

RHO-ESTIMATORS REVISITED: GENERAL THEORY AND APPLICATIONS

BY YANNICK BARAUD AND LUCIEN BIRGÉ

Université Côte d'Azur and Sorbonne Universités

Following Baraud, Birgé and Sart [*Invent. Math.* **207** (2017) 425–517], we pursue our attempt to design a robust universal estimator of the joint distribution of n independent (but not necessarily i.i.d.) observations for an Hellinger-type loss. Given such observations with an unknown joint distribution \mathbf{P} and a dominated model \mathcal{Q} for \mathbf{P} , we build an estimator $\hat{\mathbf{P}}$ based on \mathcal{Q} (a ρ -estimator) and measure its risk by an Hellinger-type distance. When \mathbf{P} does belong to the model, this risk is bounded by some quantity which relies on the local complexity of the model in a vicinity of \mathbf{P} . In most situations, this bound corresponds to the minimax risk over the model (up to a possible logarithmic factor). When \mathbf{P} does not belong to the model, its risk involves an additional bias term proportional to the distance between \mathbf{P} and \mathcal{Q} , whatever the true distribution \mathbf{P} . From this point of view, this new version of ρ -estimators improves upon the previous one described in Baraud, Birgé and Sart [*Invent. Math.* **207** (2017) 425–517] which required that \mathbf{P} be absolutely continuous with respect to some known reference measure. Further additional improvements have been brought as compared to the former construction. In particular, it provides a very general treatment of the regression framework with random design as well as a computationally tractable procedure for aggregating estimators. We also give some conditions for the maximum likelihood estimator to be a ρ -estimator. Finally, we consider the situation where the statistician has at her or his disposal many different models and we build a penalized version of the ρ -estimator for model selection and adaptation purposes. In the regression setting, this penalized estimator not only allows one to estimate the regression function but also the distribution of the errors.

1. Introduction. In a previous paper, namely Baraud, Birgé and Sart (2017), we introduced a new class of estimators that we called ρ -estimators for estimating the distribution \mathbf{P} of a random variable $\mathbf{X} = (X_1, \dots, X_n)$ with values in some measurable space $(\mathcal{X}, \mathcal{B})$ under the assumption that the X_i are independent but not necessarily i.i.d. These estimators are based on density models, a *density model* being a family of densities \mathbf{t} with respect to some reference measure $\boldsymbol{\mu}$ on \mathcal{X} . We also assumed that \mathbf{P} was absolutely continuous with respect to $\boldsymbol{\mu}$ with density \mathbf{s} and, following Le Cam (1973), we measured the performance of an estimator $\hat{\mathbf{s}}$

Received June 2016; revised November 2017.

MSC2010 subject classifications. 62G35, 62G05, 62G07, 62G08, 62C20, 62F99.

Key words and phrases. ρ -estimation, robust estimation, density estimation, regression with random design, statistical models, maximum likelihood estimators, metric dimension, VC-classes.

of \mathbf{s} in terms of $\mathbf{h}^2(\mathbf{s}, \widehat{\mathbf{s}})$, where \mathbf{h} is a Hellinger-type distance to be defined later. Originally, the motivations for this construction were to design an estimator $\widehat{\mathbf{s}}$ of \mathbf{s} with the following properties.

- Given a density model \mathbf{S} , the estimator $\widehat{\mathbf{s}}$ should be nearly optimal over \mathbf{S} from the minimax point of view, which means that it is possible to bound the risk of the estimator $\widehat{\mathbf{s}}$ over \mathbf{S} from above by some quantity $CD_n(\mathbf{S})$ which is approximately of the order of the minimax risk over \mathbf{S} .
- Since in statistics we typically have incomplete information about the true distribution of the observations, when we assume that \mathbf{s} belongs to \mathbf{S} nothing ever warrants that this is true. We may more reasonably expect that \mathbf{s} is close to \mathbf{S} which means that the model \mathbf{S} is not exact but only approximate and that the quantity $\mathbf{h}(\mathbf{s}, \mathbf{S}) = \inf_{\mathbf{t} \in \mathbf{S}} \mathbf{h}(\mathbf{s}, \mathbf{t})$ might therefore be positive. In this case, we would like the risk of $\widehat{\mathbf{s}}$ to be bounded by $C'[D_n(\mathbf{S}) + \mathbf{h}^2(\mathbf{s}, \mathbf{S})]$ for some universal constant C' . In the case of ρ -estimators, the previous bound can actually be slightly refined and expressed in the following way. It is possible to define on \mathbf{S} a positive function R such that the risk of the ρ -estimator is not larger than $R(\mathbf{s})$, with $R(\mathbf{s}) \leq CD_n(\mathbf{S})$ if \mathbf{s} belongs to the model \mathbf{S} and not larger than $C' \inf_{\bar{\mathbf{s}} \in \mathbf{S}} [R(\bar{\mathbf{s}}) + \mathbf{h}^2(\mathbf{s}, \bar{\mathbf{s}})]$ when \mathbf{s} does not belong to \mathbf{S} .

The weak sensibility of this risk bound to small deviations with respect to the Hellinger-type distance \mathbf{h} between \mathbf{s} and an element $\bar{\mathbf{s}}$ of \mathbf{S} covers some classical notions of robustness among which robustness to a possible contamination of the data and robustness to outliers, as we shall see in Section 5.

There are nevertheless some limitations to the properties of ρ -estimators as defined in Baraud, Birgé and Sart (2017).

(a) The study of random design regression required that either the distribution of the design be known or that the errors have a symmetric distribution. We want to relax these assumptions and consider the random design regression framework with greater generality.

(b) We always worked with some reference measure μ and assumed that all the probabilities we considered, including the true distribution \mathbf{P} of X , were absolutely continuous with respect to μ . This is quite natural for the probabilities that belong to our models since the models are, by assumption, dominated and, typically, defined via a reference measure μ and a family of densities with respect to μ . Nevertheless, the assumption that the true distribution \mathbf{P} of the observations be also dominated by μ is questionable. We therefore would like to get rid of it and let the true distribution be completely arbitrary, relaxing thus the assumption that the density \mathbf{s} exists. Unexpectedly, such an extension leads to subtle complications as we shall see below and this generalization is actually far from being straightforward.

(c) Our construction was necessarily restricted to countable models rather than the uncountable ones currently used in statistics.

We want here to design a method based on “probability models” rather than “density models,” which means working with dominated models \mathcal{P} consisting of probabilities rather than of densities as for \mathbf{S} . Of course, the choice of a dominating measure μ and a specific set \mathbf{S} of densities leads to a probability model \mathcal{P} . This is by the way what is actually done in statistics, but the converse is definitely not true and there exist many ways of representing a dominated probability model by a reference measure and a set of densities. It turns out (see Section 2.3) that the performance of a very familiar estimator, namely the Maximum Likelihood Estimator (MLE), can be strongly affected by the choice of a specific version of the densities. Our purpose here is to design an estimator the performance of which only depends on the probability model \mathcal{P} and not on the choice of the reference measure and the densities that are used to represent it.

In order to get rid of the above mentioned restrictions, we have to modify our original construction which leads to the new version that we present here. This new version retains all the nice properties that we proved in Baraud, Birgé and Sart (2017) and the numerous illustrations we considered there remain valid for the new version. It additionally provides a general treatment of conditional density estimation and regression, allowing the statistician to estimate both the regression function and the error distribution even when the distribution of the design is totally unknown and the errors admit no finite moments. From this point of view, our approach contrasts very much with that based on the classical least squares. An alternative point of view on the particular problem of estimating a conditional density can be found in Sart (2017).

A thorough study of the performance of the least squares estimator (or truncated versions of it) can be found in Györfi et al. (2002) and we refer the reader to the references therein. A nice feature of these results lies in the fact that they hold without any assumption on the distribution of the design. While few moment conditions on the errors are necessary to bound the \mathbb{L}_2 -integrated risk of their estimator, much stronger ones, typically boundedness of the errors, are necessary to obtain exponential deviation bounds. In contrast, in linear regression, Audibert and Catoni (2011) established exponential deviation bounds for the risk of some robust versions of the ordinary least squares estimator. Their idea is to replace the sum of squares by the sum of their truncated version in view of designing a new criterion which is less sensitive to possible outliers than the original least squares. Their way of modifying the least squares criterion shares some similarity with our way of modifying the log-likelihood criterion, as we shall see below. However, their results require some conditions on the distribution of the design as well as some (weak) moment condition on the errors while ours do not.

It is known, and we shall give an additional example below, that the MLE, which is often considered as a “universal” estimator, does not possess, in general, the properties that we require and more specifically robustness. An illustration of the lack of robustness of the MLE with respect to Hellinger deviations is provided in Baraud and Birgé (2016). Some other weaknesses of the MLE have been described

in Le Cam (1990) and Birgé (2006), among other authors, and various alternatives aimed at designing some sorts of “universal” estimators (for the problem we consider here) which would not suffer from the same weaknesses have been proposed in the past by Le Cam (1973) and (1975) followed by Birgé (1983) and (2006). The construction of ρ -estimators, as described in Baraud, Birgé and Sart (2017) was in this line. In that paper, we actually introduced ρ -estimators via a testing argument as was the case for Le Cam and Birgé for their methods. This argument remains valid for the generalized version we consider here (see Lemma D.3 of the Supplementary Material) but ρ -estimators can also be viewed as a generalization, and in fact a robustified version, of the MLE. We shall even show, in Section 6, that in favorable situations (i.i.d. observations and a convex separable set of densities as a model for the true density) the MLE is actually a ρ -estimator and, therefore, shares their properties.

