

# FROM ROBUST TESTS TO BAYES-LIKE POSTERIOR DISTRIBUTIONS

YANNICK BARAUD

ABSTRACT. In the Bayes paradigm and for a given loss function, we propose the construction of a new type of posterior distributions for estimating the law of an  $n$ -sample. The loss functions we have in mind are based on the total variation distance, the Hellinger distance as well as some  $\mathbb{L}_j$ -distances. We prove that, with a probability close to one, this new posterior distribution concentrates its mass in a neighbourhood of the law of the data, for the chosen loss function, 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. For the total variation and squared Hellinger losses, we also show that the posterior distribution keeps its concentration properties when the data are only independent, hence not necessarily i.i.d., provided that most of their marginals are close enough to some probability distribution around which the prior puts enough mass. The posterior distribution is therefore also stable with respect to the equidistribution assumption. We illustrate these results by several applications. We consider the problems of estimating a location parameter or both the location and the scale of a density in a nonparametric framework. Finally, we also tackle the problem of estimating a density, with the squared Hellinger loss, in a high-dimensional parametric model under some sparsity conditions. The results established in this paper are non-asymptotic and provide, as much as possible, explicit constants.

## 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. Our aim is to estimate  $P^*$  from the observation of  $\mathbf{X} = (X_1, \dots, X_n)$  and to evaluate the performance of an estimator with values

---

*Date:* July, 1st 2021.

*1991 Mathematics Subject Classification.* Primary 62G05, 62G35, 62F35, 62F15.

*Key words and phrases.* 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 N° 811017.

in  $\mathcal{M}$ , we introduce a loss function  $\ell$  from  $\mathcal{P} \times \mathcal{M}$  with values in  $\mathbb{R}_+$ , where  $\mathcal{P}$  denotes a suitable set of probabilities containing  $P^*$ .

The loss functions we have in mind are based on distances such as the Hellinger or the total variation one. Given two probabilities  $P, Q$  on  $(E, \mathcal{E})$ , we recall that the total variation  $\|P - Q\|$  and the squared Hellinger distance  $h^2(P, Q)$  between  $P$  and  $Q$  are respectively given by the formulas

$$(1) \quad \|P - Q\| = \frac{1}{2} \int_E \left| \frac{dP}{d\mu} - \frac{dQ}{d\mu} \right| d\mu$$

and

$$(2) \quad h^2(P, Q) = \frac{1}{2} \int_E \left( \sqrt{\frac{dP}{d\mu}} - \sqrt{\frac{dQ}{d\mu}} \right)^2 d\mu$$

where  $\mu$  dominates both  $P$  and  $Q$ , the result being independent of the choice of the reference measure  $\mu$ .

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

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

This means that with a probability close to 1, a point which is randomly drawn according to our (random) distribution  $\hat{\pi}_{\mathbf{X}}$  is likely to estimate  $P^*$  with an accuracy (with respect to the loss  $\ell$ ) not larger than  $r$ .

The classical Bayes posterior distribution has been studied at length by van der Vaart and his co-authors — see for example Ghosal, Ghosh and van der Vaart (2000). They show that it concentrates around  $P^*$  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 is unfortunately quite restrictive since such sets may be empty, and the condition therefore unsatisfied, when the probabilities in  $\mathcal{M}$  are not equivalent. The situation is even worse when the probability  $P^*$  does not belong to  $\mathcal{M}$ , even in the favourable situation where it lies close to it, since the quantity  $K(P^*, P)$  can be very large and possibly infinite even when the total variation distance between  $P^*$  and  $P$  is very small. It is actually well-known that Bayes estimators are not robust with respect to a misspecification of the model. This weakness is due to the fact that they are based on the log-likelihood function which can be unstable when an outlier belongs to the data set.

In order to overcome this issue and propose more stable posterior distributions, some authors have replaced the log-likelihood function by another

one, leading to *quasi-posterior distributions*. This is the strategy that is proposed in Baraud and Birgé (2020) (see also the references therein). The authors proposed 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 the series of papers Baraud *et al* (2017) and Baraud and Birgé (2018). The  $\rho$ -posterior distribution preserves some of the nice features of the classical Bayes one but also possesses the property of being robust with respect to the presence of outliers and contaminated data among the sample. They show that it concentrates on a Hellinger ball around  $P^*$  as soon as the prior puts enough mass around a point which is close enough to  $P^*$ . This key property confers to the  $\rho$ -posterior distribution the robustness the authors were looking for.

However, the  $\rho$ -posterior distribution is difficult to compute and therefore mainly of theoretical interest. It provides a benchmark to compare against. These difficulties are linked to the calculations of some suprema of empirical processes that are involved in the definition of the density of the  $\rho$ -posterior distribution. More deceiving is the fact that the authors do not show any improvement of their Bayes-like approach as compared to the frequentist one based on  $\rho$ -estimation. For a suitable choice of the prior, an estimator based on the  $\rho$ -posterior distribution would satisfy similar risks bound as those established for  $\rho$ -estimators. As a consequence,  $\rho$ -Bayes estimators do not seem to benefit from any gain that would result from a good choice of a prior as compared to the frequentist approach that presumes nothing.

Closer to our approach are the aggregation methods and PAC-Bayesian technics that have been popularized by Olivier Catoni in statistical learning (see Catoni (2004)). This approach has mainly been applied for the purpose of empirical risk minimization and statistical learning (see for example Alquier (2008)) and our aim is to extend it toward a versatile tool that can solve our estimation problem for various loss functions simultaneously.

The problem of designing a good estimator of  $P^*$  for a given loss function  $\ell$  was solved in a frequentist way in Baraud (2021). 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 (1973), Birgé (1983) and Birgé (2006). The aim of the present paper is to transpose this theory from the frequentist to the Bayesian paradigm.

For very general models  $\mathcal{M}$ , we prove that  $\hat{\pi}_{\mathbf{X}}$  concentrates around  $P^*$  (for the desired loss) at a rate  $r = r(n)$  which usually corresponds to the minimax one when  $\pi$  puts enough mass around  $P^*$  or, at least, around an element  $\bar{P}$  which is close enough to  $P^*$ . We therefore show that the posterior distribution  $\hat{\pi}_{\mathbf{X}}$  enjoys some optimality and robustness properties with respect to the choices of  $\mathcal{M}$  and of the prior  $\pi$ . In fact, we also show

that the posterior distribution is robust to the equidistribution assumption we started from. In particular, when  $\ell$  is the total variation or the squared Hellinger loss, the concentration properties of  $\widehat{\pi}_{\mathbf{X}}$  remain valid as long as the data are independent and most of their marginals lie close enough to  $\bar{P}$ . This result contrasts sharply with the instability of the classical Bayesian posterior distribution that we mentioned earlier.

Quite surprisingly, the concentration properties that we establish here require almost no assumption on  $P^*$  and  $\mathcal{M}$ . They mostly depend on the prior  $\pi$  and the loss function  $\ell$  that have been chosen. More precisely, for a suitable element  $P$  which belongs to the model  $\mathcal{M}$  and lies close enough to  $P^*$ , these properties depend on the minimal value of  $r$  over which the ratio  $\pi(\mathcal{B}(P, 2r))/\pi(\mathcal{B}(P, r))$  becomes large enough. This ratio was introduced earlier in Birgé (2015) for the purpose of analyzing the behaviour of the classical posterior distribution in the Bayes paradigm. In our Bayes-like paradigm, we show that the choice of the prior completely encapsulates the complexity of the model  $\mathcal{M}$ . In particular, no assumption on the VC nor on the metric dimension of the model  $\mathcal{M}$  is required. From this point of view, the results we establish here are of different nature than those obtained in the frequentist and Bayesian paradigms in Baraud (2021) and Baraud and Birgé (2020) which do require such assumptions.

Another difference with Baraud and Birgé (2020) lies in the construction of the posterior distribution. In the present paper, it does not involve any suprema of empirical processes but only integrals. It is therefore easier to compute even though the calculations of these integrals may not be necessarily easy, especially in high dimension.

The present paper is organized as follows. We present our statistical setting in Section 2. Unlike what has been described in this introduction, we actually consider independent but not necessarily i.i.d. data in order to analyse the behaviour of the posterior distribution with respect to a possible departure from equidistribution. Our main assumptions on the loss function are given and commented on in Section 3. In the remaining part of the paper, we shall mainly focus on the total variation and squared Hellinger losses. The construction of the posterior and its properties are presented in Section 4. Applications can be found in Section 5. There, we consider the problems of estimating a density in a location-scale family as well as that of a high-dimensional parameter in a parametric model under a sparsity constraint. We also show how our estimation strategy may lead to unusual rates of convergence for estimating a translation parameter in a non-regular statistical model. Finally, Section 6 is devoted to the proofs of the main theorems and Section 7 to the other proofs.

2. THE STATISTICAL SETTING

Let  $\mathbf{X} = (X_1, \dots, X_n)$  be an  $n$ -tuple of independent random variables with values in a measurable space  $(E, \mathcal{E})$  and joint distribution  $\mathbf{P}^* = \bigotimes_{i=1}^n P_i^*$ . The probabilities  $P_i^*$  are assumed to belong to a given set  $\mathcal{P}$  of probability measures on  $(E, \mathcal{E})$ . 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  $\mathbf{X}$ . To do so, we introduce a family  $\mathcal{M}$  that consists of candidate probabilities or merely finite signed measures. We endow  $\mathcal{M}$  with a  $\sigma$ -algebra  $\mathcal{A}$  and a probability measure  $\pi$ , that we call *a priori*, 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 most of the  $P_i^*$  are close enough to some point  $\bar{P}$  in  $\mathcal{M}$  around which the prior  $\pi$  puts enough mass. In order to be more specific, we introduce a loss function  $\ell$ , which is a mapping from  $(\mathcal{P} \cup \mathcal{M}) \times \mathcal{M}$  into  $\mathbb{R}_+$ , and write

$$\ell(\mathbf{P}^*, Q) = \frac{1}{n} \sum_{i=1}^n \ell(P_i^*, Q) \quad \text{for all } Q \in \mathcal{M}.$$

In order to avoid trivialities, we assume that  $\ell$  is not constantly equal to 0 on  $(\mathcal{M} \cup \mathcal{P}) \times \mathcal{M}$ . 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  $P \in \mathcal{M}$  with radius  $r > 0$  the subset of  $\mathcal{M}$  defined and denoted by  $\mathcal{B}(P, r) = \{Q \in \mathcal{M}, \ell(P, Q) \leq r\}$ . By extension, we set

$$\mathcal{B}(\mathbf{P}^*, r) = \{Q \in \mathcal{M}, \ell(\mathbf{P}^*, Q) \leq r\} \quad \text{for all } r > 0.$$

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

Throughout this paper, we use the following notations. The subsets of  $\mathbb{R}, \mathbb{R}^k$  will be equipped with their Borel  $\sigma$ -algebras. The cardinality of a set  $A$  is denoted  $|A|$  and 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  $\mathbf{x}$ . For all suitable functions  $f$  on  $(E^n, \mathcal{E}^{\otimes n})$ ,  $\mathbb{E}[f(\mathbf{X})]$  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 some measure  $S$  on  $(E, \mathcal{E})$ . For  $x \in \mathbb{R}$ ,  $(x)_+ = \max(x, 0)$  and  $(x)_- = \max\{-x, 0\}$ . 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  $\|f\|_{j, \mu} = [\int_E |f|^j d\mu]^{1/j} < +\infty$ . Finally,  $\|f\|_\infty = \sup_{x \in E} |f(x)|$  is the supremum norm of a function  $f$  on  $E$ .

### 3. OUR MAIN ASSUMPTIONS ON THE LOSS FUNCTION

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

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

is measurable and there exists a positive number  $\tau$  such that for all  $P \in \mathcal{P}$  and  $\bar{P}, Q \in \mathcal{M}$

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

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

Under such an assumption,  $\ell$ -balls are measurable, i.e. belong to  $\mathcal{A}$ . 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  $\tau = 2$ .

The construction of the posterior distribution not only depends on the prior  $\pi$  but also on the loss function  $\ell$ . For this reason, we introduce a family  $\mathcal{T}(\ell, \mathcal{M}) = \{t_{(P,Q)}, (P, Q) \in \mathcal{M}^2\}$  of functions on  $(E, \mathcal{E})$  with the following properties.

**Assumption 2.** The elements  $t_{(P,Q)}$  of  $\mathcal{T}(\ell, \mathcal{M})$  satisfy:

(i) The mapping

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

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  $P, Q \in \mathcal{M}$ ,

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

(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  $a_0 \ell(S, P) - a_1 \ell(S, P) \geq 0$ , hence that  $a_0 \geq a_1$  since  $\ell$  is not constantly equal to 0.

Many classical loss functions (among which the total variation distance, the 1-Wasserstein distance, the squared Hellinger loss, etc.) can be associated to families  $\mathcal{T}(\ell, \mathcal{M})$  satisfying our Assumption 2 — see Baraud (2021) —. Some of them possess the additional property given below.

**Assumption 3.** *Additionally to Assumption 2, there exists  $a_2 > 0$  such that*

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

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

This assumption is typically satisfied when  $\ell$  behaves as the square of a distance.

In the proposition below we provide families  $\mathcal{T}(\ell, \mathcal{M})$  that do satisfy our requirements for some loss functions  $\ell$  of interest. These results have been established in Baraud (2021) except for the squared Hellinger loss for which we refer to Baraud & Birgé (2018)[Proposition 3]. The list below is not exhaustive and other losses can also be considered, especially those that can be defined by 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. We refer to Baraud (2021) for more details on the way the families  $\mathcal{T}(\ell, \mathcal{M})$  can be obtained from the loss functions  $\ell$ .

