]^{-1} \end{aligned}$$

which together with (91) leads to the result. □

The proofs of Theorems 1 and 2 rely on suitable bounds on the Laplace transforms of sums of independent random variables and on a summation lemma. These results are presented below.

**Lemma 7** *For all  $\beta \in \mathbb{R}$  and random variable  $U$  with values in an interval of length  $l \in (0, +\infty)$ ,*

$$\log \mathbb{E} [\exp [\beta U]] \leq \beta \mathbb{E} [U] + \frac{\beta^2 l^2}{8}. \tag{92}$$

**Lemma 8** *Let  $U$  be a squared integrable random variable not larger than  $b > 0$ . For all  $\beta > 0$ ,*

$$\log \mathbb{E} [\exp [\beta U]] \leq \beta \mathbb{E} [U] + \beta^2 \mathbb{E} [U^2] \frac{\phi(\beta b)}{2}, \tag{93}$$

where  $\phi$  is defined by (24).

The proofs of Lemmas 7 and 8 can be found on pages 21 and 23 in Massart [23] (where our function  $\phi$  is defined as twice his).

**Lemma 9** *Let  $J \in \mathbb{N}$ ,  $\gamma > 0$  and  $\bar{Q} \in \mathcal{M}$ . If  $r$  satisfies  $n\beta a_1 r \geq 1$  and (12), for all  $\gamma_0 > 2\gamma$*

$$\int_{\mathcal{C}_{\mathcal{B}}(\bar{Q}, 2^J r)} \exp[-\gamma_0 n\beta a_1 \ell(\bar{Q}, P)] d\pi(P) \leq \pi(\mathcal{B}) \exp[\Xi - (\gamma_0 - 2\gamma) n\beta a_1 2^J r] \tag{94}$$

with

$$\Xi = -\gamma + \log \left[ \frac{1}{1 - \exp[-(\gamma_0 - 2\gamma)]} \right]$$

Besides,

$$\int_{\mathcal{M}} \exp[-\gamma_0 n\beta a_1 \ell(\bar{Q}, P)] d\pi(P) \leq \pi(\mathcal{B}) \exp[\Xi'] \tag{95}$$

with

$$\Xi' = \log \left[ 1 + \frac{\exp[-(\gamma_0 - \gamma)]}{1 - \exp[-(\gamma_0 - 2\gamma)]} \right].$$

**Proof** From (12), we deduce by induction that for all  $j \geq 0$

$$\begin{aligned} \pi(\mathcal{B}(\bar{Q}, 2^{j+1}r)) &\leq \exp \left[ \gamma n\beta a_1 r \sum_{k=0}^j 2^k \right] \pi(\mathcal{B}) \\ &= \exp \left[ (2^{j+1} - 1)\gamma n\beta a_1 r \right] \pi(\mathcal{B}) \end{aligned}$$

Consequently,

$$\begin{aligned} &\int_{\mathcal{C}_{\mathcal{B}}(\bar{Q}, 2^J r)} \exp[-\gamma_0 n\beta a_1 \ell(\bar{Q}, P)] d\pi(P) \\ &= \sum_{j \geq J} \int_{\mathcal{B}(\bar{Q}, 2^{j+1}r) \setminus \mathcal{B}(\bar{Q}, 2^j r)} \exp[-\gamma_0 n\beta a_1 \ell(\bar{Q}, P)] d\pi(P) \\ &\leq \pi(\mathcal{B}) \sum_{j \geq J} \frac{\pi(\mathcal{B}(\bar{Q}, 2^{j+1}r))}{\pi(\mathcal{B})} \exp[-\gamma_0 n\beta a_1 2^j r] \\ &\leq \pi(\mathcal{B}) \sum_{j \geq J} \exp \left[ \gamma n\beta a_1 (2^{j+1} - 1)r - \gamma_0 n\beta a_1 2^j r \right] \end{aligned}$$

$$\begin{aligned}
 &= \pi(\mathcal{B}) \exp[-\gamma n \beta a_1 r] \sum_{j \geq J} \exp[-(\gamma_0 - 2\gamma) n \beta a_1 2^j r] \\
 &= \pi(\mathcal{B}) \exp[-\gamma n \beta a_1 r] \sum_{j \geq 0} \exp[-(\gamma_0 - 2\gamma) n \beta a_1 2^j r].
 \end{aligned}$$

Since  $2^j \geq j + 1$  for all  $j \geq 0$  we obtain that

$$\begin{aligned}
 &\int_{c_{\mathcal{B}}(\bar{Q}, 2^J r)} \exp[-\gamma_0 n \beta a_1 \ell(\bar{Q}, P)] d\pi(P) \\
 &\leq \pi(\mathcal{B}) \exp[-\gamma n \beta a_1 r] \sum_{j \geq 0} \exp[-(\gamma_0 - 2\gamma) n \beta a_1 (j + 1) 2^j r] \\
 &\leq \pi(\mathcal{B}) \exp[-\gamma n \beta a_1 r - (\gamma_0 - 2\gamma) n \beta a_1 2^J r] \sum_{j \geq 0} \exp[-j(\gamma_0 - 2\gamma) n \beta a_1 2^J r] \\
 &= \pi(\mathcal{B}) \frac{\exp[-\gamma n \beta a_1 r]}{1 - \exp[-(\gamma_0 - 2\gamma) n \beta a_1 2^J r]} \exp[-(\gamma_0 - 2\gamma) n \beta a_1 2^J r].
 \end{aligned}$$

which leads to (94) since  $n \beta a_1 2^J r \geq n \beta a_1 r \geq 1$ . Finally, by applying this inequality with  $J = 0$  we obtain that

$$\begin{aligned}
 &\int_{\mathcal{M}} \exp[-\gamma_0 \beta n a_1 \ell(\bar{Q}, P)] d\pi(P) \\
 &= \int_{\mathcal{B}} \exp[-\gamma_0 \beta n a_1 \ell(\bar{Q}, P)] d\pi(P) + \int_{c_{\mathcal{B}}} \exp[-\gamma_0 \beta n a_1 \ell(\bar{Q}, P)] d\pi(P) \\
 &\leq \pi(\mathcal{B}) \left[ 1 + \frac{\exp[-\gamma - (\gamma_0 - 2\gamma) n \beta a_1 r]}{1 - \exp[-(\gamma_0 - 2\gamma)]} \right] \\
 &\leq \pi(\mathcal{B}) \left[ 1 + \frac{\exp[-(\gamma_0 - \gamma)]}{1 - \exp[-(\gamma_0 - 2\gamma)]} \right],
 \end{aligned}$$

which is (95). □

### 9.3 Proof of Theorem 1

For all  $i \in \{1, \dots, n\}$  and  $P, Q, Q' \in \mathcal{M}$ , let us set

$$\begin{aligned}
 U_i &= c(t_{(P, Q')}(X_i) - \mathbb{E}[t_{(P, Q')}(X_i)]) \\
 &\quad - c_1(t_{(P, Q)}(X_i) - \mathbb{E}[t_{(P, Q)}(X_i)])
 \end{aligned} \tag{96}$$

$$\begin{aligned}
 V_i &= c_2(t_{(P, Q')}(X_i) - \mathbb{E}[t_{(P, Q')}(X_i)]) \\
 &\quad - c_1(t_{(P, Q)}(X_i) - \mathbb{E}[t_{(P, Q)}(X_i)]).
 \end{aligned} \tag{97}$$

The random variables  $U_i$  are independent and under Assumption 3-(iv), they takes their values in an interval of length  $l_1 = c + c_1 = 1 + 2c$ . The  $V_i$  are also independent

and they takes their values in an interval of length  $l_2 = c_1 + c_2 = 3 + 2c$ . Applying Lemma 7, we obtain that

$$\prod_{i=1}^n \mathbb{E} [\exp [\beta U_i]] \leq \exp \left[ \frac{l_1^2 n \beta^2}{8} \right] \tag{98}$$

and

$$\prod_{i=1}^n \mathbb{E} [\exp [\beta V_i]] \leq \exp \left[ \frac{l_2^2 n \beta^2}{8} \right]. \tag{99}$$

By using Assumption 2 and the fact that  $c_0 = c_1 - ca_0/a_1 > 0$ ,

$$\begin{aligned} &c \left( a_0 \ell(\bar{P}^*, P) - a_1 \ell(\bar{P}^*, Q') \right) - c_1 \left( a_1 \ell(\bar{P}^*, P) - a_0 \ell(\bar{P}^*, Q) \right) \\ &= - (c_1 a_1 - ca_0) \ell(\bar{P}^*, P) - ca_1 \ell(\bar{P}^*, Q') + c_1 a_0 \ell(\bar{P}^*, Q) \\ &\leq -c_0 a_1 \left[ \tau^{-1} \ell(\bar{Q}, P) - \ell(\bar{P}^*, \bar{Q}) \right] - ca_1 \left[ \tau^{-1} \ell(\bar{Q}, Q') - \ell(\bar{P}^*, \bar{Q}) \right] \\ &\quad + \tau c_1 a_0 \left[ \ell(\bar{P}^*, \bar{Q}) + \ell(\bar{Q}, Q) \right] \\ &= e_0 a_1 \ell(\bar{P}^*, \bar{Q}) - \tau^{-1} c_0 a_1 \ell(\bar{Q}, P) - \tau^{-1} ca_1 \ell(\bar{Q}, Q') + \tau c_1 a_0 \ell(\bar{Q}, Q) \end{aligned} \tag{100}$$

with

$$e_0 = c_0 + c + \frac{\tau c_1 a_0}{a_1}. \tag{101}$$

It follows from (100) and Assumptions 3-(iii), more precisely its consequences (85) and (86), that

$$\begin{aligned} &n^{-1} \{ c \mathbb{E} [\mathbf{T}(X, P, Q')] - c_1 \mathbb{E} [\mathbf{T}(X, P, Q)] \} \\ &\leq c \left[ a_0 \ell(\bar{P}^*, P) - a_1 \ell(\bar{P}^*, Q') \right] - c_1 \left[ a_1 \ell(\bar{P}^*, P) - a_0 \ell(\bar{P}^*, Q) \right] \\ &\leq e_0 a_1 \ell(\bar{P}^*, \bar{Q}) - \tau^{-1} c_0 a_1 \ell(\bar{Q}, P) - \tau^{-1} ca_1 \ell(\bar{Q}, Q') + \tau c_1 a_0 \ell(\bar{Q}, Q). \end{aligned} \tag{102}$$

Since  $a_0 \geq a_1$  and  $c_2 > c_1$ ,  $c'_0 = c_2(a_0/a_1) - c_1 > 0$  and by arguing as above, we obtain similarly that

$$\begin{aligned} &n^{-1} \{ c_2 \mathbb{E} [\mathbf{T}(X, P, Q')] - c_1 \mathbb{E} [\mathbf{T}(X, P, Q)] \} \\ &\leq c_2 \left( a_0 \ell(\bar{P}^*, P) - a_1 \ell(\bar{P}^*, Q') \right) - c_1 \left( a_1 \ell(\bar{P}^*, P) - a_0 \ell(\bar{P}^*, Q) \right) \\ &= c'_0 a_1 \ell(\bar{P}^*, P) - c_2 a_1 \ell(\bar{P}^*, Q') + c_1 a_0 \ell(\bar{P}^*, Q) \\ &\leq \tau c'_0 a_1 \left[ \ell(\bar{P}^*, \bar{Q}) + \ell(\bar{Q}, P) \right] - c_2 a_1 \left[ \tau^{-1} \ell(\bar{Q}, Q') - \ell(\bar{P}^*, \bar{Q}) \right] \\ &\quad + \tau c_1 a_0 \left[ \ell(\bar{P}^*, \bar{Q}) + \ell(\bar{Q}, Q) \right] \end{aligned}$$



$$\begin{aligned} &\leq (e_1 + c_2) a_1 \ell(\overline{P}^*, \overline{Q}) + \tau c'_0 a_1 \ell(\overline{Q}, P) \\ &\quad - \tau^{-1} c_2 a_1 \ell(\overline{Q}, Q') + \tau c_1 a_0 \ell(\overline{Q}, Q), \end{aligned} \tag{103}$$

with

$$e_1 = \tau [c'_0 + c_1 a_0/a_1] = \tau [c_2(a_0/a_1) + c_1 (a_0/a_1 - 1)]. \tag{104}$$

Using (98) and (102), we deduce that for all  $P, Q, Q' \in \mathcal{M}$

$$\begin{aligned} &\mathbb{E} [\exp [\beta (c \mathbf{T}(X, P, Q') - c_1 \mathbf{T}(X, P, Q))] ] \\ &= \prod_{i=1}^n \mathbb{E} [\exp [\beta (c t_{(P, Q')}(X_i) - c_1 t_{(P, Q)}(X_i))] ] \\ &= \exp [\beta (c \mathbb{E} [\mathbf{T}(X, P, Q')] - c_1 \mathbb{E} [\mathbf{T}(X, P, Q)])] \prod_{i=1}^n \mathbb{E} [\exp [\beta U_i]] \\ &\leq \exp [n\beta [\Delta_1(P, Q) - \tau^{-1} c a_1 \ell(\overline{Q}, Q')]] \end{aligned} \tag{105}$$

with

$$\Delta_1(P, Q) = e_0 a_1 \ell(\overline{P}^*, \overline{Q}) + \tau c_1 a_0 \ell(\overline{Q}, Q) + \frac{l_1^2 \beta}{8} - \tau^{-1} c_0 a_1 \ell(\overline{Q}, P). \tag{106}$$

Using (99) and (103), we obtain similarly that for all  $P, Q, Q' \in \mathcal{M}$

$$\begin{aligned} &\mathbb{E} [\exp [\beta (c_2 \mathbf{T}(X, P, Q') - c_1 \mathbf{T}(X, P, Q))] ] \\ &\leq \exp [n\beta [\Delta_2(P, Q) - \tau^{-1} c_2 a_1 \ell(\overline{Q}, Q')]] \end{aligned} \tag{107}$$

with

$$\begin{aligned} \Delta_2(P, Q) &= (e_1 + c_2) a_1 \ell(\overline{P}^*, \overline{Q}) + \tau c'_0 a_1 \ell(\overline{Q}, P) + \tau c_1 a_0 \ell(\overline{Q}, Q) \\ &\quad + \frac{l_2^2 \beta}{8}. \end{aligned} \tag{108}$$