To explain the idea underlying the construction of ρ -estimators, let us assume that we observe an n -sample $\mathbf{X} = (X_1, \dots, X_n)$ with an unknown density q belonging to a set $\overline{\mathcal{Q}}$ of densities with respect to some reference measure μ . We may write the log-likelihood of q as $\sum_{i=1}^n \log(q(X_i))$ and the log-likelihood ratios as

$$\mathbf{L}(\mathbf{X}, q, q') = \sum_{i=1}^n \log\left(\frac{q'(X_i)}{q(X_i)}\right) = \sum_{i=1}^n \log(q'(X_i)) - \sum_{i=1}^n \log(q(X_i)),$$

so that maximizing the likelihood is equivalent to minimizing with respect to q

$$\mathbf{L}(\mathbf{X}, q) = \sup_{q' \in \overline{\mathcal{Q}}} \sum_{i=1}^n \log\left(\frac{q'(X_i)}{q(X_i)}\right) = \sup_{q' \in \overline{\mathcal{Q}}} \mathbf{L}(\mathbf{X}, q, q').$$

This happens simply because of the magic property of the logarithm which says that $\log(a/b) = \log a - \log b$. However, the use of the unbounded log function in the definition of $\mathbf{L}(\mathbf{X}, q)$ leads to various problems that are responsible for some weaknesses of the MLE. Replacing the log function by another function φ amounts to replace $\mathbf{L}(\mathbf{X}, q, q')$ by

$$(1) \quad \mathbf{T}(\mathbf{X}, q, q') = \sum_{i=1}^n \varphi\left(\frac{q'(X_i)}{q(X_i)}\right)$$

which is different from $\sum_{i=1}^n \varphi(q'(X_i)) - \sum_{i=1}^n \varphi(q(X_i))$ since φ is not the log function. We may nevertheless define the analogue of $\mathbf{L}(\mathbf{X}, q)$, namely

$$(2) \quad \mathbf{Y}(\mathbf{X}, q) = \sup_{q' \in \overline{\mathcal{Q}}} \mathbf{T}(\mathbf{X}, q, q') = \sup_{q' \in \overline{\mathcal{Q}}} \sum_{i=1}^n \varphi\left(\frac{q'(X_i)}{q(X_i)}\right)$$

and define our estimator $\widehat{q}(\mathbf{X})$ as a minimizer with respect to $q \in \overline{\mathcal{Q}}$ of the quantity $\mathbf{Y}(\mathbf{X}, q)$. The resulting estimator is an alternative to the MLE and we shall show that, for a suitable choice of a bounded function φ , it enjoys various properties, among which robustness, that are often not shared by the MLE.

To analyse the performance of this new estimator, we have to study the behaviour of the process $\mathbf{T}(X, q, q')$ when q is fixed, $q \cdot \mu$ is close to the true distribution of the X_i and q' varies in $\overline{\mathcal{Q}}$. Since the function φ is bounded, the process is similar to those considered in learning theory for the purpose of studying empirical risk minimization. As a consequence, the tools we use are also similar to those described in great detail in Koltchinskii (2006).

It is well known that working with a single model for estimating an unknown distribution is not very efficient unless one has very precise pieces of information about the true distribution, which is rarely the case. Working with many models simultaneously and performing model selection improves the situation drastically. Refining the previous construction of ρ -estimators by adding suitable penalty terms to the statistic $\mathbf{T}(X, q, q')$ allows one to work with a finite or countable family of probability models $\{\mathcal{P}_m, m \in \mathcal{M}\}$ instead of a single one, each model \mathcal{P}_m leading to a risk bound of the form $C'[D_n(\mathcal{P}_m) + \mathbf{h}^2(\mathcal{P}_m, \mathbf{P})]$, and to choose from the observations a model with approximately the best possible bound which results in a final estimator $\hat{\mathbf{P}}$ and a bound for $\mathbf{h}^2(\hat{\mathbf{P}}, \mathbf{P})$ of the form

$$C'' \inf_{m \in \mathcal{M}} [D_n(\mathcal{P}_m) + \mathbf{h}^2(\mathcal{P}_m, \mathbf{P}) + \Delta_m],$$

where the additional term Δ_m is connected to the complexity of the family of models we use.

The paper is organised as follows. We shall first make our framework, which is based on dominated families of probabilities rather than families of densities with respect to a given dominating measure, precise in Section 2. This section is devoted to the definition of models and of our new version of ρ -estimators, then to the assumptions that the function φ we use to define the statistic \mathbf{T} in (1) should satisfy. In Section 3, we define the ρ -dimension function of a model, a quantity which measures the difficulty of estimation within the model using a ρ -estimator, and present the main results, namely the performance of these new ρ -estimators. Section 4 is devoted to the extension of the construction from countable to uncountable statistical models (which are the ones currently used in statistics) under suitable assumptions. We describe the robustness properties of ρ -estimators in Section 5. In Section 6, we investigate the relationship between ρ -estimators and the MLE when the model is a convex set of densities. Section 7 provides various methods that allow one to bound the ρ -dimension functions of different types of models and indicates how these bounds are to be used to bound the risk of ρ -estimators in typical situations with applications to the minimax risk over classical statistical models. We also provide a few examples of computations of bounds for the ρ -dimension function. Many applications of our results about ρ -estimators have already been given in Baraud, Birgé and Sart (2017) and we deal here with a new one: estimation of conditional distributions in Section 8. In Section 9, we apply this additional result to the special case of random design regression when the distribution of the design is completely unknown, a situation for which not

many results are known. We provide here a complete treatment of this regression framework with simultaneous estimation of both the regression function and the density of the errors. Section 10 is devoted to estimator selection and aggregation: we show there how our procedure can be used either to select an element from a family of preliminary estimators or to aggregate them in a convex way. The supplemental article [Baraud and Birgé (2018)] contains most of the proofs as well as some additional facts and comments.

2. Our new framework and estimation strategy. As already mentioned, our method is based on statistical models which are sets of probability distributions, in opposition with more classical models which are sets of densities with respect to a given dominating measure.

2.1. *A probabilistic framework.* We observe a random variable $X = (X_1, \dots, X_n)$ defined on some probability space $(\Omega, \mathfrak{E}, \mathbb{P})$ with independent components X_i and values in the measurable product space $(\mathcal{X}, \mathfrak{B}) = (\prod_{i=1}^n \mathcal{X}_i, \otimes_{i=1}^n \mathfrak{B}_i)$. We denote by \mathcal{P} the set of all product probabilities on $(\mathcal{X}, \mathfrak{B})$ and by $\mathbf{P} = \otimes_{i=1}^n P_i \in \mathcal{P}$ the true distribution of X . We identify an element $\mathbf{Q} = \otimes_{i=1}^n Q_i$ of \mathcal{P} with the n -tuple (Q_1, \dots, Q_n) and extend this identification to the elements $\mu = \otimes_{i=1}^n \mu_i$ of the set \mathcal{M} of all σ -finite product measures on $(\mathcal{X}, \mathfrak{B})$.

When \mathbf{Q} is absolutely continuous with respect to $\mu \in \mathcal{M}$ ($\mathbf{Q} \ll \mu$) or, equivalently, μ dominates \mathbf{Q} , each Q_i , for $i = 1, \dots, n$, is absolutely continuous with respect to μ_i with density q_i so that $Q_i = q_i \cdot \mu_i$. We denote by $\mathcal{L}(\mu_i)$ the set of all densities with respect to μ_i , that is, the set of measurable functions q from \mathcal{X}_i to \mathbb{R}_+ such that $\int_{\mathcal{X}_i} q(x) d\mu_i(x) = 1$. We then write $\mathbf{Q} = \mathbf{q} \cdot \mu$ where \mathbf{q} is the n -tuple (q_1, \dots, q_n) and we say that \mathbf{q} is a density for \mathbf{Q} with respect to μ . We denote by $\mathcal{L}(\mu) = \prod_{i=1}^n \mathcal{L}(\mu_i)$ the set of such densities \mathbf{q} and by \mathcal{P}^μ the set of all those $\mathbf{P}' \in \mathcal{P}$ which are absolutely continuous with respect to μ .

Our aim is to estimate the unknown distribution $\mathbf{P} = (P_1, \dots, P_n)$ from the observation of X . In order to evaluate the performance of an estimator $\hat{\mathbf{P}}(X) \in \mathcal{P}$ of \mathbf{P} , we shall introduce, following Le Cam (1975), an Hellinger-type distance \mathbf{h} on \mathcal{P} . We recall that, given two probabilities Q and Q' on a measurable space $(\mathcal{X}, \mathfrak{B})$, the Hellinger distance and the Hellinger affinity between Q and Q' are respectively given by

$$(3) \quad \begin{cases} h^2(Q, Q') = \frac{1}{2} \int_{\mathcal{X}} \left(\sqrt{\frac{dQ}{d\mu}} - \sqrt{\frac{dQ'}{d\mu}} \right)^2 d\mu, \\ \rho(Q, Q') = \int_{\mathcal{X}} \sqrt{\frac{dQ}{d\mu} \frac{dQ'}{d\mu}} d\mu = 1 - h^2(Q, Q'), \end{cases}$$

where μ denotes any measure that dominates both Q and Q' , the result being independent of the choice of μ . The Hellinger-type distance $\mathbf{h}(\mathbf{Q}, \mathbf{Q}')$ and affinity

$\rho(\mathbf{Q}, \mathbf{Q}')$ between two elements $\mathbf{Q} = (Q_1, \dots, Q_n)$ and $\mathbf{Q}' = (Q'_1, \dots, Q'_n)$ of \mathcal{P} are then given by the formulas

$$\mathbf{h}^2(\mathbf{Q}, \mathbf{Q}') = \sum_{i=1}^n h^2(Q_i, Q'_i) = \sum_{i=1}^n [1 - \rho(Q_i, Q'_i)] = n - \rho(\mathbf{Q}, \mathbf{Q}').$$

We shall denote by \mathcal{V} the topology of the metric space $(\mathcal{P}, \mathbf{h})$.

2.2. *Models and their representations.* Let us start with this definition.

DEFINITION 1. We call *model* any dominated subset $\overline{\mathcal{Q}}$ of \mathcal{P} and we call *representation* of (the model) $\overline{\mathcal{Q}}$ a pair $\mathcal{R}(\overline{\mathcal{Q}}) = (\mu, \overline{\mathcal{Q}})$ where $\mu = (\mu_1, \dots, \mu_n)$ is a σ -finite measure which dominates $\overline{\mathcal{Q}}$ and $\overline{\mathcal{Q}}$ is a subset of $\mathcal{L}(\mu)$ such that for any \mathbf{Q} in $\overline{\mathcal{Q}}$ there exists a unique density $\mathbf{q} \in \overline{\mathcal{Q}}$ with $\mathbf{Q} = \mathbf{q} \cdot \mu$.

This means that, given a representation $(\mu, \overline{\mathcal{Q}})$ of the model $\overline{\mathcal{Q}}$, we can associate to each probability $\mathbf{Q} \in \overline{\mathcal{Q}}$ a density $\mathbf{q} \in \overline{\mathcal{Q}}$ and vice versa. Clearly, a dominated subset $\overline{\mathcal{Q}}$ has different representations depending on the choice of the dominating measure μ and the versions of the densities $q_i = dQ_i/d\mu_i$.