**Proposition 1.** *The following holds:*

(1) *Total variation. Let  $\overline{\mathcal{P}}$  be the set of all probability measures on  $(E, \mathcal{E})$ ,  $\mathcal{M}$  a subset of  $\overline{\mathcal{P}}$  dominated by some reference measure  $\mu$  and for  $P, Q \in \overline{\mathcal{P}}$ ,  $\ell(P, Q) = \|P - Q\|$  the total variation loss (TV-loss for short) between  $P$  and  $Q$ . The family  $\mathcal{T}(\ell, \mathcal{M})$  of functions  $t_{(P,Q)}$  defined for  $P = p \cdot \mu, Q = q \cdot \mu \in \mathcal{M}$  by*

$$(6) \quad t_{(P,Q)} = \frac{1}{2} [\mathbb{1}_{q>p} - Q(q > p)] - \frac{1}{2} [\mathbb{1}_{p>q} - P(p > q)]$$

*satisfies Assumption 2 with  $a_0 = 3/2$  and  $a_1 = 1/2$ . If  $\mathcal{T}(\ell, \mathcal{M})$  satisfies Assumption 2 whatever the model  $\mathcal{M}$ , it may also occasionally satisfy Assumption 3 for some specific  $\mathcal{M}$ . An example of such a model is given Section 5.2.*

(2) *Hellinger distance. Let  $\overline{\mathcal{P}}$  be the set of all probability measures on  $(E, \mathcal{E})$ ,  $\mathcal{M}$  a subset of  $\overline{\mathcal{P}}$  dominated by some reference measure  $\mu$  and for  $P, Q \in \overline{\mathcal{P}}$ ,  $\ell(P, Q) = h^2(P, Q)$  the squared Hellinger distance between  $P$  and  $Q$ . Besides, 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  $\mathcal{T}(\ell, \mathcal{M})$  of functions  $t_{(P,Q)}$  defined for  $P = p \cdot \mu, Q = q \cdot \mu \in \mathcal{M}$  by

$$(7) \quad t_{(P,Q)} = \frac{1}{2} \psi \left( \sqrt{\frac{q}{p}} \right)$$

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

(3)  $\mathbb{L}_j$ -loss with  $1 < j < +\infty$ . For  $j \in (1, +\infty)$ , let  $\overline{\mathcal{P}}_j$  be the set of finite and signed measures on  $(E, \mathcal{E})$  of the form  $P = p \cdot \mu$  with  $p \in \mathcal{L}_j(E, \mu) \cap \mathcal{L}_1(E, \mu)$  and  $\mathcal{M} = \{P = p \cdot \mu, p \in \mathcal{M}\}$  be a subset of  $\overline{\mathcal{P}}_j$ . Assume that the family  $\mathcal{M}$  of densities satisfies

$$(8) \quad \|p - q\|_\infty \leq R \|p - q\|_{\mu,j} \quad \text{for all } p, q \in \mathcal{M} \text{ and some } R > 0.$$

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

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

The family  $\mathcal{T}(\ell, \mathcal{M})$  of functions  $t_{(P,Q)}$  defined for  $P, Q \in \mathcal{M}$  by

$$(9) \quad t_{(P,Q)} = \frac{1}{2R^{j-1}} \left[ \int_E f_{(P,Q)} \frac{dP + dQ}{2} - f_{(P,Q)} \right]$$

satisfies Assumption 2 with  $a_0 = 3/(4R^{j-1})$  and  $a_1 = 1/(4R^{j-1})$  for the loss  $\ell = \ell_j$  with  $\ell_j(P, Q) = \|p - q\|_{\mu,j}$  for all  $P = p \cdot \mu$  and  $Q = q \cdot \mu$  in  $\overline{\mathcal{P}}_j$ .

When  $j = 2$ , (8) is typically satisfied when  $\mathcal{M}$  is a subset of a linear space enjoying good connections between the  $\mathbb{L}_2$  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 (1998)[Section 3] for additional examples. The property may also hold for infinite dimensional linear spaces as proven in Baraud (2021).

#### 4. CONSTRUCTION OF THE POSTERIOR DISTRIBUTION AND MAIN RESULTS

**4.1. Construction of the posterior distribution.** It relies on two positive numbers  $\beta$  and

$$(10) \quad \lambda = (1 + c)\beta \quad \text{with} \quad c > 0 \quad \text{such that} \quad c_0 = (1 + c) - c(a_0/a_1) > 0.$$

Given the family  $\mathcal{T}(\ell, \mathcal{M})$ , we set

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



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

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

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

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

and finally define the posterior distribution  $\hat{\pi}_{\mathbf{X}}$  on  $(\mathcal{M}, \mathcal{A}, \pi)$  with density

$$(11) \quad \frac{\hat{\pi}_{\mathbf{X}}}{d\pi} : P \mapsto \frac{\exp[-\beta \mathbf{T}(\mathbf{X}, P)]}{\int_{\mathcal{M}} \exp[-\beta \mathbf{T}(\mathbf{X}, P)] d\pi(P)}.$$

Our Assumption 2-(i) ensures that  $\tilde{\pi}_{\mathbf{X}}(\cdot|P)/d\pi$  is a measurable function of  $(\mathbf{X}, P, Q)$  and  $\hat{\pi}_{\mathbf{X}}/d\pi$  a measurable function of  $(\mathbf{X}, P)$ .

**4.2. The influence of the prior.** The Bayesian paradigm offers the possibility to favour some elements of  $\mathcal{M}$  as compared to others. In order to evaluate how much the prior  $\pi$  advantages or disadvantages an element  $\bar{P} \in \mathcal{M}$ , we fix some number  $\gamma > 0$ , introduce the set

$$\mathcal{R}(\beta, \bar{P}) = \left\{ r \geq \frac{1}{n\beta a_1}, \frac{\pi(\mathcal{B}(\bar{P}, 2r))}{\pi(\mathcal{B}(\bar{P}, r))} > \exp(\gamma n \beta a_1 r) \right\}$$

with the convention  $a/0 = +\infty$  for all  $a \geq 0$  and finally define

$$(12) \quad r_n(\beta, \bar{P}) = \sup \mathcal{R}(\beta, \bar{P}) \quad \text{with} \quad \sup \emptyset = 1/(n\beta a_1).$$

It follows from the definition of  $r_n(\beta, \bar{P})$  that

$$(13) \quad 0 < \pi(\mathcal{B}(\bar{P}, 2r)) \leq \exp(\gamma n \beta a_1 r) \pi(\mathcal{B}(\bar{P}, r)) \quad \text{for all } r > r_n(\beta, \bar{P}).$$

Letting  $r$  decrease to  $r_n(\beta, \bar{P})$ , we derive that (13) holds for  $r = r_n(\beta, \bar{P})$ . In particular,  $\pi(\mathcal{B}(\bar{P}, r)) > 0$  for  $r = r_n(\beta, \bar{P})$ .

We shall see below that the performance of the posterior  $\hat{\pi}_{\mathbf{X}}$  depends on those  $\bar{P} \in \mathcal{M}$  such that  $r_n(\beta, \bar{P})$  is small enough. The connection between the behaviour of the prior  $\pi$  in the vicinity of an element  $\bar{P} \in \mathcal{M}$  and the quantity  $r_n(\beta, \bar{P})$  can be made as follows. Clearly, if the prior puts no mass on the  $\ell$ -ball  $\mathcal{B}(\bar{P}, r)$  then  $r_n(\beta, \bar{P}) > r$  and  $r_n(\beta, \bar{P})$  is therefore large if  $r$  is large. In the opposite case, if the prior puts enough mass on  $\mathcal{B}(\bar{P}, r)$  in the sense that

$$(14) \quad \pi(\mathcal{B}(\bar{P}, r)) \geq \exp(-\gamma n \beta a_1 r),$$

then for all  $r' \geq r$ ,

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

which implies that  $r_n(\beta, \bar{P}) \leq r$  and  $r_n(\beta, \bar{P})$  is therefore small if  $r$  is small. Although (14) is not equivalent to (13) (it is actually stronger), the previous arguments provide a partial view on the relationship between  $\pi$  and  $r_n$  and conditions to decide whether  $\bar{P}$  is favoured by  $\pi$  or not, according to the size of  $r_n(\beta, \bar{P})$ .

**4.3. A first result on the concentration property of the posterior distribution.** Following the above discussion, when the set

$$(15) \quad \mathcal{M}(\beta) = \{\bar{P} \in \mathcal{M}, r_n(\beta, \bar{P}) \leq a_1^{-1}\beta\}$$

is non-empty, it gathers the most favoured elements of the model  $(\mathcal{M}, \pi)$  at level  $a_1^{-1}\beta$ . The set  $\mathcal{M}(\beta)$  can alternatively be defined as

$$(16) \quad \mathcal{M}(\beta) = \left\{ \bar{P} \in \mathcal{M}, \sup_{r \geq a_1^{-1}\beta} \left[ \frac{1}{\gamma n a_1 r} \log \left( \frac{\pi(\mathcal{B}(\bar{P}, 2r))}{\pi(\mathcal{B}(\bar{P}, r))} \right) \right] \leq \beta \right\}.$$

It is sometimes easier to use this latter form for the calculations. The set  $\mathcal{M}(\beta)$  will play a crucial role in our first result below.

**Theorem 1.** *Let Assumptions 1 and 2 be satisfied, fix  $\gamma < (c_0 \wedge c)/(2\tau)$  and let  $\beta \geq 1/\sqrt{n}$  be chosen in such a way that the set  $\mathcal{M}(\beta)$  defined by (15) is not empty. Then, the posterior distribution  $\hat{\pi}_{\mathbf{X}}$  defined by (11) possesses the following property. Given  $\xi > 0$ , there exists a number  $\kappa_0 \geq 2$  depending only on  $c, \tau, \gamma, \xi$  and the ratio  $a_0/a_1$  such that, for any distribution  $\mathbf{P}^*$ ,*

$$(17) \quad \mathbb{E}[\hat{\pi}_{\mathbf{X}}(\mathcal{B}(\mathbf{P}^*, \kappa_0 r))] \leq 2e^{-\xi} \quad \text{with} \quad r = \inf_{P \in \mathcal{M}(\beta)} \ell(\mathbf{P}^*, P) + a_1^{-1}\beta.$$

*In particular,*

$$\mathbb{P}[\hat{\pi}_{\mathbf{X}}(\mathcal{B}(\mathbf{P}^*, \kappa_0 r)) \geq e^{-\xi/2}] \leq 2e^{-\xi/2}.$$

A suitable choice for  $\kappa_0$  is given by (81) and it is of the form  $A + B\xi$  where  $A$  and  $B$  only depends on  $c, \tau, \gamma$  and  $a_0/a_1$ . It therefore only depends on  $\xi$  linearly and on the choice of our loss function  $\ell$  but not on the prior  $\pi$ . Hence, given a loss function  $\ell$  and a confidence level  $1 - 2e^{-\xi}$ ,  $\kappa_0$  is a numerical constant.

Our posterior distribution depends on the parameter  $\lambda = (1+c)\beta$  where  $c > 0$  satisfies (10). By Proposition 1, this constraint is satisfied for any  $c \in (0, 1/2)$  when  $\ell$  is the TV or an  $\ell_j$  loss. For the choice  $c = 1/3$ ,  $c_0 = c$  and our condition on  $\gamma$  becomes  $\gamma \leq (6\tau)^{-1}$ .

When the data are truly i.i.d. and the prior puts enough mass around their common distribution  $P^*$ , in the sense that  $P^* \in \mathcal{M}(\beta)$ , then  $r = a_1^{-1}\beta$ . 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  $\inf_{P \in \mathcal{M}(\beta)} \ell(\mathbf{P}^*, P)$ . When this quantity 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 concentration properties of the posterior distribution is therefore stable with respect to a possible misspecification of the model and a departure from the equidistribution assumption.

The value of  $r$  given by (17) depends on the choice of the parameter  $\beta$ . Since the set  $\mathcal{M}(\beta)$  is increasing with  $\beta$  (for the inclusion), the two terms  $\inf_{P \in \mathcal{M}(\beta)} \ell(\mathbf{P}^*, P)$  and  $a_1^{-1}\beta$  vary in opposite directions when  $\beta$  increases. The set  $\mathcal{M}(\beta)$  must be large enough to provide a suitable approximation of  $\mathbf{P}^*$ , and therefore include as many elements of  $\mathcal{M}$  as possible since  $\mathbf{P}^*$  is unknown, but  $\beta$  must not be too large in order to keep  $a_1^{-1}\beta$  to a reasonable size.

In order to illustrate and comment further on the choice of the parameter  $\beta$ , let us consider the following simple example.

**Example 1.** Assume that the data are truly i.i.d. with distribution  $P^*$  and that  $\mathcal{M}$  is a parametric model on which  $\pi$  behaves like the uniform distribution on a suitable bounded subset of  $\mathbb{R}^D$ . More precisely, assume that for all  $P \in \mathcal{M}$  and  $r > 0$

$$(Ar)^D \wedge 1 \leq \pi(\mathcal{B}(P, r)) \leq (Br)^D \wedge 1$$

for some positive numbers  $A \leq B$  and  $D \geq 1$ . Note that this assumption implies that  $\pi(\mathcal{B}(P, A^{-1})) = 1$  for all  $P \in \mathcal{M}$  so that the diameter of the support of  $\pi$  is bounded by  $2\tau A^{-1}$ . Then,

$$(18) \quad \frac{\pi(\mathcal{B}(P, 2r))}{\pi(\mathcal{B}(P, r))} \leq \left(\frac{2B}{A}\right)^D \quad \text{for all } P \in \mathcal{M} \quad \text{and} \quad r > 0$$

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

$$\sup_{r \geq a^{-1}\beta} \left[ \frac{1}{\gamma n a_1 r} \log \left( \frac{\pi(\mathcal{B}(P, 2r))}{\pi(\mathcal{B}(P, r))} \right) \right] \leq \frac{D}{\gamma n \beta} \log \left( \frac{2B}{A} \right)$$

and we note that the right-hand side is not larger than  $\beta \geq 1/\sqrt{n}$  for

$$(19) \quad \beta = \sqrt{\frac{D \log(2B/A)}{\gamma n}} \vee \frac{1}{\sqrt{n}}.$$

This means that for such a value of  $\beta$ ,  $P \in \mathcal{M}(\beta)$ , and since  $P$  is arbitrary we obtain that  $\mathcal{M}(\beta) = \mathcal{M}$ . We derive from Theorem 1 that the distribution  $\hat{\pi}_{\mathbf{X}}$  concentrates on an  $\ell$ -ball centered at  $\mathbf{P}^*$  with a radius of order

$$r = \inf_{P \in \mathcal{M}} \ell(\mathbf{P}^*, P) + a_1^{-1} \sqrt{\frac{D}{n}}.$$

In particular we derive, using Proposition 1, that for the TV-loss

$$r = \inf_{P \in \mathcal{M}} \left[ \frac{1}{n} \sum_{i=1}^n \|P_i^* - P\| \right] + \frac{1}{2} \sqrt{\frac{D}{n}}.$$

Provided that  $\mathcal{M}$  is of the form  $\mathcal{M} = \{P = p \cdot \mu, p \in \mathcal{M}\}$  with  $\mathcal{M}$  satisfying (8), we also derive that for the  $\ell_j$ -loss and any distribution  $\mathbf{P}^* = \otimes_{i=1}^n (p_i^* \cdot \mu)$  with  $p_1^*, \dots, p_n^* \in \mathcal{L}_j(E, \mathcal{E}, \mu)$ ,

$$r_n = \inf_{p \in \mathcal{M}} \left[ \frac{1}{n} \sum_{i=1}^n \|p_i^* - p\|_{\mu, j} \right] + 4R^{j-1} \sqrt{\frac{D}{n}}.$$

This example shows that in the context of Theorem 1, the parameter  $\beta$  must be tuned as a suitable function of  $n$ , that depends on the properties of the prior distribution, in such a way that all (or most) of the elements of  $\mathcal{M}$  belong to  $\mathcal{M}(\beta)$ .