Since  $2\gamma < \tau^{-1}c < \tau^{-1}c_2$ , we may apply Lemma 9 with  $\gamma_0 = \tau^{-1}c$  and  $\gamma_0 = \tau^{-1}c_2$  successively which leads to

$$\int_{\mathcal{M}} \exp [-\tau^{-1} c n \beta a_1 \ell(\overline{Q}, Q')] d\pi(Q') \leq \pi(\mathcal{B}) \exp [\Xi_1] \tag{109}$$

and

$$\int_{\mathcal{M}} \exp [-\tau^{-1} c_2 n \beta a_1 \ell(\overline{Q}, Q')] d\pi(Q') \leq \pi(\mathcal{B}) \exp [\Xi_1] \tag{110}$$

with

$$\begin{aligned} \Xi_1 &= \log \left[ 1 + \frac{\exp[-(\tau^{-1}c - \gamma)]}{1 - \exp[-(\tau^{-1}c - 2\gamma)]} \right] \\ &\geq \log \left[ 1 + \frac{\exp[-(\tau^{-1}c_2 - \gamma)]}{1 - \exp[-(\tau^{-1}c_2 - 2\gamma)]} \right]. \end{aligned} \tag{111}$$

Putting (107) and (110) together leads to

$$\begin{aligned} \exp[\mathbf{L}(P, Q)] &= \int_{\mathcal{M}} \mathbb{E} [\exp[\beta(c_2 \mathbf{T}(X, P, Q') - c_1 \mathbf{T}(X, P, Q))] ] d\pi(Q') \\ &\leq \exp[n\beta\Delta_2(P, Q)] \int_{\mathcal{M}} \exp[-\tau^{-1}c_2 n\beta a_1 \ell(\bar{Q}, Q')] d\pi(Q') \\ &\leq \pi(\mathcal{B}) \exp[\Xi_1 + n\beta\Delta_2(P, Q)], \end{aligned}$$

and since, for all  $(P, Q) \in \mathcal{B}^2$ , by definition (108) of  $\Delta_2(P, Q)$ ,

$$\begin{aligned} \Delta_2(P, Q) &\leq (e_1 + c_2) a_1 \ell(\bar{P}^*, \bar{Q}) + [\tau c'_0 a_1 + \tau c_1 a_0] r + \frac{l_2^2 \beta}{8} \\ &= (e_1 + c_2) a_1 \ell(\bar{P}^*, \bar{Q}) + e_1 a_1 r + \frac{l_2^2 \beta}{8} = \Delta_2 \end{aligned} \tag{112}$$

we derive that

$$\left[ \int_{\mathcal{B}^2} \exp[-\mathbf{L}(P, Q)] d\pi_{\mathcal{B}}(P) d\pi_{\mathcal{B}}(Q) \right]^{-1} \leq \pi(\mathcal{B}) \exp[\Xi_1 + n\beta\Delta_2].$$

We deduce from (84) that

$$\mathbb{P}(^cA) \leq \frac{z}{\pi(\mathcal{B})} \exp[\Xi_1 + n\beta\Delta_2].$$

In particular,  $\mathbb{P}(^cA) \leq e^{-\xi}$  for  $z$  satisfying

$$\log\left(\frac{1}{z}\right) = \xi + \log \frac{1}{\pi(\mathcal{B})} + \Xi_1 + n\beta\Delta_2. \tag{113}$$

Putting (105) and (109) together, we obtain that

$$\begin{aligned} \exp[\mathbf{M}(P, Q)] &= \int_{\mathcal{M}} \mathbb{E} [\exp[\beta(c \mathbf{T}(X, P, Q') - c_1 \mathbf{T}(X, P, Q))] ] d\pi(Q') \\ &\leq \exp[n\beta\Delta_1(P, Q)] \int_{\mathcal{M}} \exp[-\tau^{-1}cn\beta a_1 \ell(\bar{Q}, Q')] d\pi(Q') \end{aligned}$$

$$\leq \pi(\mathcal{B}) \exp[\Xi_1 + n\beta\Delta_1(P, Q)].$$

It follows from the definition (106) of  $\Delta_1(P, Q)$  that for all  $P \in \mathcal{M}$  and for all  $Q \in \mathcal{B}$ ,

$$\Delta_1(P, Q) \leq e_0 a_1 \ell(\bar{P}^*, \bar{Q}) + \tau c_1 a_0 r + \frac{l_1^2 \beta}{8} - \tau^{-1} c_0 a_1 \ell(\bar{Q}, P),$$

and consequently, for all  $P \in \mathcal{M}$  and  $Q \in \mathcal{B}$

$$\begin{aligned} & \exp[\mathbf{M}(P, Q)] \\ & \leq \pi(\mathcal{B}) \exp \left[ \Xi_1 + n\beta \left( e_0 a_1 \ell(\bar{P}^*, \bar{Q}) + \tau c_1 a_0 r + \frac{l_1^2 \beta}{8} - \tau^{-1} c_0 a_1 \ell(\bar{Q}, P) \right) \right]. \end{aligned}$$

We derive from Lemma 5 that

$$\begin{aligned} & \mathbb{E}[\exp[-\beta \mathbf{T}(X, P)]] \\ & \leq \frac{1}{\pi(\mathcal{B})} \left[ \int_{\mathcal{B}} \exp[-\mathbf{M}(P, Q)] d\pi_{\mathcal{B}}(Q) \right]^{-1} \\ & \leq \exp \left[ \Xi_1 + n\beta \left( e_0 a_1 \ell(\bar{P}^*, \bar{Q}) + \tau c_1 a_0 r + \frac{l_1^2 \beta}{8} - \tau^{-1} c_0 a_1 \ell(\bar{Q}, P) \right) \right], \end{aligned}$$

hence,

$$\begin{aligned} & \int_{c_{\mathcal{B}}(\bar{Q}, 2^J r)} \mathbb{E}[\exp[-\beta \mathbf{T}(X, P)]] d\pi(P) \\ & \leq \exp \left[ \Xi_1 + n\beta \left( e_0 a_1 \ell(\bar{P}^*, \bar{Q}) + \tau c_1 a_0 r + \frac{l_1^2 \beta}{8} \right) \right] \\ & \quad \times \int_{c_{\mathcal{B}}(\bar{Q}, 2^J r)} \exp[-\tau^{-1} c_0 n \beta a_1 \ell(\bar{Q}, P)] d\pi(P). \end{aligned} \tag{114}$$

Applying Lemma 9 with  $\gamma_0 = \tau^{-1} c_0 > 2\gamma$  and setting  $e_2 = \tau^{-1} c_0 - 2\gamma$ , we get

$$\int_{c_{\mathcal{B}}(\bar{Q}, 2^J r)} \exp[-\tau^{-1} c_0 n \beta a_1 \ell(\bar{Q}, P)] d\pi(P) \leq \pi(\mathcal{B}) \exp[\Xi_2 - e_2 n \beta a_1 2^J r]$$

with

$$\Xi_2 = -\gamma + \log \left[ \frac{1}{1 - \exp[-e_2]} \right], \tag{115}$$

which together with (114) leads to

$$\log \int_{c_{\mathcal{B}}(\bar{Q}, 2^J r)} \mathbb{E}[\exp[-\beta \mathbf{T}(X, P)]] d\pi(P)$$

$$\begin{aligned} &\leq \log [\pi (\mathcal{B})] + \Xi_1 + \Xi_2 \\ &\quad + n\beta \left[ e_0 a_1 \ell(\overline{P}^*, \overline{Q}) + \tau c_1 a_0 r + \frac{l_1^2 \beta}{8} - e_2 a_1 2^J r \right]. \end{aligned} \tag{116}$$

Using the definitions (113) of  $z$  and (112) of  $\Delta_2$  we deduce from (116) that

$$\begin{aligned} &\log \left[ \frac{1}{z} \int_{\mathcal{C}_{\mathcal{B}}(\overline{Q}, 2^J r)} \mathbb{E} [\exp [-\beta \mathbf{T}(X, P)]] d\pi(P) \right] \\ &\leq \log \left( \frac{1}{z} \right) + \log [\pi (\mathcal{B})] + \Xi_1 + \Xi_2 \\ &\quad + n\beta \left[ e_0 a_1 \ell(\overline{P}^*, \overline{Q}) + \tau c_1 a_0 r + \frac{l_1^2 \beta}{8} - e_2 a_1 2^J r \right] \\ &= \xi + \log \frac{1}{\pi(\mathcal{B})} + \Xi_1 + n\beta \Delta_2 + \log [\pi (\mathcal{B})] + \Xi_1 + \Xi_2 \\ &\quad + n\beta \left[ e_0 a_1 \ell(\overline{P}^*, \overline{Q}) + \tau c_1 a_0 r + \frac{l_1^2 \beta}{8} - e_2 a_1 2^J r \right] \\ &= n\beta \left[ (e_1 + c_2 + e_0) a_1 \ell(\overline{P}^*, \overline{Q}) + e_1 a_1 r + \frac{l_2^2 \beta}{8} + \tau c_1 a_0 r + \frac{l_1^2 \beta}{8} \right] \\ &\quad + \xi + 2\Xi_1 + \Xi_2 - e_2 n\beta a_1 2^J r \\ &= n\beta \left[ (e_0 + e_1 + c_2) a_1 \ell(\overline{P}^*, \overline{Q}) + \left[ e_1 + \frac{\tau c_1 a_0}{a_1} \right] a_1 r + \frac{(l_1^2 + l_2^2) \beta}{8} \right] \\ &\quad + \xi + 2\Xi_1 + \Xi_2 - e_2 n\beta a_1 2^J r. \end{aligned} \tag{117}$$

Setting,

$$C_1 = e_0 + e_1 + c_2 \quad \text{and} \quad C_2 = e_1 + \frac{\tau c_1 a_0}{a_1},$$

we see that the right-hand side of (117) is not larger than  $-\xi$ , provided that

$$e_2 n\beta a_1 2^J r \geq 2\xi + 2\Xi_1 + \Xi_2 + n\beta \left[ C_1 a_1 \ell(\overline{P}^*, \overline{Q}) + C_2 a_1 r + \frac{(l_1^2 + l_2^2) \beta}{8} \right]$$

or equivalently if

$$2^J \geq \frac{1}{e_2} \left[ \frac{2\xi + 2\Xi_1 + \Xi_2}{\beta n a_1 r} + \frac{C_1 \ell(\overline{P}^*, \overline{Q})}{r} + C_2 + \frac{[l_1^2 + l_2^2] \beta}{8 a_1 r} \right]. \tag{118}$$

Choosing  $\bar{Q}$  in  $\mathcal{M}(\beta)$  and using the inequalities  $a_1^{-1}\beta \geq r_n(\beta, \bar{Q}) \geq 1/(\beta na_1)$ , for

$$r = \ell(\bar{P}^*, \bar{Q}) + \frac{1}{a_1} \left( \beta + \frac{2\xi}{n\beta} \right) \geq \frac{1}{\beta na_1}$$

we obtain that the right-hand side of (118) satisfies

$$\begin{aligned} & \frac{1}{e_2} \left[ \frac{2\xi + 2\Xi_1 + \Xi_2}{\beta na_1 r} + \frac{C_1 \ell(\bar{P}^*, \bar{Q}) + C_2 r}{r} + \frac{[l_1^2 + l_2^2] \beta}{8a_1 r} \right] \\ & \leq \frac{1}{e_2} \left[ C_2 + 2\Xi_1 + \Xi_2 + \frac{C_3}{r} \left( \ell(\bar{P}^*, \bar{Q}) + \frac{1}{a_1} \left( \beta + \frac{2\xi}{n\beta} \right) \right) \right] \\ & = \frac{1}{e_2} [C_2 + 2\Xi_1 + \Xi_2 + C_3] \end{aligned}$$

with  $C_3 = \max\{1, C_1, [l_1^2 + l_2^2]/8\}$ . Inequality (118) is therefore satisfied for  $J \in \mathbb{N}$  such that

$$2^J \geq \frac{C_2 + 2\Xi_1 + \Xi_2 + C_3}{e_2} \vee 1 > 2^{J-1},$$

and we may take

$$\kappa_0 = \tau \left[ \frac{2(C_2 + 2\Xi_1 + \Xi_2 + C_3)}{e_2} \vee 1 + 1 \right] \geq \tau (2^J + 1). \tag{119}$$

We recall below, the list of constants depending on  $a_0, a_1, c, \tau$  and  $\gamma$  and we have used along the proof.

$$\begin{aligned} c_0 &= 1 + c - \frac{ca_0}{a_1}, & c_1 &= 1 + c, & c_2 &= 2 + c, \\ c'_0 &= \frac{c_2 a_0}{a_1} - c_1, & l_1 &= 1 + 2c, & l_2 &= 3 + 2c, \\ e_0 &= c_0 + c + \frac{\tau c_1 a_0}{a_1}, & e_1 &= \tau \left[ c'_0 + c_1 \frac{a_0}{a_1} \right], & e_2 &= \tau^{-1} c_0 - 2\gamma, \\ C_1 &= e_0 + e_1 + c_2, & C_2 &= e_1 + \frac{\tau c_1 a_0}{a_1}, & C_3 &= \max \left\{ 1, C_1, \frac{l_1^2 + l_2^2}{8} \right\}, \end{aligned}$$

and

$$\Xi_1 = \log \left[ 1 + \frac{\exp[-(\tau^{-1}c - \gamma)]}{1 - \exp[-(\tau^{-1}c - 2\gamma)]} \right], \quad \Xi_2 = -\gamma + \log \left[ \frac{1}{1 - \exp[-e_2]} \right].$$