Our estimation strategy is based on specific dominated subsets of \mathcal{P} that we call ρ -models.

DEFINITION 2. A ρ -model is a *countable* (which in this paper always means either finite or infinite and countable) subset \mathcal{Q} of \mathcal{P} .

A ρ -model \mathcal{Q} being countable, it is necessarily dominated. One should think of it as a probability set to which the true distribution is believed to be close (with respect to the Hellinger-type distance \mathbf{h}).

2.3. *Construction of a ρ -estimator on a model \mathcal{Q} .* Given the model \mathcal{Q} , our estimator is defined as a random element of $\text{Cl}(\mathcal{Q})$, where $\text{Cl}(\mathcal{R})$ denotes the closure of the subset \mathcal{R} of \mathcal{P} in the metric space $(\mathcal{P}, \mathbf{h})$, and its construction relies on a particular representation $\mathcal{R}(\mathcal{Q})$ of the model \mathcal{Q} . It actually depends on three elements with specific properties to be made precise below:

(i) A function ψ (which will serve as a substitute for the logarithm to derive an alternative to the MLE) with the following properties.

ASSUMPTION 1. The function ψ is nondecreasing from $[0, +\infty]$ to $[-1, 1]$, Lipschitz and satisfies

$$(4) \quad \psi(x) = -\psi(1/x) \quad \text{for all } x \in [0, +\infty), \text{ hence } \psi(1) = 0.$$

Throughout this paper, we shall only consider, without further notice, functions ψ satisfying Assumption 1.

(ii) A model $\mathcal{Q} \subset \mathcal{P}$ (in most cases a ρ -model) with a representation $\mathcal{R}(\mathcal{Q}) = (\boldsymbol{\mu}, \mathcal{Q})$.

(iii) A *penalty function* “**pen**” mapping \mathcal{Q} to \mathbb{R} , the role of which will be explained later in Section 3. We may, at first reading, assume that this penalty function is identically 0.

It is essential to note that the dominating measure $\boldsymbol{\mu}$ is chosen by the statistician and that there is no reason that the true distribution \mathbf{P} of X be absolutely continuous with respect to $\boldsymbol{\mu}$. On the contrary, all probabilities \mathbf{P}' on \mathcal{X} belonging to $\text{Cl}(\mathcal{Q})$ are absolutely continuous with respect to $\boldsymbol{\mu}$.

Given the function ψ and the representation $\mathcal{R}(\mathcal{Q})$, we define the real-valued function \mathbf{T} on $\mathcal{X} \times \mathcal{Q} \times \mathcal{Q}$ by

$$(5) \quad \mathbf{T}(\mathbf{x}, \mathbf{q}, \mathbf{q}') = \sum_{i=1}^n \psi \left(\sqrt{\frac{q'_i(x_i)}{q_i(x_i)}} \right) \quad \text{for } \mathbf{x} = (x_1, \dots, x_n) \in \mathcal{X} \text{ and } \mathbf{q}, \mathbf{q}' \in \mathcal{Q},$$

with the conventions $0/0 = 1$ and $a/0 = +\infty$ for all $a > 0$. We then set (with $\mathbf{Q} = \mathbf{q} \cdot \boldsymbol{\mu}$ and $\mathbf{Q}' = \mathbf{q}' \cdot \boldsymbol{\mu}$)

$$(6) \quad \Upsilon(X, \mathbf{q}) = \sup_{\mathbf{q}' \in \mathcal{Q}} [\mathbf{T}(X, \mathbf{q}, \mathbf{q}') - \mathbf{pen}(\mathbf{Q}')] + \mathbf{pen}(\mathbf{Q}) \quad \text{for all } \mathbf{q} \in \mathcal{Q}.$$

DEFINITION 3 (ρ -estimators). Let $\mathcal{E}(\psi, X)$ be the (nonvoid) set

$$(7) \quad \mathcal{E}(\psi, X) = \left\{ \mathbf{Q} = \mathbf{q} \cdot \boldsymbol{\mu}, \mathbf{q} \in \mathcal{Q} \mid \Upsilon(X, \mathbf{q}) < \inf_{\mathbf{q}' \in \mathcal{Q}} \Upsilon(X, \mathbf{q}') + \frac{\kappa}{25} \right\},$$

where the positive constant κ is given by (19) below. A ρ -estimator $\widehat{\mathbf{P}} = \widehat{\mathbf{P}}(X)$ relative to $(\mathcal{R}(\mathcal{Q}), \mathbf{pen})$ is any (measurable) element of $\text{Cl}(\mathcal{E}(\psi, X))$.

Since $\widehat{\mathbf{P}}$ belongs to $\text{Cl}(\mathcal{E}(\psi, X))$, the elements of which are dominated by $\boldsymbol{\mu}$, there exists a random density $\widehat{\mathbf{p}} = (\widehat{p}_1, \dots, \widehat{p}_n)$ with $\widehat{p}_i \in \mathcal{L}(\mu_i)$ for $i = 1, \dots, n$ such that $\widehat{\mathbf{P}} = \widehat{\mathbf{p}} \cdot \boldsymbol{\mu}$. Note that $\widehat{\mathbf{P}}$ might not belong to \mathcal{Q} .

As an immediate consequence of Assumption 1 and the convention $1/0 = +\infty$, $\psi(+\infty) = -\psi(0)$ and

$$(8) \quad \mathbf{T}(X, \mathbf{q}, \mathbf{q}') = -\mathbf{T}(X, \mathbf{q}', \mathbf{q}) \quad \text{for all } \mathbf{q}, \mathbf{q}' \in \mathcal{Q}.$$

Moreover,

$$\Upsilon(X, \mathbf{q}) \geq [\mathbf{T}(X, \mathbf{q}, \mathbf{q}) - \mathbf{pen}(\mathbf{Q})] + \mathbf{pen}(\mathbf{Q}) = \mathbf{T}(X, \mathbf{q}, \mathbf{q}) = n\psi(1) = 0$$

for all $\mathbf{q} \in \mathcal{Q}$, which implies that any element $\widehat{\mathbf{P}} = \widehat{\mathbf{p}} \cdot \boldsymbol{\mu}$ in \mathcal{Q} such that $\Upsilon(X, \widehat{\mathbf{p}}) < \kappa/25$ is a ρ -estimator. In particular, when $\mathbf{pen}(\mathbf{Q}) = 0$ for all $\mathbf{Q} \in \mathcal{Q}$ (which we shall write in the sequel $\mathbf{pen} = \mathbf{0}$) and $\Upsilon(X, \widehat{\mathbf{p}}) = 0$, it follows from (6) that

$$\mathbf{T}(X, \widehat{\mathbf{p}}, \mathbf{q}) \leq \Upsilon(X, \widehat{\mathbf{p}}) = 0 = \mathbf{T}(X, \widehat{\mathbf{p}}, \widehat{\mathbf{p}}) \leq -\mathbf{T}(X, \widehat{\mathbf{p}}, \mathbf{q}) = \mathbf{T}(X, \mathbf{q}, \widehat{\mathbf{p}})$$

for all $\mathbf{q} \in \mathcal{Q}$. This means that, in this case, $(\hat{\mathbf{p}}, \hat{\mathbf{p}})$ is a saddle point of the map $(\mathbf{q}, \mathbf{q}') \mapsto \mathbf{T}(X, \mathbf{q}, \mathbf{q}')$.

A ρ -estimator \mathbf{P} depends on the chosen representation $\mathcal{R}(\mathcal{Q})$ of \mathcal{Q} and there are different versions of the ρ -estimators associated to \mathcal{Q} , even though most of the time \mathcal{Q} will directly be given by a specific representation, that is a family \mathcal{Q} of densities with respect to some reference measure μ . Here is the important point, to be proven in Section 3: when \mathcal{Q} is a ρ -model, the risk bounds we shall derive only depend on \mathcal{Q} and the penalty function but not on the chosen representation of \mathcal{Q} , which allows us to choose the more convenient one for the construction. In contrast, the performances of many classical estimators are sensitive to the representation of the model \mathcal{Q} and this is in particular the case of the MLE as shown by the following example.

PROPOSITION 1. *Let us consider a sequence of i.i.d. random variables $(X_k)_{k \geq 1}$ defined on a measurable space $(\Omega, \mathcal{A}, \mathbb{P})$ with normal distribution $P_\theta = \mathcal{N}(\theta, 1)$ for some unknown $\theta \in \mathbb{R}$. We choose for reference measure $\mu = \mathcal{N}(0, 1)$ and for the version of $dP_\theta/d\mu$, $\theta \in \mathbb{R}$, the function*

$$(9) \quad p_\theta(x) = \exp\left[\theta x - (\theta^2/2) + (\theta^2/2) \exp(x^2) \mathbb{1}_\theta(x) \mathbb{1}_{(0, +\infty)}(\theta)\right].$$

Whatever the value of the true parameter θ , on a set of probability tending to 1 when n goes to infinity, the MLE is given by $X_{(n)} = \max\{X_1, \dots, X_n\}$ and is therefore inconsistent.

The proof of Proposition 1 is given in Section D.1 of the Supplementary Material [Baraud and Birgé (2018)]. Note that the usual choice for p_θ : $x \mapsto \exp[-x\theta + (\theta^2/2)]$ for $dP_\theta/d\mu$ is purely conventional. Mathematically speaking, our choice (9) is perfectly correct but leads to an inconsistent MLE. Also note that the usual tools that are used to prove consistency of the MLE, like bracketing entropy [see, for instance, Theorem 7.4 of van de Geer (2000)] are not stable with respect to changes of versions of the densities in the family. The same is true for arguments based on VC-classes that we used in Baraud, Birgé and Sart (2017). Choosing a convenient set of densities to work with is well grounded as long as the reference measure μ not only dominates the model but also the true distribution \mathbf{P} . If not, sets of null measure with respect to μ might have a positive probability under \mathbf{P} and it becomes unclear how the choice of this set of densities influences the performance of the estimator.