Quite surprisingly, the situation changes drastically when Assumption 3 is met as we shall see in the next section.

#### 4.4. The concentration property of the posterior distribution under Assumption 3.

Let us define the mapping

$$(20) \quad \phi : \begin{cases} (0, +\infty) & \longrightarrow \mathbb{R}_+ \\ z & \longmapsto \phi(z) = \frac{2(e^z - 1 - z)}{z^2}. \end{cases}$$

The function  $\phi$  is increasing on  $(0, +\infty)$  and tends to 1 when  $z$  tends to 0.

**Theorem 2.** *Assume that Assumptions 1 and 3 hold and define*

$$(21) \quad c_1 = c_0 - \beta a_2 a_1^{-1} \tau^2 \phi[\beta(1+2c)](1+2c(1+c));$$

$$(22) \quad c_2 = c - \beta a_2 a_1^{-1} \tau^2 \phi[\beta(1+2c)]c^2;$$

$$(23) \quad c_3 = (2+c) - \beta a_2 a_1^{-1} \tau^2 \phi[\beta(3+2c)](2+c)^2.$$

Let  $\gamma < (c_1 \wedge c_2 \wedge c_3)/(2\tau)$  and  $\beta_0$  be the value of  $\beta$  for which  $c_1 \wedge c_2 \wedge c_3 = 0$ . Then, for  $\beta \in (0, \beta_0)$  and  $n \geq 1/(\beta a_1)$ , the posterior distribution  $\hat{\pi}_{\mathbf{X}}$  defined by (11) satisfies the following property. Given  $\xi > 0$ , there exists a number  $\kappa_0$  depending only on  $a_0, a_1, a_2, c, \tau, \beta, \gamma$  and  $\xi$  such that, for any distribution  $\mathbf{P}^*$ ,

$$(24) \quad \mathbb{E}[\hat{\pi}_{\mathbf{X}}({}^c\mathcal{B}(\mathbf{P}^*, \kappa_0 r))] \leq 2e^{-\xi} \text{ with } r = \inf_{P \in \mathcal{M}} [\ell(\mathbf{P}^*, P) + r_n(\beta, P)].$$

In particular,

$$\mathbb{P}[\hat{\pi}_{\mathbf{X}}({}^c\mathcal{B}(\mathbf{P}^*, \kappa_0 r)) \geq e^{-\xi/2}] \leq 2e^{-\xi/2}.$$

It follows from the proof that one may take  $\kappa_0$  given by (93), which is of the form  $(A' + B'\xi)r$  with  $A'$  and  $B'$  depending only on  $a_0, a_1, a_2, c, \tau, \beta$  and  $\gamma$ . Note that the constraints on  $c, \beta$  and  $\gamma$  that are required in our Theorem 2 only depend on  $a_0, a_1$  and  $a_2$ , hence on the choice of the loss function  $\ell$ , but not on the model  $(\mathcal{M}, \pi)$ . In particular, unlike Theorem 1, the value of  $\beta$  can be chosen as a universal constant for a given loss function. For example,

when  $\ell = h^2$ , we know from Proposition 1 that  $a_0 = 2$ ,  $a_1 = 3/16$  and  $a_2 = 3\sqrt{2}/4$ , and we may take  $c = 0.05$ ,  $\beta = 0.01$ ,  $\gamma = 0.01$  and  $\tau = 2$  in (3).

In order to illustrate this new result and compare it with that of Theorem 1, let us go back to the framework of Example 1.

**Example 2** (Example 1 continued). Assume that  $\ell = h^2$  is the squared Hellinger loss and that the quantities  $c, \beta$  and  $\gamma$  have been chosen as above.

We see that that the right-hand side of (18) is not larger than  $\exp(\gamma n a_1 \beta r)$  provided that

$$r \geq \frac{D \log(2B/A)}{\gamma n a_1 \beta}.$$

For  $\gamma \leq D \log(2B/A)$ , the right-hand side of this latter inequality is not smaller than  $1/(n\beta a_1)$  and we derive from the definition (12) of  $r_n(\beta, P)$  that

$$r_n(\beta, P) \leq \frac{D \log(2B/A)}{\gamma n a_1 \beta} \quad \text{for all } P \in \mathcal{M}.$$

By applying Theorem 2 we obtain that the posterior distribution  $\hat{\pi}_{\mathbf{X}}$  concentrates on an  $\ell$ -ball with radius of order  $\inf_{P \in \mathcal{M}} \ell(\mathbf{P}^*, P) + D/n$ , hence at rate  $1/n$  when the data are i.i.d. with distribution  $P^* \in \mathcal{M}$ .

Applying our Theorem 1 under the only Assumption 2 and ignoring the fact that the loss  $\ell = h^2$  additionally satisfies Assumption 3, would lead, by arguing as for the TV-loss, to the weaker result that the posterior distribution concentrates on an  $\ell$ -ball with radius of order  $\inf_{P \in \mathcal{M}} \ell(\mathbf{P}^*, P) + \sqrt{D/n}$ .

We conclude that Theorem 2 leads to a stronger result on the concentration properties of  $\hat{\pi}_{\mathbf{X}}$  as compared to Theorem 1 when the loss function  $\ell$  satisfies Assumption 3 on the model  $\mathcal{M}$ .

## 5. APPLICATIONS

**5.1. How big is the set  $\mathcal{M}(\beta)$  in a translation model?** In this section, we consider the translation model  $\mathcal{M} = \{P_\theta = p(\cdot - \theta) \cdot \mu, \theta \in \mathbb{R}\}$  associated to a density  $p$  on  $\mathbb{R}$  with respect to the Lebesgue measure  $\mu$ . Given a density  $q$  on the real line and a scale parameter  $\sigma$ , we estimate the location parameter  $\theta$  by choosing the prior  $\nu_\sigma$  with density  $q_\sigma : \theta \mapsto \sigma^{-1}q(\theta/\sigma)$  with respect to  $\mu$ . The prior  $\pi$  on  $\mathcal{M}$  is the image of  $\nu_\sigma$  by the mapping  $\theta \mapsto P_\theta$  and we use the total variation distance to measure the quality of an estimator of  $P_\theta$ .

By Theorem 1, we know that the concentration properties of the posterior  $\hat{\pi}_{\mathbf{X}}$  depend on the size of the set  $\mathcal{M}(\beta)$  given by (16). Given a compact interval  $I \subset \mathbb{R}$ , our aim is to find a value of  $\beta \geq 1/\sqrt{n}$  for which  $\mathcal{M}(\beta)$  contains the subset  $\{P_\theta, \theta \in I\} \subset \mathcal{M}$ . We assume the following.

**Assumption 4.** *The density  $q$  is positive, symmetric and decreasing on  $\mathbb{R}_+$ . Besides, there exists  $L \in (0, +\infty]$  such that the mapping  $H$  defined by*

$$(25) \quad H : \begin{cases} [0, L) & \longrightarrow [0, 1) \\ t & \longmapsto \|P_t - P_0\| \end{cases}$$

*is bijective.*

Under Assumption 4,  $H$  is necessarily increasing on  $[0, L)$  and we may define its inverse  $G : [0, 1) \rightarrow [0, L)$ . We set

$$(26) \quad \bar{\Gamma} = \max \left\{ \left[ \sup_{0 < r \leq 1/4} \frac{G(2r)}{G(r)} \right] q(0), \frac{1}{2G(1/4)} \right\}$$

and assume that this quantity is finite. Note that  $\bar{\Gamma}$  only depends on the choices of  $p$  and  $q$ . For example, when  $p(x) = (1/2)e^{-|x|}$ ,  $L = +\infty$ ,  $H : t \mapsto 1 - \exp[-t/2]$ ,  $G : r \mapsto -2 \log(1-r)$  and since the mapping  $r \mapsto [G(2r)/G(r)]$  is increasing,

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

If  $p : x \mapsto (\alpha/2)(1 - |x|)^{-1+\alpha} \mathbb{1}_{|x| < 1}$  with  $\alpha > 0$ ,  $L = 2$ ,  $H : t \mapsto 1 - (1 - t/2)^\alpha$ ,  $G : r \mapsto 2[1 - (1 - r)^{1/\alpha}]$ . Since  $G(r) \sim 2r/\alpha$  in a neighbourhood of 0, the mapping  $r \mapsto G(2r)/G(r)$  is continuous on  $[0, 1/4]$  and therefore bounded. Given  $q(0)$ ,  $\bar{\Gamma}$  is therefore a finite number.

The following result is proven in Section (7.1).

**Proposition 2.** *Let Assumption 4 hold,  $\bar{\Gamma}$  the quantity defined by (26),  $\gamma \leq \log 4$  and  $t$  a positive number such that  $\nu_1([t, +\infty)) \leq 1/4$ . The set  $\mathcal{M}(\beta)$  contains the subset  $\{P_\theta, \theta \in [-\sigma t, \sigma t]\}$  if*

$$(27) \quad \beta \geq \sqrt{\frac{1}{n\gamma} \max \left\{ \log \left( \frac{\bar{\Gamma}(\sigma \vee 1)}{q(2t)} \right), \log 4 \right\}}.$$

Note that the interval  $I = [-\sigma t, \sigma t]$  can be enlarged by increasing the value of  $\sigma$  or that of  $t$ . In the first case, increasing  $\sigma$  makes the prior  $\nu_\sigma$  flatter and for a fixed value of  $t > 0$ , the right-hand side of (27) increases as  $\sqrt{\log \sigma}$  when  $\sigma$  becomes larger than 1. In the other case, for a fixed value of  $\sigma$ , the right-hand side of (27) increases as  $\sqrt{\log(1/q(2t))}$ . 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.

**5.2. Fast rates.** We go back to the statistical framework described in Section 5.1 in the special case where  $p$  is the density  $x \mapsto \alpha x^{\alpha-1} \mathbb{1}_{(0,1]}$  with  $\alpha \in (0, 1]$ . As before, we choose the TV-loss. Since  $\ell$  is a distance we may take  $\tau = 1$  in Assumption 1. Besides, we have seen in Proposition 1 that the family  $\mathcal{T}(\mathcal{M}, \ell)$  given by (6) satisfies our Assumption 2 with  $a_0 = 3/2$

and  $a_1 = 1/2$ . In fact, it turns out that in the specific situation we consider here, the TV-loss also satisfies Assumption 3 with  $a_2 = 1$ . This means that from a more statistical point of view, the TV-loss rather behaves here as the square of a distance. In addition, some simple calculations show that

$$(28) \quad \|P_\theta - P_{\theta'}\| = |\theta - \theta'|^\alpha \wedge 1 \quad \text{for all } \theta, \theta' \in \mathbb{R}.$$

These facts are proven in Baraud (2021) [Examples 5 and 6].

Since Assumption 3 holds true, we may apply Theorem 2 and the reader can check that the constants  $c = 0.3$ ,  $\beta = 0.1$  and  $\gamma = 0.01$  satisfy the requirements of this theorem.

To estimate the location parameter  $\theta$ , we choose a prior  $\nu_\sigma = \sigma^{-1}q(\cdot/\sigma) \cdot \mu$  associated to a density  $q$  that satisfies the requirements of Assumption 4 so that, by (28), this assumption holds true with  $L = 1$ ,  $G : r \rightarrow r^{1/\alpha}$  and

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

We prove in Section 7.2 the following result.

**Proposition 3.** *Let  $t_0$  be the third quartile of  $\nu_1$ . If the density  $q$  is positive, symmetric and decreasing on  $[0, +\infty)$ , for all  $\theta \in \mathbb{R}$ ,*

$$r_n(\beta, P_\theta) \leq \frac{2000}{n} \max \left\{ \log(\bar{\Gamma}(\sigma \vee 1)) - \log \left( q \left[ 2 \left( \frac{|\theta|}{\sigma} \vee t_0 \right) \right] \right), \log 4 \right\}.$$

By applying Theorem 2, we conclude that for all  $\xi > 0$ , with a probability at least  $1 - 2e^{-\xi/2}$ , the posterior distribution satisfies

$$\hat{\pi}_{\mathbf{X}}(\mathcal{B}(\mathbf{P}^*, \kappa'_0 r)) \geq 1 - e^{-\xi/2}$$

for some  $\kappa'_0 \geq 2$  only depending on  $\xi, \alpha, q(0)$  and for

$$r = \inf_{\theta \in \mathbb{R}} \left[ \frac{1}{n} \sum_{i=1}^n \|P_i^* - P_\theta\| + \frac{1}{n} \left[ 1 - \log \left( q \left[ 2 \left( \frac{|\theta|}{\sigma} \vee t_0 \right) \right] \right) + \log(\sigma \vee 1) \right] \right].$$

In particular, when the data are i.i.d. with distribution  $P_{\theta^*}$ , with probability close to 1, an element  $P_{\hat{\theta}}$  drawn randomly according to the posterior distribution  $\hat{\pi}_{\mathbf{X}}$  satisfies with a probability close to 1,

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

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

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

The parameter is therefore estimated at rate  $n^{-1/\alpha}$  which is much faster than the usual  $1/\sqrt{n}$ -parametric one that is reached by an estimator based on a moment method for instance. Since the densities are unbounded, note that the maximum likelihood estimator does not exist and is therefore useless. It is well-known, mainly from the work of Le Cam, that is impossible to

estimate a distribution in a translation model at a rate faster than  $1/n$  for the TV-loss. Because of (28), the rate we get is not only optimal for estimating  $P_{\theta^*}$  but also for estimating  $\theta^*$  with respect to the Euclidean distance. An alternative 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 optimal and robust.

