# On Estimating a Density Using Hellinger Distance and Some Other Strange Facts

Lucien Birgé*

U.E.R. de Sciences Economiques, Université Paris X – Nanterre, 200 Av. de la République,
F-92001 Nanterre-Cedex, France, and
Mathematical Sciences Research Institute, Berkeley, California, USA

**Résumé.** On s'intéresse ici aux possibles vitesses d'estimation d'une densité à support compact dans $\mathbb{R}^m$ sous des hypothèses de régularité, lorsque la perte est mesurée par le carré de la distance de Hellinger (on regardera aussi le cas connu des normes $\mathbb{L}^q$ pour $1 \leq q \leq 2$) et le risque est le risque minimax sur la famille. On donne une méthode générale permettant de traiter les problèmes dans le cadre de la théorie de l'approximation sous des conditions concernant l'entropie métrique et l' $\varepsilon$-capacité des familles à estimer. Les rapports entre régularité et entropie métrique étant bien connus, nous pourrons aussi traiter les cas classiques et d'autres qui le sont moins. Sous des conditions de bornes inférieures les vitesses sont celles observées pour la norme $\mathbb{L}^1$ mais elles diffèrent dans le cas général. On montre aussi que les restrictions sur la compacité du support ou la régularité sont indispensables et que leur absence mène à l'impossibilité d'obtenir une estimation raisonnable en ce sens que n'importe quelle suite d'estimateurs sera arbitrairement mauvaise en un point au moins. Un résultat analogue est vrai sous des conditions de régularité.

## I. Introduction

In recent years, numerous papers have been devoted to studying the problem of global estimation of a density function on $\mathbb{R}^m$, especially if this density is known to belong to a family of functions with a given smoothness, the loss being measured by a power of the distance associated to the $\mathbb{L}^q$ norm for $1 \leq q \leq +\infty$. Pioneering work was done by Bretagnolle and Huber [5] and a very extensive study may be found in Ibragimov and Khas'minskii [8] and [9]. The best obtainable rates of convergence are now well known for the classical smooth families. Some connected results appeared in slightly different directions: in [5], some results are given in particular cases using the Hellinger

---

distance as a loss function. In [18] and [19] the classical results are extended to the estimation of derivatives of a density. In his recent paper [7] Devroye proved that estimation is impossible unless there exist uniform bounds for both the smoothness of the functions in the family and the compactness of their support. Actually a similar result was already known from Ibragimov and Khas'minskii (see [8] and [9]) but Devroye's theorem is stronger. Related facts are also found in [4]. The study of more general rates of convergence than the usual polynomial ones is carried out in [2] or [3] where it is proved that all "reasonable" rates may occur. Those special facts concerning the limits of the possibility of estimating a function without enough "a priori" restrictions, the fact that the speeds which are found when one measures the loss using Hellinger distance are different from those one finds when one works with $\mathbb{L}^1$ or $\mathbb{L}^2$ norms and the lack of unification in the theory suggested the following study with those three purposes:

i) Try to get a general way of finding the best rates of convergence, even for more general families than the usual ones, by putting the problem in the framework which is given in [2] or [3] and links the rates to the dimensional properties of the family. This allows us to express the assumptions in terms of approximation theory and then use all the well-known results from this theory to apply it to any particular case of density estimation. With these tools we shall be able to treat families with special moduli of continuity which are not found in the literature.

ii) Simultaneously give an extensive treatment for the case of Hellinger distance $h$.

iii) Indicate useful and general tools for finding lower bounds and use them in order to strengthen some known results in this direction and especially complement the results of Devroye [7].

We shall also recover a number of classical results from this theory, mainly in Sect. 3. Most of this section consists in a summary of previous results and is included here for sake of completeness as an illustration of the relations between the rates of convergence and the metric structure of the parameter space.

Actually, there is no fundamental difference in the treatment between $\mathbb{L}^1$ and Hellinger; the fact that different rates may occur when one uses Hellinger distance is related to this very simple remark: when you compute $h(f; f+g)$ where $f$ is a density and $g$ a small perturbation of order $\varepsilon$ for the sup norm, if $f$ is not small, this distance will be of order $\varepsilon$ just like the $\mathbb{L}^1$ distance, but if $f$ is also of order $\varepsilon$ then $h$ will be of order $\sqrt{\varepsilon}$, $h$ and $\mathbb{L}^1$ being no longer comparable. This will lead us to distinguish in the sequel between two cases, $\alpha$) and $\beta$): either $\alpha$) our family $\Theta$ has a uniform lower bound and then $h$ just behaves as the variation distance and so does not deserve a special treatment (that was the assumption used in [2] and [3]). Or $\beta$) there is no such lower bound, the densities in $\Theta$ may be arbitrarily small on some interval and we then get different estimates for the rates, such as the ones which were found in [5]. When dealing with the $\mathbb{L}^q$ norm, only case $\alpha$) will be relevant.

In Sect. 2 we shall set up our assumptions using a dimensional point of view (see [3, 13] and [14]); they will allow us to give a simplified and unified

presentation of the results. Translating smoothness assumptions into those is straightforward most of the time and is based on known results in approximation theory which we shall use freely. Section 3 will show how those assumptions give us upper bounds for the risk and Sect. 4 will be devoted to the tools necessary to get the corresponding lower bounds. Also a result in the spirit of [7] Theorem 1, ii) will be given for the usual smooth families. This will be Corollary 4.5 and will lead naturally to the problem of Sect. 5. In [7], Devroye proved that if we consider all density functions on [0; 1] which are bounded by 2, for any sequence $\{c_n\}$ converging to 0 and any sequence $\{\hat{f}_n\}$ of estimates there exists a density $f$ such that

$$\limsup_n c_n^{-1} \, \mathbb{E}_f[\int |\hat{f}_n(x) - f(x)| \, dx] \geq 1. \tag{1.3}$$

Notice the "*lim sup*" here which suggests the two questions: is this result a truly "*asymptotic*" one or not; what is to be expected about the "*lim inf*"? The proof of the analogous Corollary 4.5 suggests that it might be really "*asymptotic*" and, as L. Le Cam pointed out to me, in the finite dimensional case (ordinary parametric case), it is always possible to design a sequence of estimates for which the "*lim inf*" is uniformly small at each point. It is not difficult to see that in non-parametric cases with compact parameter spaces having finite $\varepsilon$-entropy for all positive $\varepsilon$, the same will be true. On the contrary, in the case studied by Devroye, the parameter space is not compact and contains $\varepsilon$-separated subsets of arbitrarily large cardinality. Using such subsets it is then possible to see that (1.3) also holds with "$\inf_n$" instead of "$\limsup_n$" provided that the sequence $\{c_n\}$ has a suitable upper bound. We shall also prove the analogous result for Hellinger distance. See [4] for similar non-asymptotic results for Hellinger balls.

Before proceeding, let us set up some notations and recall a few results. For two probability measures $P$ and $Q$, with densities $f$ and $g$ with respect to $\mu$, we may consider the following distances

$$h(P; Q) = \tfrac{1}{2} \int (\sqrt{f} - \sqrt{g})^2 \, d\mu;$$
$$d_q(P; Q) = \|f - g\|_q = [\int |f - g|^q \, d\mu]^{1/q};$$
$$D(P; Q) = \tfrac{1}{2} \int |f - g| \, d\mu = \tfrac{1}{2} d_1(P; Q);$$

which are respectively Hellinger distance, $\mathbb{L}^q$-distance and total variation distance. Some other quantities are also of interest which are the testing affinity, the Hellinger affinity and the Kullback information:

$$\pi(P; Q) = \int (f \wedge g) \, d\mu; \quad \rho(P; Q) = \int \sqrt{fg} \, d\mu; \quad K(P; Q) = \int \log \frac{f}{g} \, f \, d\mu.$$

The following relations are well-known (see for example [6, 13] or [14]):

$$\pi = 1 - D; \quad h^2 = 1 - \rho; \quad h^2 \leq D \leq h\sqrt{2}; \quad \pi^2 \leq \rho^2 \leq \pi(2 - \pi), \tag{1.1}$$

$$\rho(P^n, Q^n) = \rho^n(P, Q); \quad K(P^n, Q^n) = nK(P, Q). \tag{1.2}$$

If $\Theta$ is our parameter space assumed to be a set of densities with respect to the measure $\mu$, the points in $\Theta$ will be called alternatively $\theta$, $f_\theta$ (density) or $P_\theta$ (probability) and by convention for any distance $d$ on $\Theta$ $d(\theta; \theta') = d(f_\theta; f_{\theta'}) = d(P_\theta; P_{\theta'})$. Given $n$ i.i.d. observations following the unknown law $P_\theta$, the minimax risk will be defined by

$$R_M(d^q, n) = \inf_{T_n} \sup_{\theta \in \Theta} \mathbb{E}_\theta[d^q(P_\theta; T_n)]$$

where $T_n$ is any estimate (considered as a random probability) depending on $n$ observations. Sometimes, when $d = d_q$ and $q > 1$ we shall also use the normalized risk, which is

$$R'_M(d_q^q, n) = \inf_{T_n} \sup_{\theta \in \Theta} \mathbb{E}_\theta\left[\frac{d_q^q(P_\theta; T_n)}{\|f_\theta\|_q^q}\right].$$

We shall also need a few dimensional concepts concerning metric spaces. $\mathcal{B}(x, r)$ will denote the open ball of center $x$ and radius $r$. If $S$ is a subset of a metric space, $N$ is called an $\varepsilon$-net for $S$ if any point in $S$ is within a distance $\varepsilon$ of some point in $N$; $S$ is said to be $\varepsilon$-separated if any two points in $S$ are at a distance greater than $\varepsilon$. For any totally bounded subset $\Theta$ of a metric space we may define its $\varepsilon$-entropy $N_\varepsilon$. It is the logarithm (base 2) of the minimum cardinality of an $\varepsilon$-net. Its $\varepsilon$-capacity $\mathscr{C}_\varepsilon$ is the logarithm of the maximum cardinality of an $\varepsilon$-separated set. The definitions and properties of those quantities and the way of computing them in many cases are given in the remarkable paper [12] by Kolmogorov and Tikhomirov.