2.4. Notation and conventions. Throughout this paper, given a representation $\mathcal{R}(\overline{\mathcal{Q}}) = (\mu, \overline{\mathcal{Q}})$ of a model $\overline{\mathcal{Q}}$, we shall use lower case letters $\mathbf{q}, \mathbf{q}', \dots$ and q_i, q'_i, \dots for denoting the chosen densities of $\mathbf{Q}, \mathbf{Q}', \dots$ and Q_i, Q'_i, \dots with respect to the reference measures μ and μ_i , respectively, for all $i = 1, \dots, n$. We set $\log_+(x) = \max\{\log x, 0\}$ for all $x > 0$; $|A|$ denotes the cardinality of the

set A ; $\mathcal{B}(\mathbf{P}, r) = \{\mathbf{Q} \in \mathcal{P} \mid \mathbf{h}(\mathbf{P}, \mathbf{Q}) \leq r\}$ is the closed Hellinger-type ball in \mathcal{P} with center \mathbf{Q} and radius r . Given a set E , a nonnegative function ℓ on $E \times E$, $x \in E$ and $A \subset E$, we set $\ell(x, A) = \inf_{y \in A} \ell(x, y)$. In particular, for $\mathcal{R} \subset \mathcal{P}$, $\mathbf{h}(\mathbf{P}, \mathcal{R}) = \inf_{\mathbf{R} \in \mathcal{R}} \mathbf{h}(\mathbf{P}, \mathbf{R})$. We set $x \vee y$ and $x \wedge y$ for $\max\{x, y\}$ and $\min\{x, y\}$, respectively. By convention $\sup_{\emptyset} = 0$, the ratio $u/0$ equals $+\infty$ for $u > 0$, $-\infty$ for $u < 0$ and 1 for $u = 0$.

2.5. *Our assumptions.* Given the ρ -model \mathcal{Q} , let us now indicate what properties the function ψ (satisfying Assumption 1) are required in view of controlling the risk of the resulting ρ -estimators.

ASSUMPTION 2. Let \mathcal{Q} be the ρ -model to be used for the construction of ρ -estimators. There exist three positive constants a_0, a_1, a_2 with $a_0 \geq 1 \geq a_1$ and $a_2^2 \geq 1 \vee (6a_1)$ such that, whatever the representation $\mathcal{R}(\mathcal{Q}) = (\boldsymbol{\mu}, \mathcal{Q})$ of \mathcal{Q} , the densities $\mathbf{q}, \mathbf{q}' \in \mathcal{Q}$, the probability $\mathbf{R} \in \mathcal{P}$ and $i \in \{1, \dots, n\}$,

$$(10) \quad \int_{\mathcal{X}_i} \psi \left(\sqrt{\frac{q'_i}{q_i}} \right) dR_i \leq a_0 h^2(R_i, Q_i) - a_1 h^2(R_i, Q'_i)$$

and

$$(11) \quad \int_{\mathcal{X}_i} \psi^2 \left(\sqrt{\frac{q'_i}{q_i}} \right) dR_i \leq a_2^2 [h^2(R_i, Q_i) + h^2(R_i, Q'_i)].$$

Note that the left-hand sides of (10) and (11) depend on the choices of the reference measures μ_i and versions of the densities $q_i = dQ_i/d\mu_i$ and $q'_i = dQ'_i/d\mu_i$ while the corresponding right-hand sides do not.

Given ψ that satisfies Assumption 2, the values of a_0, a_1 and a_2 are clearly not uniquely defined but, in the sequel, when we shall say that Assumption 2 holds, this will mean that the function ψ satisfies (10) and (11) with given values of these constants which will therefore be considered as fixed once ψ has been chosen. When we shall say that some quantity depends on ψ , it will implicitly mean that it depends on these chosen values of a_0, a_1 and a_2 .

An important consequence of (8), (10) and (11) is the fact that, for all \mathbf{Q}, \mathbf{Q}' in \mathcal{Q} and $\mathbf{P} \in \mathcal{P}$,

$$(12) \quad a_1 \mathbf{h}^2(\mathbf{P}, \mathbf{Q}) - a_0 \mathbf{h}^2(\mathbf{P}, \mathbf{Q}') \leq \mathbb{E}[\mathbf{T}(\mathbf{X}, \mathbf{q}, \mathbf{q}')] \leq a_0 \mathbf{h}^2(\mathbf{P}, \mathbf{Q}) - a_1 \mathbf{h}^2(\mathbf{P}, \mathbf{Q}').$$

These inequalities follow by summing the inequalities (10) with respect to i with $\mathbf{R} = \mathbf{P}$, then exchanging the roles of \mathbf{Q} and \mathbf{Q}' and applying (8). They imply that the sign of $\mathbb{E}[\mathbf{T}(\mathbf{X}, \mathbf{q}, \mathbf{q}')]$ tells us which of the two distributions \mathbf{Q} and \mathbf{Q}' is closer to the true one when the ratio between the distances $\mathbf{h}(\mathbf{P}, \mathbf{Q})$ and $\mathbf{h}(\mathbf{P}, \mathbf{Q}')$ is far enough from one.

In view of checking that a given function ψ satisfies Assumption 2, the next result to be proved in Section D.3 of the Supplementary Material [Baraud and Birgé (2018)] is useful.

PROPOSITION 2. *If, for a particular representation $\mathcal{R}(\mathcal{Q}) = (\boldsymbol{\mu}, \mathcal{Q})$ of the ρ -model \mathcal{Q} and any probability $\mathbf{R} \in \mathcal{P}^\mu$, the function ψ satisfies (10) and (11) for positive constants $a_0 > 2$, $a_1 \leq [(a_0 - 2)/2] \wedge 1$ and $a_2^2 \geq 1 \vee (6a_1)$, then it satisfies Assumption 2 with the same constants a_0 , a_1 and a_2 .*

This proposition means that, up to a possible adjustment of the constants a_0 and a_1 , it is actually enough to check that (10) and (11) hold true for a given representation $(\boldsymbol{\mu}, \mathcal{Q})$ of \mathcal{Q} and all probabilities $\mathbf{R} \ll \boldsymbol{\mu}$.

Let us now introduce two functions ψ which do satisfy Assumption 2.

PROPOSITION 3. *Let ψ_1 and ψ_2 be the functions taking the value 1 at $+\infty$ and defined for $x \in \mathbb{R}_+$ by*

$$\psi_1(x) = \frac{x - 1}{\sqrt{x^2 + 1}} \quad \text{and} \quad \psi_2(x) = \frac{x - 1}{x + 1}.$$

These two functions are continuously increasing from $[0, +\infty]$ to $[-1, 1]$, Lipschitz (with respective Lipschitz constants 1.143 and 2) and satisfy Assumption 2 for all ρ -models \mathcal{Q} with $a_0 = 4.97$, $a_1 = 0.083$, $a_2^2 = 3 + 2\sqrt{2}$ for ψ_1 and $a_0 = 4$, $a_1 = 3/8$, $a_2^2 = 3\sqrt{2}$ for ψ_2 .

Both functions can therefore be used everywhere in the applications of the present paper. Nevertheless, we prefer ψ_2 because it leads to better constants in the risk bounds of the estimator. Proposition 3 is proved in Section D.4 of the Supplementary Material [Baraud and Birgé (2018)]. Some comments on Assumption 2 can be found in Section D.2 of the Supplementary Material [Baraud and Birgé (2018)]. When the ρ -model reduces to two elements, our selection procedure can be interpreted as a robust test between two simple hypotheses. Upper bounds on the errors of the first and second kinds are established in Section D.10 of the Supplementary Material [Baraud and Birgé (2018)].

3. The performance of ρ -estimators on ρ -models.

3.1. *The ρ -dimension function.* The deviation $\mathbf{h}(\mathbf{P}, \widehat{\mathbf{P}})$ between the true distribution \mathbf{P} and a ρ -estimator $\widehat{\mathbf{P}}$ built on the ρ -model \mathcal{Q} is controlled by two terms which are the analogue of the classical bias and variance terms and we shall first introduce a function that replaces here the variance.

Let $y > 0$, $\mathbf{P}, \bar{\mathbf{P}} \in \mathcal{P}$ and \mathcal{D}_0 be an arbitrary subset of \mathcal{P} ; we define

$$\mathcal{B}^{\mathcal{D}_0}(\mathbf{P}, \bar{\mathbf{P}}, y) = \left\{ \mathbf{Q} \in \mathcal{D}_0 \mid \mathbf{h}^2(\mathbf{P}, \bar{\mathbf{P}}) + \mathbf{h}^2(\mathbf{P}, \mathbf{Q}) < y^2 \right\}$$

and for measurable nonnegative functions \mathbf{q}, \mathbf{q}' on $(\mathcal{X}, \mathcal{B})$, we set

$$(13) \quad \mathbf{Z}(X, \mathbf{q}, \mathbf{q}') = \mathbf{T}(X, \mathbf{q}, \mathbf{q}') - \mathbb{E}[\mathbf{T}(X, \mathbf{q}, \mathbf{q}')].$$

Given a representation $\mathcal{R} = (\boldsymbol{\mu}, \mathcal{Q})$ of $\mathcal{Z} \cup \{\bar{\mathbf{P}}\}$, we define

$$(14) \quad w(\mathcal{R}, \mathcal{Z}, \mathbf{P}, \bar{\mathbf{P}}, y) = \mathbb{E} \left[\sup_{\mathbf{Q} \in \mathcal{B}^{\mathcal{Z}}(\mathbf{P}, \bar{\mathbf{P}}, y)} |\mathbf{Z}(X, \bar{\mathbf{p}}, \mathbf{q})| \right],$$

where, for $\mathbf{Q} \in \mathcal{B}^{\mathcal{Z}}(\mathbf{P}, \bar{\mathbf{P}}, y) \subset \mathcal{Z}$, \mathbf{q} denotes the (unique) element of \mathcal{Q} such as $\mathbf{Q} = \mathbf{q} \cdot \boldsymbol{\mu}$ and $\bar{\mathbf{p}}$ denotes the element of \mathcal{Q} such that $\bar{\mathbf{P}} = \bar{\mathbf{p}} \cdot \boldsymbol{\mu}$. We recall that we use the convention $\sup_{\emptyset} = 0$. Since \mathcal{Z} is countable, so is $\mathcal{B}^{\mathcal{Z}}(\mathbf{P}, \bar{\mathbf{P}}, y) \subset \mathcal{Z}$. Therefore, the supremum of $|\mathbf{Z}(X, \bar{\mathbf{p}}, \cdot)|$ over $\mathcal{B}^{\mathcal{Z}}(\mathbf{P}, \bar{\mathbf{P}}, y)$ is measurable and the right-hand side of (14) is well defined. Also note that, since $\mathbf{T}(X, \bar{\mathbf{p}}, \bar{\mathbf{p}}) = n\psi(1) = 0$,

$$\mathbb{E} \left[\sup_{\mathbf{Q} \in \mathcal{B}^{\mathcal{Z}}(\mathbf{P}, \bar{\mathbf{P}}, y)} |\mathbf{Z}(X, \bar{\mathbf{p}}, \mathbf{q})| \right] = \mathbb{E} \left[\sup_{\mathbf{Q} \in \mathcal{B}^{\mathcal{Z} \cup \{\bar{\mathbf{P}}\}}(\mathbf{P}, \bar{\mathbf{P}}, y)} |\mathbf{Z}(X, \bar{\mathbf{p}}, \mathbf{q})| \right].$$

Hence $w(\mathcal{R}, \mathcal{Z}, \mathbf{P}, \bar{\mathbf{P}}, y) = w(\mathcal{R}, \mathcal{Z} \cup \{\bar{\mathbf{P}}\}, \mathbf{P}, \bar{\mathbf{P}}, y)$.