It is also interesting to see how the critical radius  $r$  behaves under a misspecification of the prior  $\nu_\sigma$ , i.e. when the size of the parameter  $\theta^*$  is large compared to  $\sigma$ . When  $q$  is Gaussian,  $r$  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.

**5.3. A general result under entropy.** In this section,  $(E, \mathcal{E}) = (\mathbb{R}^k, \mathcal{B}(\mathbb{R}^k))$  that we equip with the Lebesgue measure  $\mu$  and the norm  $|\cdot|_\infty$ . We consider the TV-loss  $\ell$  and the location-scale family

$$(29) \quad \mathcal{M} = \left\{ P_{(p, \mathbf{m}, \sigma)} = \frac{1}{\sigma^k} p \left( \frac{\cdot - \mathbf{m}}{\sigma} \right) \cdot \mu, p \in \mathcal{M}_0, \mathbf{m} \in \mathbb{R}^k, \sigma > 0 \right\},$$

where  $\mathcal{M}_0$  is a set of densities on  $(\mathbb{R}^k, \mathcal{B}(\mathbb{R}^k), \mu)$  that satisfies the following entropy condition:

**Assumption 5.** *Let  $\tilde{D}$  be a continuous non-increasing 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*

$$(30) \quad |\mathcal{M}_0[\eta]| \leq \exp \left[ \tilde{D}(\eta) \right]$$

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

$$(31) \quad \ell(P_{(p, \mathbf{0}, 1)}, P_{(\bar{p}, \mathbf{0}, 1)}) = \frac{1}{2} \int_{\mathbb{R}^k} |p - \bar{p}| d\mu \leq \eta.$$

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

$$(32) \quad \ell(P_{(p, \mathbf{0}, 1)}, P_{(p, \mathbf{m}, \sigma)}) \leq \left[ A \left( \left( \left| \frac{\mathbf{m}}{\sigma} \right|_\infty \right)^\alpha + \left( 1 - \frac{1}{\sigma} \right)^\alpha \right) \right] \wedge 1.$$

For  $\eta, \delta > 0$ , we define

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

and for  $\theta = \theta(\bar{p}, j_0, \mathbf{j}) \in \Theta[\eta, \delta]$ , set

$$(33) \quad L_\theta = (k + 1)L + \log |\mathcal{M}_0[\eta]| + 2 \sum_{i=0}^k \log(1 + |j_i|)$$



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 prior  $\pi$  defined as

$$(34) \quad \pi(\{P_\theta\}) = e^{-L\theta} \quad \text{for all } \theta \in \Theta[\eta, \delta].$$

**Corollary 1.** *Let  $\xi > 0$ ,  $K > 1$ ,  $\ell$  be the TV-loss and  $\mathcal{M}$  the family of probabilities given by (29) where  $\mathcal{M}_0$  satisfies Assumption 5. Consider the parameters*

$$(35) \quad \eta_n = \inf \mathcal{D}_n \quad \text{with} \quad \mathcal{D}_n = \left\{ \eta > 0, \tilde{D}(\eta) \leq \frac{n\eta^2}{24} \right\}$$

$$(36) \quad \delta = \delta_n = \left( \frac{\eta_n}{2A} \right)^{1/\alpha},$$

$$(37) \quad \beta = \beta_n = \frac{1}{2} \left[ K\eta_n + 2\sqrt{\frac{18.6(k+1)}{n}} \right]$$

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

$$(38) \quad |\log \sigma| \vee \left| \frac{\mathbf{m}}{\sigma} \right|_\infty \leq \Lambda_n = \exp \left[ \frac{(K^2 - 1)n\eta_n^2}{48(k+1)} + \log \log(1 + \delta_n) \right].$$

Then, the posterior distribution  $\hat{\pi}_{\mathbf{X}}$  associated to the value  $\lambda = 4\beta/3$ , the prior  $\pi$  given by (34) and the family  $\mathcal{T}(\ell, \mathcal{M})$  given by (6) possesses the following property: there exists  $\kappa_0 \geq 2$  only depending on  $\xi$  such that

$$(39) \quad \mathbb{E} [\hat{\pi}_{\mathbf{X}} (\mathcal{C}\mathcal{B}(\mathbf{P}^*, \kappa_0 r))] \leq 2e^{-\xi}$$

with

$$(40) \quad r = \inf_{P \in \mathcal{M}_n(K)} \ell(\mathbf{P}^*, P) + K\eta_n + \sqrt{\frac{k+1}{n}}.$$

To comment on Condition (32), let us assume that the set  $\mathcal{M}_0$  consists of densities  $p$  that are supported on  $[0, 1]^k$ , satisfies  $\sup_{p \in \mathcal{M}_0} \|p\|_\infty \leq L_0$  and

$$(41) \quad \sup_{p \in \mathcal{M}_0} |p(\mathbf{x}) - p(\mathbf{x}')| \leq L_1 |\mathbf{x} - \mathbf{x}'|^\alpha \quad \text{for all } \mathbf{x}, \mathbf{x}' \in \mathbb{R}^k$$

for some constants  $L_0, L_1 > 0$  and  $\alpha \in (0, 1]$ . 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 (41), we deduce that

for all  $p \in \mathcal{M}_0$ ,  $\sigma \geq 1$  and  $\mathbf{m} \in \mathbb{R}^k$ ,

$$\begin{aligned}
& \ell(P_{(p, \mathbf{0}, 1)}, P_{(p, \mathbf{m}, \sigma)}) \\
& \leq \ell(P_{(p, \mathbf{0}, 1)}, P_{(p, \mathbf{0}, \sigma)}) + \ell(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)^\alpha |\mathbf{x}|^\alpha d\mathbf{x} \\
& \quad + \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|^\alpha d\mathbf{x} \\
& \leq \frac{1}{2} \left(1 - \frac{1}{\sigma^k}\right) + \frac{L_1 k^{\alpha/2}}{2\sigma^k} \left(1 - \frac{1}{\sigma}\right)^\alpha + \frac{L_0}{2} \left(1 - \frac{1}{\sigma^k}\right) + L_1 \left|\frac{\mathbf{m}}{\sigma}\right|^\alpha \\
& \leq \frac{1}{2} [1 + L_1 k^{\alpha/2} + L_0] \left(1 - \frac{1}{\sigma}\right)^\alpha + L_1 \left|\frac{\mathbf{m}}{\sigma}\right|^\alpha
\end{aligned}$$

and (32) is therefore satisfied with  $A = L_1 \vee [(1 + L_1 k^{\alpha/2} + L_0)/2]$ .

In Lemma 1 below, we show that (32) may also be satisfied in a situation where the densities in  $\mathcal{M}_0$  are irregular and possibly discontinuous. It makes it possible to consider the following example.

**Example 3.** We consider here the situation where  $k = 1$  and  $\mathcal{M}_0$  is the set of all non-increasing densities on  $[0, 1]$  that are bounded by  $B > 1$ . Then,  $\mathcal{M}$  consists of all the probabilities whose densities are non-increasing and supported on intervals  $I$  with positive lengths and which are bounded by  $B/\mu(I)$ . Birman and Solomjak (1967) proved that  $\mathcal{M}_0$  satisfies Assumption 5 with  $\widetilde{D}(\eta)$  of order  $(1/\eta) \vee 1$  (up to some constant that depends on  $B$ ) and consequently  $\eta_n$  is of order  $n^{-1/3}$ . Besides, it follows from Lemma 1 below that (32) is satisfied with  $A = B$  and  $\alpha = 1$ . We may therefore apply Corollary 1. For a value of  $K$  large enough compared to 1,  $\Lambda_n$  defined by (38) is larger than  $\exp[CK^2 n^{1/3}]$  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

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

the posterior distribution  $\hat{\pi}_{\mathbf{X}}$  satisfies for all  $\xi > 0$ , with a probability at least  $1 - 2e^{-\xi/2}$ ,  $\hat{\pi}_{\mathbf{X}} [\mathcal{B}(P^*, C'n^{-1/3})] \geq 1 - e^{-\xi/2}$  where the constant  $C' > 0$  only depends on  $\xi, K, B$  but not on  $m^*$  and  $\sigma^*$ . This means that the concentration properties of  $\hat{\pi}_{\mathbf{X}}$  hold true over a huge range of translation and scale parameters when  $n$  is large enough.

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

$$(42) \quad \frac{1}{2} \int_{\mathbb{R}} \left| \frac{1}{\sigma} p\left(\frac{x}{\sigma}\right) - p(x) \right| dx \leq \left(1 - \frac{1}{\sigma}\right).$$

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

$$(43) \quad \frac{1}{2} \int_{\mathbb{R}} |p(x) - p(x - m)| dx \leq (|m|B) \wedge 1.$$

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

$$(44) \quad \frac{1}{2} \int_{\mathbb{R}} \left| \frac{1}{\sigma} p\left(\frac{x - m}{\sigma}\right) - p(x) \right| dx \leq \left[ B \left| \frac{m}{\sigma} \right| + \left(1 - \frac{1}{\sigma}\right) \right] \wedge 1.$$

**5.4. Estimating a parameter under sparsity.** Let us consider a family of distributions  $\mathcal{M} = \{P_{\theta} = p_{\theta} \cdot \mu, \theta \in \mathbb{R}^k\}$  that are parametrized by  $\mathbb{R}^k$  where the dimension  $k$  is large. We presume, even though this might not be true, that the data are i.i.d. with a distribution  $P_{\theta^*} \in \mathcal{M}$  associated to a parameter  $\theta^*$  the coordinates of which are all zero except, maybe, a small number of these. For  $m \subset \{1, \dots, k\}$ ,  $m \neq \emptyset$ , we introduce the sub-family  $\mathcal{M}_m$  that gathers the 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  $\mathcal{M}_m = \{P_{\theta}, \theta \in \Theta_m\}$  for all  $m \subset \{1, \dots, k\}$ ,  $m \neq \emptyset$ .

Throughout this section we consider the squared Hellinger loss and, given some  $R > 0$ , we assume that there exist  $\alpha \in (0, 1]$  and a positive number  $B_k = B_k(R)$  possibly depending on  $k$  and  $R$ , although we do not specify the dependency with respect to  $R$ , such that for all  $\theta, \theta' \in \mathbb{R}^k$  with  $|\theta|_{\infty} \vee |\theta'|_{\infty} \leq R$

$$(45) \quad h(P_{\theta}, P_{\theta'}) \leq \sqrt{B_k} |\theta - \theta'|_{\infty}^{\alpha}.$$

As a consequence, the mapping  $\theta \mapsto P_{\theta}$  is continuous. We endow  $\mathcal{M}$  with the Borel  $\sigma$ -algebra  $\mathcal{A}$  associated to the Hellinger distance so that, for all probabilities  $P$  on  $(E, \mathcal{E})$ , the mapping  $Q \mapsto \ell(P, Q) = h^2(P, Q)$  is continuous, hence measurable, on  $(\mathcal{M}, \mathcal{A})$ , and Assumption 1 is satisfied with  $\tau = 2$ .

Given a nonempty subset  $m$  of  $\{1, \dots, k\}$ , we endow  $\Theta_m$  with the uniform distribution  $\nu_m$  on the cube  $\Theta_m(R) = \{\theta \in \Theta_m, |\theta|_{\infty} \leq R\}$ . This leads

to a prior  $\pi_m = \pi_m(R)$  on  $\mathcal{M}_m$  defined as the image of  $\nu_m$  by the mapping  $\boldsymbol{\theta} \mapsto P_{\boldsymbol{\theta}}$ . For  $m = \emptyset$ , we set  $\Theta_{\emptyset} = \Theta_{\emptyset}(R) = \{\mathbf{0}\}$  and we endow it with the Dirac mass at  $\mathbf{0}$  so that  $\mathcal{M}_{\emptyset} = \{P_{\mathbf{0}}\}$  and  $\pi_{\emptyset}$  is the Dirac mass at  $P_{\mathbf{0}}$ . We finally define our prior  $\pi = \pi(R)$  on  $(\mathcal{M}, \mathcal{A})$  as

$$(46) \quad \pi = \sum_{m \subset \{1, \dots, k\}} e^{-L_m} \pi_m \quad \text{with} \quad L_m = |m| \log k + k \log \left(1 + \frac{1}{k}\right).$$

It is not difficult to check that  $\sum_{m \subset \{1, \dots, k\}} e^{-L_m} = 1$ , hence that  $\pi$  is a genuine probability on  $(\mathcal{M}, \mathcal{A})$ .

The following result holds.

**Corollary 2.** *Let  $\ell = h^2$  be the squared Hellinger distance,  $\mathcal{M} = \{P_{\boldsymbol{\theta}} = p_{\boldsymbol{\theta}} \cdot \mu, \boldsymbol{\theta} \in \mathbb{R}^k\}$  a dominated statistical model satisfying (45) and the assumption that the mapping*

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

is measurable. We assume that  $RB_k^{1/(2\alpha)} \geq 1$  and endow  $\mathcal{M}$  with the prior  $\pi = \pi(R)$  defined by (46) and define the posterior distribution  $\hat{\pi}_{\mathbf{X}}$  by (11) with  $\beta = 0.01$ ,  $\lambda = 1.05\beta$  and the family  $\mathcal{F}(\ell, \mathcal{M})$  given by (7). Then there exists  $\kappa_0 \geq 2$ , only depending on the value of  $\xi > 0$ , such that

$$\mathbb{E} [\hat{\pi}_{\mathbf{X}}({}^c\mathcal{B}(\mathbf{P}^*, \kappa_0 r))] \leq 2e^{-\xi}$$

where

$$(47) \quad r = \inf_{m \subset \{1, \dots, k\}} \left[ \inf_{\boldsymbol{\theta} \in \Theta_m(R)} \ell(\mathbf{P}^*, P_{\boldsymbol{\theta}}) + \frac{|m| \log \left(2kR(nB_k)^{1/(2\alpha)} + 1\right)}{n} \right].$$

Let us now comment on and illustrate this result.