Throughout this paper we shall denote by $A_i$, $C$ or $C_i$ different constants but the $A_i$'s will always have the same meaning (issued from our assumptions) while the $C_i$'s will be generic constants (often depending on the former ones). For two sequences $\{u_n\}$ and $\{v_n\}$, $u_n \asymp v_n$ will mean that $C_1 u_n \leq v_n \leq C_2 u_n$ for some positive constants $C_1, C_2$. In this case we shall say that $\{u_n\}$ and $\{v_n\}$ are of the same order.

## II. Smoothness and Dimension

As was mentioned in the introduction a smoothness assumption for a family of densities generally implies some dimensional properties. As a very simple illustration we may consider the case of all density functions on $[0; 1]$ satisfying a uniform Lipschitz condition such as

$$|f(x) - f(y)| \leq C|x - y| \qquad x, y \in [0; 1].$$

Then the $\varepsilon$-entropy $N_\varepsilon$ of this set is of order $\varepsilon^{-1}$. Using the techniques introduced in [2] and [3] we immediately find, by solving the equation $n\varepsilon^2 = C_1 \varepsilon^{-1}$ that an upper bound for the speed of estimation will be $C_2 n^{-1/3}$ when the risk is measured by $d_1$ for example. Analogous considerations using

results concerning $\varepsilon$-capacity of this set give us the lower bounds. Since the relations between smoothness and dimension have been known for a long time by workers in approximation theory (see [12] and more recent results in [15] and [16] for example), we may express our assumptions in terms of dimensional properties of the density families. This will allow us to give a unified treatment of different problems. Our A1 reduces to an assumption concerning the metric $\varepsilon$-entropy of $\Theta$, A2 is closely related to the $\varepsilon$-capacity of $\Theta$ (see [12] for details). We shall also need two other technical assumptions, ULB being used to distinguish between the two possible behaviors of Hellinger distance and CS to ensure the possibility of estimating (see counterexamples in [7] and Sect. 4).

**CS.** *All densities $f$ in $\Theta$ have the same compact support.*

In this case we may always assume that $\mu$ is Lebesgue probability on the common support. This rescaling will not change the index of smoothness of the family but only the constants which are irrelevant here.

**ULB.** *There exists some constant $L > 0$ such that $f \geq L$ for $f$ in $\Theta$.*

The two following assume that $k$ is a decreasing function on $[0; 1]$.

**A1$(d, k)$.** *The set $\Theta$ is metricized by some distance $d$ and for each $\varepsilon$ in $]0; 1]$ there exists an $\varepsilon$-net $N(\varepsilon)$ for $\Theta$ the cardinal of which is smaller than $8^{A_0 k(\varepsilon)}$. Also $A_0 k(1) \geq 1/2$.*

**A2$(k)$.** *For any $\varepsilon$ in $]0; 1]$ there exists some function $g_\varepsilon$ with support on $I$, $\mu(I) = v$ and $\int_I g_\varepsilon d\mu = 0$ with*

$$\|g_\varepsilon\|_\infty \leq \varepsilon \quad and \quad \mu\{|g_\varepsilon| \geq A_1 \varepsilon\} \geq A_2 v. \tag{2.1}$$

*There exist $r$ translates $I + x_i$ of the support of $g_\varepsilon$ which do not intersect and are such that the family*

$$\left\{ f + \sum_{i=1}^r \delta_i g_\varepsilon(x - x_i) \right\} \tag{2.2}$$

*belongs to $\Theta$ with either*

$$\text{i) } \delta_i = \pm 1 \qquad \text{ii) } \delta_i = 0 \text{ or } 1,$$

*$f$ being constant and equal to $c$ on $\bigcup_{i=1}^r I + x_i$. Moreover the following inequalities hold*

$$v \leq A_3 k^{-1}(\varepsilon); \quad 1 \geq rv \geq A_4 > 0. \tag{2.3}$$

*and either*

$$\alpha) \ 2 \geq c \geq A_5 > 0 \quad or \quad \beta) \ A_1 \varepsilon \leq c \leq A_5 \varepsilon. \tag{2.4}$$

It is clear that under ULB $\alpha)$ only is possible, $\beta)$ is useful only in the case that $d = h$ and ULB does not hold.

*Remark.* 1) A1 is actually much stronger that what is needed in order to apply Theorem 3.1 but since it is easier to check and is satisfied in all subsequent applications, there is no need for the refined version given in [3].

2) In Sect. 3 we shall always assume that A1 holds with $d = d_\infty$. This implies that $\Theta$ is compact for $d_\infty$ and therefore uniformly bounded. In this case all $d_p$'s for $1 < p \leq 2$ are bounded by some multiple of $h$ as noticed in [3].

In the simplest cases the function $k$ will be a power function of $\varepsilon$ which leads us to the following

*Definition 2.1.* If a family $(\Theta, d)$ satisfies A1 and A2 with $k(\varepsilon) = \varepsilon^{-\delta}$, $\delta > 0$ it will be said to have an exponent of dimension $\delta$.

This exponent is easily computed for the classical smooth families. For example assume that for $f$ in $\Theta$

$$|f^{(p)}(x) - f^{(p)}(y)| \leq C |x - y|^\alpha \qquad p \geq 0, \ 0 < \alpha \leq 1, \ x, y \in [0; 1]. \tag{2.5}$$

Then the family will be said to have an exponent of smoothness $s = p + \alpha$ and it is known that for such a family $\delta = s^{-1}$. For densities on some compact convex subset of dimension $m$ in $\mathbb{R}^m$, if (2.5) holds ($f^{(p)}$ being now any partial derivative of order $p$) we shall find $\delta = \dfrac{m}{s}$. We may extend this to anisotropic smoothness on $\mathbb{R}^m$, which means that, as functions of coordinate $x_i$, $1 \leq i \leq m$, the densities have smoothness $s_i$, but $s_i$ depends on $i$. In this case we may define an exponent of global smoothness $s$ by $s^{-1} = m^{-1} \sum_{i=1}^{m} s_i^{-1}$ and the following will hold:

**Proposition 2.2.** *If $\Theta$ is a family of densities on $\mathbb{R}^m$ with common compact convex support of dimension $m$ and exponent of smoothness $s > 0$, then $A1(d_\infty, k)$ and $A2(k)$ hold with an exponent of dimension $\delta = \dfrac{m}{s}$.*

Much more general families can be used. In [2, 3] we consider families with a given concave modulus of continuity, i.e. the set of all functions satisfying

$$|f(x) - f(y)| \leq C \omega(|x - y|) \tag{2.6}$$

$\omega$ being a concave non-decreasing function and $C$ a constant such that $\omega(\varepsilon) \geq \varepsilon$. The smooth families with $s \leq 1$ are of this type. We may also consider generalizations of (2.5) of the form

$$|f^{(p)}(x) - f^{(p)}(y)| \leq C \omega(|x - y|) \tag{2.7}$$

and multidimensional analogues (isotropic or not). For the general case on $\mathbb{R}^m$ with different conditions of type (2.7) for each coordinate, denoting by $\eta = \psi_i(\varepsilon)$ the reciprocal function of $\eta \mapsto \varepsilon = \eta^{p_i} \omega_i(\eta)$ we would find that

$$\mathrm{Log}\, k(\varepsilon) = - \sum_{i=1}^{m} \mathrm{Log}\, \psi_i(\varepsilon)$$

which is a simple extension of the ordinary case when $\omega_i(\eta) = \eta^{\alpha_i}$. We shall not insist on this because those properties are only related to approximation theory and may be checked by the arguments used in [12]. Actually other types of families could be dealt using assumptions A1 and A2 and more general functions $k$ may occur (see [2] and [3] for illustrations of this fact). Properties of type A1 may be obtained directly for distances other than $d_\infty$, depending upon the conditions imposed on the family (examples may be found in [16]). Here we shall always assume $A1(d_\infty, k)$ and then derive $A1(d_q, k)$ from the fact that $d_q \leq d_\infty$ when the dominating measure is a probability, which is assumed here if CS holds. For $h$, the problem is different. It is known from [2, 3] that under some ULB condition we have $h \leq C d_\infty$; in this case there will not be any problem $\left(\text{just change } h \text{ into } C^{-1}h \text{ which is not}\right.$ important for the type of result we look for, or change $k(\varepsilon)$ into $\left. k\left(\frac{\varepsilon}{C}\right)\right)$. But without ULB condition we only know that $2h^2 \leq d_1 \leq d_\infty$ so that $A1(d_\infty, k)$ implies $A1(h, k_1)$ with $k_1(\varepsilon) = k(\varepsilon^2)$. As a consequence if a family has an exponent of dimension $\delta$, the rates of convergence will be functions of $\delta$ for $d_q$ or $h$ under ULB and functions of $2\delta$ for $h$ without ULB, assuming that case $\beta$) of A2 holds. But before we come to this, let us give the following useful but simple lemma:

**Lemma 2.3.** *Assume that $g_\varepsilon$ is a function with support on $I$, $\mu(I) = v$, that (2.1) holds and that $f$ is constant and equal to $c \geq \varepsilon$ on $I$. Then*

$$A_1^q A_2 \varepsilon^q v \leq d_q^q(f; f+g) = 2^{-q} d_q^q(f-g; f+g) \leq \varepsilon^q v, \tag{2.8}$$

$$\frac{A_1^2 A_2}{12} \frac{\varepsilon^2 v}{c} \leq h^2(f; f \pm g) \leq \frac{1}{2} \frac{\varepsilon^2 v}{c}, \tag{2.9}$$

$$\frac{A_1^2 A_2}{2} \frac{\varepsilon^2 v}{c} \leq h^2(f+g; f-g) \leq \frac{\varepsilon^2 v}{c}. \tag{2.10}$$

## III. Upper Bounds for the Risk

The idea of relating the speed of convergence with the dimensional properties of the parameter space and of the construction of "universal" dimensional estimates goes back to Le Cam [13]. We shall use here the following result which may be found in [2, 3].

**Theorem 3.1.** *Suppose we observe $n$ i.i.d. random variables from an unknown distribution $P_\theta$ in some parameter space $(\Theta, d)$ with metric $d$. Suppose also that the following properties are true for some constants $B, D$:*

  i) *$d(\theta; \theta') \leq B h(P_\theta; P_{\theta'})$     $P_\theta, P_{\theta'} \in \Theta$,*
  ii) *for some $\varepsilon$ with $n\varepsilon^2 \geq D \geq \frac{1}{2}$, there exists an $\varepsilon$-net $N$ in $(\Theta, d)$ such that*

$$\text{Card} \left[ N \cap \mathscr{B}(\theta, 2^j \varepsilon) \right] \leq 2^{jD} \quad \forall \theta \in \Theta, \ j \geq 3.$$

*Then there exists an estimate $\hat{\theta}_n$ of $\theta$ such that*

$$\sup_{\theta \in \Theta} \mathbb{E}_\theta[d^t(\hat{\theta}_n; \theta)] \leq C(B, t)\, \varepsilon^t. \tag{3.1}$$

Let us now apply this theorem to our particular case. If we restrict ourselves to $d_1$ or $h$, i) is certainly true because of (1.1). It is also true for $d_q$ if $1 < q \leq 2$ and $A1(d_\infty, k)$ holds because $\Theta$ is then uniformly bounded (see [3]). For $q > 2$ this becomes false and the treatment of such cases does not naturally follow from this theory. That is why we shall restrict ourselves to the case $1 \leq q \leq 2$ throughout this section.

Let us now assume that $A1(d_\infty, k)$ holds. Then we may distinguish between two cases:

$\alpha$) $d = d_q$, $1 \leq q \leq 2$ or $d = Ch$ and ULB holds.

$\beta$) $d = h$ and ULB does not hold.

In case $\alpha$) $A1(d, k)$ holds; in case $\beta$) $A1(d, k_1)$, $k_1(\varepsilon) = k(\varepsilon^2)$. But in order to apply Theorem 3.1 we only need to check that

$$n\varepsilon^2 \geq D \geq \tfrac{1}{2} \quad \text{with} \quad \alpha)\ D = A_0 k(\varepsilon) \quad \text{or} \quad \beta)\ D = A_0 k(\varepsilon^2),$$

that is

$$\alpha)\ n\varepsilon^2 \geq A_0 k(\varepsilon) \qquad \beta)\ n\varepsilon^2 \geq A_0 k(\varepsilon^2). \tag{3.2}$$

If $\varepsilon$ satisfies (3.2) we shall immediately get that

$$R_M(d^t, n) \leq C\varepsilon^t \tag{3.3}$$

which leads to the following corollary using Proposition 2.2.

**Corollary 3.2.** *Assume that $\Theta$ is a family of densities for which $A1(d_\infty, \varepsilon^{-\delta})$ holds with a distance $d$ satisfying either $\alpha$) or $\beta$), then there exist estimates $\hat{\theta}_n$ of $\theta$ depending on $n$ observations such that*

$$\sup_{\theta \in \Theta} \mathbb{E}_\theta[d^t(\hat{\theta}_n; \theta)] \leq \begin{cases} C_1\, n^{\frac{-t}{2+\delta}} & \text{case } \alpha \\ C_2\, n^{\frac{-t}{2(1+\delta)}} & \text{case } \beta \end{cases}.$$

*If $\Theta$ is a family of smooth densities on some compact convex subset of $\mathbb{R}^m$ of dimension $m$ with an exponent of global smoothness $s$ then*

$$R_M(d^t, n) \leq C_3\, n^{\frac{-ts}{2s+m}} \qquad \text{case } \alpha,$$

$$R_M(h^t, n) \leq C_4\, n^{\frac{-ts}{2(s+m)}} \qquad \text{case } \beta.$$

These are the classical results of [2, 3, 5, 8, 9, 18], but we may also deal with more special families using for example the moduli of continuity. Let us just consider two simple cases using Theorem 3.1.

*Example 1.* Let us look at densities on $[0; 1]$ satisfying (2.6) with $\omega(x) = -x^\alpha \log x$; $0 < \alpha \leq 1$; $x \leq \tfrac{1}{3}$. Then, following the same computations as in [2, 3], we easily find

*Case* α: $R_M(d^t, n) \leq C_5 n^{\frac{-t\alpha}{2\alpha+1}} \log n^{\frac{t}{2\alpha+1}}$,

*Case* β: $R_M(h^t, n) \leq C_6 n^{\frac{-t\alpha}{2\alpha+2}} \log n^{\frac{t}{2\alpha+2}}$.

*Example* 2. Same case as Example 1 but $\omega(x) = \dfrac{-1}{\log x}$ for $x \leq 0.1$. Then

*Case* α: $R_M(d^t, n) \leq C_7 (\log n)^{-t}$,

*Case* β: $R_M(h^t, n) \leq C_8 (\log n)^{-t/2}$,

which shows that very slow rates may occur (those bounds being optimal as we shall see).

Other cases including classes of analytic functions could also be treated but to get lower bounds for those classes one generally needs assumptions that are more complex than A2. This would make this presentation more obscure. For examples of such cases see [2] or [3] and [8].

## IV. Lower Bounds for the Risk

Let us begin with two general theorems which could be used in many circumstances to find lower bounds. They are Fano's Lemma and Assouad's Lemma. The first one has been applied very often in [11, 8, 9, 4] (see [10] or [3] for a proof) and is in a sense more general because it applies in more general situations. It could also replace Assouad's Lemma in almost any practical case. Assouad's Lemma is more specific but gives better constants and is also easier to apply. It is in the spirit of Proposition 2.1 of [5] or [18] and may be found in [1].

**Fano's Lemma** (see [3], Lemma 2.7). *Let S be a finite set of probabilities S* $= \{P_1; \dots; P_r\}$ *such that*

$$K(P_i, P_j) \leq K \quad \forall P_i, P_j \in S.$$

*Then for any estimate* $\phi$ *with values in* $\{1; 2; \dots; r\}$ *we have*

$$r^{-1} \sum_{i=1}^{r} P_i[\phi \neq i] \geq \gamma \tag{4.1}$$

*if*

$$K + \log 2 \leq (1 - \gamma) \log (r - 1). \tag{4.2}$$

In the case of $n$ i.i.d. observations, we have $P_i = Q_i^n$ and $K(P_i, P_j) = nK(Q_i, Q_j)$ so that (4.2) is equivalent to

$$n[\sup_{i, j} K(Q_i, Q_j)] + \log 2 \leq (1 - \gamma) \log (r - 1). \tag{4.3}$$