DEFINITION 4 (ρ -dimension function). Let \mathcal{Z} be a ρ -model and ψ some function satisfying Assumption 2 with constants a_0, a_1 and a_2 . The ρ -dimension function $D^{\mathcal{Z}}$ of \mathcal{Z} is the mapping from $\mathcal{P} \times \mathcal{P}$ to $[1, +\infty)$ given by

$$(15) \quad D^{\mathcal{Z}}(\mathbf{P}, \bar{\mathbf{P}}) = \left[\beta^2 \sup \left\{ y^2 \mid \mathbf{w}^{\mathcal{Z}}(\mathbf{P}, \bar{\mathbf{P}}, y) > \frac{a_1 y^2}{8} \right\} \right] \vee 1$$

with $\beta = a_1/(4a_2)$ and

$$\mathbf{w}^{\mathcal{Z}}(\mathbf{P}, \bar{\mathbf{P}}, y) = \inf_{\mathcal{R}} w(\mathcal{R}, \mathcal{Z}, \mathbf{P}, \bar{\mathbf{P}}, y) \quad \text{for all } y > 0,$$

where the infimum runs over all the representations $\mathcal{R} = (\boldsymbol{\mu}, \mathcal{Q})$ of $\mathcal{Z} \cup \{\bar{\mathbf{P}}\}$.

Note that the ρ -dimension function of \mathcal{Z} depends on the choice of the function ψ and not on the choice of the representations of $\mathcal{Z} \cup \{\bar{\mathbf{P}}\}$. Since it measures the local fluctuations of the centred empirical process $\mathbf{Z}(X, \bar{\mathbf{p}}, \mathbf{q})$ indexed by $\mathbf{q} \in \mathcal{Q}$, it is quite similar to the local Rademacher complexity introduced in Koltchinskii (2006) for the purpose of studying empirical risk minimization. Its importance comes from the following property.

PROPOSITION 4. Let \mathcal{Z} be a ρ -model, $\bar{\mathbf{P}} \in \mathcal{P}$ and $\mathcal{R} = (\boldsymbol{\mu}, \mathcal{Q})$, an arbitrary representation of $\mathcal{Z} \cup \{\bar{\mathbf{P}}\}$. Whatever $\mathbf{P} \in \mathcal{P}$,

$$(16) \quad w(\mathcal{R}, \mathcal{Z}, \mathbf{P}, \bar{\mathbf{P}}, y) \leq \mathbf{w}^{\mathcal{Z}}(\mathbf{P}, \bar{\mathbf{P}}, y) + 8\mathbf{h}^2(\mathbf{P}, \bar{\mathbf{P}}) \quad \text{for all } y > 0,$$

hence, for all $y > \beta^{-1} \sqrt{D^{\mathcal{Z}}(\mathbf{P}, \bar{\mathbf{P}})}$

$$(17) \quad w(\mathcal{R}, \mathcal{Z}, \mathbf{P}, \bar{\mathbf{P}}, y) \leq (a_1 y^2 / 8) + 8\mathbf{h}^2(\mathbf{P}, \bar{\mathbf{P}}).$$

The proof is provided in Section D.5 of the Supplementary Material [Baraud and Birgé (2018)].

3.2. *Exponential deviation bounds.* Our first theorem, to be proven in Section A.2 of the Supplementary Material [Baraud and Birgé (2018)], deals with the situation of a null penalty function $\text{pen} = \mathbf{0}$.

THEOREM 1. *Let \mathbf{P} be an arbitrary distribution in \mathcal{P} , \mathcal{Q} a ρ -model and ψ a function satisfying Assumption 2. Whatever the representation \mathcal{R} of \mathcal{Q} , a ρ -estimator $\widehat{\mathbf{P}}$ relative to $(\mathcal{R}, \mathbf{0})$ as defined in Section 2.3 satisfies, for all $\overline{\mathbf{P}} \in \mathcal{Q}$ and $\xi > 0$,*

$$(18) \quad \mathbb{P}\left[\mathbf{h}^2(\mathbf{P}, \widehat{\mathbf{P}}) \leq \gamma \mathbf{h}^2(\mathbf{P}, \overline{\mathbf{P}}) + \frac{4\kappa}{a_1} \left(\frac{D^{\mathcal{Q}}(\mathbf{P}, \overline{\mathbf{P}})}{4.7} + 1.49 + \xi \right)\right] \geq 1 - e^{-\xi},$$

with

$$(19) \quad \gamma = \frac{4(a_0 + 8)}{a_1} + 2 + \frac{84}{a_2^2} \quad \text{and} \quad \kappa = \frac{35a_2^2}{a_1} + 74, \quad \text{hence} \quad \frac{\kappa}{25} \geq 11.36.$$

In particular, if the ρ -dimension function $D^{\mathcal{Q}}$ is bounded on $\mathcal{P} \times \mathcal{Q}$ by $D_n \geq 1$, then

$$(20) \quad \mathbb{P}\left[\mathbf{C}\mathbf{h}^2(\mathbf{P}, \widehat{\mathbf{P}}) \leq \mathbf{h}^2(\mathbf{P}, \mathcal{Q}) + D_n + \xi\right] \geq 1 - e^{-\xi} \quad \text{for all } \xi > 0$$

and for some constant $C > 0$ which only depends on the choice of ψ .

None of the quantities involved in (18) depends on the chosen representation \mathcal{R} of \mathcal{Q} , which means that the performance of $\widehat{\mathbf{P}}$ does not depend on \mathcal{R} although its construction depends on it. We shall therefore (abusively) refer to $\widehat{\mathbf{P}}$ as a ρ -estimator on \mathcal{Q} omitting to mention what representation is used for its construction.

Introducing a nontrivial penalty function allows one to favour some probabilities as compared to others in \mathcal{Q} and gives thus a Bayesian flavour to our estimation procedure. We shall mainly use it when we have at our disposal not only one single ρ -model for \mathbf{P} but rather a countable collection $\{\mathcal{Q}_m, m \in \mathcal{M}\}$ of candidate ones, in which case $\mathcal{Q} = \bigcup_{m \in \mathcal{M}} \mathcal{Q}_m$ is still a ρ -model that we call the *reference ρ -model*. The penalty function may not only be used for estimating \mathbf{P} but also for performing model selection among the family $\{\mathcal{Q}_m, m \in \mathcal{M}\}$ by deciding that the procedure selects the ρ -model $\mathcal{Q}_{\widehat{m}}$ if the resulting estimator $\widehat{\mathbf{P}}$ belongs to $\mathcal{Q}_{\widehat{m}}$. Since $\widehat{\mathbf{P}}$ may belong to several ρ -models, this selection procedure may result in a (random) set of possible ρ -models for \mathbf{P} and a common way of selecting one is to choose that with the smallest *complexity* in a suitable sense. In the present paper, the complexity of a ρ -model \mathcal{Q}_m will be measured by means of a nonnegative weight function Δ mapping \mathcal{M} into \mathbb{R}_+ and which satisfies

$$(21) \quad \sum_{m \in \mathcal{M}} e^{-\Delta(m)} \leq 1,$$

where the number “1” is chosen for convenience. When equality holds in (21), $e^{-\Delta(\cdot)}$ can be viewed as a prior distribution on the family of ρ -models $\{\mathcal{Q}_m, m \in \mathcal{M}\}$.

In such a context, we shall describe how our penalty term should depend on this weight function Δ in view of selecting a suitable ρ -model for \mathbf{P} . The next theorem is proved in Section A.3 of the Supplementary Material [Baraud and Birgé (2018)].

THEOREM 2. *Let \mathbf{P} be an arbitrary distribution in \mathcal{P} , $\{\mathcal{Q}_m, m \in \mathcal{M}\}$ be a countable collection of ρ -models, Δ a weight function satisfying (21), $\mathcal{R}(\mathcal{Q})$ a representation of $\mathcal{Q} = \bigcup_{m \in \mathcal{M}} \mathcal{Q}_m$, ψ a function satisfying Assumption 2 and κ be given by (19). Assume that there exists a mapping $D_n : \mathcal{M} \rightarrow \mathbb{R}_+$ and a number $K \geq 0$ such that, whatever $m \in \mathcal{M}$,*

$$(22) \quad D^{\mathcal{Q}_m}(\mathbf{P}, \bar{\mathbf{P}}) \leq D_n(m) + K D_n(m') \quad \text{for all } (\mathbf{P}, \bar{\mathbf{P}}) \in \mathcal{P} \times \mathcal{Q}_{m'}.$$

Let the penalty function satisfy, for some constant $\kappa_1 \in \mathbb{R}$,

$$(23) \quad \text{pen}(\mathbf{Q}) = \kappa_1 + \kappa \inf_{\{m \in \mathcal{M} \mid \mathcal{Q}_m \ni \mathbf{Q}\}} \left[\frac{D_n(m)}{4.7} + \Delta(m) \right] \quad \text{for all } \mathbf{Q} \in \mathcal{Q}.$$