**9.4 Proof of Theorem 2**

The proof follows the same lines as that of Theorem 1. Under Assumption 3-(iv), the random variables  $U_i$  and  $V_i$  defined by (96) and (97) are not larger than with  $b = c + c_1 = l_1$  and  $b = c_2 + c_1 = l_2$  respectively. Since under Assumption 4, more precisely its consequence (87), that

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \mathbb{E} [U_i^2] &\leq 2 \left[ \frac{c^2}{n} \sum_{i=1}^n \text{Var} [t_{(P, Q')}(X_i)] + \frac{c_1^2}{n} \sum_{i=1}^n \text{Var} [t_{(P, Q)}(X_i)] \right] \\ &\leq 2a_2 \left[ (c^2 + c_1^2)\ell(\bar{P}^*, P) + c^2\ell(\bar{P}^*, Q') + c_1^2\ell(\bar{P}^*, Q) \right] \end{aligned}$$

and

$$\frac{1}{n} \sum_{i=1}^n \mathbb{E} [V_i^2] \leq 2a_2 \left[ (c_2^2 + c_1^2)\ell(\bar{P}^*, P) + c_2^2\ell(\bar{P}^*, Q') + c_1^2\ell(\bar{P}^*, Q) \right]$$

we may apply Lemma 8 and using the notation  $\Lambda_1 = \tau\phi(\beta l_1)$ ,  $\Lambda_2 = \tau\phi(\beta l_2)$  and Assumption 1, we get

$$\begin{aligned} &\frac{1}{n\beta} \log \left[ \prod_{i=1}^n \mathbb{E} [\exp [\beta U_i]] \right] \\ &\leq \phi(\beta l_1)\beta a_2 \left[ (c^2 + c_1^2)\ell(\bar{P}^*, P) + c^2\ell(\bar{P}^*, Q') + c_1^2\ell(\bar{P}^*, Q) \right] \\ &\leq 2\Lambda_1\beta a_2 \left[ c^2 + c_1^2 \right] \ell(\bar{P}^*, \bar{Q}) \\ &\quad + \Lambda_1\beta a_2 \left[ (c^2 + c_1^2)\ell(\bar{Q}, P) + c^2\ell(\bar{Q}, Q') + c_1^2\ell(\bar{Q}, Q) \right] \end{aligned} \tag{120}$$

and similarly

$$\begin{aligned} &\frac{1}{n\beta} \log \left[ \prod_{i=1}^n \mathbb{E} [\exp [\beta V_i]] \right] \\ &\leq 2\Lambda_2\beta a_2 \left[ c_2^2 + c_1^2 \right] \ell(\bar{P}^*, \bar{Q}) \\ &\quad + \Lambda_2\beta a_2 \left[ (c_2^2 + c_1^2)\ell(\bar{Q}, P) + c_2^2\ell(\bar{Q}, Q') + c_1^2\ell(\bar{Q}, Q) \right]. \end{aligned} \tag{121}$$

It follows from (102) that

$$\begin{aligned} E_1 &= n^{-1} \{ c\mathbb{E} [\mathbf{T}(X, P, Q')] - c_1\mathbb{E} [\mathbf{T}(X, P, Q)] \} \\ &\quad + 2\Lambda_1\beta a_2 \left[ c^2 + c_1^2 \right] \ell(\bar{P}^*, \bar{Q}) \\ &\quad + \Lambda_1\beta a_2 \left[ (c^2 + c_1^2)\ell(\bar{Q}, P) + c^2\ell(\bar{Q}, Q') + c_1^2\ell(\bar{Q}, Q) \right] \end{aligned}$$

$$\begin{aligned} &\leq \left[ e_0 a_1 + 2\Lambda_1 \beta a_2 (c^2 + c_1^2) \right] \ell(\bar{P}^*, \bar{Q}) \\ &\quad - \left[ \tau^{-1} c_0 a_1 - \Lambda_1 \beta a_2 (c^2 + c_1^2) \right] \ell(\bar{Q}, P) \\ &\quad - \left[ \tau^{-1} c a_1 - \Lambda_1 \beta a_2 c^2 \right] \ell(\bar{Q}, Q') \\ &\quad + \left[ \tau c_1 a_0 + \Lambda_1 \beta a_2 c_1^2 \right] \ell(\bar{Q}, Q). \end{aligned}$$

Using the definitions (25) of  $\bar{c}_1$  and (26) of  $\bar{c}_2$ , that is,

$$\bar{c}_1 = c_0 - \tau \Lambda_1 \beta a_2 a_1^{-1} (c^2 + c_1^2) \quad \text{and} \quad \bar{c}_2 = c - \tau \Lambda_1 \beta a_2 a_1^{-1} c^2$$

and setting

$$\begin{aligned} e_3 &= e_0 + 2\Lambda_1 \beta \frac{a_2 (c^2 + c_1^2)}{a_1} \\ e_4 &= \frac{1}{a_1} \left[ \tau c_1 a_0 + \Lambda_1 \beta a_2 c_1^2 \right] \end{aligned}$$

and arguing as in the proof of inequality (105), we deduce from (120) that

$$\begin{aligned} &\log \mathbb{E} \left[ \exp \left[ \beta (c \mathbf{T}(X, P, Q') - c_1 \mathbf{T}(X, P, Q)) \right] \right] \\ &\leq n \beta E_1 \\ &\leq n \beta a_1 \left[ e_3 \ell(\bar{P}^*, \bar{Q}) - \tau^{-1} [\bar{c}_1 \ell(\bar{Q}, P) + \bar{c}_2 \ell(\bar{Q}, Q')] + e_4 \ell(\bar{Q}, Q) \right]. \end{aligned} \tag{122}$$

It follows from (103) that

$$\begin{aligned} E_2 &= n^{-1} \left\{ c_2 \mathbb{E} [\mathbf{T}(X, P, Q')] - c_1 \mathbb{E} [\mathbf{T}(X, P, Q)] \right\} \\ &\quad + 2\Lambda_2 \beta a_2 \left[ c_2^2 + c_1^2 \right] \ell(\bar{P}^*, \bar{Q}) \\ &\quad + \Lambda_2 \beta a_2 \left[ (c_2^2 + c_1^2) \ell(\bar{Q}, P) + c_2^2 \ell(\bar{Q}, Q') + c_1^2 \ell(\bar{Q}, Q) \right] \\ &\leq (e_1 + c_2) a_1 \ell(\bar{P}^*, \bar{Q}) + \tau c'_0 a_1 \ell(\bar{Q}, P) - \tau^{-1} c_2 a_1 \ell(\bar{Q}, Q') \\ &\quad + \tau c_1 a_0 \ell(\bar{Q}, Q) + 2\Lambda_2 \beta a_2 \left[ c_2^2 + c_1^2 \right] \ell(\bar{P}^*, \bar{Q}) \\ &\quad + \Lambda_2 \beta a_2 \left[ (c_2^2 + c_1^2) \ell(\bar{Q}, P) + c_2^2 \ell(\bar{Q}, Q') + c_1^2 \ell(\bar{Q}, Q) \right] \\ &= \left[ (e_1 + c_2) a_1 + 2\Lambda_2 \beta a_2 (c_2^2 + c_1^2) \right] \ell(\bar{P}^*, \bar{Q}) \\ &\quad + \left[ \tau c'_0 a_1 + \Lambda_2 \beta a_2 (c_2^2 + c_1^2) \right] \ell(\bar{Q}, P) \\ &\quad - \left[ \tau^{-1} c_2 a_1 - \Lambda_2 \beta a_2 c_2^2 \right] \ell(\bar{Q}, Q') \\ &\quad + \left[ \tau c_1 a_0 + \Lambda_2 \beta a_2 c_1^2 \right] \ell(\bar{Q}, Q). \end{aligned}$$

Using the definition (27) of  $\bar{c}_3$ , that is,

$$\bar{c}_3 = c_2 - \tau \Lambda_2 \beta a_2 a_1^{-1} c_2^2,$$

and setting

$$e_5 = e_1 + c_2 + 2\Lambda_2 \beta \frac{a_2 (c_2^2 + c_1^2)}{a_1}, \quad e_6 = \tau c'_0 + \Lambda_2 \beta \frac{a_2 (c_2^2 + c_1^2)}{a_1}$$

$$e_7 = \frac{1}{a_1} \left[ \tau c_1 a_0 + \Lambda_2 \beta a_2 c_1^2 \right],$$

and arguing as in the proof of (107), we deduce from (121) that

$$\log \mathbb{E} \left[ \exp \left[ \beta (c_2 \mathbf{T}(X, P, Q') - c_1 \mathbf{T}(X, P, Q)) \right] \right] \leq n\beta E_2$$

$$= n\beta a_1 \left( e_5 \ell(\bar{P}^*, \bar{Q}) + e_6 \ell(\bar{Q}, P) - \tau^{-1} \bar{c}_3 \ell(\bar{Q}, Q') + e_7 \ell(\bar{Q}, Q) \right). \quad (123)$$

Under our assumption on  $\beta$ , we know that the quantities  $\bar{c}_2$  and  $\bar{c}_3$  are positive and that  $2\gamma < \tau^{-1} (\bar{c}_2 \wedge \bar{c}_3)$ . We may therefore apply Lemma 9 with  $\gamma_0 = \tau^{-1} \bar{c}_2$  and  $\gamma_0 = \tau^{-1} \bar{c}_3$  successively and get

$$\int_{\mathcal{M}} \exp \left[ -\tau^{-1} \bar{c}_2 n\beta a_1 \ell(\bar{Q}, Q') \right] d\pi(Q') \leq \pi(\mathcal{B}) \exp \left[ \bar{\Xi}_1 \right] \quad (124)$$

and

$$\int_{\mathcal{M}} \exp \left[ -\tau^{-1} \bar{c}_3 n\beta a_1 \ell(\bar{Q}, Q') \right] d\pi(Q') \leq \pi(\mathcal{B}) \exp \left[ \bar{\Xi}_1 \right] \quad (125)$$

with

$$\bar{\Xi}_1 = \log \left[ 1 + \frac{\exp \left[ -(\tau^{-1} (\bar{c}_2 \wedge \bar{c}_3) - \gamma) \right]}{1 - \exp \left[ -(\tau^{-1} (\bar{c}_2 \wedge \bar{c}_3) - 2\gamma) \right]} \right]. \quad (126)$$

Putting (123) and (125) together, we obtain that for all  $(P, Q) \in \mathcal{B}^2$

$$\exp \left[ \mathbf{L}(P, Q) \right] = \int_{\mathcal{M}} \mathbb{E} \left[ \exp \left[ \beta (c_2 \mathbf{T}(X, P, Q') - c_1 \mathbf{T}(X, P, Q)) \right] \right] d\pi(Q')$$

$$\leq \exp \left[ n\beta a_1 \left( e_5 \ell(\bar{P}^*, \bar{Q}) + e_6 \ell(\bar{Q}, P) + e_7 \ell(\bar{Q}, Q) \right) \right]$$

$$\times \int_{\mathcal{M}} \exp \left[ -\tau^{-1} \bar{c}_3 n\beta a_1 \ell(\bar{Q}, Q') \right] d\pi(Q')$$

$$\leq \pi(\mathcal{B}) \exp \left[ \bar{\Xi}_1 + n\beta a_1 \left( e_5 \ell(\bar{P}^*, \bar{Q}) + e_6 \ell(\bar{Q}, P) + e_7 \ell(\bar{Q}, Q) \right) \right]$$

$$\leq \pi(\mathcal{B}) \exp \left[ \bar{\Xi}_1 + n\beta a_1 \left( e_5 \ell(\bar{P}^*, \bar{Q}) + (e_6 + e_7)r \right) \right].$$



Consequently,

$$\begin{aligned} & \left[ \int_{\mathcal{B}^2} \exp[-\mathbf{L}(P, Q)] d\pi_{\mathcal{B}}(P)d\pi_{\mathcal{B}}(Q) \right]^{-1} \\ & \leq \pi(\mathcal{B}) \exp \left[ \bar{\Xi}_1 + n\beta a_1 \left( e_5 \ell(\bar{P}^*, \bar{Q}) + (e_6 + e_7)r \right) \right]. \end{aligned}$$

We deduce from (84) that

$$\mathbb{P}^{(CA)} \leq \frac{z}{\pi(\mathcal{B})} \exp \left[ \bar{\Xi}_1 + n\beta a_1 \left( e_5 \ell(\bar{P}^*, \bar{Q}) + (e_6 + e_7)r \right) \right].$$

In particular,  $\mathbb{P}^{(CA)} \leq e^{-\xi}$  for  $z$  satisfying

$$\log \left( \frac{1}{z} \right) = \xi + \log \frac{1}{\pi(\mathcal{B})} + \bar{\Xi}_1 + n\beta a_1 \left[ e_5 \ell(\bar{P}^*, \bar{Q}) + (e_6 + e_7)r \right]. \tag{127}$$

Putting (122) and (124) together, we obtain that for all  $Q \in \mathcal{B}$

$$\begin{aligned} & \exp[\mathbf{M}(P, Q)] \\ & = \int_{\mathcal{M}} \mathbb{E} \left[ \exp \left[ \beta \left( c\mathbf{T}(\mathbf{X}, P, Q') - c_1\mathbf{T}(\mathbf{X}, P, Q) \right) \right] \right] d\pi(Q') \\ & \leq \exp \left[ n\beta a_1 \left( e_3 \ell(\bar{P}^*, \bar{Q}) - \tau^{-1} \bar{c}_1 \ell(\bar{Q}, P) + e_4 \ell(\bar{Q}, Q) \right) \right] \\ & \quad \times \int_{\mathcal{M}} \exp \left[ -\tau^{-1} \bar{c}_2 n\beta a_1 \ell(\bar{Q}, Q') \right] d\pi(Q') \\ & \leq \pi(\mathcal{B}) \exp \left[ \bar{\Xi}_1 + n\beta a_1 \left( e_3 \ell(\bar{P}^*, \bar{Q}) + e_4 r - \tau^{-1} \bar{c}_1 \ell(\bar{Q}, P) \right) \right]. \end{aligned}$$