We shall now prove a suitable version of Assouad's Lemma following the original treatment of [1]. The hypercube $\mathscr{C}(r)$ of dimension $r$ will be the set $\{0; 1\}^r$ with points $x = \{x_1; \dots; x_r\}$, $x_i = 0$ or 1 and distance $\Delta$:

$$\Delta(x, x') = \sum_{i=1}^{r} |x_i - x_i'|.$$

**Assouad's Lemma.** *Assume we are given $n$ i.i.d. observations from some unknown distribution in some parameter space $\Theta$. Assume that some finite subset $\Theta_\mathscr{C}$ of $\Theta$ is in one-to-one correspondence with $\mathscr{C}(r)$:*

$$x \in \mathscr{C}(r) \leftrightarrow P_x \in \Theta_\mathscr{C}$$

*with the following properties*

$$h^2(P_x; P_{x'}) \leq \beta_i \leq 1 \quad \text{if} \quad x_j = x'_j, \, j \neq i. \tag{4.4}$$

*$\ell$ is a loss function defined on the space of probabilities and satisfying for all $P$:*

$$\ell(P; P_x) \geq \sum_{i=1}^{r} [x_i \ell_i(P) + (1 - x_i) \ell'_i(P)];$$
$$\ell_i(P), \ell'_i(P) \geq 0; \quad \ell_i(P) + \ell'_i(P) \geq \alpha > 0 \quad i = 1, \ldots, r. \tag{4.5}$$

*Let $\hat{P}_n$ be any estimate and let $R_B$ be the Bayes risk of $\hat{P}_n$ with respect to the uniform probability on $\Theta_\mathscr{C}$:*

$$R_B = 2^{-r} \sum_{x \in \mathscr{C}} \int \ell(\hat{P}_n, P_x) \, dP_x^n.$$

*Then if $\beta = r^{-1} \sum_{i=1}^{r} \beta_i$,*

$$R_B \geq \frac{r\alpha}{2} \max \left[1 - (2n\beta)^{1/2}; \tfrac{1}{2}(1 - \beta)^{2n}\right]. \tag{4.6}$$

*Proof.* Using (4.5) we may write

$$R_B \geq 2^{-r} \sum_{i=1}^{r} \left[ \sum_{x \mid x_i = 1} \int \ell_i(\hat{P}_n) \, dP_x^n + \sum_{x \mid x_i = 0} \int \ell'_i(\hat{P}_n) \, dP_x^n \right]$$

$$= \tfrac{1}{2} \sum_{i=1}^{r} \left[ \int \ell_i(\hat{P}_n) \, d\left[2^{-r+1} \sum_{x \mid x_i = 1} dP_x^n\right] + \int \ell'_i(\hat{P}_n) \, d\left[2^{-r+1} \sum_{x \mid x_i = 0} dP_x^n\right] \right]$$

$$\geq \frac{\alpha}{2} \sum_{i=1}^{r} \pi \left[2^{-r+1} \sum_{x \mid x_i = 1} P_x^n; \, 2^{-r+1} \sum_{x \mid x_i = 0} P_x^n\right].$$

$\operatorname{Inf}(x, y)$ being a concave function of $(x, y)$, Jensen's inequality entails that

$$R_B \geq \frac{\alpha}{2} \sum_{i=1}^{r} \left[2^{1-r} \sum_{J_i} \pi(P_x^n, P_{x'}^n)\right],$$
$$J_i = \{(x, x') \mid x_i = 0, \, x'_i = 1, \, x_j = x'_j \text{ for } j \neq i\}.$$

For the pairs $(P_x, P_{x'})$ such that $(x, x')$ belongs to $J_i$, (4.4) holds. Using (1.1) and (1.2) this leads to $\rho^2(P_x^n, P_{x'}^n) \geq (1 - \beta_i)^{2n}$ and

$$\pi(P_x^n, P_{x'}^n) \geq 1 - [1 - (1 - \beta_i)^{2n}]^{1/2} \geq \max \left[1 - (2n\beta_i)^{1/2}; \tfrac{1}{2}(1 - \beta_i)^{2n}\right]$$

which implies

$$R_B \geq \frac{r\alpha}{2} r^{-1} \sum_{i=1}^{r} \max\left[1-(2n\beta_i)^{1/2}; \tfrac{1}{2}(1-\beta_i)^{2n}\right];$$

the conclusion follows from a convexity argument. □

In the sequel we shall always use this lemma with all $\beta_i$'s equal to $\beta$ but this version may prove useful in other cases.

The use of (4.6) in conjunction with assumption A2 will then become fairly obvious:

**Corollary 4.1.** *Let $\Theta_\mathscr{C}$ be a subset of $\Theta$ of the following type*

$$\Theta_\mathscr{C} = \left\{ f + \sum_{i=1}^{r} \delta_i g_i \right\} \quad \text{with i) } \delta_i = \pm 1 \text{ or ii) } \delta_i = 0 \text{ or } 1.$$

*$f$ is a density which is constant on the sum of the supports of the $g_i$'s which supports are disjoint and the $g_i$'s are translates of the same function $g$ with $\int g(x)\,d\mu(x)=0$. Assume that $\ell(P,Q)=d_q^q(P,Q)$ or $h^2(P,Q)$ (in this case we take $q=2$ in (4.9)) and*

$$h^2(f+g_1; f-g_1) \leq \beta \text{ (case i)} \quad \text{or} \quad h^2(f; f+g_1) \leq \beta \text{ (case ii)} \tag{4.7}$$

$$\ell(f+g_1; f-g_1) \geq \alpha \text{ (case i)} \quad \text{or} \quad \ell(f; f+g_1) \geq \alpha \text{ (case ii)}. \tag{4.8}$$

*Then the (uniform) Bayes risk on $\Theta_\mathscr{C}$ for any estimate depending on $n$ i.i.d. variables satisfies*

$$R_B \geq 2^{-q} r\alpha \max\left[1-(2n\beta)^{1/2}; \tfrac{1}{2}(1-\beta)^{2n}\right]. \tag{4.9}$$

*Proof.* We shall just look at case i) for $\ell = d_q^q$ and apply Assouad's Lemma to $\Theta_\mathscr{C}$ which obviously satisfies (4.4) since all $g_i$'s are translates of the same function with disjoint supports. Because of the definition of $\ell$ we easily see that

$$\ell\left(P, f + \sum_{i=1}^{r} \delta_i g_i\right) \geq \sum_{i=1}^{r} \int_{I_i} \left|\frac{dP}{d\mu} - (f+\delta_i g_i)\right|^q d\mu$$

$I_i$ being the support of $g_i$, which implies (4.5) with

$$\ell_i = \int_{I_i} \left|\frac{dP}{d\mu} - f - g_i\right|^q d\mu \qquad \ell_i' = \int_{I_i} \left|\frac{dP}{d\mu} - f + g_i\right|^q d\mu$$

and because $(a+b)^q \leq 2^{q-1}(a^q+b^q)$

$$\ell_i + \ell_i' \geq 2^{1-q} \int_{I_i} \left(\left|\frac{dP}{dx} - f - g_i\right| + \left|\frac{dP}{dx} - f + g_i\right|\right)^q d\mu$$

$$\geq 2^{1-q} \int_{I_i} |2g_i|^q \, dx \geq 2^{1-q} \alpha,$$

hence the result. □

*Remark.* The assumption that $f$ is constant wasn't used in this case and could be removed but is useful for case ii) or when $\ell = h^2$; also in what follows, the constancy of $f$ could be removed when $\delta_i = \pm 1$ and $\ell = d_q^q$.

**Proposition 4.2.** *Assume that* A2(k) *holds and* $\varepsilon$ *satisfies* $n\varepsilon^2 \leqq A_6 k(\varepsilon^2)$ *if* $d=h$ *and* $\beta$) *holds,* $n\varepsilon^2 \leqq A_6 k(\varepsilon)$ *in all other cases. Then*

$$R_M(d^q, n) \geqq C\varepsilon^q. \tag{4.10}$$

*Proof.* Let us prove it when $d=h$ in case $\beta$, the other cases being analogous. Then A2(k) holds with (2.4) $\beta$) and (2.2) i) (say). The hypercube $\Theta_\mathscr{C}$ of Corollary 4.1 is then given by the family (2.2) with $g_i(x)=g_{\varepsilon^2}(x-x_i)$. Using (2.10) and (2.4) we may choose $\beta=\varepsilon^2 v A_1^{-1}$ in (4.7) and $\alpha=A_1^2 A_2 (2A_5)^{-1} \varepsilon^2 v$ in (4.8). Putting these values in (4.9) and using (2.3) gives the result. $\square$

We may now use our whole set of assumptions to put together the upper and lower bounds using Theorem 3.1 and Proposition 4.2 to get

**Proposition 4.3.** *Assume that* A1($d_\infty$, k) *and* A2(k) *hold and that* $\varepsilon_n$ *is a sequence such that*

$$n\varepsilon_n^2 \asymp k(\varepsilon_n) \text{ case } \alpha); \qquad n\varepsilon_n^2 \asymp k(\varepsilon_n^2) \text{ case } \beta)$$

*then*

$$R_M(d^q, n) \asymp \varepsilon_n^q.$$

*Proof.* We just consider case $\alpha$. Then by assumption $n\varepsilon_n^2 \leqq Ck(\varepsilon_n)$; choosing $\varepsilon = [(A_0/3C)^{1/2} \vee 1]\, \varepsilon_n$, (3.2) will hold and (3.3) will give an upper bound of order $\varepsilon^q$ which is also of order $\varepsilon_n^q$. The lower bound case is similar using the fact that (4.10) holds if $n\varepsilon^2 \leqq A_6 k(\varepsilon)$. $\square$

As a consequence we find the well-known fact that the rates of convergence for smooth densities with exponent $s$ in $\mathbb{R}^m$ given by Corollary 3.2 are the optimal ones. The same holds for the examples of Sect. 3.