Then any ρ -estimator $\hat{\mathbf{P}}$ relative to $(\mathcal{R}(\mathcal{Q}), \text{pen})$ satisfies, for all $\xi > 0$ with probability at least $1 - e^{-\xi}$ and with γ given by (19),

$$(24) \quad \begin{aligned} \mathbf{h}^2(\mathbf{P}, \hat{\mathbf{P}}) &\leq \inf_{m \in \mathcal{M}} \left[\gamma \mathbf{h}^2(\mathbf{P}, \mathcal{Q}_m) + \frac{4\kappa}{a_1} \left(\frac{K+1}{4.7} D_n(m) + \Delta(m) \right) \right] \\ &\quad + \frac{4\kappa}{a_1} (1.49 + \xi). \end{aligned}$$

3.3. The case of density estimation. Of special interest is the situation where the X_i are assumed to be i.i.d. with values in a measurable set $(\mathcal{X}, \mathcal{B})$ in which case $\mathcal{X} = \mathcal{X}^n$, $\mathcal{B} = \mathcal{B}^{\otimes n}$, \mathcal{P} and \mathcal{M} denote respectively the set of all probability distributions and all positive σ -finite measures on $(\mathcal{X}, \mathcal{B})$ and \mathbf{P} is expected (although this is not necessarily true) to belong to $\mathcal{P}^n = \{P^{\otimes n}, P \in \mathcal{P}\}$. Note that the Hellinger distance $h(\cdot, \cdot)$ on \mathcal{P} is related to the Hellinger-type distance $\mathbf{h}(\cdot, \cdot)$ on \mathcal{P}^n in the following way:

$$\mathbf{h}^2(\mathbf{Q}, \mathbf{Q}') = n h^2(Q, Q') \quad \text{for all } Q, Q' \in \mathcal{P} \text{ with } \mathbf{Q} = Q^{\otimes n}, \mathbf{Q}' = (Q')^{\otimes n}.$$

If $\mathbf{P} = P^{\otimes n} \in \mathcal{P}^n$, estimating \mathbf{P} then amounts to estimating the marginal distribution P and we model the probability P rather than \mathbf{P} .

DEFINITION 5. We call *density ρ -model* any countable subset \mathcal{Q} of \mathcal{P} .

Given such a density ρ -model \mathcal{Q} for P with representation (μ, \mathcal{Q}) (which implies that the mapping $q \mapsto Q = q \cdot \mu$ is one to one), the corresponding ρ -model for \mathbf{P} is simply $\mathcal{Q} = \{\mathbf{Q} = Q^{\otimes n}, Q \in \mathcal{Q}\}$ with representation (μ, \mathcal{Q}) , $\mu = \mu^{\otimes n}$

and $\mathcal{Q} = \{\mathbf{q} : (x_1, \dots, x_n) \mapsto (q(x_1) \dots q(x_n)), q \in \mathcal{Q}\}$. In this case, for simplicity, we write $\mathbf{T}(X, q, q')$ and $\Upsilon(X, q)$ for $\mathbf{T}(X, \mathbf{q}, \mathbf{q}')$ and $\Upsilon(X, \mathbf{q})$, respectively. Examples involving density estimation will be considered in Sections 5, 6, 8 and 9 below.

We may also work with several density ρ -models $\{\mathcal{Q}_m, m \in \mathcal{M}\}$ for P simultaneously, in which case $\mathcal{Q} = \bigcup_{m \in \mathcal{M}} \mathcal{Q}_m$ is also a density ρ -model. A penalty function pen on \mathcal{Q} leads to a penalty function \mathbf{pen} on $\mathcal{Q} = \bigcup_{m \in \mathcal{M}} \mathcal{Q}_m$ defined by $\mathbf{pen}(\mathbf{Q}) = \mathbf{pen}(Q^{\otimes n}) = \text{pen}(Q)$ for all $Q \in \mathcal{Q}$. Any ρ -estimator $\hat{\mathbf{P}}$ relative to $((\mu, \mathcal{Q}), \mathbf{pen})$ is of the form $\hat{\mathbf{P}} = \hat{P}^{\otimes n}$ with $\hat{P} \in \text{Cl}(\mathcal{Q})$ and \hat{P} will be called a *density ρ -estimator* for P relative to $((\mu, \mathcal{Q}), \mathbf{pen})$.

We deduce that, under the assumptions of Theorem 1, if \mathbf{P} is truly of the form $\mathbf{P} = P^{\otimes n}$, for all $\bar{P} \in \mathcal{Q}$,

$$\mathbb{P} \left[Ch^2(P, \hat{P}) \leq h^2(P, \bar{P}) + \frac{D^{\mathcal{Q}}(\mathbf{P}, \bar{P}^{\otimes n})}{n} + \frac{\xi}{n} \right] \geq 1 - e^{-\xi} \quad \text{for all } \xi > 0.$$

Under the assumptions of Theorem 2, for all $\xi > 0$ and a positive constant C depending only on ψ ,

$$\mathbb{P} \left[Ch^2(P, \hat{P}) \leq \inf_{m \in \mathcal{M}} \left[h^2(P, \mathcal{Q}_m) + \frac{D_n(m) + \Delta(m)}{n} \right] + \frac{\xi}{n} \right] \geq 1 - e^{-\xi}.$$

4. From ρ -models to uncountable statistical models. The previous results apply to statistical models \mathcal{Q} that are countable, which is not the common case in statistics. The aim of this section is to explain how our general theory on ρ -models can be used to solve estimation problems on models that are possibly uncountable. Hereafter, we shall denote by \mathcal{Q} a *general statistical model*, that is, an arbitrary subset of \mathcal{P} .

4.1. *Working with nets.* Let us first recall this classical definition.

DEFINITION 6. Given $\eta \geq 0$, a subset \mathcal{Q} of $\overline{\mathcal{Q}}$ such that $\mathbf{h}(\mathbf{Q}, \mathcal{Q}) \leq \eta$ for all $\mathbf{Q} \in \overline{\mathcal{Q}}$ is called an η -net of $\overline{\mathcal{Q}}$. The case $\eta = 0$ corresponds to the situation where \mathcal{Q} is \mathcal{V} -dense in $\overline{\mathcal{Q}}$.

If there exists a countable η -net \mathcal{Q} for $\overline{\mathcal{Q}}$, it is a ρ -model. If its ρ -dimension $D^{\mathcal{Q}}$ is bounded by $D_n = D_n(\eta) \geq 1$ on $\mathcal{P} \times \mathcal{Q}$, we deduce from Theorem 1 and the inequality $\mathbf{h}(\mathbf{P}, \mathcal{Q}) \leq \mathbf{h}(\mathbf{P}, \overline{\mathcal{Q}}) + \eta$ that any ρ -estimator on \mathcal{Q} satisfies

$$(25) \quad \mathbb{P} \left[Ch^2(\mathbf{P}, \hat{\mathbf{P}}) \leq \mathbf{h}^2(\mathbf{P}, \overline{\mathcal{Q}}) + D_n(\eta) + \eta^2 + \xi \right] \geq 1 - e^{-\xi} \quad \text{for all } \xi > 0,$$

hence

$$(26) \quad \mathbb{E} \left[\mathbf{h}^2(\mathbf{P}, \hat{\mathbf{P}}) \right] \leq C' \left[\mathbf{h}^2(\mathbf{P}, \overline{\mathcal{Q}}) + D_n(\eta) + \eta^2 \right]$$

for some constants $C, C' > 0$ depending on ψ only. Most of the statistical models $\overline{\mathcal{Q}}$ that are used in statistics possess η -nets for all values of $\eta \geq 0$. Since the ρ -dimension function $D^{\mathcal{Q}}$ can only increase with inclusion, choosing for each $\eta \geq 0$ an η -net with the smallest possible cardinality and then the value η^* of η that minimizes $D_n(\eta) + \eta^2$ leads to a ρ -estimator $\widehat{\mathbf{P}}$ with the smallest possible risk bound in (26). This risk bound turns out to be minimax (up to possible extra logarithmic factors) in all cases we know; see Section 7.1.

4.2. *Models that are universally separable.* Following Pollard (1984), we shall say that a class of densities $\overline{\mathcal{Q}} \subset \mathcal{L}(\mu)$ is *universally separable* if one can find a countable subset $\mathcal{Q} \subset \overline{\mathcal{Q}}$ such that, for each $\mathbf{q} \in \overline{\mathcal{Q}}$, there exists a sequence $(\mathbf{q}^{(j)})_{j \geq 1}$ in \mathcal{Q} which converges towards \mathbf{q} pointwise, that is,

$$(27) \quad q_i^{(j)}(x) \xrightarrow{j \rightarrow +\infty} q_i(x) \quad \text{for } 1 \leq i \leq n \text{ and all } x \in \mathcal{X}_i.$$

We shall then say that \mathcal{Q} is \mathcal{T} -dense in $\overline{\mathcal{Q}}$. Note that if $(\mathbf{q}^{(j)})_{j \geq 1}$ converges towards \mathbf{q} pointwise, by Scheffé’s lemma, the sequence of probabilities $\mathbf{Q}_j = \mathbf{q}^{(j)} \cdot \mu$ converges in total variation, hence in Hellinger distance, towards $\mathbf{Q} = \mathbf{q} \cdot \mu$. This implies that if \mathcal{Q} is \mathcal{T} -dense in $\overline{\mathcal{Q}}$, the set of probabilities $\mathcal{Q} = \{\mathbf{q} \cdot \mu, \mathbf{q} \in \mathcal{Q}\}$ is \mathcal{V} -dense in $\overline{\mathcal{Q}} = \{\mathbf{q} \cdot \mu, \mathbf{q} \in \overline{\mathcal{Q}}\}$.

We shall work here within the following framework. For some $\mu \in \mathcal{M}$, let $\{\overline{\mathcal{Q}}_m, m \in \mathcal{M}\}$ be a countable family of universally separable subsets of $\mathcal{L}(\mu)$ with $\mathcal{Q}_m \subset \overline{\mathcal{Q}}_m$ a countable and \mathcal{T} -dense subset of $\overline{\mathcal{Q}}_m$. We set $\overline{\mathcal{Q}} = \bigcup_{m \in \mathcal{M}} \overline{\mathcal{Q}}_m$, $\mathcal{Q} = \bigcup_{m \in \mathcal{M}} \mathcal{Q}_m$, $\overline{\mathcal{Q}}_m = \{\mathbf{q} \cdot \mu, \mathbf{q} \in \overline{\mathcal{Q}}_m\}$ for all $m \in \mathcal{M}$, $\overline{\mathcal{Q}} = \{\mathbf{q} \cdot \mu, \mathbf{q} \in \overline{\mathcal{Q}}\}$ and $\mathcal{Q} = \{\mathbf{q} \cdot \mu, \mathbf{q} \in \mathcal{Q}\}$. Note that \mathcal{Q} is a ρ -model since \mathcal{Q} is countable and that \mathcal{Q} is \mathcal{V} -dense in $\overline{\mathcal{Q}}$ since \mathcal{Q} is \mathcal{T} -dense in $\overline{\mathcal{Q}}$. Let now $\overline{\text{pen}}$ be some penalty function on $\overline{\mathcal{Q}}$ with the following property.