We derive from Lemma 5 that

$$\begin{aligned} & \mathbb{E} \left[ \exp[-\beta \mathbf{T}(\mathbf{X}, P)] \right] \\ & \leq \frac{1}{\pi(\mathcal{B})} \left[ \int_{\mathcal{B}} \exp[-\mathbf{M}(P, Q)] d\pi_{\mathcal{B}}(Q) \right]^{-1} \\ & \leq \exp \left[ \bar{\Xi}_1 + n\beta a_1 \left( e_3 \ell(\bar{P}^*, \bar{Q}) + e_4 r - \tau^{-1} \bar{c}_1 \ell(\bar{Q}, P) \right) \right], \end{aligned}$$

and consequently,

$$\begin{aligned} & \int_{c_{\mathcal{B}}(\bar{Q}, 2J_r)} \mathbb{E} \left[ \exp[-\beta \mathbf{T}(\mathbf{X}, P)] \right] d\pi(P) \\ & \leq \exp \left[ \bar{\Xi}_1 + n\beta a_1 \left( e_3 \ell(\bar{P}^*, \bar{Q}) + e_4 r \right) \right] \\ & \quad \times \int_{c_{\mathcal{B}}(\bar{Q}, 2J_r)} \exp \left[ -\tau^{-1} \bar{c}_1 n\beta a_1 \ell(\bar{Q}, P) \right] d\pi(P). \tag{128} \end{aligned}$$

Since under our assumptions,  $\bar{c}_1 > 0$  and  $2\gamma < \tau^{-1}\bar{c}_1$  we may apply Lemma 9 with  $\gamma_0 = \tau^{-1}\bar{c}_1$ , and setting  $e_8 = \tau^{-1}\bar{c}_1 - 2\gamma$  which leads to

$$\int_{\mathcal{C}_{\mathcal{B}}(\bar{Q}, 2^J r)} \exp \left[ -\tau^{-1}\bar{c}_1 n\beta a_1 \ell(\bar{Q}, P) \right] d\pi(P) \leq \pi(\mathcal{B}) \exp \left[ \bar{\Xi}_2 - e_8 n\beta a_1 2^J r \right].$$

with

$$\bar{\Xi}_2 = -\gamma + \log \left[ \frac{1}{1 - \exp[-e_8]} \right], \tag{129}$$

which together with (128) leads to

$$\begin{aligned} & \int_{\mathcal{C}_{\mathcal{B}}(\bar{Q}, 2^J r)} \mathbb{E} \left[ \exp[-\beta \mathbf{T}(X, P)] \right] d\pi(P) \\ & \leq \pi(\mathcal{B}) \exp \left[ \bar{\Xi}_1 + \bar{\Xi}_2 + n\beta a_1 \left( e_3 \ell(\bar{P}^*, \bar{Q}) + e_4 r - e_8 2^J r \right) \right]. \end{aligned} \tag{130}$$

Using the definition (127) of  $z$ , we deduce that

$$\begin{aligned} & \log \left[ \frac{1}{z} \int_{\mathcal{C}_{\mathcal{B}}(\bar{Q}, 2^J r)} \mathbb{E} \left[ \exp[-\beta \mathbf{T}(X, P)] \right] d\pi(P) \right] \\ & \leq \log \left( \frac{1}{z} \right) + \log \pi(\mathcal{B}) + \bar{\Xi}_1 + \bar{\Xi}_2 + n\beta a_1 \left( e_3 \ell(\bar{P}^*, \bar{Q}) + e_4 r - e_8 2^J r \right) \\ & = \xi + \log \frac{1}{\pi(\mathcal{B})} + \bar{\Xi}_1 + n\beta a_1 \left[ e_5 \ell(\bar{P}^*, \bar{Q}) + (e_6 + e_7)r \right] \\ & \quad + \log \pi(\mathcal{B}) + \bar{\Xi}_1 + \bar{\Xi}_2 + n\beta a_1 \left( e_3 \ell(\bar{P}^*, \bar{Q}) + e_4 r - e_8 2^J r \right) \\ & = \xi + 2\bar{\Xi}_1 + \bar{\Xi}_2 + n\beta a_1 \left[ (e_3 + e_5) \ell(\bar{P}^*, \bar{Q}) + (e_4 + e_6 + e_7)r \right] \\ & \quad - e_8 n\beta a_1 2^J r. \end{aligned}$$

The right-hand side is not larger than  $-\xi$  provided that

$$2^J \geq \frac{1}{e_8} \left[ \frac{2\xi + 2\bar{\Xi}_1 + \bar{\Xi}_2}{n\beta a_1 r} + \left[ (e_3 + e_5) \frac{\ell(\bar{P}^*, \bar{Q})}{r} + e_4 + e_6 + e_7 \right] \right]. \tag{131}$$

Using the fact that  $r_n(\beta, \bar{Q}) \geq 1/(n\beta a_1)$ , with the choice

$$r = \ell(\bar{P}^*, \bar{Q}) + r_n(\beta, \bar{Q}) + \frac{2\xi}{n\beta a_1} \geq \ell(\bar{P}^*, \bar{Q}) + \frac{1 + 2\xi}{n\beta a_1} \geq \frac{1}{n\beta a_1},$$

the right-hand side of (131) satisfies

$$\frac{1}{e_8} \left[ \frac{2\xi + 2\bar{\Xi}_1 + \bar{\Xi}_2}{n\beta a_1 r} + \left[ (e_3 + e_5) \frac{\ell(\bar{P}^*, \bar{Q})}{r} + e_4 + e_6 + e_7 \right] \right]$$

$$\begin{aligned} &\leq \frac{1}{e_8} \left[ 2\bar{\Xi}_1 + \bar{\Xi}_2 + e_4 + e_6 + e_7 + \frac{(e_3 + e_5) \vee 1}{r} \left( \ell(\bar{P}^*, \bar{Q}) + \frac{2\xi}{n\beta a_1} \right) \right] \\ &\leq \frac{2\bar{\Xi}_1 + \bar{\Xi}_2 + e_4 + e_6 + e_7 + (e_3 + e_5) \vee 1}{e_8}. \end{aligned}$$

Inequality (131) holds for  $J \in \mathbb{N}$  such that

$$2^J \geq \frac{2\bar{\Xi}_1 + \bar{\Xi}_2 + e_4 + e_6 + e_7 + (e_3 + e_5) \vee 1}{e_8} \vee 1 > 2^{J-1},$$

and we may take

$$\begin{aligned} \kappa_0 &= \tau \left[ \frac{2 \left[ 2\bar{\Xi}_1 + \bar{\Xi}_2 + e_4 + e_6 + e_7 + (e_3 + e_5) \vee 1 \right]}{e_8} \vee 1 + 1 \right] \\ &\geq \tau \left( 2^J + 1 \right). \end{aligned} \tag{132}$$

In complements to constants listed at the end of the proof of Theorem 1, we recall that

$$\Lambda_1 = \tau\phi(\beta l_1), \quad \Lambda_2 = \tau\phi(\beta l_2)$$

$$\bar{c}_1 = c_0 - \tau\Lambda_1\beta \frac{a_2(c^2 + c_1^2)}{a_1}, \quad \bar{c}_2 = c - \tau\Lambda_1\beta \frac{a_2c^2}{a_1}, \quad \bar{c}_3 = c_2 - \tau\Lambda_2\beta \frac{a_2c_2^2}{a_1},$$

$$\begin{aligned} e_3 &= e_0 + 2\Lambda_1\beta \frac{a_2(c^2 + c_1^2)}{a_1}, & e_4 &= \frac{1}{a_1} \left[ \tau c_1 a_0 + \Lambda_1\beta a_2 c_1^2 \right], \\ e_5 &= e_1 + c_2 + 2\Lambda_2\beta \frac{a_2(c_2^2 + c_1^2)}{a_1}, & e_6 &= \tau c'_0 + \Lambda_2\beta \frac{a_2(c_2^2 + c_1^2)}{a_1}, \\ e_7 &= \frac{1}{a_1} \left[ \tau c_1 a_0 + \Lambda_2\beta a_2 c_1^2 \right], & e_8 &= \tau^{-1}\bar{c}_1 - 2\gamma, \end{aligned}$$

and

$$\begin{aligned} \bar{\Xi}_1 &= \log \left[ 1 + \frac{\exp [ - (\tau^{-1}(\bar{c}_2 \wedge \bar{c}_3) - \gamma) ]}{1 - \exp [ - (\tau^{-1}(\bar{c}_2 \wedge \bar{c}_3) - 2\gamma) ]} \right], \\ \bar{\Xi}_2 &= -\gamma + \log \left[ \frac{1}{1 - \exp [ -e_8 ]} \right]. \end{aligned}$$

### 9.5 Proof of Theorem 3

Let us take  $r \geq \varepsilon$  and set  $\varpi = 2\xi + 1$  so that

$$\pi \left( {}^c\mathcal{B} \right) \leq \pi \left( {}^c\mathcal{B}(\overline{Q}, \varepsilon) \right) \leq e^{-\varpi} \pi \left( \mathcal{B}(\overline{Q}, \varepsilon) \right) \leq e^{-\varpi} \pi \left( \mathcal{B} \right).$$

In order to prove the first part, let us go back to the proof of Theorem 1. Clearly,

$$\int_{\mathcal{M}} \exp \left[ -\tau^{-1} c n \beta a_1 \ell(\overline{Q}, Q') \right] d\pi(Q') \leq 1 = \pi \left( \mathcal{B} \right) + \pi \left( {}^c\mathcal{B} \right) \leq \pi \left( \mathcal{B} \right) (1 + e^{-\varpi})$$

and similarly,

$$\int_{\mathcal{M}} \exp \left[ -\tau^{-1} c_2 n \beta a_1 \ell(\overline{Q}, Q') \right] d\pi(Q') \leq \pi \left( \mathcal{B} \right) (1 + e^{-\varpi}).$$

Inequalities (109) and (110) are therefore satisfied with  $\Xi_1 = \log(1 + e^{-1})$ . Moreover,

$$\begin{aligned} & \int_{{}^c\mathcal{B}(\overline{Q}, 2^J r)} \exp \left[ -\tau^{-1} c_0 n \beta a_1 \ell(\overline{Q}, P) \right] d\pi(P) \\ & \leq \exp \left[ -\tau^{-1} c_0 n \beta a_1 2^J r \right] \pi \left( {}^c\mathcal{B} \right) \leq \pi \left( \mathcal{B} \right) \exp \left[ -\varpi - \tau^{-1} c_0 n \beta a_1 2^J r \right]. \end{aligned}$$

We deduce from (114) that

$$\begin{aligned} & \int_{{}^c\mathcal{B}(\overline{Q}, 2^J r)} \mathbb{E} \left[ \exp \left[ -\beta \mathbf{T}(X, P) \right] \right] d\pi(P) \\ & \leq \exp \left[ \Xi_1 + n\beta \left( e_0 a_1 \ell(\overline{P}^*, \overline{Q}) + \tau c_1 a_0 r + \frac{l_1^2 \beta}{8} \right) \right] \\ & \quad \times \pi \left( \mathcal{B} \right) \exp \left[ -\varpi - \tau^{-1} c_0 n \beta a_1 2^J r \right], \end{aligned}$$

and consequently,

$$\begin{aligned} & \log \int_{{}^c\mathcal{B}(\overline{Q}, 2^J r)} \mathbb{E} \left[ \exp \left[ -\beta \mathbf{T}(X, P) \right] \right] d\pi(P) \\ & \leq \log \pi \left( \mathcal{B} \right) + \Xi_1 - \varpi \\ & \quad + n\beta \left[ e_0 a_1 \ell(\overline{P}^*, \overline{Q}) + \tau c_1 a_0 r + \frac{l_1^2 \beta}{8} - \tau^{-1} c_0 a_1 2^J r \right]. \end{aligned}$$

Using the definitions (113) of  $z$  and (112) of  $\Delta_2$ , we deduce that

$$\log \left[ \frac{1}{z} \int_{{}^c\mathcal{B}(\overline{Q}, 2^J r)} \mathbb{E} \left[ \exp \left[ -\beta \mathbf{T}(X, P) \right] \right] d\pi(P) \right]$$

$$\begin{aligned}
 &\leq \xi + \log \frac{1}{\pi(\mathcal{B})} + \Xi_1 + n\beta\Delta_2 + \log \pi(\mathcal{B}) + \Xi_1 - \varpi \\
 &\quad + n\beta \left[ e_0 a_1 \ell(\bar{P}^*, \bar{Q}) + \tau c_1 a_0 r + \frac{l_1^2 \beta}{8} - \tau^{-1} c_0 a_1 2^J r \right] \\
 &= \xi + 2\Xi_1 + n\beta \left[ (e_1 + c_2) a_1 \ell(\bar{P}^*, \bar{Q}) + e_1 a_1 r + \frac{l_2^2 \beta}{8} \right] - \varpi \\
 &\quad + n\beta \left[ e_0 a_1 \ell(\bar{P}^*, \bar{Q}) + \tau c_1 a_0 r + \frac{l_1^2 \beta}{8} - \tau^{-1} c_0 a_1 2^J r \right] \\
 &= \xi + 2\Xi_1 - \varpi \\
 &\quad + n\beta a_1 \left[ C_1 \ell(\bar{P}^*, \bar{Q}) + C_2 r + \frac{(l_1^2 + l_2^2) \beta}{8 a_1} - \tau^{-1} c_0 2^J r \right],
 \end{aligned}$$

where the constants  $C_1$  and  $C_2$  are the same as those defined in the proof of Theorem 1. If we choose  $r = \ell(\bar{P}^*, \bar{Q}) \vee (\beta/a_1) \vee \varepsilon$  and  $J$  such that  $\tau^{-1} c_0 2^J \geq C_1 + C_2 + (l_1^2 + l_2^2)/8$ , we obtain that

$$\log \left[ \frac{1}{z} \int_{c_{\mathcal{B}}(\bar{Q}, 2^J r)} \mathbb{E} [\exp [-\beta \mathbf{T}(\mathbf{X}, P)]] d\pi(P) \right] \leq \xi + 2\Xi_1 - \varpi \leq -\xi$$

since  $\varpi = 2\xi + 1 \geq 2(\xi + \Xi_1)$ . We conclude as in the proof of Theorem 1.