We shall now prove an analogue of Theorem 1, ii) of [7] in the case of a given convergence to zero of the risk. Devroye's result could be proved in the same way but since we shall get an improved version in the next section we shall not consider his case but the following one, deliberately restricted for the sake of simplicity and brevity: $\Theta$ is assumed to be a family of densities with support on $[0; 1]$, $d=h$ and ULB holds. It could easily be generalized to $\mathbb{R}^m$, $d_q$ or case $\beta$. We shall need a slight modification of A2 allowing us to use simultaneously different $g_\varepsilon$'s corresponding to different values of $\varepsilon$. It is practically no more restrictive than A2.

**A'2(k).** *Suppose that the* $\varepsilon_i$'s, $i \geqq 1$ *are small enough* ($\varepsilon_i \leqq \eta_i$, *say*), *that they form a decreasing sequence and that the* $g_i$'s *$(g_i=g_{\varepsilon_i})$ are chosen according to A2 with (2.1). Then for some numbers* $r_i$ *the family*

$$\Theta_\infty = \left\{ 1 + \sum_{i \geqq 1} \sum_{j=1}^{r_i} \delta_{ij} g_i(x-x_{ij}) \right\}, \qquad \delta_{ij}=0 \text{ or } 1,$$

*is a subset of* $\Theta$, *the functions* $g_i(x-x_{ij})$ *have disjoint supports and*

$$A_4 2^{-i} \leqq r_i v_i \leqq 2^{-i}; \qquad v_i \leqq A_3 k^{-1}(\varepsilon_i). \tag{4.11}$$

**Proposition 4.4.** *Let us consider two decreasing sequences* $\{a_n\}$ *and* $\{c_n\}$, *converging to zero with*

$$C_1 \, n a_n^2 \leqq k(a_n) \leqq C_2 \, n a_n^2. \tag{4.12}$$

*If A'2(k) holds and $\{\hat{f}_n\}$ is a sequence of estimates, there exists at least one point θ in $\Theta_\infty$ such that*

$$\limsup_n c_n^{-1} a_n^{-2} \, \mathbb{E}_\theta[h^2(f_\theta; \hat{f}_n)] \geqq 1. \tag{4.13}$$

Notice that under assumptions A1 and A2 (4.12) implies that $a_n^2$ is the speed of convergence of $R_M(h^2, n)$. Statement (4.13) shows that it cannot be improved for a given sequence of estimates. For the smooth densities it implies the following

**Corollary 4.5.** $\Theta$ *being a family of smooth densities on* $[0; 1]$ *with exponent s and* $\{c_n\}$ *a sequence decreasing to zero, for any sequence* $\{\hat{f}_n\}$ *of estimates there exists a fixed f in* $\Theta$ *such that*

$$\limsup_n c_n^{-1} \, n^{\frac{2s}{2s+1}} \, \mathbb{E}_f[h^2(\hat{f}_n; f)] \geqq 1.$$

*Proof of Proposition 4.4.* Let us denote by $\delta_i$ the vector $\{\delta_{ij}\}_{j=1,\dots,r_i}$ and call the $\delta_i$'s the coordinates of the points of $\Theta_\infty$. $\Theta_k$ is the subset of points with $\delta_i = 0$ for $i > k$. Choose a sequence $\{m_i\}$ such that $c_{m_i} \leqq C_3 \, C_4^2 \, 2^{-i-2}$ and a sequence $n_i$ satisfying

$$n_i \geqq m_i; \qquad C_4 \, a_{n_i} \leqq \eta_i, \tag{4.14}$$

$$a_{n_{i+1}} \leqq \frac{C_4 \, C_3}{5} \, n_i^{-1/2} \, 2^{-i/2} \, a_{n_i}^2, \tag{4.15}$$

with

$$C_3 = \frac{A_1^2 \, A_2 \, A_4}{100}; \qquad C_4 = \left(\frac{C_1}{4A_3}\right)^{1/2} \wedge 1.$$

This is always possible inductively since $a_n \to 0$ when $n \to +\infty$. Fix $\varepsilon_i = C_4 \, a_{n_i}$ and notice that if $f_0$ and $f_1$ in $\Theta_\infty$ just differ by one coordinate $\delta_i$ and only one $\delta_{ij}$ then by (2.9)

$$\frac{A_1^2 \, A_2 \, v_i \, \varepsilon_i^2}{12} \leqq h^2(f_0; f_1) \leqq \frac{v_i \, \varepsilon_i^2}{2}, \tag{4.16}$$

and if $f_0$ is in $\Theta_k$ and $f_1$ has the same $k$ first coordinates

$$h^2(f_0; f_1) \leqq \sum_{i>k} r_i \, \frac{v_i \, \varepsilon_i^2}{2} \leqq \sum_{i>k} 2^{-i-1} \, \varepsilon_i^2 \leqq 2^{-k-1} \, \varepsilon_{k+1}^2.$$

This implies using (1.1), (1.2), (4.15) that

$$\rho^{n_k}(f_0; f_1) \geqq \exp\left[n_k \log\left(1 - \frac{\varepsilon_{k+1}^2}{2^{k+1}}\right)\right] \geqq 1 - n_k \, \frac{\varepsilon_{k+1}^2}{2^{k+1}},$$

$$D(P_0^{n_k}; P_1^{n_k}) \leqq n_k^{1/2} \, \varepsilon_{k+1} \, 2^{-k/2} \leqq \frac{C_3}{5} \, \varepsilon_k^2 \, 2^{-k}, \tag{4.17}$$

and also for any probability $P$

$$h^2(P_1; P) \geq \tfrac{1}{2} h^2(P_0; P) - h^2(P_0; P_1) \geq \tfrac{1}{2} h^2(P_0; P) - \frac{C_3}{20} 2^{-k} \varepsilon_k^2. \qquad (4.18)$$

Now consider the subset of $\Theta_k$ formed by the points having their $k-1$ first coordinates fixed, $\delta_k$ being free and some estimate $\hat{P}$ depending on $n_k$ observations. We are then in a situation to apply Corollary 4.1 which gives using (4.16) and (4.11)

$$R_B \geq \tfrac{1}{4} r_k \frac{A_1^2 A_2}{12} v_k \varepsilon_k^2 [1 - (n_k v_k \varepsilon_k^2)^{1/2}] \geq 2^{-k} 2 C_3 \varepsilon_k^2 [1 - (n_k A_3 \varepsilon_k^2 k^{-1}(\varepsilon_k))^{1/2}].$$

But since $\varepsilon_k \leq a_{n_k}$ we have $n_k \varepsilon_k^2 k^{-1}(\varepsilon_k) \leq n_k C_4^2 k^{-1}(a_{n_k}) a_{n_k}^2 \leq \dfrac{1}{4 A_3}$ and then there exists one point $P_0$ in $\Theta_k$ with those first $k-1$ coordinates and

$$\mathbb{E}_{P_0^{n_k}}[h^2(\hat{P}; P_0)] \geq 2^{-k} C_3 \varepsilon_k^2.$$

Because of (4.18), (4.17) for any $P_1$ having the same $k$ first coordinates as $P_0$

$$\mathbb{E}_{P_1^{n_k}}[h^2(\hat{P}; P_1)] \geq \left[ \frac{1}{2} 2^{-k} C_3 \varepsilon_k^2 - \frac{C_3}{20} 2^{-k} \varepsilon_k^2 \right] - \frac{C_3}{5} \varepsilon_k^2 2^{-k} \geq \frac{C_3}{4} 2^{-k} \varepsilon_k^2.$$

Now by induction on $k$, we shall find that for any $k$ there exists $\delta_1, \dots, \delta_k$ such that any point $f_\theta$ in $\Theta_\infty$ with those first coordinates will satisfy for all $i \leq k$,

$$\mathbb{E}_\theta[h^2(f_\theta; \hat{f}_{n_i})] \geq \frac{C_3}{4} 2^{-i} \varepsilon_i^2 \geq \frac{C_3 C_4^2}{4} 2^{-i} a_{n_i}^2 \geq a_{n_i}^2 c_{n_i}. \qquad (4.19)$$

We deduce (4.13) by a simple limiting argument since $\Theta_\infty$ is compact and the $\theta$'s satisfying (4.19) up to order $k$ form a decreasing sequence of closed subsets of $\Theta_\infty$. $\square$

## V. Bad Problems

We shall consider here what may happen when the families are not nice. Let us begin with the case considered by Devroye [7] of the densities bounded by 2 on $[0; 1]^m$. Actually we may take $m = 1$, the extension to $\mathbb{R}^m$ being obvious.