ASSUMPTION 3. There exists a function $p : \mathcal{M} \rightarrow \mathbb{R}$ such that

$$(28) \quad \overline{\text{pen}}(\mathbf{Q}) = \inf_{m \in \mathcal{M}, \overline{\mathcal{Q}}_m \ni \mathbf{Q}} p(m) \quad \text{for all } \mathbf{Q} \in \overline{\mathcal{Q}}$$

and, for any $\mathbf{Q} \in \overline{\mathcal{Q}}$, there exists some $m_{\mathbf{Q}} \in \mathcal{M}$ such that $\mathbf{Q} \in \overline{\mathcal{Q}}_{m_{\mathbf{Q}}}$ and $\overline{\text{pen}}(\mathbf{Q}) = p(m_{\mathbf{Q}})$.

Note that this assumption holds in particular in the case of a single model with $\overline{\text{pen}} = \mathbf{0}$. Within this framework, we can prove the following result.

THEOREM 3. Let $\{\overline{\mathcal{Q}}_m, ; m \in \mathcal{M}\}$ be a countable family of universally separable subsets of $\mathcal{L}(\mu)$ and $\overline{\text{pen}}$ a penalty function on $\overline{\mathcal{Q}}$ that satisfies Assumption 3. Any ρ -estimator $\widehat{\mathbf{P}}$ on $\overline{\mathcal{Q}}$ relative to $((\mu, \overline{\mathcal{Q}}), \overline{\text{pen}})$ is also a ρ -estimator on the ρ -model \mathcal{Q} relative to $((\mu, \mathcal{Q}), \text{pen})$ where pen is the restriction of $\overline{\text{pen}}$ to \mathcal{Q} .

The proof is postponed to Section A.5 of the Supplementary Material [Baraud and Birgé (2018)].

This result says that, provided that the penalty function satisfies (28), which is consistent with (23), the construction of a ρ -estimator on the possibly uncountable set $\overline{\mathcal{Q}}$ with representation $(\mu, \overline{\mathcal{Q}})$ actually results in a ρ -estimator based on the ρ -model \mathcal{Q} .

As soon as we can control the ρ -dimension function of \mathcal{Q} by some features of $\overline{\mathcal{Q}}$, in the case of a single model, or the ρ -dimension functions of the ρ -models $\mathcal{Q}_m = \{\mathbf{q} \cdot \mu, \mathbf{q} \in \mathcal{Q}_m\}$ by the features of the models $\overline{\mathcal{Q}}_m$, in the general case, we are able to bound the risk of the ρ -estimator relative to $((\mu, \overline{\mathcal{Q}}), \overline{\mathbf{pen}})$ using the results of Theorems 1 and 2.

For illustration, let us mention a few examples of density sets that are universally separable:

- (a) the set $\overline{\mathcal{H}}_D$ of right-continuous histograms on \mathbb{R} with at most $D \geq 1$ pieces;
- (b) for $L > 0$ and $\alpha = r + \beta$ with $r \in \mathbb{N}, \beta \in (0, 1]$, the set $\overline{\mathcal{H}}_\alpha(L)$ of functions f on $[0, 1]$ that are r -times differentiable and satisfy

$$|f^{(r)}(x) - f^{(r)}(y)| \leq L|x - y|^\beta \quad \text{for all } x, y \in [0, 1];$$

- (c) the set $\overline{\mathcal{H}}_\downarrow$ of nonincreasing and right-continuous densities on $(0, +\infty)$.

The set $\overline{\mathcal{H}}_\alpha(L)$ is universally separable because the larger set consisting of continuous functions on $[0, 1]$ is separable for the topology induced by the norm of the uniform convergence; hence all its subsets are separable with respect to this topology which implies pointwise convergence. We prove that the sets $\overline{\mathcal{H}}_D$ and $\overline{\mathcal{H}}_\downarrow$ are universally separable in Section B of the Supplementary Material [Baraud and Birgé (2018)]. We shall see in Section 6 that the MLE on the convex density sets $\overline{\mathcal{H}}_\alpha(L)$ and $\overline{\mathcal{H}}_\downarrow$ is actually a ρ -estimator.

5. Why is a ρ -estimator robust? The aim of this section is to analyse the robustness properties of ρ -estimators. For the sake of simplicity, we shall restrict ourselves to the particular case of density estimation as described in Section 3.3.

5.1. *Misspecification and contamination.* We assume here that we work with a single ρ -model \mathcal{Q} (so that Theorem 1 applies) for which $D^\mathcal{Q}(\mathbf{P}, \overline{\mathbf{P}})$ is bounded from above independently of $\mathbf{P} \in \mathcal{P}$ and $\overline{\mathbf{P}} \in \mathcal{Q}$ by some number $D_n(\mathcal{Q}) \geq 1$ depending on the marginal model \mathcal{Q} and the number n of marginals. Examples of such situations will be provided in Section 7.

When $\mathbf{P} = P^{\otimes n}$, that is when the data are truly i.i.d. with marginal distribution P , (18) becomes

$$(29) \quad \mathbb{P}\left[Ch^2(P, \widehat{P}) \leq h^2(P, \mathcal{Q}) + n^{-1}[D_n(\mathcal{Q}) + \xi]\right] \geq 1 - e^{-\xi} \quad \text{for all } \xi > 0,$$

where C is a positive constant only depending on ψ .

The bias term in (29), namely $h^2(P, \mathcal{Q})$, accounts for the robustness property of the ρ -estimator with respect to the Hellinger distance and measures the additional loss we get as compared to the case when P belongs to \mathcal{Q} . If this quantity is small, the performance of the ρ -estimator will not deteriorate too much as compared to the ideal situation where P does belong to \mathcal{Q} . In fact, if there exists some probability $\bar{P} \in \mathcal{Q}$ such that $h^2(P, \mathcal{Q}) = h^2(P, \bar{P})$ is small as compared to $D_n(\mathcal{Q})/n$, everything is almost as if the ρ -estimator \hat{P} were built from an i.i.d. sample with distribution \bar{P} . The ρ -estimators under P and \bar{P} would therefore look the same. This includes the following situations.

Misspecification. The true distribution P of the observations does not belong to \mathcal{Q} but is close to \mathcal{Q} . For example, let \mathcal{Q} be countable and \mathcal{V} -dense in the set of all Gaussian distributions on \mathbb{R}^k with identity covariance matrix and mean vector belonging to a linear subspace $\bar{S} \subset \mathbb{R}^k$. Assume that the true distribution P has the same form except for the fact that its mean does not belong to \bar{S} but is at Euclidean distance $\varepsilon > 0$ from \bar{S} . Then it follows from classical formulas that

$$h^2(P, \mathcal{Q}) = 1 - e^{-\varepsilon^2/8} \leq \varepsilon^2/8.$$

Contamination. The true distribution P is of the form $(1 - \varepsilon)\bar{P} + \varepsilon R$ with $\bar{P} \in \mathcal{Q}$ and $R \neq \bar{P}$ but otherwise arbitrary. This situation arises when a proportion $\varepsilon \in (0, 1)$ of the sample X_1, \dots, X_n is contaminated by another sample. It follows from the convexity property of the Hellinger distance that

$$h^2(P, \mathcal{Q}) \leq h^2(P, \bar{P}) \leq \varepsilon h^2(R, \bar{P}) \leq \varepsilon,$$

and this bound holds whatever the contaminating distribution R . From a more practical point of view, one can see the contaminated case as follows: for each i , one decides between no contamination with a probability $1 - \varepsilon$ and contamination with a probability ε and draws X_i accordingly with distribution either \bar{P} or R . If it were possible to extract from the sample X_1, \dots, X_n these N data, with $N \sim \mathcal{B}(n, 1 - \varepsilon)$, which are really distributed according to the distribution $\bar{P} \in \mathcal{Q}$, we would build a ρ -estimator \tilde{P} on these data. The robustness property ensures that the ρ -estimator \hat{P} based on the whole data set remains close to \tilde{P} . Everything works almost as if the ρ -estimator \hat{P} only considered the noncontaminated subsample and ignored the other data, at least when ε is small enough.

5.2. More robustness. There is an additional aspect of robustness that is not apparent in (29). Our general result about the performance of ρ -estimators, as stated in (18), actually allows that our observations be independent but not necessarily i.i.d., in which case the joint distribution \mathbf{P} of (X_1, \dots, X_n) is actually of the form $\otimes_{i=1}^n P_i$ but not necessarily of the form $P^{\otimes n}$. Of course we do not know whether \mathbf{P} is of the first form or the second and, proceeding as if X_1, \dots, X_n were

i.i.d., we build a ρ -estimator $\widehat{P} \in \text{Cl}(\mathcal{Q})$ of the presumed common density P and make a mistake which is no longer $h^2(P, \widehat{P})$ but

$$\frac{1}{n} \mathbf{h}^2(\mathbf{P}, \widehat{\mathbf{P}}) \quad \text{with } \widehat{\mathbf{P}} = \widehat{P}^{\otimes n} \quad \text{and} \quad \mathbf{h}^2(\mathbf{P}, \widehat{\mathbf{P}}) = \sum_{i=1}^n h^2(P_i, \widehat{P}),$$

which is consistent with the i.i.d. case $P_i = P$ for all i . In this context, we actually get the following analogue of (29): for all $\xi > 0$,

$$(30) \quad \mathbb{P} \left[\frac{C}{n} \mathbf{h}^2(\mathbf{P}, \widehat{\mathbf{P}}) \leq \inf_{Q \in \mathcal{Q}} \left(\frac{1}{n} \sum_{i=1}^n h^2(P_i, Q) \right) + \frac{D_n(\mathcal{Q}) + \xi}{n} \right] \geq 1 - e^{-\xi}.$$