In order to prove the second part of Theorem 3, we go back to the proof of Theorem 2. The arguments are similar. As before,

$$\int_{\mathcal{M}} \exp \left[ -\tau^{-1} \bar{c}_2 n \beta a_1 \ell(\bar{Q}, Q') \right] d\pi(Q') \leq \pi(\mathcal{B}) (1 + e^{-\varpi})$$

and

$$\int_{\mathcal{M}} \exp \left[ -\tau^{-1} \bar{c}_3 n \beta a_1 \ell(\bar{Q}, Q') \right] d\pi(Q') \leq \pi(\mathcal{B}) (1 + e^{-\varpi}).$$

Inequalities (124) and (125) are therefore both satisfied with  $\bar{\Xi}_1 = \log(1 + e^{-1})$ . Moreover

$$\begin{aligned}
 &\int_{c_{\mathcal{B}}(\bar{Q}, 2^J r)} \exp \left[ -\tau^{-1} \bar{c}_1 n \beta a_1 \ell(\bar{Q}, P) \right] d\pi(P) \\
 &\leq \pi(\mathcal{B}) \exp \left[ -\varpi - \tau^{-1} \bar{c}_1 n \beta a_1 2^J r \right],
 \end{aligned}$$

and we deduce from (128) that

$$\int_{c_{\mathcal{B}}(\bar{Q}, 2^J r)} \mathbb{E} [\exp [-\beta \mathbf{T}(\mathbf{X}, P)]] d\pi(P)$$

$$\begin{aligned} &\leq \exp \left[ \bar{\Xi}_1 + n\beta a_1 \left( e_3 \ell(\bar{P}^*, \bar{Q}) + e_4 r \right) \right] \\ &\quad \times \int_{\mathcal{C}_{\mathcal{B}}(\bar{Q}, 2^J r)} \exp \left[ -\tau^{-1} \bar{c}_1 n\beta a_1 \ell(\bar{Q}, P) \right] d\pi(P) \\ &\leq \pi(\mathcal{B}) \exp \left[ \bar{\Xi}_1 + n\beta a_1 \left[ e_3 \ell(\bar{P}^*, \bar{Q}) + e_4 r - \tau^{-1} \bar{c}_1 2^J r \right] - \varpi \right]. \end{aligned}$$

Using the definition (127) of  $z$ , we deduce that

$$\begin{aligned} &\log \left[ \frac{1}{z} \int_{\mathcal{C}_{\mathcal{B}}(\bar{Q}, 2^J r)} \mathbb{E} \left[ \exp \left[ -\beta \mathbf{T}(X, P) \right] \right] d\pi(P) \right] \\ &\leq \xi + \log \frac{1}{\pi(\mathcal{B})} + \bar{\Xi}_1 + n\beta a_1 \left[ e_5 \ell(\bar{P}^*, \bar{Q}) + (e_6 + e_7)r \right] \\ &\quad + \log \pi(\mathcal{B}) + \bar{\Xi}_1 + n\beta a_1 \left[ e_3 \ell(\bar{P}^*, \bar{Q}) + e_4 r - \tau^{-1} \bar{c}_1 2^J r \right] - \varpi \\ &= \xi + 2\bar{\Xi}_1 - \varpi \\ &\quad + n\beta a_1 \left[ (e_3 + e_5) \ell(\bar{P}^*, \bar{Q}) + (e_4 + e_6 + e_7)r - \tau^{-1} \bar{c}_1 2^J r \right]. \end{aligned}$$

Taking  $r = \ell(\bar{P}^*, \bar{Q}) \vee \varepsilon \geq \varepsilon$  and  $J \geq 0$  such that

$$\tau^{-1} \bar{c}_1 2^J \geq e_3 + e_5 + e_4 + e_6 + e_7$$

we obtain that

$$\log \left[ \frac{1}{z} \int_{\mathcal{C}_{\mathcal{B}}(\bar{Q}, 2^J r)} \mathbb{E} \left[ \exp \left[ -\beta \mathbf{T}(X, P) \right] \right] d\pi(P) \right] \leq -\xi$$

and we conclude as before.

## 10 Other proofs

### 10.1 Proof of Lemma 1

Let  $Y$  be a random variable with gamma distribution  $\gamma(s, 1)$ . Since  $\sigma Y \sim \gamma(s, \sigma)$ , it is sufficient to prove the result for  $\sigma = 1$ . Using the inequality  $\log(1 - x) \geq -x/(1 - x)$  which holds for all  $x \in [0, 1)$ , we obtain that

$$\log \mathbb{E} \left[ e^{\beta(Y-s)} \right] = -s \left[ \log(1 - \beta) + \beta \right] \leq \frac{s\beta^2}{1 - \beta} \quad \text{for all } \beta \in [0, 1).$$

Applying Lemma 8.2 in Birgé [10] with  $a = \sqrt{s}$  and  $b = 1$ , we obtain that

$$\mathbb{P} \left[ Y \geq s + 2\sqrt{s\xi} + \xi \right] \leq e^{-\xi} \quad \text{for all } \xi \geq 0$$

which proves (44). Let us now turn to the lower bound. For  $x \geq 0$ , let us set

$$g(x) = x - \log(1 + x) \leq \left(\frac{x^2}{2}\right) \wedge x.$$

For all  $t, u \geq 0$ ,

$$\begin{aligned} \int_{t+u}^{+\infty} x^t e^{-x} dx &= \int_u^{+\infty} (t + y)^t e^{-t-y} dy = t^t e^{-t} \int_u^{+\infty} e^{-tg(y/t)} dy \\ &\geq t^t e^{-t} \left( \int_u^{+\infty} e^{-y^2/(2t)} dy \vee \int_u^{+\infty} e^{-y} dy \right) \\ &= t^t e^{-t} \left[ \left( \sqrt{2\pi t} \bar{F} \left( \frac{u}{\sqrt{t}} \right) \right) \vee e^{-u} \right], \end{aligned}$$

where  $\bar{F}(z) = \mathbb{P}[\mathcal{N}(0, 1) \geq z]$  for all  $z \in \mathbb{R}$ . Using the the following inequalities

$$t^{t-1/2} e^{-t} \sqrt{2\pi} \leq \Gamma(t) \leq t^{t-1/2} e^{-t} \sqrt{2\pi} \exp[1/(12t)], \tag{133}$$

that can be found in Whittaker and Watson [25, p. 253], with  $t = s - 1 > 0$ , we deduce that

$$\begin{aligned} \mathbb{P}[Y \geq t + u] &= \frac{1}{\Gamma(t + 1)} \int_{t+u}^{+\infty} x^t e^{-x} dx = \frac{1}{t\Gamma(t)} \int_{t+u}^{+\infty} x^t e^{-x} dx \\ &\geq \left[ \bar{F} \left( \frac{u}{\sqrt{t}} \right) e^{-1/(12t)} \right] \vee \left[ \frac{e^{-u-1/(12t)}}{\sqrt{2\pi t}} \right]. \end{aligned}$$

Using the fact that  $\bar{F}(\Phi^{-1}(z)) = e^{-z}$  for all  $z \geq 0$ , we obtain that for the choice

$$u = \left[ \sqrt{t} \Phi^{-1} \left( \xi - \frac{1}{12t} \right) \right] \vee \log \left( \frac{e^{\xi-1/(12t)}}{\sqrt{2\pi t}} \right),$$

which is nonnegative for  $\xi \geq \log 2 + 1/(12t)$ , the quantity  $\mathbb{P}[Y \geq t + u]$  is at least  $e^{-\xi}$ , which proves (45).

### 10.2 Proof of Theorem 4

Throughout this proof,  $a_0 = 2, a_1 = 3/16, \beta = 2\gamma = 1/500$  and  $\kappa$  denotes a positive numerical constant that may vary from line to line. It follows from Corollary 4 that for  $n$  large enough,  $r_n(\beta, P_{\theta^*}) \leq r_n^* = \kappa k/n$ . Applying our Corollary 2 with  $\ell = h^2$  (and  $2\xi$  in place of  $\xi$ ), we obtain that for  $n$  large enough, with a probability at least  $1 - 2e^{-\xi}$ ,

$$1 - e^{-\xi} \leq \widehat{v}_X^h \left( \left\{ \theta \in \Theta, h^2(\theta, \theta^*) \leq r_n(\xi) \right\} \right) \text{ with } r_n(\xi) = \frac{\kappa(k + \xi)}{n}.$$

We know by Proposition 9 that under the assumptions of Corollary 4, Assumption 9-(i) is satisfied with  $s = 2$ ,  $|\cdot|_*$  given by (81) and  $\varepsilon = 1/2$ . This implies that for  $n$  large

$$\left\{ \theta \in \Theta, h^2(\theta, \theta^*) \leq r_n(\xi) \right\} \subset \left\{ \theta \in \Theta, |\theta - \theta^*|_*^2 \leq 2r_n(\xi) \right\},$$

which leads to (46).

### 10.3 Proof of Proposition 4

Let us denote by  $F_\sigma$  the distribution function of  $\nu_\sigma$ . Throughout this proof, we fix some  $\theta^* \in [-\sigma t, \sigma t]$ . Our aim is to prove that  $P_{\theta^*}$  belongs to  $\mathcal{M}(\bar{\beta})$ .

Since the total variation distance is translation invariant,  $\|P_\theta - P_{\theta^*}\| = \|P_{\theta - \theta^*} - P_0\| = \|P_{\theta^* - \theta} - P_0\|$  and consequently, for all  $r \in [0, 1)$ ,

$$\{\theta \in \Theta, \|P_\theta - P_{\theta^*}\| \leq r\} = \{\theta \in \Theta, |\theta^* - \theta| \leq \varphi(r)\} \quad \text{for all } r \in [0, 1)$$

while for  $r \geq 1$ ,  $\{\theta \in \Theta, \|P_\theta - P_{\theta^*}\| \leq r\} = \Theta = \mathbb{R}$ .

We set  $r_0 = \sup\{r > 0, \varphi(r) \leq \sigma t\}$  and distinguish between two cases.

**Case 1** Assume  $r_0 \leq 1/4$ . For all  $r < r_0$ ,  $\varphi(r) < \sigma t$ ,  $2r < 1$ , and since  $q$  is symmetric, positive and decreasing on  $\mathbb{R}_+$ ,

$$\begin{aligned} \frac{\pi(\mathcal{B}(P_{\theta^*}, 2r))}{\pi(\mathcal{B}(P_{\theta^*}, r))} &= \frac{\nu_\sigma(\{\theta \in \mathbb{R}, \|P_\theta - P_{\theta^*}\| \leq 2r\})}{\nu_\sigma(\{\theta \in \mathbb{R}, \|P_\theta - P_{\theta^*}\| \leq r\})} \\ &= \frac{\nu_\sigma(\{\theta \in \mathbb{R}, |\theta - \theta^*| \leq \varphi(2r)\})}{\nu_\sigma(\{\theta \in \mathbb{R}, |\theta - \theta^*| \leq \varphi(r)\})} \leq \frac{2q_\sigma(0)\varphi(2r)}{2q_\sigma(|\theta^*| + \varphi(r))\varphi(r)} \\ &\leq \frac{q_\sigma(0)\varphi(2r)}{q_\sigma(|\theta^*| + \sigma t)\varphi(r)} \leq \frac{q_\sigma(0)\varphi(2r)}{q_\sigma(2\sigma t)\varphi(r)} \\ &= \frac{q(0)\varphi(2r)}{q(2t)\varphi(r)} \leq \frac{\bar{\Gamma}}{q(2t)}. \end{aligned}$$

For all  $r_0 < r < 1$ ,  $|\theta^*| \leq \sigma t < \varphi(r)$ , hence  $F_\sigma(|\theta^*| - \varphi(r)) \leq F_\sigma(0) = 1/2$  and  $F_\sigma(|\theta^*| + \varphi(r)) \geq F_\sigma(\varphi(r)) \geq F_\sigma(\sigma t) = F_1(t) \geq 3/4$  under our assumption on  $t$ . Consequently,

$$\begin{aligned} \frac{\pi(\mathcal{B}(P_{\theta^*}, 2r))}{\pi(\mathcal{B}(P_{\theta^*}, r))} &\leq \frac{1}{\nu_\sigma(\{\theta \in \mathbb{R}, |\theta - \theta^*| \leq \varphi(r)\})} \\ &= \frac{1}{F_\sigma(|\theta^*| + \varphi(r)) - F(|\theta^*| - \varphi(r))} \\ &\leq \frac{1}{3/4 - 1/2} = 4. \end{aligned}$$

Note that the result also holds for  $r = r_0$  by letting  $r$  decrease to  $r_0$ .