Let us start with some $\mathscr{C}^\infty$ function $\tilde{g}$ in $[0; 1]$ with the following properties

$$\tilde{g}(x + \tfrac{1}{2}) = -\tilde{g}(\tfrac{1}{2} - x) \qquad 0 \leq x \leq \tfrac{1}{2} \qquad \text{so that } \int_0^1 \tilde{g}(x)\, dx = 0$$

$$|\tilde{g}| \leq 1 \quad \text{and} \quad \tilde{g} = 1 \text{ on } \left[ \frac{\varepsilon}{4}; \frac{1}{2} - \frac{\varepsilon}{4} \right]; \quad \tilde{g}^{(p)}(0) = 0 \ \forall p.$$

$$(5.1)$$

Now fix an integer $r$ and put $g(x) = \tilde{g}(xr)$; then it is immediate to check that

$$h^2(1 + g; 1 - g) = \alpha, \qquad \frac{1 - \varepsilon}{r} \leq \alpha \leq \frac{1}{r}; \qquad d_q^q(1 + g; 1 - g) \geq 2^q \frac{1 - \varepsilon}{r}.$$

Let us consider the family $\Theta_{\mathscr{C}} = \left\{ 1 + \sum_{i=1}^{r} \delta_i g\left( x - \frac{i-1}{r} \right) \right\}$, $\delta_i = \pm 1$. Applying Assouad's Lemma we get from Corollary 4.1,

$$R_M(h^2, n) \geqq \frac{1-\varepsilon}{4} (1 - (2nr^{-1})^{1/2}); \qquad R_M(d_q^q, n) \geqq (1-\varepsilon)(1 - (2nr^{-1})^{1/2}).$$

Letting $r \to +\infty$ and $\varepsilon \to 0$ we find that for the set $\Theta$ of $\mathscr{C}^\infty$ functions on $[0; 1]$ bounded by two

$$R_M(h^2, n) \geqq \tfrac{1}{4} \quad \text{and} \quad R_M(d_q^q, n) \geqq 1.$$

To get Theorem 1, i) of [7] we just need a proper normalization by the $\mathbb{L}^q$ norm of the true density. But in $\Theta_{\mathscr{C}}$, all elements have the same $\mathbb{L}^q$ norm, which is approximately $2^{q-1}$; dividing by this quantity we find for the normalized risk $R'_M(d_q^q) \geqq 2^{-q+1}$.

Let us now come back to any smooth family but just assume that its elements may have arbitrary compact support (i.e. CS is no more true) and consider the following transformation $\Lambda$ from $\mathscr{C}^\infty[0; 1]$ to $\mathscr{C}^\infty(\mathbb{R})$ given by $\Lambda f(x) = \lambda f(\lambda x)$, $\lambda > 0$. It transforms a density into a density and has the following property that

$$h^2(\Lambda f, \Lambda g) = h^2(f; g); \quad d_q^q(\Lambda f; \Lambda g) = \lambda^{q-1} d_q^q(f; g); \quad \frac{d_q^q(\Lambda f; \Lambda g)}{\|\Lambda f\|_q^q} = \frac{d_q^q(f; g)}{\|f\|_q^q}.$$

Because of this property the risks $R_M(h^2, n)$ or $R'_M(d_q^q, h)$ will remain constant if we transform $\Theta_{\mathscr{C}}$ using $\Lambda$. But for given $r$, even very large, there exists some $\lambda \leqq 1$ such that $\Lambda(\Theta_{\mathscr{C}})$ will become arbitrarily smooth with compact support because all elements of $\Theta_{\mathscr{C}}$ have their $p^{\text{th}}$ derivative uniformly bounded. That means that any family of arbitrary smooth densities with compact (not the same) support contains some $\Lambda\Theta_{\mathscr{C}}$ and that for those families

$$R_M(h^2, n) \geqq \tfrac{1}{4}; \qquad R'_M(d_q^q, n) \geqq 2^{-q+1}.$$

*Remark.* This is the analogue of Devroye's result for the family $\mathscr{L}(g)$. We could get exactly the same result as his using a modified version of the above proof and applying case ii) of Corollary 4.1.

We shall now improve the second part of Devroye's theorem. We shall need some preliminary considerations and first recall this known lemma:

**Lemma 5.1.** *Assume that $\Theta_{\mathscr{C}}$ is in one-to-one correspondence with some hypercube of dimension $r$ and that $d$ is a distance such that*

$$d^q(P_x; P_{x'}) \geqq \alpha\Delta(x, x').$$

*For any $\lambda$ with $0 < \lambda < 1$ there exists a subset $S$ of $\Theta_{\mathscr{C}}$ such that*

$$\text{Card } S \geqq \exp\left[ \frac{r\lambda^2}{2} \right]; \qquad d^q(P; Q) > \alpha r \frac{1-\lambda}{2} \quad \forall P, Q \in S, \ P \neq Q.$$

*Remark.* In conjunction with Fano's Lemma this leads to an equivalent of Assouad's Lemma which was used in [2, 3] or [8] to get lower bounds and could have been used here, but the constants would not have been the same and Assouad's Lemma is simpler.

First we must describe the families of densities in $\mathscr{C}^\infty[0; 1]$ that we shall use. Let us come back to the function $\tilde{g}$ of (5.1) and consider two sequences $\{a_i\}, \{r_i\}$ for $i \geq 1$ with

$$a_i > 0; \quad \sum_{i=1}^{+\infty} a_i = 1; \quad r_i \to +\infty \quad \text{when} \quad i \to +\infty, \ r_i \in \mathbb{N}.$$

Define

$$A_0 = 0; \quad A_i = \sum_{j=1}^{i} a_j \quad \text{if} \ i \geq 1; \quad \lambda_i = r_i^{-1} a_i; \quad I_i = [A_{i-1}; A_i];$$

$$f_{ij}(x) = (1-\varepsilon^2)^{1/2} \left[ \tilde{g} \left( \frac{x - A_{i-1} - (j-1)\lambda_i}{\lambda_i} \right) \right] \quad i \geq 1; \ j = 1, 2, \ldots, r_i$$

and consider

$$\Theta'_1 = \left\{ 1 + \sum_{i=1}^{+\infty} \left[ \sum_{j=1}^{r_i} \delta_{ij} f_{ij}(x) \right] \right\}, \quad \delta_{ij} = \pm 1.$$

$\Theta'_1$ is a family of $\mathscr{C}^\infty$ densities on $[0; 1]$ bounded by 2. Easy computations with $K = 2 \log(2\varepsilon^{-1})$ give

$$K(1 + \delta_{ij} f_{ij}; 1 - \delta_{ij} f_{ij}) \leq K \lambda_i; \quad h^2(1 + f_{ij}; 1 - f_{ij}) \leq \lambda_i, \tag{5.2}$$

$$h^2(1 + f_{ij}; 1 - f_{ij}) \geq \lambda_i (1-\varepsilon)^2; \quad d_q^q(1 + f_{ij}; 1 - f_{ij}) \geq 2^q \lambda_i (1-\varepsilon)(1-\varepsilon^2)^{q/2}. \tag{5.3}$$

From now on we shall restrict our discussion to the case of $h$, the other one being similar up to multiplication by $2^q$ which cancels in (5.11). $\Theta'_1$ may be considered as a product of hypercubes of dimension $r_i$: $\Theta'_1 = \prod_{i=1}^{\infty} \mathscr{C}'_i$. Applying Lemma 5.1 to $\mathscr{C}'_i$ with $\lambda = 10^{-1/2}$ and $d$ replaced by $h$ restricted to the set $I_i$, we shall find a subset $\mathscr{C}_i$ of $\mathscr{C}'_i$, with cardinal larger than $\exp\left[\frac{r_i}{20}\right]$, having the following property: for any two points $x, x'$ in $\mathscr{C}_i$, (using (5.3)) and suitable choice of $\varepsilon$,

$$h^2(P_x 1_{I_i}; P_{x'} 1_{I_i}) \geq \lambda_i r_i (1-\varepsilon)^2 \frac{1 - 10^{-1/2}}{2} \geq \frac{a_i}{3}. \tag{5.4}$$

Since the set $\mathscr{C}_i$ will only depend on the part of the measures supported by $I_i$, i.e. on coordinate $i$, the subset $\Theta_1 = \prod_{i=1}^{+\infty} \mathscr{C}_i$ will have the following property that if $\theta$ and $\theta'$ are two points in $\Theta_1$, $\theta = \{\theta_i\}_{i \geq 1}$, $\theta_i \in \mathscr{C}_i$ and $\Delta(\theta_i, \theta'_i) = 0$ if $\theta_i = \theta'_i$, 1 in the other case:

$$h^2(\theta, \theta') = \sum_{i=1}^{+\infty} \Delta(\theta_i; \theta'_i) h_i^2; \quad \frac{a_i}{3} \leq h_i^2 \leq a_i, \tag{5.5}$$

$$K(\theta, \theta') \leq K \sum_{i=1}^{+\infty} \Delta(\theta_i; \theta_i') \, a_i \qquad (5.6)$$

because of (5.2) and (5.4). Moreover we may assume that

$$\exp\left(\frac{r_i}{20}\right) \leq \operatorname{Card} \mathscr{C}_i = U_i \qquad (5.7)$$

and notice that $\Theta_1$ is totally bounded and closed, hence compact by (5.5).