When  $B_k$  does not increase faster than a power of  $k$ ,  $r$  only depends logarithmically on the dimension  $k$ , as expected.

Since the mapping

$$R \mapsto \sup \left\{ \frac{h(P_{\boldsymbol{\theta}}, P_{\boldsymbol{\theta}'})}{|\boldsymbol{\theta} - \boldsymbol{\theta}'|_{\infty}^{\alpha}} \mid \boldsymbol{\theta}, \boldsymbol{\theta}' \in \mathbb{R}^k, \boldsymbol{\theta} \neq \boldsymbol{\theta}', |\boldsymbol{\theta}|_{\infty} \vee |\boldsymbol{\theta}'|_{\infty} \leq R \right\}$$

is non-decreasing, our condition  $RB_k^{1/(2\alpha)} \geq 1$  holds for  $R$  large enough.

For illustration, assume that  $P_{\boldsymbol{\theta}}$  is the Gaussian distribution with mean  $\boldsymbol{\theta} \in \mathbb{R}^k$  and covariance matrix  $\sigma^2 I_k$  where  $I_k$  denotes the  $k \times k$  identity matrix. Then,

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

and we obtain that (45) is satisfied with  $B_k = k/(8\sigma^2)$  and  $\alpha = 1$ . In particular, our condition  $RB_k^{1/(2\alpha)} \geq 1$  is equivalent to  $R \geq 2\sigma\sqrt{(2/k)}$ . In this case, we obtain that the value of  $r$  given by (47) is of order

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

More generally, when  $\mathcal{M} = \{P_{\boldsymbol{\theta}}, \boldsymbol{\theta} \in \mathbb{R}^k\}$  is a regular statistical model with Fisher information  $\mathbf{I}(\boldsymbol{\theta})$  at  $\boldsymbol{\theta}$ , we know from the book of Ibragimov and Has'minskiĭ (1981)[Theorem 7.1 p.81] that for all  $\boldsymbol{\theta}, \boldsymbol{\theta}' \in \mathbb{R}^k$  such that  $|\boldsymbol{\theta}|_{\infty} \vee |\boldsymbol{\theta}'|_{\infty} \leq R$

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

Then, Assumption (45) holds with  $\alpha = 1$  and

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

where  $\varrho(\mathbf{I}(\boldsymbol{\theta}''))$  denotes the largest eigenvalue of the matrix  $\mathbf{I}(\boldsymbol{\theta}'')$ . It is well known that this value is independent of  $\boldsymbol{\theta}''$  when  $\mathcal{M}$  is a translation model.

## 6. PROOFS OF THEOREMS 1 AND 2

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

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

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

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

**Lemma 2.** *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 (2011) [Lemma 4.2, P. 28].

**Lemma 3.** *For  $P, Q \in \mathcal{M}$ , we set*

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

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

$$(48) \quad \mathbb{E}[\exp[-\beta \mathbf{T}(\mathbf{X}, P)]] \leq \frac{1}{\pi(\mathcal{B})} \left[ \int_{\mathcal{B}} \exp[-\mathbf{M}(P, Q)] d\pi_{\mathcal{B}}(Q) \right]^{-1}.$$

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

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

Then,

$$\begin{aligned} & \mathbb{E}[\exp[-I(\mathbf{X}, P, Q)]] \\ &= \mathbb{E} \left[ \exp \left[ -c_1 \beta \mathbf{T}(\mathbf{X}, P, Q) + \log \int_{\mathcal{M}} \exp[c\beta \mathbf{T}(\mathbf{X}, P, Q')] d\pi(Q') \right] \right] \\ &= \mathbb{E} \left[ \int_{\mathcal{M}} \exp[c\beta \mathbf{T}(\mathbf{X}, P, Q') - c_1 \beta \mathbf{T}(\mathbf{X}, P, Q)] d\pi(Q') \right] \\ (49) \quad &= \exp[\mathbf{M}(P, Q)]. \end{aligned}$$

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

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

Hence,

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

Applying Lemma 2 with  $U = \mathbf{X}$ ,  $V = Q$  with distribution  $\pi_{\mathcal{B}}$ , and  $h(U, V) = -I(\mathbf{X}, P, Q)$ , we obtain that

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

and (48) follows from (49).  $\square$

**Lemma 4.** For  $P, Q \in \mathcal{M}$ , we set

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

For all  $r \in \mathcal{V}(\pi, \bar{P})$ ,

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

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

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

Then,

$$\begin{aligned} & \mathbb{E} [\exp[-H(\mathbf{X}, P, Q)]] \\ & = \mathbb{E} \left[ \exp[-\beta c_1 \mathbf{T}(\mathbf{X}, P, Q)] \int_{\mathcal{M}} \exp[c_2 \beta \mathbf{T}(\mathbf{X}, P, Q')] d\pi(Q') \right] \\ & = \mathbb{E} \left[ \int_{\mathcal{M}} \exp[\beta (c_2 \mathbf{T}(\mathbf{X}, P, Q') - c_1 \mathbf{T}(\mathbf{X}, P, Q))] d\pi(Q') \right] \\ (50) \quad & = \exp[\mathbf{L}(P, Q)]. \end{aligned}$$

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

$$\begin{aligned} \mathbb{E} [\exp[\beta \mathbf{T}(\mathbf{X}, P)]] & = \mathbb{E} \left[ \exp \left[ \int_{\mathcal{M}} [\beta \mathbf{T}(\mathbf{X}, P, Q)] d\tilde{\pi}_{\mathbf{X}}(Q) \right] \right] \\ & \leq \mathbb{E} \left[ \int_{\mathcal{M}} \exp[\beta \mathbf{T}(\mathbf{X}, P, Q)] d\tilde{\pi}_{\mathbf{X}}(Q) \right] \\ & = \mathbb{E} \left[ \frac{\int_{\mathcal{M}} \exp[c_2 \beta \mathbf{T}(\mathbf{X}, P, Q)] d\pi(Q)}{\int_{\mathcal{M}} \exp[c_1 \beta \mathbf{T}(\mathbf{X}, P, Q)] d\pi(Q)} \right] \\ & = \mathbb{E} \left[ \frac{1}{\int_{\mathcal{M}} \exp[H(\mathbf{X}, P, Q)] d\pi(Q)} \right]. \end{aligned}$$

Applying Lemma 2 with  $U = \mathbf{X}$  and  $V = Q$  with distribution  $\pi$  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 (50) that for all  $P \in \mathcal{M}$

$$\begin{aligned} \mathbb{E} [\exp[\beta \mathbf{T}(\mathbf{X}, P)]] & \leq \left[ \int_{\mathcal{M}} \exp[-\mathbf{L}(P, Q)] d\pi(Q) \right]^{-1} \\ (51) \quad & \leq \frac{1}{\pi(\mathcal{B})} \left[ \int_{\mathcal{B}} \exp[-\mathbf{L}(P, Q)] d\pi_{\mathcal{B}}(Q) \right]^{-1}. \end{aligned}$$

Applying Lemma 2 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 (51) leads to the result.  $\square$

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

$$(52) \quad \log \mathbb{E}[\exp[\beta U]] \leq \beta \mathbb{E}[U] + \frac{\beta^2 l^2}{8}.$$

**Lemma 6.** *Let  $U$  be a squared integrable random variable not larger than  $b > 0$ . For all  $\beta > 0$ ,*

$$(53) \quad \log \mathbb{E}[\exp[\beta U]] \leq \beta \mathbb{E}[U] + \beta^2 \mathbb{E}[U^2] \phi(\beta b),$$

where  $\phi$  is defined by (20).

The proofs of Lemmas 5 and 6 can be found on pages 21 and 23 in Massart (2007).

**Lemma 7.** *Let  $J \in \mathbb{N}$ ,  $\gamma > 0$  and  $\bar{P} \in \mathcal{M}$ . If  $r$  satisfies  $n\beta a_1 r \geq 1$  and (13), for all  $\gamma_0 > 2\gamma$*

$$(54) \quad \begin{aligned} &\int_{\varepsilon_{\mathcal{B}}(\bar{P}, 2^J r)} \exp[-\gamma_0 n\beta a_1 \ell(\bar{P}, P)] d\pi(P) \\ &\leq \pi(\mathcal{B}) \exp\left[\Xi - (\gamma_0 - 2\gamma) n\beta a_1 2^J r\right] \end{aligned}$$

with

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

Besides,

$$(55) \quad \int_{\mathcal{M}} \exp[-\gamma_0 n\beta a_1 \ell(\bar{P}, P)] d\pi(P) \leq \pi(\mathcal{B}) \exp[\Xi']$$

with

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



*Proof.* From (13), we deduce by induction that for all  $j \geq 0$

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

Consequently,

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

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

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

which leads to (54) 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

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

which is (55). □

**6.2. Main parts of the proofs of Theorems 1 and 2.** Throughout the proofs of these two theorems, we fix some arbitrary element  $\bar{P} \in \mathcal{M}$  and  $r \geq r_n(\beta, \bar{P})$ . It follows from the definition of  $r_n(\beta, \bar{P})$  that  $r$  satisfies both  $n\beta a_1 r \geq 1$  and inequality (13). For a positive number  $z$ , that will be chosen later as well, we set

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

It follows from the definition (11) of  $\hat{\pi}_{\mathbf{X}}$  that, given  $J \in \mathbb{N}$ ,

$$\begin{aligned} \mathbb{E} \left[ \hat{\pi}_{\mathbf{X}} \left( {}^c\mathcal{B}(\bar{P}, 2^J r) \right) \right] &= \mathbb{E} \left[ \hat{\pi}_{\mathbf{X}} \left( {}^c\mathcal{B}(\bar{P}, 2^J r) \right) \mathbb{1}_{cA} \right] + \mathbb{E} \left[ \hat{\pi}_{\mathbf{X}} \left( {}^c\mathcal{B}(\bar{P}, 2^J r) \right) \mathbb{1}_A \right] \\ &\leq \mathbb{P}(cA) + \frac{1}{z} \mathbb{E} \left[ \int_{{}^c\mathcal{B}(\bar{P}, 2^J r)} \exp[-\beta \mathbf{T}(\mathbf{X}, P)] d\pi(P) \right] \\ (56) \quad &= \mathbb{P}(cA) + \frac{1}{z} \int_{{}^c\mathcal{B}(\bar{P}, 2^J r)} \mathbb{E} [\exp[-\beta \mathbf{T}(\mathbf{X}, P)]] d\pi(P). \end{aligned}$$

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 (56) is not larger than  $e^{-\xi}$ . To achieve this goal, we bound the first term of the right-hand side of (56) by first applying Markov's inequality

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

and then by using Lemma 4. This leads to

$$(57) \quad \mathbb{P}(cA) \leq \frac{z}{\pi^2(\mathcal{B})} \left[ \int_{\mathcal{B}^2} \exp[-\mathbf{L}(P, Q)] d\pi_{\mathcal{B}}(P) d\pi_{\mathcal{B}}(Q) \right]^{-1}.$$

To show that the first term in the right hand-side of (56) is not larger than  $e^{-\xi}$  we therefore prove that this is the case of the right-hand side of (57) for  $z$  small enough. We bound the second term of (56) by using Lemma 3.

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  $\bar{P} \in \mathcal{M}(\beta)$ ,

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

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

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

We obtain (17) by monotone convergence, taking a sequence  $(\bar{P}_N)_{N \geq 0} \subset \mathcal{M}(\beta)$  such that  $\ell(\mathbf{P}^*, \bar{P}_N)$  is non-increasing to  $\inf_{P \in \mathcal{M}(\beta)} \ell(\mathbf{P}^*, P)$ , so that

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

and (17) 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  $\bar{P} \in \mathcal{M}$ ,

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

and we get (24) by arguing similarly.

**6.3. Proof of Theorem 1.** For all  $i \in \{1, \dots, n\}$  and  $P, Q, Q' \in \mathcal{M}$ , let us set

$$(58) \quad \begin{aligned} U_i &= c \left( t_{(P, Q')}(X_i) - \mathbb{E} \left[ t_{(P, Q')}(X_i) \right] \right) \\ &\quad - c_1 \left( t_{(P, Q)}(X_i) - \mathbb{E} \left[ t_{(P, Q)}(X_i) \right] \right) \end{aligned}$$

$$(59) \quad \begin{aligned} V_i &= c_2 \left( t_{(P, Q')}(X_i) - \mathbb{E} \left[ t_{(P, Q')}(X_i) \right] \right) \\ &\quad - c_1 \left( t_{(P, Q)}(X_i) - \mathbb{E} \left[ t_{(P, Q)}(X_i) \right] \right). \end{aligned}$$

The random variables  $U_i$  are independent and under Assumption 2-(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 5, we obtain that

$$(60) \quad \prod_{i=1}^n \mathbb{E} [\exp [\beta U_i]] \leq \exp \left[ \frac{l_1^2 n \beta^2}{8} \right]$$

and

$$(61) \quad \prod_{i=1}^n \mathbb{E} [\exp [\beta V_i]] \leq \exp \left[ \frac{l_2^2 n \beta^2}{8} \right].$$