**Case 2** Assume that  $r_0 > 1/4$ . Then  $\varphi(1/4) \leq \sigma t$  and arguing as before, we obtain that for all  $r \leq 1/4 < r_0$ ,

$$\begin{aligned} \frac{\pi(\mathcal{B}(P_{\theta^*}, 2r))}{\pi(\mathcal{B}(P_{\theta^*}, r))} &\leq \frac{2q_\sigma(0)\varphi(2r)}{2q_\sigma(|\theta^*| + \varphi(r))\varphi(r)} = \frac{q_\sigma(0)\varphi(2r)}{q_\sigma(|\theta^*| + \varphi(1/4))\varphi(r)} \\ &\leq \frac{q_\sigma(0)\varphi(2r)}{q_\sigma(2\sigma t)\varphi(r)} \leq \frac{\bar{\Gamma}}{q(2t)}. \end{aligned}$$

For all  $r \in (1/4, 1)$ ,  $\varphi(r) \geq \varphi(1/4)$  and

$$\begin{aligned} \frac{\pi(\mathcal{B}(P_{\theta^*}, 2r))}{\pi(\mathcal{B}(P_{\theta^*}, r))} &\leq \frac{1}{v_\sigma(\{\theta \in \mathbb{R}, |\theta - \theta^*| \leq \varphi(r)\})} \\ &\leq \frac{1}{v_\sigma(\{\theta \in \mathbb{R}, |\theta - \theta^*| \leq \varphi(1/4)\})} \\ &\leq \frac{1}{2q_\sigma(|\theta^*| + \varphi(1/4))\varphi(1/4)} \\ &\leq \frac{1}{2q_\sigma(2\sigma t)\varphi(1/4)} \leq \frac{\bar{\Gamma}\sigma}{q(2t)}. \end{aligned}$$

We obtain that in any case, for all  $r \in (0, 1)$  and  $\theta^* \in [-\sigma t, \sigma t]$ ,

$$\log\left(\frac{\pi(\mathcal{B}(P_{\theta^*}, 2r))}{\pi(\mathcal{B}(P_{\theta^*}, r))}\right) \leq \max\left\{\log\left(\frac{\bar{\Gamma}(\sigma \vee 1)}{q(2t)}\right), \log 4\right\}. \tag{134}$$

The inequality is also clearly true for  $r \geq 1$  since then  $\pi(\mathcal{B}(P_{\theta^*}, 2r)) = \pi(\mathcal{B}(P_{\theta^*}, r)) = 1$ . Hence, for all  $r \geq a_1^{-1}\beta$

$$\begin{aligned} \frac{1}{n\gamma a_1 r} \log\left(\frac{\pi(\mathcal{B}(P_{\theta^*}, 2r))}{\pi(\mathcal{B}(P_{\theta^*}, r))}\right) &\leq \frac{1}{n\gamma\beta} \sup_{r>0} \log\left(\frac{\pi(\mathcal{B}(P_{\theta^*}, 2r))}{\pi(\mathcal{B}(P_{\theta^*}, r))}\right) \\ &\leq \frac{1}{n\gamma\beta} \max\left\{\log\left(\frac{\bar{\Gamma}(\sigma \vee 1)}{q(2t)}\right), \log 4\right\}. \end{aligned}$$

The right-hand side is not larger than  $\beta$  provided that it satisfies (50) and this lower bound is not smaller than  $1/\sqrt{n}$  since  $\gamma \leq 1$ . We conclude by using (15).

**10.4 Proof of Proposition 5**

Under our assumption on  $q$ , Assumption 6 is satisfied and

$$\bar{\Gamma} = 2^{1/s} \max\left\{q(0), 2^{(1/s)-1}\right\}.$$

Let  $t = (|\theta|/\sigma) \vee t_0$ . Then,  $\theta \in [-\sigma t, \sigma t]$ ,  $\nu_1([t, +\infty)) \leq 1/4$  and inequality (134) holds true. We deduce from (11) that

$$r_n(\beta, P_\theta) \leq \frac{1}{\gamma n a_1 \beta} \max \left\{ \log \left( \frac{\bar{\Gamma}(\sigma \vee 1)}{q(2t)} \right), \log 4 \right\}$$

and the result follows from our specific choices of  $a_1$ ,  $\gamma$  and  $\beta$ .

### 10.5 Proof of Corollary 3

We set for short  $\Theta = \Theta[\eta, \delta]$  with the parameters  $\eta$  and  $\delta$  defined by (60) and (61) respectively and also define

$$J_n = \exp \left[ \frac{(K^2 - 1)\gamma \tau^4 a_1^2 n \eta_n^2}{2(k + 1)} \right] \tag{135}$$

so that  $\mathcal{M}_n(K)$  contains the elements  $P = P_{(p, \mathbf{m}, \sigma)}$  of  $\mathcal{M}$  such that

$$|\log \sigma| \vee \left| \frac{\mathbf{m}}{\sigma} \right|_\infty \leq \log(1 + \delta) J_n.$$

Hereafter we fix  $P = P_{(p, \mathbf{m}, \sigma)} \in \mathcal{M}_n(K)$ . There exist  $\theta = \theta(P) = (\bar{Q}, \bar{\mathbf{m}}, \bar{\sigma}) \in \Theta$  with  $\bar{\sigma} = (1 + \delta)^{j_0}$ ,  $\bar{\mathbf{m}} = \bar{\sigma} \delta \mathbf{j}$ ,  $(j_0, \mathbf{j}) \in \mathbb{Z} \times \mathbb{Z}^k$  such that

$$\frac{\bar{\sigma}}{(1 + \delta)} \leq \sigma < \bar{\sigma} \quad \text{and} \quad \bar{m}_i = j_i \bar{\sigma} \delta \leq m_i < \bar{m}_i + \bar{\sigma} \delta, \tag{136}$$

for all  $i \in \{1, \dots, k\}$ . Consequently,

$$0 \leq \left(1 - \frac{\sigma}{\bar{\sigma}}\right) \leq \frac{\delta}{1 + \delta} < \delta \quad \text{and} \quad \left| \frac{\mathbf{m} - \bar{\mathbf{m}}}{\bar{\sigma}} \right|_\infty \leq \delta, \tag{137}$$

and we infer from (56) and (57) and the fact that the total variation loss is translation and scale invariant that  $P_\theta$  satisfies

$$\begin{aligned} \ell(P_{(p, \mathbf{m}, \sigma)}, P_\theta) &\leq \ell(P_{(p, \mathbf{m}, \sigma)}, P_{(\bar{Q}, \mathbf{m}, \sigma)}) + \ell(P_{(\bar{Q}, \mathbf{m}, \sigma)}, P_{(\bar{Q}, \bar{\mathbf{m}}, \bar{\sigma})}) \\ &\leq \ell(P_{(p, \mathbf{0}, 1)}, P_{(\bar{Q}, \mathbf{0}, 1)}) + \ell(P_{(\bar{Q}, \mathbf{0}, 1)}, P_{(\bar{Q}, \frac{\bar{\mathbf{m}} - \mathbf{m}}{\bar{\sigma}}, \frac{\bar{\sigma}}{\sigma})}) \\ &\leq \eta + \left[ A \left( \left| \frac{\mathbf{m} - \bar{\mathbf{m}}}{\bar{\sigma}} \right|_\infty^s + \left(1 - \frac{\sigma}{\bar{\sigma}}\right)^s \right) \right] \wedge 1 \\ &\leq \eta + 2A\delta^s = 2\eta. \end{aligned}$$

Besides, the parameters  $(j_0, \mathbf{j}) \in \mathbb{Z} \times \mathbb{Z}^k$  can be controlled in the following way. Using that  $\sigma \leq \bar{\sigma}$ , the inequality  $\log(1 + \delta) \leq \delta$  and (137), we obtain that for all

$i \in \{1, \dots, k\}$ ,

$$|j_i| = \left| \frac{\bar{m}_i}{\bar{\sigma}\delta} \right| = \frac{1}{\bar{\sigma}\delta} |\bar{m}_i - m_i + m_i| \leq \frac{1}{\bar{\sigma}\delta} \left[ \bar{\sigma}\delta + \sigma \left| \frac{m_i}{\sigma} \right| \right] \leq 1 + \frac{1}{\log(1 + \delta)} \left| \frac{m_i}{\sigma} \right|.$$

Besides,

$$\begin{aligned} j_0 &= \frac{\log \bar{\sigma}}{\log(1 + \delta)} = \frac{1}{\log(1 + \delta)} \left[ -\log \left( 1 + \frac{\sigma}{\bar{\sigma}} - 1 \right) + \log \sigma \right] \\ &\leq \frac{1}{\log(1 + \delta)} \left[ -\log \left( 1 - \frac{\delta}{1 + \delta} \right) + |\log \sigma| \right] \\ &= \frac{1}{\log(1 + \delta)} \left[ \log(1 + \delta) + |\log \sigma| \right] \leq 1 + \frac{|\log \sigma|}{\log(1 + \delta)} \end{aligned}$$

and using the inequality  $\log(1 + 2x) \leq 2 \log(1 + x)$ , which holds for all  $x \geq 0$ , we obtain that

$$j_0 \geq \frac{\log \sigma}{\log(1 + \delta)} \geq -\frac{|\log \sigma|}{\log(1 + \delta)} \geq -\left[ 1 + \frac{|\log \sigma|}{\log(1 + \delta)} \right].$$

Putting these inequalities together and using the fact that  $P \in \mathcal{M}_n(K)$ , we get

$$|(j_0, \mathbf{j})|_\infty \leq 1 + \frac{1}{\log(1 + \delta)} \left[ |\log \sigma| \vee \left| \frac{\mathbf{m}}{\sigma} \right|_\infty \right] \leq 1 + J_n. \tag{138}$$

For all  $r > 0$ ,  $e^{-L\theta} \leq \pi(\mathcal{B}(P_\theta, r)) \leq 1$  and these two inequalities together with the definition (60) of  $\eta$  and Assumption 7 imply that for all  $r > 0$

$$\begin{aligned} \frac{\pi(\mathcal{B}(P_\theta, 2r))}{\pi(\mathcal{B}(P_\theta, r))} &\leq \exp[L_\theta] \leq \exp \left[ \tilde{D}(\eta) + 2 \sum_{i=0}^k \left[ \frac{L}{2} + \log(1 + |j_i|) \right] \right] \\ &\leq \exp \left[ \gamma \tau^4 a_1^2 n \eta^2 + (k + 1) [L + 2 \log(1 + |(j_0, \mathbf{j})|_\infty)] \right]. \end{aligned}$$

Using (138), the definition (135) of  $J_n$  and the fact that  $\log(2 + x) \leq \log 3 + \log x$  for all  $x \geq 1$ , we derive that

$$\begin{aligned} \frac{\pi(\mathcal{B}(P_\theta, 2r))}{\pi(\mathcal{B}(P_\theta, r))} &\leq \exp \left[ \gamma \tau^4 a_1^2 n \eta^2 + (k + 1)L + 2(k + 1) \log(2 + J_n) \right], \\ &\leq \exp \left[ K^2 \gamma \tau^4 a_1^2 n \eta^2 + (k + 1)(L + \log 9) \right] \end{aligned}$$

and since  $\gamma = 1/6 \leq L' = L + \log 9 < 3.1$ ,

$$\frac{1}{n\beta a_1} \leq r_n(\beta, P_\theta) \leq \frac{1}{\gamma n\beta a_1} \left[ K^2 \gamma \tau^4 a_1^2 n \eta^2 + (k + 1)L' \right]$$

$$= \frac{1}{a_1\beta} \left[ K^2\tau^4 a_1^2 \eta^2 + \frac{(k+1)L'}{\gamma n} \right].$$

For the choice of  $\beta = \beta_n$  given by (62),

$$\beta \geq \sqrt{K^2\tau^4 a_1^2 \eta^2 + \frac{(k+1)L'}{\gamma n}} \geq \sqrt{\frac{k+1}{n}} \vee \frac{K\eta}{2}$$

hence,  $r_n(\beta, P_\theta) \leq a_1^{-1}\beta$  and  $P_\theta \in \mathcal{M}(\beta)$ . This implies that

$$\begin{aligned} \inf_{P' \in \mathcal{M}(\beta)} \ell(\bar{P}^*, P') + a_1^{-1}\beta &\leq \ell(\bar{P}^*, P_\theta) + a_1^{-1}\beta \\ &\leq \ell(\bar{P}^*, P) + \ell(P, P_\theta) + a_1^{-1}\beta \\ &\leq \ell(\bar{P}^*, P) + 2\eta + \left[ K\tau^2\eta + \frac{1}{a_1} \sqrt{\frac{(k+1)L'}{\gamma n}} \right], \end{aligned}$$

and the result follows by applying Corollary 1 and by using the fact that  $P$  is arbitrary in  $\mathcal{M}_n(K)$ .

### 10.6 Proof of Lemma 2

For all  $p \in \mathcal{M}_0$ ,  $\sigma \geq 1$  and  $\mathbf{m} \in \mathbb{R}^k$ , the supports of the functions  $\mathbf{x} \mapsto p(\mathbf{x}/\sigma)$  and  $\mathbf{x} \mapsto p((\mathbf{x} - \mathbf{m})/\sigma)$  are included in the set  $\mathcal{K} = [0, \sigma]^k \cup \{\mathbf{m} + \mathbf{x}, \mathbf{x} \in [0, \sigma]^k\}$  the Lebesgue measure of which is not larger than  $2\sigma^k$ . Consequently, using (66), we deduce that for all  $p \in \mathcal{M}_0$ ,  $\sigma \geq 1$  and  $\mathbf{m} \in \mathbb{R}^k$ ,

$$\begin{aligned} &\|P_{(p, \mathbf{0}, 1)} - P_{(p, \mathbf{m}, \sigma)}\| \\ &\leq \|P_{(p, \mathbf{0}, 1)} - P_{(p, \mathbf{0}, \sigma)}\| + \|P_{(p, \mathbf{0}, \sigma)} - P_{(p, \mathbf{m}, \sigma)}\| \\ &= \frac{1}{2} \int_{\mathbb{R}^k} \left| p(\mathbf{x}) - \frac{1}{\sigma^k} p\left(\frac{\mathbf{x}}{\sigma}\right) \right| d\mathbf{x} + \frac{1}{2\sigma^k} \int_{\mathbb{R}^k} \left| p\left(\frac{\mathbf{x}}{\sigma}\right) - p\left(\frac{\mathbf{x} - \mathbf{m}}{\sigma}\right) \right| d\mathbf{x} \\ &\leq \frac{1}{2} \int_{\mathbb{R}^k} \left| p(\mathbf{x}) - \frac{1}{\sigma^k} p(\mathbf{x}) \right| d\mathbf{x} + \frac{1}{2\sigma^k} \int_{\mathbb{R}^k} \left| p(\mathbf{x}) - p\left(\frac{\mathbf{x}}{\sigma}\right) \right| d\mathbf{x} \\ &\quad + \frac{1}{2\sigma^k} \int_{\mathbb{R}^k} \left| p\left(\frac{\mathbf{x}}{\sigma}\right) - p\left(\frac{\mathbf{x} - \mathbf{m}}{\sigma}\right) \right| d\mathbf{x} \\ &\leq \frac{1}{2} \int_{\mathbb{R}^k} \left| p(\mathbf{x}) - \frac{1}{\sigma^k} p(\mathbf{x}) \right| d\mathbf{x} + \frac{1}{2\sigma^k} \int_{[0, 1]^k} \left| p(\mathbf{x}) - p\left(\frac{\mathbf{x}}{\sigma}\right) \right| d\mathbf{x} \\ &\quad + \frac{1}{2\sigma^k} \int_{[0, \sigma]^k \setminus [0, 1]^k} \left| p\left(\frac{\mathbf{x}}{\sigma}\right) \right| d\mathbf{x} + \frac{1}{2\sigma^k} \int_{\mathcal{K}} \left| p\left(\frac{\mathbf{x}}{\sigma}\right) - p\left(\frac{\mathbf{x} - \mathbf{m}}{\sigma}\right) \right| d\mathbf{x} \\ &\leq \frac{1}{2} \left( 1 - \frac{1}{\sigma^k} \right) + \frac{1}{2\sigma^k} \int_{[0, 1]^k} L_1 \left( 1 - \frac{1}{\sigma} \right)^s |\mathbf{x}|^s d\mathbf{x} \end{aligned}$$