Given a subset $T_j$ of $\Theta_1$ of the following form

$$T_j = \{\theta_1\} \times \ldots \times \{\theta_{j-1}\} \times S_j \times \{\theta_{j+1}\} \times \ldots, \qquad S_j \subset \mathscr{C}_j \qquad (5.8)$$

with $\operatorname{Card} T_j = \operatorname{Card} S_j = V_j \leq U_j$ and an estimate $\hat{\theta}_j$ of $\theta_j$, with values in $S_j$, depending on $n$ observations, we may apply Fano's Lemma and using (5.6) find that if

$$n K a_j + \log 2 \leq \tfrac{1}{13} \log(V_j - 1) \qquad (5.9)$$

then $V_j^{-1} \sum_{\theta_j \in S_j} P_{\theta_j}[\hat{\theta}_j \neq \theta_j] \geq \tfrac{12}{13}$, which implies because of (5.5)

$$V_j^{-1} \sum_{\theta_j \in S_j} \mathbb{E}_{\theta_j}[h^2(P_{\theta_j}; P_{\hat{\theta}_j})] \geq \frac{4 a_j}{13}.$$

By (5.5) $h^2(P_{\theta_j}; P_{\hat{\theta}_j})$ is smaller that $a_j$; then there exists a subset $S_j'$ of $S_j$ with $\operatorname{Card} S_j' \geq \tfrac{2}{11} V_j$ and

$$\inf_{\theta_j \in S_j'} \mathbb{E}_{\theta_j}[h^2(P_{\theta_j}; P_{\hat{\theta}_j})] \geq \frac{2 a_j}{13}. \qquad (5.10)$$

Given an arbitrary probability $P$ let us call $\tilde{\theta}_j$ any point in $S_j$ which minimizes $h(P 1_{I_j}; P_{\theta_j} 1_{I_j})$. It is easy to see that we always have (with $\theta$, $\tilde{\theta}$ in $T_j$)

$$h^2(P; P_\theta) \geq \tfrac{1}{4} h^2(P_{\tilde{\theta}}; P_\theta).$$

This remark and (5.10) give the following lemma:

**Lemma 5.2.** *Given a subset $T_j$ of $\Theta_1$ defined by (5.8) with $\operatorname{Card} S_j = V_j$ satisfying (5.9) and an estimate $\hat{P}$ depending on $n$ observations, there exists a subset $T_j'$ of $T_j$ with*

$$\inf_{\theta \in T_j'} \mathbb{E}_\theta[h^2(P_\theta; \hat{P})] \geq \frac{a_j}{26}; \qquad \operatorname{Card} T_j' \geq \frac{2}{11} V_j. \qquad (5.11)$$

$$\left(\text{or } \inf_{\theta \in T_j'} \mathbb{E}_\theta[d_q^q(P_\theta; \hat{P})] \geq \frac{2 a_j}{13}\right).$$

We are now in a position to prove our theorem. Let any sequence $\{c_i\}_{i \geq 1}$ converging to zero be given such that $\tfrac{1}{78} < \sup_i c_i < \tfrac{1}{26}$. Fix $a_1 = 26 \sup_i c_i$ and $a_j = 2^{1-j}(1 - a_1)$ for $j \geq 2$. Then $\{a_j\}_{j \geq 1}$ is a strictly decreasing sequence. Define $m_0 = 0$ and

$$m_j = \left[ m_{j-1} \vee \sup \left\{ i \mid c_i > \frac{a_j}{26} \right\} \right] + 1, \qquad j \geq 1.$$

Then $m_1 = 1$ and $\{m_j\}_{j \geq 1}$ is a strictly increasing sequence such that $i \geq m_j$ implies $c_i \leq \dfrac{a_j}{26}$. For convenience let us set the following notations:

$$M_j = m_{j+1} - m_j; \quad m'_j = m_{j+1} - 1; \quad B_j = (2/11)^{M_j};$$
$$C_j = (2/11)^{-m'_j + 1}; \quad D_j = 1 + 2^{13} \exp[13 K m'_j a_j].$$

Notice that $C_1 = B_1^{-1}$ and $C_{j+1} = C_j B_{j+1}^{-1}$. We may now choose the $r_j$'s inductively in order that the $U_i$ satisfy (5.7) together with

$$U_1 \geq C_1 D_1; \quad U_j \geq C_j D_j [\prod_{i \leq j-1} U_i^{j-i-1}], \quad j \geq 2; \tag{5.12}$$

(in the case of $d_q$ just replace in those formulas 26 by 6.5). We shall first prove the following

**Lemma 5.3.** *For any $T_j$ of the form* (5.8) *such that*

$$m'_j K a_j + \log 2 \leq \tfrac{1}{13} \log[B_j V_j - 1] \tag{5.13}$$

*and any sequence of $M_j$ estimates $\{\hat{P}_i \mid i = m_j; \; m_j + 1; \ldots; m'_j\}$ such that $\hat{P}_i$ is a function of $i$ variables there exists a subset $T''_j$ of $T_j$ with*

$$\inf_{\theta \in T''_j} \; \inf_{i = m_j, \ldots, m'_j} c_i^{-1} \, \mathbb{E}_\theta[h^2(\hat{P}_i; P_\theta)] \geq 1; \quad \operatorname{Card} T''_j \geq B_j V_j.$$

*Proof.* Assuming that (5.9) is true with $n = m_j$ we may use Lemma 5.2 to get some $T'_j$, $\operatorname{Card} T'_j \geq \tfrac{2}{11} V_j$ and satisfying (5.11) which gives us since $\dfrac{a_j}{26} \geq c_i$ for $i \geq m_j$

$$\inf_{\theta \in T'_j} c_{m_j}^{-1} \, \mathbb{E}_\theta[h^2(\hat{P}_{m_j}; P_\theta)] \geq 1.$$

But we may now apply the same lemma to $T'_j$ if (5.9) is still valid with $n = m_j + 1$ and $V_j$ is replaced by $\tfrac{2}{11} V_j$, and so on. The cardinal of the last set $T''_j$ will be at least $B_j V_j$. Since during this process the right member of (5.9) is decreasing and the left one is increasing, all conditions will be true if the last one is true which amounts to (5.13).   $\square$

Suppose we are given a sequence $\{\hat{P}_n\}_{n \geq 1}$ of estimates, $\hat{P}_n$ depending on $n$ variables. In order to apply Lemma 5.3, we may notice that (5.13) is equivalent to $B_j V_j \geq D_j$ which will always be satisfied because of (5.12) if we have

$$V_1 = U_1; \quad V_j \geq C_{j-1}^{-1} [\prod_{i \leq j-1} U_i^{j-i-1}]^{-1} U_j, \quad j \geq 2. \tag{5.14}$$

We start the construction with $T_1$ given by (5.8), and $S_1 = \mathscr{C}_1$ so that $V_1 = U_1$ and we may apply Lemma 5.3 getting some $T''_1$ of the form

$$T''_1 = S''_1 \times \{\theta_2\} \times \ldots; \quad \operatorname{Card} S''_1 \geq B_1 V_1,$$

where $S''_1$ obviously depends on $\theta_2, \theta_3, \ldots$. Keeping $\theta_3, \theta_4, \ldots$ fixed we shall now vary $\theta_2$ and write $S''_1(\theta_2)$. When $\theta_2$ varies in $\mathscr{C}_2$ we shall consider

$\sum\limits_{\theta_2 \in \mathscr{C}_2} \mathrm{Card}\, S_1''(\theta_2)$. Since $S_1''(\theta_2)$ is a subset of $\mathscr{C}_1$ we may see that if each $\theta_1$ in $\mathscr{C}_1$ belongs to no more than $V_2$ distinct $S_1''(\theta_2)$ then

$$V_2\, U_1 \geqq \sum_{\theta_2 \in \mathscr{C}_2} \mathrm{Card}\, S_1''(\theta_2) \geqq U_2\, B_1\, U_1 \tag{5.15}$$

which implies that $V_2 \geqq B_1\, U_2$ and proves the existence of some $S_2$ with $\mathrm{Card}\, S_2 = V_2 \geqq B_1\, U_2 = U_2\, C_1^{-1}$ and some $\theta_1 \in \bigcap\limits_{\theta_2 \in S_2} S_1''(\theta_2)$ such that with $T_2 = \{\theta_1\} \times S_2 \times \{\theta_3\} \times \dots$,

$$\inf_{\theta \in T_2}\ \inf_{i=1,\dots,m_1'}\ c_i^{-1}\, \mathbb{E}_\theta[h^2(\hat{P}_i; P_\theta)] \geqq 1.$$

Let us now proceed by induction. Assume that for any sequence $\theta_{j+1}, \theta_{j+2},\dots$ there exists $\theta_1, \theta_2, \dots, \theta_{j-1}$ and some subset $S_j$ of $\mathscr{C}_j$ such that (with $T_j$ as in (5.8))

$$\inf_{\theta \in T_j}\ \inf_{i=1,\dots,m_{j-1}'}\ c_i^{-1}\, \mathbb{E}_\theta[h^2(\hat{P}_i; P_\theta)] \geqq 1,$$
$$\mathrm{Card}\, S_j = V_j \geqq U_j\, C_{j-1}^{-1}\,\Big[\prod_{i \leqq j-1} U_i^{j-i-1}\Big]^{-1}.$$