By using Assumption 1 and the fact that  $c_0 = c_1 - ca_0/a_1 > 0$ ,

$$\begin{aligned} &c(a_0 \ell(\mathbf{P}^*, P) - a_1 \ell(\mathbf{P}^*, Q')) - c_1(a_1 \ell(\mathbf{P}^*, P) - a_0 \ell(\mathbf{P}^*, Q)) \\ &= -(c_1 a_1 - ca_0) \ell(\mathbf{P}^*, P) - ca_1 \ell(\mathbf{P}^*, Q') + c_1 a_0 \ell(\mathbf{P}^*, Q) \\ &\leq -c_0 a_1 [\tau^{-1} \ell(\bar{P}, P) - \ell(\mathbf{P}^*, \bar{P})] - ca_1 [\tau^{-1} \ell(\bar{P}, Q') - \ell(\mathbf{P}^*, \bar{P})] \\ &\quad + \tau c_1 a_0 [\ell(\mathbf{P}^*, \bar{P}) + \ell(\bar{P}, Q)] \\ (62) \quad &= e_0 a_1 \ell(\mathbf{P}^*, \bar{P}) - \tau^{-1} c_0 a_1 \ell(\bar{P}, P) - \tau^{-1} ca_1 \ell(\bar{P}, Q') + \tau c_1 a_0 \ell(\bar{P}, Q) \end{aligned}$$

with

$$(63) \quad e_0 = c_0 + c + \frac{\tau c_1 a_0}{a_1}.$$

It follows from Assumptions 2-(iii) and (62) that

$$(64) \quad \begin{aligned} & n^{-1} \{c\mathbb{E} [\mathbf{T}(\mathbf{X}, P, Q')] - c_1\mathbb{E} [\mathbf{T}(\mathbf{X}, P, Q)]\} \\ & \leq c [a_0\ell(\mathbf{P}^*, P) - a_1\ell(\mathbf{P}^*, Q')] - c_1 [a_1\ell(\mathbf{P}^*, P) - a_0\ell(\mathbf{P}^*, Q)] \\ & \leq e_0 a_1 \ell(\mathbf{P}^*, \bar{P}) - \tau^{-1} c_0 a_1 \ell(\bar{P}, P) - \tau^{-1} c a_1 \ell(\bar{P}, Q') + \tau c_1 a_0 \ell(\bar{P}, Q). \end{aligned}$$

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

$$(65) \quad \begin{aligned} & n^{-1} \{c_2\mathbb{E} [\mathbf{T}(\mathbf{X}, P, Q')] - c_1\mathbb{E} [\mathbf{T}(\mathbf{X}, P, Q)]\} \\ & \leq c_2 (a_0\ell(\mathbf{P}^*, P) - a_1\ell(\mathbf{P}^*, Q')) - c_1 (a_1\ell(\mathbf{P}^*, P) - a_0\ell(\mathbf{P}^*, Q)) \\ & = c'_0 a_1 \ell(\mathbf{P}^*, P) - c_2 a_1 \ell(\mathbf{P}^*, Q') + c_1 a_0 \ell(\mathbf{P}^*, Q) \\ & \leq \tau c'_0 a_1 [\ell(\mathbf{P}^*, \bar{P}) + \ell(\bar{P}, P)] - c_2 a_1 [\tau^{-1} \ell(\bar{P}, Q') - \ell(\mathbf{P}^*, \bar{P})] \\ & \quad + \tau c_1 a_0 [\ell(\mathbf{P}^*, \bar{P}) + \ell(\bar{P}, Q)] \\ & \leq (e_1 + c_2) a_1 \ell(\mathbf{P}^*, \bar{P}) + \tau c'_0 a_1 \ell(\bar{P}, P) \\ & \quad - \tau^{-1} c_2 a_1 \ell(\bar{P}, Q') + \tau c_1 a_0 \ell(\bar{P}, Q), \end{aligned}$$

with

$$(66) \quad e_1 = \tau [c'_0 + c_1 a_0 / a_1] = \tau [c_2 (a_0 / a_1) + c_1 (a_0 / a_1 - 1)].$$

Using (60) and (64), we deduce that for all  $P, Q, Q' \in \mathcal{M}$

$$(67) \quad \begin{aligned} & \mathbb{E} [\exp [\beta (c\mathbf{T}(\mathbf{X}, P, Q') - c_1\mathbf{T}(\mathbf{X}, P, Q))]] \\ & = \prod_{i=1}^n \mathbb{E} [\exp [\beta (ct_{(P, Q')}(X_i) - c_1 t_{(P, Q)}(X_i))] ] \\ & = \exp [\beta (c\mathbb{E} [\mathbf{T}(\mathbf{X}, P, Q')] - c_1\mathbb{E} [\mathbf{T}(\mathbf{X}, P, Q)])] \prod_{i=1}^n \mathbb{E} [\exp [\beta U_i]] \\ & \leq \exp [n\beta [\Delta_1(P, Q) - \tau^{-1} c a_1 \ell(\bar{P}, Q')]] \end{aligned}$$

with

$$(68) \quad \Delta_1(P, Q) = e_0 a_1 \ell(\mathbf{P}^*, \bar{P}) + \tau c_1 a_0 \ell(\bar{P}, Q) + \frac{l_1^2 \beta}{8} - \tau^{-1} c_0 a_1 \ell(\bar{P}, P).$$

Using (61) and (65), we obtain similarly that for all  $P, Q, Q' \in \mathcal{M}$

$$(69) \quad \begin{aligned} & \mathbb{E} [\exp [\beta (c_2\mathbf{T}(\mathbf{X}, P, Q') - c_1\mathbf{T}(\mathbf{X}, P, Q))]] \\ & \leq \exp [n\beta [\Delta_2(P, Q) - \tau^{-1} c_2 a_1 \ell(\bar{P}, Q')]] \end{aligned}$$

with

$$(70) \quad \begin{aligned} \Delta_2(P, Q) &= (e_1 + c_2) a_1 \ell(\mathbf{P}^*, \bar{P}) + \tau c'_0 a_1 \ell(\bar{P}, P) + \tau c_1 a_0 \ell(\bar{P}, Q) \\ &\quad + \frac{l_2^2 \beta}{8}. \end{aligned}$$

Since  $2\gamma < \tau^{-1}c < \tau^{-1}c_2$ , we may apply Lemma 7 with  $\gamma_0 = \tau^{-1}c$  and  $\gamma_0 = \tau^{-1}c_2$  successively which leads to

$$(71) \quad \int_{\mathcal{M}} \exp[-\tau^{-1}cn\beta a_1 \ell(\bar{P}, Q')] d\pi(Q') \leq \pi(\mathcal{B}) \exp[\Xi_1]$$

and

$$(72) \quad \int_{\mathcal{M}} \exp[-\tau^{-1}c_2n\beta a_1 \ell(\bar{P}, Q')] d\pi(Q') \leq \pi(\mathcal{B}) \exp[\Xi_1]$$

with

$$(73) \quad \begin{aligned} \Xi_1 &= \log \left[ 1 + \frac{\exp[-(\tau^{-1}c - \gamma)]}{1 - \exp[-(\tau^{-1}c - 2\gamma)]} \right] \\ &\geq \log \left[ 1 + \frac{\exp[-(\tau^{-1}c_2 - \gamma)]}{1 - \exp[-(\tau^{-1}c_2 - 2\gamma)]} \right]. \end{aligned}$$

Putting (69) and (72) together leads to

$$\begin{aligned} \exp[\mathbf{L}(P, Q)] &= \int_{\mathcal{M}} \mathbb{E}[\exp[\beta(c_2 \mathbf{T}(\mathbf{X}, P, Q') - c_1 \mathbf{T}(\mathbf{X}, P, Q))] ] d\pi(Q') \\ &\leq \exp[n\beta \Delta_2(P, Q)] \int_{\mathcal{M}} \exp[-\tau^{-1}c_2n\beta a_1 \ell(\bar{P}, Q')] d\pi(Q') \\ &\leq \pi(\mathcal{B}) \exp[\Xi_1 + n\beta \Delta_2(P, Q)], \end{aligned}$$

and since, for all  $(P, Q) \in \mathcal{B}^2$ , by definition (70) of  $\Delta_2(P, Q)$ ,

$$(74) \quad \begin{aligned} \Delta_2(P, Q) &\leq (e_1 + c_2) a_1 \ell(\mathbf{P}^*, \bar{P}) + [\tau c'_0 a_1 + \tau c_1 a_0] r + \frac{l_2^2 \beta}{8} \\ &= (e_1 + c_2) a_1 \ell(\mathbf{P}^*, \bar{P}) + e_1 a_1 r + \frac{l_2^2 \beta}{8} = \Delta_2 \end{aligned}$$

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 (57) that

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

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

$$(75) \quad \log\left(\frac{1}{z}\right) = \xi + \log \frac{1}{\pi(\mathcal{B})} + \Xi_1 + n\beta \Delta_2.$$

Putting (67) and (71) together, we obtain that

$$\begin{aligned}
& \exp [\mathbf{M}(P, Q)] \\
&= \int_{\mathcal{M}} \mathbb{E} \left[ \exp \left[ \beta \left( c\mathbf{T}(\mathbf{X}, P, Q') - c_1\mathbf{T}(\mathbf{X}, P, Q) \right) \right] \right] d\pi(Q') \\
&\leq \exp [n\beta\Delta_1(P, Q)] \int_{\mathcal{M}} \exp \left[ -\tau^{-1}cn\beta a_1\ell(\bar{P}, Q') \right] d\pi(Q') \\
&\leq \pi(\mathcal{B}) \exp [\Xi_1 + n\beta\Delta_1(P, Q)].
\end{aligned}$$

It follows from the definition (68) of  $\Delta_1(P, Q)$  that for all  $P \in \mathcal{M}$  and for all  $Q \in \mathcal{B}$ ,

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

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

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

We derive from Lemma 3 that

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

hence,

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

Applying Lemma 7 with  $\gamma_0 = \tau^{-1} c_0 > 2\gamma$  and setting  $e_2 = \tau^{-1} c_0 - 2\gamma$ , we get

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

with

$$(77) \quad \Xi_2 = -\gamma + \log \left[ \frac{1}{1 - \exp [-e_2]} \right],$$

which together with (76) leads to

$$\begin{aligned}
 & \log \int_{\mathcal{C}_{\mathcal{B}}(\bar{P}, 2^J r)} \mathbb{E} [\exp [-\beta \mathbf{T}(\mathbf{X}, P)]] d\pi(P) \\
 & \leq \log [\pi(\mathcal{B})] + \Xi_1 + \Xi_2 \\
 (78) \quad & + n\beta \left[ e_0 a_1 \ell(\mathbf{P}^*, \bar{P}) + \tau c_1 a_0 r + \frac{l_1^2 \beta}{8} - e_2 a_1 2^J r \right].
 \end{aligned}$$

Using the definitions (75) of  $z$  and (74) of  $\Delta_2$  we deduce from (78) that

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

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 (79) is not larger than  $-\xi$ , provided that

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

or equivalently if

$$(80) \quad 2^J \geq \frac{1}{e_2} \left[ \frac{2\xi + 2\Xi_1 + \Xi_2}{\beta n a_1 r} + \frac{C_1 \ell(\mathbf{P}^*, \bar{P}) + C_2 r}{r} + \frac{[l_1^2 + l_2^2] \beta}{8 a_1 r} \right].$$

For  $\bar{P} \in \mathcal{M}(\beta)$  and the choice  $r = \ell(\mathbf{P}^*, \bar{P}) + a_1^{-1} \beta \geq r_n(\beta, \bar{P}) \geq 1/(n a_1 \beta)$ , (80) is satisfied if

$$2^J \geq \frac{1}{e_2} \left( 2\xi + 2\Xi_1 + \Xi_2 + C_1 + C_2 + \frac{[l_1^2 + l_2^2]}{8} \right)$$

and the requirement  $\kappa_0 \geq \tau(2^J + 1)$  if

$$(81) \quad \kappa_0 = \tau \left[ 1 + \frac{2}{e_2} \left( 2\xi + 2\Xi_1 + \Xi_2 + C_1 + C_2 + \frac{[l_1^2 + l_2^2]}{8} \right) \right].$$

**6.4. Proof of Theorem 2.** The proof follows the same lines as that of Theorem 1. Under Assumption 2-(iv), the random variables  $U_i$  and  $V_i$  defined by (58) and (59) are not larger than with  $b = c + c_1 = l_1$  and  $b = c_2 + c_1 = l_2$  respectively. Since under Assumption 3, for all  $i \in \{1, \dots, n\}$

$$\begin{aligned} \mathbb{E} [U_i^2] &\leq 2 [c^2 \text{Var} [t_{(P,Q')}(X_i)] + c_1^2 \text{Var} [t_{(P,Q)}(X_i)]] \\ &\leq 2a_2 [(c^2 + c_1^2)\ell(P_i^*, P) + c^2\ell(P_i^*, Q') + c_1^2\ell(P_i^*, Q)] \end{aligned}$$

and

$$\mathbb{E} [V_i^2] \leq 2a_2 [(c_2^2 + c_1^2)\ell(P_i^*, P) + c_2^2\ell(P_i^*, Q') + c_1^2\ell(P_i^*, Q)]$$

we may apply Lemma 6 and using the notations  $\Lambda_1 = \tau\phi(\beta l_1)$ ,  $\Lambda_2 = \tau\phi(\beta l_2)$  as well as Assumption 1, we get

$$\begin{aligned} &\frac{1}{n\beta} \log \left[ \prod_{i=1}^n \mathbb{E} [\exp [\beta U_i]] \right] \\ &\leq 2\phi(\beta l_1)\beta a_2 [(c^2 + c_1^2)\ell(\mathbf{P}^*, P) + c^2\ell(\mathbf{P}^*, Q') + c_1^2\ell(\mathbf{P}^*, Q)] \\ &\leq 2\Lambda_1\beta a_2 [c^2 + c_1^2] \ell(\mathbf{P}^*, \bar{P}) \\ (82) \quad &+ \Lambda_1\beta a_2 [(c^2 + c_1^2)\ell(\bar{P}, P) + c^2\ell(\bar{P}, Q') + c_1^2\ell(\bar{P}, Q)] \end{aligned}$$

and

$$\begin{aligned} &\frac{1}{n\beta} \log \left[ \prod_{i=1}^n \mathbb{E} [\exp [\beta V_i]] \right] \\ &\leq 2\Lambda_2\beta a_2 [c_2^2 + c_1^2] \ell(\mathbf{P}^*, \bar{P}) \\ (83) \quad &+ \Lambda_2\beta a_2 [(c_2^2 + c_1^2)\ell(\bar{P}, P) + c_2^2\ell(\bar{P}, Q') + c_1^2\ell(\bar{P}, Q)]. \end{aligned}$$