$$\begin{aligned}
 & + \frac{1}{2} \int_{[0,1]^k \setminus [0,1/\sigma]^k} |p(\mathbf{x})| d\mathbf{x} + \frac{L_1}{2\sigma^k} \int_{\mathcal{K}} \left| \frac{\mathbf{m}}{\sigma} \right|^s d\mathbf{x} \\
 \leq & \frac{1}{2} \left( 1 - \frac{1}{\sigma^k} \right) + \frac{L_1 k^{s/2}}{2\sigma^k} \left( 1 - \frac{1}{\sigma} \right)^s + \frac{L_0}{2} \left( 1 - \frac{1}{\sigma^k} \right) + L_1 \left| \frac{\mathbf{m}}{\sigma} \right|^s \\
 \leq & \frac{1}{2} \left[ 1 + L_1 k^{s/2} + L_0 \right] \left( 1 - \frac{1}{\sigma} \right)^s + L_1 \left| \frac{\mathbf{m}}{\sigma} \right|^s
 \end{aligned}$$

and (57) is therefore satisfied with  $A = L_1 \vee [(1 + L_1 k^{s/2} + L_0)/2]$ .

### 10.7 Proof of Lemma 3

By doing the change of variables  $u = x - m$  in (68) if ever necessary, we may assume with no loss of generality that  $m > 0$ . Then, since  $p$  is nonincreasing in  $(0, +\infty)$  and vanishes elsewhere  $p(x - m) \geq p(x)$  for all  $x \geq m$  and  $p(x) \geq p(x - m) = 0$  for all  $x \in (0, m)$ . Consequently,

$$\begin{aligned}
 \int_{\mathbb{R}} |p(x) - p(x - m)| dx & = \int_0^m p(x) dx + \int_m^{+\infty} [p(x - m) - p(x)] dx \\
 & = 2 \int_0^m p(x) dx + \int_m^{+\infty} p(x - m) dx - \int_0^{+\infty} p(x) dx \\
 & \leq 2mB + 1 - 1,
 \end{aligned}$$

and we obtain (68).

Since  $\sigma \geq 1$ ,  $p(x/\sigma) \geq p(x)$  and  $p(x)/\sigma \leq p(x)$  for all  $x > 0$ . Hence,

$$\begin{aligned}
 & \int_{\mathbb{R}} \left| \frac{1}{\sigma} p\left(\frac{x}{\sigma}\right) - p(x) \right| dx \\
 & \leq \int_{\mathbb{R}} \left| \frac{1}{\sigma} p\left(\frac{x}{\sigma}\right) - \frac{1}{\sigma} p(x) \right| dx + \int_{\mathbb{R}} \left| \frac{1}{\sigma} p(x) - p(x) \right| dx \\
 & = \frac{1}{\sigma} \int_{\mathbb{R}} \left( p\left(\frac{x}{\sigma}\right) - p(x) \right) dx + \int_{\mathbb{R}} \left( p(x) - \frac{1}{\sigma} p(x) \right) dx \\
 & = 2 \left( 1 - \frac{1}{\sigma} \right),
 \end{aligned}$$

which leads to (67).

Finally, by combining (68) and (67) we deduce that for all  $m \in \mathbb{R}$  and  $\sigma \geq 1$

$$\begin{aligned}
 & \frac{1}{2} \int_{\mathbb{R}} \left| \frac{1}{\sigma} p\left(\frac{x - m}{\sigma}\right) - p(x) \right| dx \\
 & = \frac{1}{2} \int_{\mathbb{R}} \left| \frac{1}{\sigma} p\left(\frac{x - m}{\sigma}\right) - \frac{1}{\sigma} p\left(\frac{x}{\sigma}\right) \right| dx + \frac{1}{2} \int_{\mathbb{R}} \left| \frac{1}{\sigma} p\left(\frac{x}{\sigma}\right) - p(x) \right| dx \\
 & = \frac{1}{2} \int_{\mathbb{R}} \left| p\left(u - \frac{m}{\sigma}\right) - p(u) \right| du + \frac{1}{2} \int_{\mathbb{R}} \left| \frac{1}{\sigma} p\left(\frac{x}{\sigma}\right) - p(x) \right| dx
 \end{aligned}$$

$$\leq B \left| \frac{m}{\sigma} \right| + \left( 1 - \frac{1}{\sigma} \right)$$

which yields to (69).

### 10.8 Proof of Proposition 6

This proposition is a consequence of Corollary 2. Let us first check that the assumptions of this corollary are satisfied. For all  $S \in \mathcal{S}$ , the mapping  $\theta \mapsto h(S, P_\theta)$  is continuous because of (71). It is therefore measurable and it follows from the definition of the algebra  $\mathcal{A}$  that Assumption 1 is satisfied. Since the mapping  $(x, \theta) \mapsto p(x, \theta)$  is measurable, so are the mappings

$$\mathbf{p} : \begin{cases} (\mathbb{R}^k \times E \times E \longrightarrow \mathbb{R}_+ \\ (x, \theta, \theta') \longmapsto (p_\theta(x), p_{\theta'}(x)). \end{cases}$$

and  $(x, \theta, \theta') \mapsto \psi(\sqrt{p_{\theta'}(x)/p_\theta(x)})$ , since  $\psi$  is measurable. We deduce that  $(x, P, P') \mapsto t_{(P, P')}(x)$  is measurable on  $(E \times \mathcal{M} \times \mathcal{M}, \mathcal{E} \otimes \mathcal{A} \otimes \mathcal{A})$  which proves that Assumption 3-(i) holds true. The requirements of Corollary 2 are therefore satisfied and we may apply it. In order to evaluate the quantity  $r_n(\beta, P_\theta)$  for  $\theta \in \mathbb{R}^k$ , we use the following lemma the proof of which is postponed to Sect. 10.9.

**Lemma 10** *Let  $\theta \in [-R, R]^k$ . For all  $m \subset \{1, \dots, k\}$  and  $r > 0$*

$$v_m \left( \left\{ \theta' \in \mathbb{R}^k, |\theta' - \theta|_\infty \leq r \right\} \right) = \begin{cases} \frac{1}{2^{|m|}} \prod_{i \in m} \left[ \left( 1 - \frac{|\theta_i|}{R} \right) \wedge \frac{r}{R} + \left( 1 + \frac{|\theta_i|}{R} \right) \wedge \frac{r}{R} \right] & \text{if } |\theta_i| \leq r \text{ for all } i \notin m \\ 0 & \text{otherwise,} \end{cases}$$

with the convention  $\prod_{\emptyset} = 1$ . In particular, if  $\theta \in \Theta_m(R)$  and

$$v_m \left( \left\{ \theta' \in \mathbb{R}^k, |\theta' - \theta|_\infty \leq r \right\} \right) \geq \frac{1}{2^{|m|}} \left( \frac{r}{R} \wedge 1 \right)^{|m|} \tag{139}$$

and for all  $K > 1$

$$\frac{v_m \left( \left\{ \theta' \in \mathbb{R}^k, |\theta' - \theta|_\infty \leq Kr \right\} \right)}{v_m \left( \left\{ \theta' \in \mathbb{R}^k, |\theta' - \theta|_\infty \leq r \right\} \right)} \leq K^{|m|}. \tag{140}$$

Let us set  $B = B_k$  for short and define  $m^*$  as the subset of  $\{1, \dots, k\}$  that minimizes over those  $m \subset \{1, \dots, k\}$  the mapping

$$m \mapsto \inf_{\theta \in \Theta_m(R)} \ell(\overline{P}^*, P_\theta) + \frac{|m| \log(2kR(nB)^{1/s}) + 1}{\gamma n \beta a_1}.$$

Finally, let  $\theta^*$  for some arbitrary element of  $\Theta_{m^*}(R)$ . It follows from (71) and (139) that for all  $r > 0$ ,

$$\begin{aligned}
 1 &\geq \pi_m(\mathcal{B}(P_{\theta^*}, r)) \\
 &= v_m\left(\left\{\theta \in \mathbb{R}^k, h^2(P_{\theta^*}, P_{\theta}) \leq r\right\}\right) \\
 &\geq v_m\left(\left\{\theta \in \mathbb{R}^k, \|\theta - \theta^*\|_{\infty} \leq (r/B)^{1/s}\right\}\right) \\
 &\geq \frac{1}{2^{|m|}} \left(\frac{(r/B)^{1/s}}{R} \wedge 1\right)^{|m|} \geq \frac{1}{2^{|m|}} \left(\frac{(r \wedge 1)^{1/s}}{RB^{1/s}}\right)^{|m|}, \tag{141}
 \end{aligned}$$

where the last inequality holds true under the assumption that  $RB^{1/s} \geq 1$ .

We deduce from (141) that for all  $r > 0$

$$\begin{aligned}
 \frac{\pi(\mathcal{B}(P_{\theta^*}, 2r))}{\pi(\mathcal{B}(P_{\theta^*}, r))} &\leq \frac{1}{\pi(\mathcal{B}(P_{\theta^*}, r))} \\
 &\leq \frac{1}{\sum_{m \subset \{1, \dots, k\}} e^{-L_m} v_m(\{\theta \in \mathbb{R}^k, \|\theta - \theta^*\|_{\infty} \leq (r/B)^{1/s}\})} \\
 &\leq \frac{e^{L_{m^*}}}{v_{m^*}(\{\theta \in \Theta_{m^*}, \|\theta - \theta^*\|_{\infty} \leq (r/B)^{1/s}\})} \\
 &\leq \exp\left[L_{m^*} + |m^*| \log\left(\frac{2RB^{1/s}}{(r \wedge 1)^{1/s}}\right)\right] \\
 &= \exp\left[|m^*| \log(2kRB^{1/s}) + k \log\left(1 + \frac{1}{k}\right) + \frac{|m^*|}{s} \log\left(\frac{1}{r} \vee 1\right)\right]. \tag{142}
 \end{aligned}$$

Provided that

$$r \geq \frac{|m^*| \log(2kR(nB)^{1/s}) + 1}{\gamma n \beta a_1} \geq \frac{1}{n},$$

we obtain

$$\begin{aligned}
 &|m^*| \log(2kRB^{1/s}) + k \log\left(1 + \frac{1}{k}\right) + \frac{|m^*|}{s} \log\left(\frac{1}{r} \vee 1\right) \\
 &\leq |m^*| \log(2kRB^{1/s}) + k \log\left(1 + \frac{1}{k}\right) + |m^*| \log(n^{1/s}) \\
 &\leq |m^*| \log(2kR(nB)^{1/s}) + 1 \leq \gamma n \beta a_1 r
 \end{aligned}$$

and deduce from (142) that  $r_n(\beta, P_{\theta^*})$  defined by (11) satisfies

$$\frac{1}{n\beta a_1} \leq r_n(\beta, P_{\theta^*}) \leq \frac{|m^*| \log(2kR(nB)^{1/s}) + 1}{\gamma n \beta a_1}.$$

Applying Corollary 2, we obtain that for some numerical constant  $\kappa'_0 > 0$ ,

$$\mathbb{E} \left[ \widehat{\pi}_X \left( \mathcal{C}_{\mathcal{B}}(\overline{P}^*, \kappa'_0 r(m^*, \theta^*)) \right) \right] \leq 2e^{-\xi}$$

with

$$r(m^*, \theta^*) = \ell(\overline{P}^*, P_{\theta^*}) + \frac{|m^*| \log(2kR(nB)^{1/s}) + \xi}{\gamma n \beta a_1}.$$

Finally, the conclusion follows from the definition of  $m^*$  and the fact that  $\theta^*$  is arbitrary in  $\Theta_{m^*}(R)$ .