Letting all $\theta_i$ fixed except for $\theta_j$, we may apply Lemma 5.3 to $T_j$ since (5.14) is satisfied and for some $T_j''$ and $S_j''$ with $\mathrm{Card}\, S_j'' = V_j' \geqq B_j\, V_j$

$$\inf_{\theta \in T_j''}\ \inf_{i=1,\dots,m_j'}\ c_i^{-1}\, \mathbb{E}_\theta[h^2(\hat{P}_i; P_\theta)] \geqq 1.$$

To make the induction work we keep $\theta_{j+2}, \dots$ fixed and let $\theta_{j+1}$ change; then $\theta_1, \dots, \theta_{j-1}$ and $S_j''$ are functions of $\theta_{j+1}$. But the possible choices for $\theta_1, \dots, \theta_{j-1}$ are at most $\prod\limits_{i=1}^{j-1} U_i$. Since because of (5.12) $U_{j+1}\big[\prod\limits_{i \leqq j-1} U_i\big]^{-1} \geqq 1$, we may restrict $\theta_{j+1}$ to vary within a subset $S_{j+1}'$ of $\mathscr{C}_{j+1}$ of cardinal $U_{j+1}' \geqq U_{j+1}\big[\prod\limits_{i \leqq j-1} U_i\big]^{-1}$ and keep $\theta_1, \dots, \theta_{j-1}$ fixed for all the $\theta_{j+1}$ in $S_{j+1}'$. Now $S_j''$ is still depending on $\theta_{j+1}$ but we may also fix it using the same argument that we used for (5.15) and get some fixed $\theta_j$ for a subset $S_{j+1}$ of $S_{j+1}'$ of cardinal $V_{j+1}$ and

$$V_{j+1}\, U_j \geqq \sum_{\theta_{j+1} \in S_{j+1}'} \mathrm{Card}\, S_j''(\theta_{j+1}) \geqq U_{j+1}'\, V_j'$$

which implies

$$V_{j+1} \geqq \frac{U_{j+1}'\, V_j'}{U_j} \geqq B_j\, \frac{V_j}{U_j}\,\Big[\prod_{i \leqq j-1} U_i\Big]^{-1}\, U_{j+1}. \tag{5.16}$$

Using these $\theta_1, \dots, \theta_{j-1}, \theta_j, S_{j+1}$, to define $T_{j+1}$, we get

$$\inf_{\theta \in T_{j+1}}\ \inf_{i=1,\dots,m_j'}\ c_i^{-1}\, \mathbb{E}_\theta[h^2(\hat{P}_i; P_\theta)] \geqq 1$$

and from (5.12) and (5.16)

$$V_{j+1} \geqq \frac{U_{j+1}}{C_j}\,\Big[\prod_{i \leqq j} U_i^{j-i}\Big]^{-1},$$

which is our induction assumption. We have then proved the following theorem in case i) for $n < +\infty$. Case ii) is absolutely analogous because all points in $\Theta_1$ have the same $\mathbb{L}^q$ norm and this norm is approximately $2^{q-1}$ when $\varepsilon$ is small. The case $n = +\infty$ follows from the fact that the subset of the $\theta$'s satisfying (5.17) for some $n$ is closed then compact, non-void and decreasing with $n$.

**Theorem 5.4.** *For any given sequence* $\{c_n\}_{n \geq 1}$ *converging to zero such that*

i) $\frac{1}{78} < \sup\limits_i c_i < \frac{1}{26}$ *if* $d = h$,        ii) $\frac{2}{39} < \sup\limits_i c_i < \frac{2}{13}$ *if* $d = d_q$,

*there exists a compact subset* $\Theta$ *of densities in* $\mathscr{C}^\infty[0; 1]$ *bounded by two and for any sequence* $\{\hat{P}_n\}_{n \geq 1}$ *of estimates depending on n observations there exist infinite subsets* $\Theta'_n$ *with*

$$\inf_{\theta \in \Theta'_n} \inf_{k < n} c_k^{-1} \, \mathbb{E}_\theta[h^2(\hat{P}_k; P_\theta)] \geq 1,$$

*or*

$$\inf_{\theta \in \Theta'_n} \inf_{k < n} c_k^{-1} \, \mathbb{E}_\theta[d_q^q(\hat{P}_k; P_\theta)] \geq 1, \tag{5.17}$$

*or*

$$\inf_{\theta \in \Theta'_n} \inf_{k < n} c_k^{-1} \, \mathbb{E}_\theta\left[\frac{d_q^q(\hat{P}_k; P_\theta)}{\|P_\theta\|_q^q}\right] \geq 2^{1-q}$$

*and there is at least one point in* $\Theta$ *satisfying* (5.17) *with* $n = +\infty$.

*Remarks.* 1) The constants which we used here ($\frac{1}{26}$, etc. ...) are not optimal at all and were chosen for convenience. The same proof with longer computations could lead to $\frac{1}{8}$ which is still not optimal.

2) The same method also gives a similar result for the second family $\mathscr{L}(g)$ considered by Devroye for $d_1$ or $h$ (not for $d_q$ if $q > 1$). We just need minor modifications and the transformation $\Lambda$ mentioned at the beginning of the section. We should use a family $\{\Lambda_i\}_{i \geq 1}$ of such applications with different values of the $\lambda_i$'s: $\lambda_i = r_i^{-1} a_i$ and apply $\Lambda_i$ to the parts of the measures which have a support on $I_i$. This will transform $I_i$ into an interval of length $r_i$ and each $f_{ij}$ will have a support of length 1. The details are easy. It does not work for $q > 1$ just as in Devroye's paper because there is no longer a suitable normalization by the $\mathbb{L}^q$ norm.

## References

1. Assouad, P.: Deux remarques sur l'estimation. C.R. Acad. Sci. Paris **296**, Sér. I, 1021–1024 (1983)
2. Birgé, L.: Thèse, 3ᵉ-partie. Université Paris VII (1980)
3. Birgé, L.: Approximation dans les espaces métriques et théorie de l'estimation. Z. Wahrscheinlichkeitstheor. Verw. Geb. **65**, 181–237 (1983)
4. Birgé, L.: Non-asymptotic minimax risk for Hellinger balls. Probability and Math. Statistics. To appear

5. Bretagnolle, J., Huber, C.: Estimation des densités: risque minimax. Z. Wahrscheinlich-keitstheor. Verw. Geb. **47**, 119–137 (1979)
6. Dacunha-Castelle, D.: École d'Eté de Probabilités de Saint-Flour VII. Lecture Notes in Mathematics no. **678**. Berlin-Heidelberg-New York: Springer 1977
7. Devroye, L.: On arbitrarily slow rates of convergence in density estimation. Z. Wahrscheinlich-keitstheor. Verw. Gebiete **62**, 475–483 (1983)
8. Ibragimov, I.A., Khas'minskii, R.Z.: Estimation of distribution density. Zap. Nauchn. Semin. LOMI **98**, 61–85 (1980) in russian, and J. Sov. Math. **21**, 40–57 (1983)
9. Ibragimov, I.A., Khas'minskii, R.Z.: On the non-parametric density function. Zap. Nauchn. Semin. LOMI **108**, 73–89, in russian (1981)
10. Ibragimov, I.A., Khas'minskii, R.Z.: Statistical Estimation, Asymptotic Theory. Berlin-Heidel-berg-New York: Springer 1981
11. Khas'minskii, R.Z.: A lower bound on the risks of non-parametric estimates of densities in the uniform metric. Theory Probab. Appl. **23**, 794–796 (1978)
12. Kolmogorov, A.N., Tikhomirov, V.M.: $\varepsilon$-entropy and $\varepsilon$-capacity of sets in function spaces. Am. Math. Soc. Transl. (2) **17**, 277–364 (1961)
13. Le Cam, L.: Convergence of estimates under dimensionality restrictions. Ann. Statist. **1**, 38–53 (1973)
14. Le Cam, L.: Asymptotic methods in statistical decision theory. To be published
15. Lorentz, G.G.: Metric entropy and approximation. Bull. Amer. Math. Soc. **72**, 903–937 (1966)
16. Lorentz, G.G.: Approximation of Functions. New York: Holt, Rinehart, Winston 1966
17. Müller, H.G., Gasser, T.: Optimal convergence properties of kernel estimates of derivatives of a density function. In: Smoothing Techniques for Curve Estimation (T. Gasser and M. Rosenblatt eds.). Lecture Notes in Mathematics no. **757**, pp. 144–154. Berlin-Heidelberg-New York: Springer 1979
18. Stone, C.J.: Optimal rates of convergence for nonparametric estimators. Ann. Statist. **8**, 1348–1360 (1980)
19. Stone, C.J.: Optimal uniform rate of convergence for nonparametric estimators of a density function or its derivatives. University of California Berkeley: preprint (1983)