It follows from (64) that

$$\begin{aligned} E_1 &= n^{-1} \{ c\mathbb{E} [\mathbf{T}(\mathbf{X}, P, Q')] - c_1\mathbb{E} [\mathbf{T}(\mathbf{X}, P, Q)] \} \\ &\quad + 2\Lambda_1\beta a_2 [c^2 + c_1^2] \ell(\mathbf{P}^*, \bar{P}) \\ &\quad + \Lambda_1\beta a_2 [(c^2 + c_1^2)\ell(\bar{P}, P) + c^2\ell(\bar{P}, Q') + c_1^2\ell(\bar{P}, Q)] \\ &= [e_0 a_1 + 2\Lambda_1\beta a_2 (c^2 + c_1^2)] \ell(\mathbf{P}^*, \bar{P}) \\ &\quad - [\tau^{-1} c_0 a_1 - \Lambda_1\beta a_2 (c^2 + c_1^2)] \ell(\bar{P}, P) \\ &\quad - [\tau^{-1} c a_1 - \Lambda_1\beta a_2 c^2] \ell(\bar{P}, Q') \\ &\quad + [\tau c_1 a_0 + \Lambda_1\beta a_2 c_1^2] \ell(\bar{P}, Q). \end{aligned}$$



Using the definitions (21) of  $\tau^{-1}c_1$  and (22) of  $\tau^{-1}c_2$  and setting

$$\begin{aligned} e_3 &= e_0 + 2\Lambda_1\beta \frac{a_2(c_2^2 + c_1^2)}{a_1} \\ e_4 &= \frac{1}{a_1} [\tau c_1 a_0 + \Lambda_1\beta a_2 c_1^2] \end{aligned}$$

and arguing as in the proof of inequality (67), we deduce from (82) that

$$\begin{aligned} &\mathbb{E} [\exp [\beta (c\mathbf{T}(\mathbf{X}, P, Q') - c_1\mathbf{T}(\mathbf{X}, P, Q))]] \\ &\leq \exp [n\beta E_1] \\ (84) \quad &= \exp [n\beta a_1 (e_3\ell(\mathbf{P}^*, \bar{P}) - \tau^{-1} [c_1\ell(\bar{P}, P) + c_2\ell(\bar{P}, Q')] + e_4\ell(\bar{P}, Q))] \end{aligned}$$

It follows from (65) that

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

Using the definition (23) of  $\tau^{-1}c_3$ , setting

$$\begin{aligned} e_5 &= e_1 + c_2 + 2\Lambda_2\beta \frac{a_2(c_2^2 + c_1^2)}{a_1} \\ e_6 &= \tau c'_0 + \Lambda_2\beta \frac{a_2(c_2^2 + c_1^2)}{a_1} \\ e_7 &= \frac{1}{a_1} [\tau c_1 a_0 + \Lambda_2\beta a_2 c_1^2], \end{aligned}$$

and arguing as in the proof of (69), we deduce from (83) that

$$\begin{aligned} &\mathbb{E} [\exp [\beta (c_2\mathbf{T}(\mathbf{X}, P, Q') - c_1\mathbf{T}(\mathbf{X}, P, Q))]] \\ &\leq \exp [n\beta E_2] \\ (85) \quad &= \exp [n\beta a_1 (e_5\ell(\mathbf{P}^*, \bar{P}) + e_6\ell(\bar{P}, P) - \tau^{-1}c_3\ell(\bar{P}, Q') + e_7\ell(\bar{P}, Q))] . \end{aligned}$$

Under our assumption on  $\beta$ , we know that the quantities  $\tau^{-1}c_2$  and  $\tau^{-1}c_3$  are positive and that  $2\gamma < \tau^{-1}(c_2 \wedge c_3)$ . We may therefore apply Lemma 7

with  $\gamma_0 = \tau^{-1}c_2$  and  $\gamma_0 = \tau^{-1}c_3$  successively and get

$$(86) \quad \int_{\mathcal{M}} \exp[-\tau^{-1}c_2 n \beta a_1 \ell(\bar{P}, Q')] d\pi(Q') \leq \pi(\mathcal{B}) \exp[\bar{\Xi}_1]$$

and

$$(87) \quad \int_{\mathcal{M}} \exp[-\tau^{-1}c_3 n \beta a_1 \ell(\bar{P}, Q')] d\pi(Q') \leq \pi(\mathcal{B}) \exp[\bar{\Xi}_1]$$

with

$$(88) \quad \bar{\Xi}_1 = \log \left[ 1 + \frac{\exp[-(\tau^{-1}(c_2 \wedge c_3) - \gamma)]}{1 - \exp[-(\tau^{-1}(c_2 \wedge c_3) - 2\gamma)]} \right].$$

Putting (85) and (87) together, we obtain that for all  $(P, Q) \in \mathcal{B}^2$

$$\begin{aligned} & \exp[\mathbf{L}(P, Q)] \\ &= \int_{\mathcal{M}} \mathbb{E} [\exp[\beta(c_2 \mathbf{T}(\mathbf{X}, P, Q') - c_1 \mathbf{T}(\mathbf{X}, P, Q))] ] d\pi(Q') \\ &\leq \exp[n\beta a_1 (e_5 \ell(\mathbf{P}^*, \bar{P}) + e_6 \ell(\bar{P}, P) + e_7 \ell(\bar{P}, Q))] \\ &\quad \times \int_{\mathcal{M}} \exp[-\tau^{-1}c_3 n \beta a_1 \ell(\bar{P}, Q')] d\pi(Q') \\ &\leq \pi(\mathcal{B}) \exp[\bar{\Xi}_1 + n\beta a_1 (e_5 \ell(\mathbf{P}^*, \bar{P}) + e_6 \ell(\bar{P}, P) + e_7 \ell(\bar{P}, Q))] \\ &\leq \pi(\mathcal{B}) \exp[\bar{\Xi}_1 + n\beta a_1 (e_5 \ell(\mathbf{P}^*, \bar{P}) + (e_6 + e_7)r)]. \end{aligned}$$

Consequently,

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

We deduce from (57) that

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

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

$$(89) \quad \log\left(\frac{1}{z}\right) = \xi + \log \frac{1}{\pi(\mathcal{B})} + \bar{\Xi}_1 + n\beta a_1 [e_5 \ell(\mathbf{P}^*, \bar{P}) + (e_6 + e_7)r].$$

Putting (84) and (86) together, we obtain that for all  $Q \in \mathcal{B}$

$$\begin{aligned} & \exp[\mathbf{M}(P, Q)] \\ &= \int_{\mathcal{M}} \mathbb{E} [\exp[\beta(c \mathbf{T}(\mathbf{X}, P, Q') - c_1 \mathbf{T}(\mathbf{X}, P, Q))] ] d\pi(Q') \\ &\leq \exp[n\beta a_1 (e_3 \ell(\mathbf{P}^*, \bar{P}) - \tau^{-1}c_1 \ell(\bar{P}, P) + e_4 \ell(\bar{P}, Q))] \\ &\quad \times \int_{\mathcal{M}} \exp[-\tau^{-1}c_2 n \beta a_1 \ell(\bar{P}, Q')] d\pi(Q') \\ &\leq \pi(\mathcal{B}) \exp[\bar{\Xi}_1 + n\beta a_1 (e_3 \ell(\mathbf{P}^*, \bar{P}) + e_4 r - \tau^{-1}c_1 \ell(\bar{P}, P))]. \end{aligned}$$

We derive from Lemma 3 that

$$\begin{aligned} & \mathbb{E} [\exp [-\beta \mathbf{T}(\mathbf{X}, P)]] \\ & \leq \frac{1}{\pi(\mathcal{B})} \left[ \int_{\mathcal{B}} \exp [-\mathbf{M}(P, Q)] d\pi_{\mathcal{B}}(Q) \right]^{-1} \\ & \leq \exp [\bar{\Xi}_1 + n\beta a_1 (e_3 \ell(\mathbf{P}^*, \bar{P}) + e_4 r - \tau^{-1} c_1 \ell(\bar{P}, P))], \end{aligned}$$

and consequently,

$$\begin{aligned} & \int_{c_{\mathcal{B}}(\bar{P}, 2^J r)} \mathbb{E} [\exp [-\beta \mathbf{T}(\mathbf{X}, P)]] d\pi(P) \\ & \leq \exp [\bar{\Xi}_1 + n\beta a_1 (e_3 \ell(\mathbf{P}^*, \bar{P}) + e_4 r)] \\ (90) \quad & \times \int_{c_{\mathcal{B}}(\bar{P}, 2^J r)} \exp [-\tau^{-1} c_1 n\beta a_1 \ell(\bar{P}, P)] d\pi(P). \end{aligned}$$

Since under our assumptions,  $\tau^{-1} c_1 > 0$  and  $2\gamma < \tau^{-1} c_1$  we may apply Lemma 7 with  $\gamma_0 = \tau^{-1} c_1$ , which leads to

$$\int_{c_{\mathcal{B}}(\bar{P}, 2^J r)} \exp [-\tau^{-1} c_1 n\beta a_1 \ell(\bar{P}, P)] d\pi(P) \leq \pi(\mathcal{B}) \exp [\bar{\Xi}_2 - (\tau^{-1} c_1 - 2\gamma) n\beta a_1 2^J r].$$

with

$$(91) \quad \bar{\Xi}_2 = -\gamma + \log \left[ \frac{1}{1 - \exp [-(\tau^{-1} c_1 - 2\gamma)]} \right],$$

which together with (90) leads to

$$\begin{aligned} & \int_{c_{\mathcal{B}}(\bar{P}, 2^J r)} \mathbb{E} [\exp [-\beta \mathbf{T}(\mathbf{X}, P)]] d\pi(P) \\ (92) \quad & \leq \pi(\mathcal{B}) \exp [\bar{\Xi}_1 + \bar{\Xi}_2 + n\beta a_1 (e_3 \ell(\mathbf{P}^*, \bar{P}) + e_4 r - (\tau^{-1} c_1 - 2\gamma) 2^J r)]. \end{aligned}$$

Using the definition (89) of  $z$ , we deduce that

$$\begin{aligned} & \log \left[ \frac{1}{z} \int_{c_{\mathcal{B}}(\bar{P}, 2^J r)} \mathbb{E} [\exp [-\beta \mathbf{T}(\mathbf{X}, P)]] d\pi(P) \right] \\ & \leq \log \left( \frac{1}{z} \right) + \log \pi(\mathcal{B}) + \bar{\Xi}_1 + \bar{\Xi}_2 + n\beta a_1 (e_3 \ell(\mathbf{P}^*, \bar{P}) + e_4 r - (\tau^{-1} c_1 - 2\gamma) 2^J r) \\ & = \xi + \log \frac{1}{\pi(\mathcal{B})} + \bar{\Xi}_1 + n\beta a_1 [e_5 \ell(\mathbf{P}^*, \bar{P}) + (e_6 + e_7) r] \\ & \quad + \log \pi(\mathcal{B}) + \bar{\Xi}_1 + \bar{\Xi}_2 + n\beta a_1 (e_3 \ell(\mathbf{P}^*, \bar{P}) + e_4 r - (\tau^{-1} c_1 - 2\gamma) 2^J r) \\ & = \xi + 2\bar{\Xi}_1 + \bar{\Xi}_2 + n\beta a_1 [(e_3 + e_5) \ell(\mathbf{P}^*, \bar{P}) + (e_4 + e_6 + e_7) r] \\ & \quad - (\tau^{-1} c_1 - 2\gamma) n\beta a_1 2^J r. \end{aligned}$$

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

$$2^J \geq \frac{1}{\tau^{-1} c_1 - 2\gamma} \left[ \frac{2\xi + 2\bar{\Xi}_1 + \bar{\Xi}_2}{n\beta a_1 r} + \left[ (e_3 + e_5) \frac{\ell(\mathbf{P}^*, \bar{P})}{r} + e_4 + e_6 + e_7 \right] \right].$$

Choosing  $r = \ell(\mathbf{P}^*, \bar{P}) + r_n(\beta, \bar{P}) \geq 1/(n\beta a_1)$ , this inequality is satisfied as soon as

$$2^J \geq \frac{1}{\tau^{-1}c_1 - 2\gamma} [2\xi + 2\bar{\Xi}_1 + \bar{\Xi}_2 + e_3 + e_5 + e_4 + e_6 + e_7]$$

and the requirement  $\kappa_0 \geq \tau(2^J + 1)$  if

$$(93) \quad \kappa_0 = \tau \left[ 1 + \frac{2}{\tau^{-1}c_1 - 2\gamma} (2\xi + 2\bar{\Xi}_1 + \bar{\Xi}_2 + e_3 + e_5 + e_4 + e_6 + e_7) \right].$$

## 7. OTHER PROOFS

**7.1. Proof of Proposition 2.** Let  $|\bar{\theta}| \leq \sigma t$  and  $F_\sigma$  be the distribution function of  $\nu_\sigma$ . For conveniency, when  $L$  is finite, we define  $H(x) = 1$  for all  $x \geq L$ . Since the total variation distance is translation invariant,

$$\|P_\theta - P_{\bar{\theta}}\| = \|P_{\theta - \bar{\theta}} - P_0\| = \|P_{\bar{\theta} - \theta} - P_0\| = H(|\bar{\theta} - \theta|)$$

for all  $\theta, \bar{\theta} \in \mathbb{R}$ . We distinguish between two cases

**Case 1:**  $r_0 = H(\sigma t) \leq 1/4$ . Since  $q$  is symmetric, positive and decreasing on  $\mathbb{R}_+$ , for all  $r \leq r_0$ ,  $G(r) \leq \sigma t$  and

$$\begin{aligned} \frac{\pi(\mathcal{B}(P_{\bar{\theta}}, 2r))}{\pi(\mathcal{B}(P_{\bar{\theta}}, r))} &= \frac{\nu_\sigma(\{\theta \in \mathbb{R}, \|P_\theta - P_{\bar{\theta}}\| \leq 2r\})}{\nu_\sigma(\{\theta \in \mathbb{R}, \|P_\theta - P_{\bar{\theta}}\| \leq r\})} = \frac{\nu_\sigma(\{\theta \in \mathbb{R}, H(|\theta - \bar{\theta}|) \leq 2r\})}{\nu_\sigma(\{\theta \in \mathbb{R}, H(|\theta - \bar{\theta}|) \leq r\})} \\ &= \frac{\nu_\sigma(\{\theta \in \mathbb{R}, |\theta - \bar{\theta}| \leq G(2r)\})}{\nu_\sigma(\{\theta \in \mathbb{R}, |\theta - \bar{\theta}| \leq G(r)\})} \leq \frac{2q_\sigma(0)G(2r)}{2q_\sigma(|\theta| + G(r))G(r)} \\ &\leq \frac{q_\sigma(0)G(2r)}{q_\sigma(|\theta| + G(r_0))G(r)} \leq \frac{q_\sigma(0)G(2r)}{q_\sigma(|\theta| + \sigma t)G(r)} \leq \frac{q_\sigma(0)G(2r)}{q_\sigma(2\sigma t)G(r)} \\ &= \frac{q(0)G(2r)}{q(2t)G(r)} \leq \frac{\bar{\Gamma}}{q(2t)}. \end{aligned}$$