**10.9 Proof of Lemma 10**

Let  $\theta \in \mathbb{R}$  and  $\nu$  be the uniform distribution on  $[-R, R]$ . For all  $\theta \in [-R, R]$  and  $r > 0$ ,

$$\begin{aligned} \nu([\theta - r, \theta + r]) &= \frac{1}{2R} [(\theta + r) \wedge R - (\theta - r) \vee (-R)]_+ \\ &= \frac{1}{2R} [(r + \theta) \wedge R + (r - \theta) \wedge R]_+ \\ &= \frac{1}{2R} [(r + |\theta|) \wedge R + (r - |\theta|) \wedge R]_+ \\ &= \frac{1}{2} \left[ \left(1 - \frac{|\theta|}{R}\right) \wedge \frac{r}{R} + \left(1 + \frac{|\theta|}{R}\right) \wedge \frac{r}{R} \right]. \end{aligned}$$

Let now  $\theta \in \mathbb{R}^k$  such that  $|\theta|_\infty \leq R$ . For all  $m \subset \{1, \dots, k\}$ ,  $m \neq \emptyset$ ,

$$\nu_m(\{\theta' \in \Theta_m, |\theta' - \theta|_\infty \leq r\}) = 0$$

if there exists  $i \notin m$  such that  $|\theta_i| > r$ . Otherwise

$$\begin{aligned} \nu_m(\{\theta' \in \mathbb{R}^k, |\theta' - \theta|_\infty \leq r\}) &= \nu_m\left(\{\theta' \in \Theta_m, \max_{i \in m} |\theta'_i - \theta_i| \leq r\}\right) \\ &= \prod_{i \in m} \nu([\theta_i - r, \theta_i + r]) \\ &= \frac{1}{2^{|m|}} \prod_{i \in m} \left[ \left(1 - \frac{|\theta_i|}{R}\right) \wedge \frac{r}{R} + \left(1 + \frac{|\theta_i|}{R}\right) \wedge \frac{r}{R} \right]. \end{aligned}$$

If  $m = \emptyset$ ,

$$\nu_\emptyset(\{\theta' \in \mathbb{R}^k, |\theta' - \theta|_\infty \leq r\}) = \mathbb{1}_{|\theta|_\infty \leq r}.$$



Let us now turn to the proof of (140). Since  $\theta \in \Theta_m(R)$ , for all  $K' \in \{1, K\}$

$$\begin{aligned} & \nu_m \left( \left\{ \theta' \in \mathbb{R}^k, \|\theta' - \theta\|_\infty \leq K'r \right\} \right) \\ &= \nu_m \left( \left\{ \theta' \in \Theta_m, \max_{i \in m} |\theta'_i - \theta_i| \leq K'r \right\} \right) \\ &= \prod_{i \in m} \nu([\theta_i - K'r, \theta_i + K'r]), \end{aligned}$$

It is therefore enough to show that for all  $r > 0$  and  $\theta \in [0, R]$

$$\Delta(r) = \frac{\nu([\theta - Kr, \theta + Kr])}{\nu([\theta - r, \theta + r])} \leq K.$$

This is what we do now by distinguishing between several cases.

When  $\theta + Kr \leq R, \theta - Kr \geq 2\theta - R \geq -R$  and consequently,  $\Delta(r) = K$ . When  $\theta + Kr > R$  and  $-R \leq \theta - Kr$ ,

$$\Delta(r) = \frac{R - (\theta - Kr)}{(\theta + r) \wedge R - (\theta - r)} = \begin{cases} \frac{R - \theta + Kr}{R - \theta + r} & \text{when } \theta + r > R \\ \frac{R - \theta + Kr}{2r} & \text{when } \theta + r \leq R, \end{cases}$$

and the conclusion follows from the facts that  $0 \leq R - \theta \leq Kr$ . When  $\theta + Kr > R$  and  $\theta - Kr < -R, r \geq (\theta + R)/K \geq R/K$ , hence  $R + r - \theta \geq 2R/K$  and  $R \leq Kr$ . Consequently,

$$\begin{aligned} \Delta(r) &= \frac{2R}{(\theta + r) \wedge R - (\theta - r) \vee (-R)} \\ &= \begin{cases} \frac{2R}{2R} = 1 & \text{when } \theta + r > R \text{ and } \theta - r < -R \\ \frac{2R}{R + r - \theta} \leq K & \text{when } \theta + r > R \text{ and } \theta - r \geq -R \\ \frac{2R}{2r} \leq K & \text{when } \theta + r \leq R, \end{cases} \end{aligned}$$

which concludes the proof.

### 10.10 Proof of Proposition 8

Let  $\varepsilon$  be a small enough positive number. Since  $q$  is continuous and positive at  $\theta^*$  and since  $\mathcal{K}$  has a nonempty interior, there exists  $z^* > 0$  such that  $\Theta^* = \mathcal{B}_*(\theta^*, z^*) \subset \mathcal{K}$ ,

$$0 < \underline{b}^* \leq q(\theta) \leq \bar{b}^* \quad \text{with } \bar{b}^*/\underline{b}^* \leq 1 + \varepsilon, \tag{143}$$

for all  $\theta \in \Theta^*$  and

$$(1 - \varepsilon)|\theta - \theta^*|_*^s \leq \ell(\theta, \theta^*) \leq (1 + \varepsilon)|\theta - \theta^*|_*^s. \tag{144}$$

In particular,  $\nu(\Theta^*) > 0$  and we may define the distribution  $\nu^* = \nu(\cdot \cap \Theta^*)/\nu(\Theta^*)$  on  $\Theta^*$  with density  $q^* = q\mathbb{1}_{\Theta^*}/\nu(\Theta^*)$ . Let  $\mathcal{M}^* = \{P_\theta, \theta \in \Theta^*\}$  and  $\pi^*$  be the prior on  $\mathcal{M}^*$  associated with  $\nu^*$ . The parameter space  $\Theta^*$  is convex and it follows from (144) that  $(\Theta^*, \theta^*, \ell, \nu^*)$  satisfy Assumption 8-(i) with  $\bar{a} = 1 + \varepsilon, \underline{a} = 1 - \varepsilon$ . Besides, it follows from (143) that the density  $q^*$  satisfies condition (77) on  $\Theta^*$ . We may apply Proposition 7 and deduce that for the model  $(\mathcal{M}^*, \pi^*), r_n^* = r_n^*(\beta, P_{\theta^*})$  is not larger than  $\kappa_0^*k/(\beta n)$  with

$$\kappa_0^* = \frac{1}{a_1\gamma} \left\{ \left[ 1 + \frac{\log [2(1 + \varepsilon)/(1 - \varepsilon)]}{s \log 2} \right] \log (2(1 + \varepsilon)) \right\} \vee 1 < \frac{(1 + s^{-1})}{a_1\gamma}$$

for  $\varepsilon$  small enough. Consequently, by definition of  $r_n^*$ , for all  $r \geq r_n^*$

$$\begin{aligned} \pi^*(\mathcal{B}(P_{\theta^*}, 2r)) &= \frac{1}{\nu(\Theta^*)} \nu(\{\theta \in \Theta, \ell(\theta, \theta^*) \leq 2r\} \cap \Theta^*) \\ &\leq \frac{\exp(\gamma n \beta a_1 r)}{\nu(\Theta^*)} \nu(\{\theta \in \Theta, \ell(\theta, \theta^*) \leq r\} \cap \Theta^*) \\ &\leq \frac{\exp(\gamma n \beta a_1 r)}{\nu(\Theta^*)} \nu(\{\theta \in \Theta, \ell(\theta, \theta^*) \leq r\}). \end{aligned} \tag{145}$$

Let  $r_1 = [(z^*)^s \underline{a}_K] \wedge \eta/2$ . If  $r \in (0, r_1)$  and the parameter  $\theta \in \Theta$  satisfies  $\ell(\theta, \theta^*) \leq 2r$ , then  $\ell(\theta, \theta^*) < \eta$  and  $\theta$  necessarily belongs to  $\mathcal{K}$  under Assumption 9-(ii). Applying (79) we deduce that for such a parameter  $\theta \in \Theta$

$$\underline{a}_K |\theta - \theta^*|_*^s \leq \ell(\theta, \theta^*) \leq 2r < 2r_1 \leq \underline{a}_K (z^*)^s,$$

which implies that  $\theta \in \Theta^*$ . For  $n$  large enough,  $r_n^* = \kappa_0^*k/n < r_1$  and for  $r \in (r_n^*, r_1)$  we may therefore write, using (145),

$$\begin{aligned} \pi(\mathcal{B}(P_{\theta^*}, 2r)) &= \nu(\{\theta \in \Theta, \ell(\theta, \theta^*) \leq 2r\}) \\ &= \nu(\{\theta \in \Theta, \ell(\theta, \theta^*) \leq 2r\} \cap \Theta^*) \\ &\leq \exp(\gamma n \beta a_1 r) \nu(\{\theta \in \Theta, \ell(\theta, \theta^*) \leq r\}) \\ &= \exp(\gamma n \beta a_1 r) \pi(\mathcal{B}(P_{\theta^*}, r)). \end{aligned} \tag{146}$$

Since  $q$  is bounded away from 0 in a neighbourhood of  $\theta^*$ ,  $\pi(\mathcal{B}(P_{\theta^*}, r_1)) > 0$  and we may also write that for  $r \geq r_1$  and  $n$  large enough

$$\begin{aligned} \pi(\mathcal{B}(P_{\theta^*}, r)) &\geq \pi(\mathcal{B}(P_{\theta^*}, r_1)) \\ &= \exp[\log \pi(\mathcal{B}(P_{\theta^*}, r_1)) + \gamma n \beta a_1 r_1 - \gamma n \beta a_1 r] \\ &\geq \exp[-\gamma n \beta a_1 r_1] \geq \exp[-\gamma n \beta a_1 r] \end{aligned}$$

$$\geq \exp[-\gamma n \beta a_1 r] \pi(\mathcal{B}(P_{\theta^*}, 2r)). \quad (147)$$

Putting (146) and (147) together we obtain that for  $n$  large enough

$$\pi(\mathcal{B}(P_{\theta^*}, 2r)) \leq \exp(\gamma n \beta a_1 r) \pi(\mathcal{B}(P_{\theta^*}, r)) \quad \text{for all } r \geq r_n^*$$

and consequently that  $r_n(\beta, P_{\theta^*}) \leq r_n^* = \kappa_0^* k/n$ .

**Acknowledgements** The author is grateful to Lucien Birgé and the two anonymous referees for their support and suggestions, which contributed to improving a previous version of the present paper.

## References

1. Alquier, P.: PAC-Bayesian bounds for randomized empirical risk minimizers. *Math. Methods Stat.* **17**(4), 279–304 (2008)
2. Atchadé, Y.A.: On the contraction properties of some high-dimensional quasi-posterior distributions. *Ann. Stat.* **45**(5), 2248–2273 (2017)
3. Audibert, J.-Y., Catoni, O.: Linear regression through PAC-Bayesian truncation. [arXiv:1010.0072](https://arxiv.org/abs/1010.0072) (2011)
4. Baraud, Y.: Tests and estimation strategies associated to some loss functions. *Probab. Theory Relat. Fields* **180**(3), 799–846 (2021)
5. Baraud, Y., Birgé, L.: Rho-estimators revisited: general theory and applications. *Ann. Stat.* **46**(6B), 3767–3804 (2018)
6. Baraud, Y., Birgé, L.: Robust Bayes-like estimation: Rho-Bayes estimation. *Ann. Stat.* **48**(6), 3699–3720 (2020)
7. Baraud, Y., Birgé, L., Sart, M.: A new method for estimation and model selection:  $\rho$ -estimation. *Invent. Math.* **207**(2), 425–517 (2017)
8. Bhattacharya, A., Pati, D., Yang, Y.: Bayesian fractional posteriors. *Ann. Stat.* **47**(1), 39–66 (2019)
9. Birgé, L.: Approximation dans les espaces métriques et théorie de l'estimation. *Z. Wahrsch. Verw. Gebiete* **65**(2), 181–237 (1983)
10. Birgé, L.: An alternative point of view on Lepski's method. In: *State of the Art in Probability and Statistics (Leiden, 1999)*, Volume 36 of IMS Lecture Notes Monograph Series, pp. 113–133. Institute of Mathematical Statistics, Beachwood (2001)
11. Birgé, L.: Model selection via testing: an alternative to (penalized) maximum likelihood estimators. *Ann. Inst. H. Poincaré Probab. Stat.* **42**(3), 273–325 (2006)
12. Birgé, L.: About the non-asymptotic behaviour of Bayes estimators. *J. Stat. Plan. Inference* **166**, 67–77 (2015)
13. Birgé, L.: About the non-asymptotic behaviour of Bayes estimators. *J. Stat. Plan. Inference* **166**, 67–77 (2015)
14. Birgé, L., Massart, P.: Minimum contrast estimators on sieves: exponential bounds and rates of convergence. *Bernoulli* **4**(3), 329–375 (1998)
15. Birman, Mv., Solomjak, M.Z.: Piecewise polynomial approximations of functions of classes  $W_p^\alpha$ . *Mat. Sb. (N.S.)* **73**(115), 331–355 (1967)
16. Bissiri, P.G., Holmes, C.C., Walker, S.G.: A general framework for updating belief distributions. *J. R. Stat. Soc. Ser. B. Stat. Methodol.* **78**(5), 1103–1130 (2016)
17. Catoni, O.: Statistical learning theory and stochastic optimization. In: *Lecture Notes from the 31st Summer School on Probability Theory Held in Saint-Flour, July 8–25, 2001*. Springer, Berlin (2004)
18. Chernozhukov, V., Hong, H.: An MCMC approach to classical estimation. *J. Econom.* **115**(2), 293–346 (2003)
19. Ghosal, S., Ghosh, J.K., van der Vaart, A.W.: Convergence rates of posterior distributions. *Ann. Stat.* **28**(2), 500–531 (2000)
20. Ibragimov, I.A., Has'minskiĭ, R.Z.: *Statistical Estimation. Asymptotic Theory*, vol. 16. Springer, New York (1981)

21. Jiang, W., Tanner, M.A.: Gibbs posterior for variable selection in high-dimensional classification and data mining. *Ann. Stat.* **36**(5), 2207–2231 (2008)
22. Le Cam, L.: Convergence of estimates under dimensionality restrictions. *Ann. Stat.* **1**, 38–53 (1973)
23. Massart, P.: Concentration Inequalities and Model Selection, Volume 1896 of Lecture Notes in Mathematics. Springer, Berlin. Lectures from the 33rd Summer School on Probability Theory held in Saint-Flour, July 6–23, 2003 (2007)
24. van der Vaart, A.W.: Asymptotic Statistics, Volume 3 of Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, Cambridge (1998)
25. Whittaker, E.T., Watson, G.N.: A Course of Modern Analysis. Cambridge Mathematical Library. Cambridge University Press, Cambridge. An introduction to the general theory of infinite processes and of analytic functions; with an account of the principal transcendental functions, Reprint of the fourth (1927) edition (1996)

**Publisher's Note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Springer Nature or its licensor (e.g. a society or other partner) holds exclusive rights to this article under a publishing agreement with the author(s) or other rightsholder(s); author self-archiving of the accepted manuscript version of this article is solely governed by the terms of such publishing agreement and applicable law.