For all  $r_0 < r < 1$ ,  $|\bar{\theta}| \leq \sigma t = G(r_0) \leq G(r)$ , hence  $F_\sigma(|\bar{\theta}| - G(r)) \leq F_\sigma(0) = 1/2$  and  $F_\sigma(|\bar{\theta}| + G(r)) \geq F_\sigma(G(r)) \geq F_\sigma(\sigma t) = F_1(t) \geq 3/4$ . Consequently,

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

**Case 2:**  $r_0 > 1/4$ . Arguing as before, we obtain that for all  $r \leq 1/4 < r_0$ ,

$$\begin{aligned} \frac{\pi(\mathcal{B}(P_{\bar{\theta}}, 2r))}{\pi(\mathcal{B}(P_{\bar{\theta}}, r))} &\leq \frac{2q_\sigma(0)G(2r)}{2q_\sigma(|\theta| + G(r))G(r)} = \frac{q_\sigma(0)G(2r)}{q_\sigma(|\theta| + G(r_0))G(r)} \\ &\leq \frac{q_\sigma(0)G(2r)}{q_\sigma(2\sigma t)G(r)} \leq \frac{\bar{\Gamma}}{q(2t)}. \end{aligned}$$

For all  $r \in (1/4, 1)$ ,  $G(r) \geq G(1/4)$  and  $G(1/4) \leq \sigma t$

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

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

$$(94) \quad \log \left( \frac{\pi(\mathcal{B}(P_{\bar{\theta}}, 2r))}{\pi(\mathcal{B}(P_{\bar{\theta}}, r))} \right) \leq \max \left\{ \log \left( \frac{\bar{\Gamma}(\sigma \vee 1)}{q(2t)} \right), \log 4 \right\},$$

hence, for all  $r \geq a_1^{-1}\beta$

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

The right-hand side is not larger than  $\beta$  provided that it satisfies (27) and this lower bound is not smaller than  $1/\sqrt{n}$  under the assumption  $\eta \leq \log 4$ . We conclude by using (16).

**7.2. Proof of Proposition 3.** Let us take  $t = (|\bar{\theta}|/\sigma) \vee t_0$ . For such a value of  $t$ ,  $\bar{\theta} \in [-\sigma t, \sigma t]$  and  $\nu_1([t, +\infty)) \leq 1/4$ . Since Assumption 4 is satisfied, (94) holds true and we deduce from (12) that

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

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

**7.3. Proof of Corollary 1.** Throughout this proof,  $a_0 = 3/2$ ,  $a_1 = 1/2$ ,  $\tau = 1$ ,  $\mathcal{A}$  the  $\sigma$ -algebra generated by all the subsets of  $\mathcal{M}$  so that, Assumptions 1 and 2 hold. The constants  $c = 1/3$ , hence  $\lambda = 4\beta/3$ , and  $\gamma = 1/6$  satisfy the constraints of Theorem 1. We may therefore apply it with these choices. Besides, we set  $\delta = \delta_n$ ,  $\eta = \eta_n$ ,  $\beta = \beta_n$  and  $\Theta = \Theta[\eta, \delta]$  for short and

$$(95) \quad J_n = \exp \left[ \frac{(K^2 - 1)\gamma\tau^4 a_1^2 n \eta_n^2}{2(k+1)} \right]$$

so that  $\mathcal{M}_n(K)$  gathers 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  $\theta = \theta(P) = (\bar{p}, \bar{\mathbf{m}}, \bar{\sigma}) \in \Theta$  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

$$(96) \quad \frac{\bar{\sigma}}{(1 + \delta)} \leq \sigma < \bar{\sigma} \quad \text{and} \quad \bar{m}_i = j_i \bar{\sigma} \delta \leq m_i < \bar{m}_i + \bar{\sigma} \delta,$$

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

$$(97) \quad 0 \leq \left(1 - \frac{\sigma}{\bar{\sigma}}\right) \leq \frac{\delta}{1 + \delta} < \delta \quad \text{and} \quad \left| \frac{\mathbf{m} - \bar{\mathbf{m}}}{\bar{\sigma}} \right|_{\infty} \leq \delta,$$

and we infer from (31) and (32) and the fact that the total variation loss is translation and scale invariant that  $P_{\theta}$  satisfies

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

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 (97), 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,

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

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

$$(98) \quad |(j_0, \mathbf{j})|_{\infty} \leq 1 + \frac{1}{\log(1 + \delta)} \left[ |\log \sigma| \vee \left| \frac{\mathbf{m}}{\sigma} \right|_{\infty} \right] \leq 1 + J_n.$$

For all  $r > 0$ ,  $e^{-L\theta} \leq \pi(\mathcal{B}(P_\theta, r)) \leq 1$  and these two inequalities together with the definition (35) of  $\eta$  and Assumption 5 imply that for all  $r > 0$

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

Using (98), the definition (95) of  $J_n$  and the fact that  $\log(2+x) \leq \log 3 + \log x$  for all  $x \geq 1$ , we derive that

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

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

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

For the choice of  $\beta = \beta_n$  given by (37),

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

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

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

and the result follows from Theorem 1 and the fact that  $P$  is arbitrary in  $\mathcal{M}_n(K)$ .

**7.4. Proof of Lemma 1.** By doing the change of variables  $u = x - m$  in (43) if ever necessary, we may assume with no loss of generality that  $m > 0$ . Then, since  $p$  is non-increasing 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 (43).

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 (42).

Finally, by combining (43) and (42) we deduce that for all  $m \in \mathbb{R}$  and  $\sigma \geq 1$

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

which yields to (44).

**7.5. Proof of Corollary 2.** It follows from Proposition 1 and our condition on  $p$  that the family  $\mathcal{T}(\ell, \mathcal{M})$  satisfies Assumption 2 with  $a_0 = 2$ ,  $a_1 = 3/16$  and  $a_2 = 3\sqrt{2}/4$  for the loss  $\ell = h^2$ . Besides, Assumption 1 holds true with  $\tau = 2$  and the constants  $\gamma = 0.01$ ,  $\beta = 0.01$ ,  $\lambda = (1+c)\beta$  with  $c = 0.05$  satisfy the constraints of Theorem 2.

We also use the following lemma the proof of which is postponed to Section 7.6.

**Lemma 8.** *Let  $\theta \in \mathbb{R}^k$  be such that  $|\theta|_{\infty} \leq R$ . For all  $m \subset \{1, \dots, k\}$  and  $r > 0$*

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

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

$$(99) \quad \nu_m \left( \left\{ \theta' \in \mathbb{R}^k, |\theta' - \theta|_{\infty} \leq r \right\} \right) \geq \frac{1}{2^{|m|}} \left( \frac{r}{R} \wedge 1 \right)^{|m|}$$



and for all  $K > 1$

$$(100) \quad \frac{\nu_m(\{\boldsymbol{\theta}' \in \mathbb{R}^k, |\boldsymbol{\theta}' - \boldsymbol{\theta}|_\infty \leq Kr\})}{\nu_m(\{\boldsymbol{\theta}' \in \mathbb{R}^k, |\boldsymbol{\theta}' - \boldsymbol{\theta}|_\infty \leq r\})} \leq K^{|m|}.$$

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

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

Finally, let  $\bar{\boldsymbol{\theta}}$  for some arbitrary element of  $\Theta_{\bar{m}}(R)$ .

It follows from (45) and (99) that for all  $r > 0$ ,

$$(101) \quad \begin{aligned} 1 &\geq \pi_m(\mathcal{B}(P_{\bar{\boldsymbol{\theta}}}, r)) \\ &= \nu_m(\{\boldsymbol{\theta} \in \mathbb{R}^k, h^2(P_{\bar{\boldsymbol{\theta}}}, P_{\boldsymbol{\theta}}) \leq r\}) \\ &\geq \nu_m(\{\boldsymbol{\theta} \in \mathbb{R}^k, |\boldsymbol{\theta} - \bar{\boldsymbol{\theta}}|_\infty \leq (r/B)^{1/(2\alpha)}\}) \\ &\geq \frac{1}{2^{|m|}} \left( \frac{(r/B)^{1/(2\alpha)}}{R} \wedge 1 \right)^{|m|} \geq \frac{1}{2^{|m|}} \left( \frac{(r \wedge 1)^{1/(2\alpha)}}{RB^{1/(2\alpha)}} \right)^{|m|}, \end{aligned}$$

where the last inequality holds true under the assumption that  $RB^{1/(2\alpha)} \geq 1$ .

We deduce from (101) that for all  $r > 0$

$$(102) \quad \begin{aligned} \frac{\pi(\mathcal{B}(P_{\bar{\boldsymbol{\theta}}}, 2r))}{\pi(\mathcal{B}(P_{\bar{\boldsymbol{\theta}}}, r))} &\leq \frac{1}{\pi(\mathcal{B}(P_{\bar{\boldsymbol{\theta}}}, r))} \\ &\leq \frac{1}{\sum_{m \subset \{1, \dots, k\}} e^{-L_m} \nu_m(\{\boldsymbol{\theta} \in \mathbb{R}^k, |\boldsymbol{\theta} - \bar{\boldsymbol{\theta}}|_\infty \leq (r/B)^{1/(2\alpha)}\})} \\ &\leq \frac{e^{L_{\bar{m}}}}{\nu_{\bar{m}}(\{\boldsymbol{\theta} \in \Theta_{\bar{m}}, |\boldsymbol{\theta} - \bar{\boldsymbol{\theta}}|_\infty \leq (r/B)^{1/(2\alpha)}\})} \\ &\leq \exp \left[ L_{\bar{m}} + |\bar{m}| \log \left( \frac{2RB^{1/(2\alpha)}}{(r \wedge 1)^{1/(2\alpha)}} \right) \right] \\ &= \exp \left[ |\bar{m}| \log(2kRB^{1/(2\alpha)}) + k \log \left( 1 + \frac{1}{k} \right) + \frac{|\bar{m}|}{2\alpha} \log \left( \frac{1}{r} \vee 1 \right) \right]. \end{aligned}$$

Provided that

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

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

and we deduce from (102) that  $r_n(\beta, P_{\bar{\theta}})$  defined by (12) satisfies

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

By applying Theorem 2, we conclude that (24) holds with

$$r \leq \ell(\mathbf{P}^*, P_{\bar{\theta}}) + \frac{|\bar{m}| \log \left( 2kR(nB)^{1/(2\alpha)} \right) + 1}{\gamma n \beta a_1}$$

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

**7.6. Proof of Lemma 8.** 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  $\boldsymbol{\theta} \in \mathbb{R}^k$  such that  $|\boldsymbol{\theta}|_\infty \leq R$ . For all  $m \subset \{1, \dots, k\}$ ,  $m \neq \emptyset$ ,

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

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

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

If  $m = \emptyset$ ,

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

Let us now turn to the proof of (100). Since  $\boldsymbol{\theta} \in \Theta_m(R)$ , for all  $K' \in \{1, K\}$

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

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.

## REFERENCES

- Alquier, P. (2008). PAC-Bayesian bounds for randomized empirical risk minimizers. *Math. Methods Statist.*, 17(4):279–304.
- Audibert, J.-Y. and Catoni, O. (2011). Linear regression through PAC-Bayesian truncation. *arXiv:1010.0072*.
- Baraud, Y. (2021). Tests and estimation strategies associated to some loss functions. *Probab. Theory Relat. Fields*, 180(3):799–846.

- Baraud, Y. and Birgé, L. (2018). Rho-estimators revisited: General theory and applications. *Ann. Statist.*, 46(6B):3767–3804.
- Baraud, Y. and Birgé, L. (2020). Robust bayes-like estimation: Rho-bayes estimation. *Ann. Statist.*, 48(6):3699–3720.
- Baraud, Y., Birgé, L., and Sart, M. (2017). A new method for estimation and model selection:  $\rho$ -estimation. *Invent. Math.*, 207(2):425–517.
- Birgé, L. (1983). Approximation dans les espaces métriques et théorie de l'estimation. *Z. Wahrsch. Verw. Gebiete*, 65(2):181–237.
- Birgé, L. (2006). Model selection via testing: an alternative to (penalized) maximum likelihood estimators. *Ann. Inst. H. Poincaré Probab. Statist.*, 42(3):273–325.
- Birgé, L. (2015). About the non-asymptotic behaviour of Bayes estimators. *J. Statist. Plann. Inference*, 166:67–77.
- Birgé, L. and Massart, P. (1998). Minimum contrast estimators on sieves: exponential bounds and rates of convergence. *Bernoulli*, 4(3):329–375.
- Birman, M. v. and Solomjak, M. Z. (1967). Piecewise polynomial approximations of functions of classes  $W_p^\alpha$ . *Mat. Sb. (N.S.)*, 73 (115):331–355.
- Catoni, O. (2004). 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-Verlag, Berlin.
- Ghosal, S., Ghosh, J. K., and van der Vaart, A. W. (2000). Convergence rates of posterior distributions. *Ann. Statist.*, 28(2):500–531.
- Ibragimov, I. A. and Has'minskiĭ, R. Z. (1981). *Statistical Estimation. Asymptotic Theory*, volume 16. Springer-Verlag, New York.
- Le Cam, L. (1973). Convergence of estimates under dimensionality restrictions. *Ann. Statist.*, 1:38–53.
- Massart, P. (2007). *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.

DEPARTMENT OF MATHEMATICS,  
 UNIVERSITY OF LUXEMBOURG  
 MAISON DU NOMBRE  
 6 AVENUE DE LA FONTE  
 L-4364 ESCH-SUR-ALZETTE  
 GRAND DUCHY OF LUXEMBOURG  
 Email address: yannick.baraud@uni.lu