This allows many more possibilities of deviations between \mathbf{P} and the statistical model $\{Q^{\otimes n}, Q \in \mathcal{Q}\}$. For instance, we may have $h(P_i, \bar{P}) \leq \varepsilon$ for some $\bar{P} \in \mathcal{Q}$ and all i , $P_i \neq P_{i'}$ for all $i \neq i'$, and nevertheless

$$\inf_{Q \in \mathcal{Q}} \left(\frac{1}{n} \sum_{i=1}^n h^2(P_i, Q) \right) \leq \varepsilon^2.$$

An alternative situation corresponds to a small number of “outliers”, namely, $P_i = P$ except on a subset $J \subset \{1, \dots, n\}$ of indices of small cardinality and, for $i \in J$, P_i is completely arbitrary, for instance a Dirac measure. In such a case, for any probability Q ,

$$\left(1 - \frac{|J|}{n} \right) h^2(P, Q) \leq \frac{1}{n} \sum_{i=1}^n h^2(P_i, Q) \leq \left(1 - \frac{|J|}{n} \right) h^2(P, Q) + \frac{|J|}{n},$$

and we deduce from (30) that, on a set of probability at least $1 - e^{-\xi}$,

$$\begin{aligned} \frac{C(n - |J|)}{n} h^2(P, \widehat{P}) \leq C \frac{\mathbf{h}^2(\mathbf{P}, \widehat{\mathbf{P}})}{n} &\leq \left[\left(\frac{n - |J|}{n} \right) h^2(P, \mathcal{Q}) + \frac{|J|}{n} \right] \\ &+ \frac{D_n(\mathcal{Q}) + \xi}{n}. \end{aligned}$$

Finally,

$$\mathbb{P} \left[C h^2(P, \widehat{P}) \leq h^2(P, \mathcal{Q}) + \frac{|J| + D_n(\mathcal{Q}) + \xi}{n - |J|} \right] \geq 1 - e^{-\xi} \quad \text{for all } \xi > 0.$$

When $|J|/n$ is small enough, this bound appears to be a slight modification of what we would get from (29) if \mathbf{P} were of the form $P^{\otimes n}$. This means that the ρ -estimator \widehat{P} is also robust with respect to a possible departure from the assumption that the X_i are i.i.d.

6. The ρ -estimators and the MLE. As mentioned in the [Introduction](#), there are some deep connexions between the MLE and ρ -estimators which are mostly due to the similarities in the neighbourhood of 1 between the logarithm and the functions ψ of Proposition 3. A nice result in this direction was communicated to the authors by Weijie Su in October 2016. It concerns the case of density estimation, as described in Section 3.3 with a single density model $\overline{\mathcal{Q}} = \{q \cdot \mu, q \in \overline{\mathcal{Q}}\}$ where $\overline{\mathcal{Q}}$ is universally separable as defined in Section 4.2.

ASSUMPTION 4. The function $x \mapsto \varphi(x) = \psi(\sqrt{x})$, where ψ is the function used to define the statistic \mathbf{T} in (5), satisfies $\varphi(1) = 0$, is concave and admits a positive derivative at 1.

PROPOSITION 5 [Weijie Su, private communication (2016)]. *Let Assumption 4 hold, $\overline{\mathcal{Q}}$ be a convex set of densities on the measured space $(\mathcal{X}, \mathcal{B}, \mu)$ and the likelihood be not identically equal to 0 on $\overline{\mathcal{Q}}$. The maximum likelihood estimator $\widehat{Q} = \widehat{q} \cdot \mu$ on the density model $\overline{\mathcal{Q}} = \{q \cdot \mu, q \in \overline{\mathcal{Q}}\}$, when it exists, satisfies*

$$\Upsilon(\mathbf{X}, \widehat{q}) = \sup_{q' \in \overline{\mathcal{Q}}} \mathbf{T}(\mathbf{X}, \widehat{q}, q') = 0 = \inf_{q \in \overline{\mathcal{Q}}} \sup_{q' \in \overline{\mathcal{Q}}} \mathbf{T}(\mathbf{X}, q, q') = \inf_{q \in \overline{\mathcal{Q}}} \Upsilon(\mathbf{X}, q)$$

and is therefore a ρ -estimator relative to $(\overline{\mathcal{Q}}, 0)$.

PROOF. Given the data X_1, \dots, X_n , if the maximum likelihood \widehat{q} exists, it is unique since the logarithm is strictly concave. Moreover, $\widehat{q}(X_i) > 0$ for all $i \in \{1, \dots, n\}$. Since $\Upsilon(\mathbf{X}, \widehat{q}) \geq \mathbf{T}(\mathbf{X}, \widehat{q}, \widehat{q}) = 0$, it suffices to prove that

$$L(q) = \mathbf{T}(\mathbf{x}, \widehat{q}, q) = \sum_{i=1}^n \varphi\left(\frac{q(X_i)}{\widehat{q}(X_i)}\right) \leq 0 \quad \text{for all } q \in \overline{\mathcal{Q}}.$$

For $q \in \overline{\mathcal{Q}}$ and $\varepsilon \in [0, 1]$, $(1 - \varepsilon)\widehat{q} + \varepsilon q \in \overline{\mathcal{Q}}$ and, when $\varepsilon \rightarrow 0$,

$$\begin{aligned} L((1 - \varepsilon)\widehat{q} + \varepsilon q) &= n\varphi(1) + \varepsilon \left[\varphi'(1) \sum_{i=1}^n \frac{q(X_i)}{\widehat{q}(X_i)} + o(1) \right] \\ (31) \qquad \qquad \qquad &= \varepsilon \left[\varphi'(1) \sum_{i=1}^n \frac{q(X_i)}{\widehat{q}(X_i)} + o(1) \right] \end{aligned}$$

since $\varphi(1) = 0$. When φ is the logarithm and $\varepsilon > 0$, the right-hand side of (31) is negative since \widehat{q} is the unique MLE. Letting ε go to 0 we derive that

$$(32) \qquad \sum_{i=1}^n \frac{q(X_i)}{\widehat{q}(X_i)} \leq 0 \quad \text{for all } q \in \overline{\mathcal{Q}}.$$

Moreover, the concavity of φ implies that for all $\varepsilon \in [0, 1]$

$$\varphi\left(\frac{(1 - \varepsilon)\widehat{q}(X_i) + \varepsilon q(X_i)}{\widehat{q}(X_i)}\right) \leq \varphi(1) + \varepsilon \frac{q(X_i)}{\widehat{q}(X_i)}\varphi'(1) = \varepsilon\varphi'(1) \frac{q(X_i)}{\widehat{q}(X_i)}$$

so that, for all $q \in \overline{\mathcal{Q}}$, $L((1 - \varepsilon)\widehat{q} + \varepsilon q) \leq \varepsilon\varphi'(1) \sum_{i=1}^n q(X_i)/\widehat{q}(X_i)$ and

$$\mathbf{T}(\mathbf{x}, \widehat{q}, q) = L(q) \leq \varphi'(1) \sum_{i=1}^n \frac{q(X_i)}{\widehat{q}(X_i)} \leq 0 \quad \text{for all } q \in \overline{\mathcal{Q}}$$

by (32), which completes the proof. \square

Note that both functions ψ_1 and ψ_2 of Proposition 3 satisfy Assumption 4.

We may now derive the following relationship between the MLE and ρ -estimators, the proof of which immediately follows from Theorem 3 and Su’s proposition.

COROLLARY 1. *Let $\overline{\mathcal{Q}}$ be a convex set of densities on the measured space $(\mathcal{X}, \mathcal{B}, \mu)$ which is universally separable on \mathcal{X} with countable and \mathcal{T} -dense subset \mathcal{Q} and ψ satisfy Assumptions 1 and 4. The maximum likelihood estimator $\widehat{Q} = \widehat{q} \cdot \mu$ on the density model $\mathcal{Q} = \{q \cdot \mu, q \in \overline{\mathcal{Q}}\}$, when it exists, is a ρ -estimator on the ρ -density model $\mathcal{Q} = \{q \cdot \mu, q \in \mathcal{Q}\}$.*

For illustration, the set $\overline{\mathcal{Q}} = \overline{\mathcal{H}}_{\mathcal{I}}$ of right-continuous histograms based on a fixed partition \mathcal{I} of $[0, 1)$ into $D \geq 1$ intervals is convex and obviously universally separable. The usual histogram \widehat{p} based on \mathcal{I} , which corresponds to the MLE on $\overline{\mathcal{H}}_{\mathcal{I}}$, can be viewed as a ρ -estimator on a countable subset of $\overline{\mathcal{H}}_{\mathcal{I}}$. Taking back some of the examples of convex and universally separable density sets given in Section 4.2, we deduce that the MLE on $\overline{\mathcal{H}}_{\downarrow}$, that is, the Grenander estimator, or on the set $\overline{\mathcal{H}}_{\alpha}(L)$ are also ρ -estimators.

7. Bounding the ρ -dimension function of a ρ -model with applications to the risk of ρ -estimators. It clearly follows from the results of Section 3 that bounding the risk of ρ -estimators amounts to bounding the ρ -dimension of ρ -models which we shall now do under various assumptions. Throughout this section, we fix the function ψ satisfying Assumption 2 (typically ψ_1 or ψ_2) and when we shall say that some quantity depends on ψ , this will mean that it actually depends on a_1 and a_2 .

In view of (20), of special interest is the situation where the ρ -dimension function $(\mathbf{P}, \overline{\mathbf{P}}) \mapsto D^{\mathcal{Q}}(\mathbf{P}, \overline{\mathbf{P}})$ of the ρ -model \mathcal{Q} is uniformly bounded from above on $\mathcal{P} \times \mathcal{Q}$ by some constant $D_n \geq 1$. Let us begin by a few elementary considerations. If one can find a representation $\mathcal{R} = (\mu, \mathcal{Q})$ of $\mathcal{Q} \cup \{\overline{\mathbf{P}}\}$ such that $w(\mathcal{R}, \mathcal{Q}, \mathbf{P}, \overline{\mathbf{P}}, y) \leq a_1 y^2/8$ for all $y \geq \beta^{-1}\sqrt{D}$, we immediately derive from the definition of $D^{\mathcal{Q}}$ that

$$(33) \quad D^{\mathcal{Q}}(\mathbf{P}, \overline{\mathbf{P}}) \leq D \vee 1 \quad \text{for } (\mathbf{P}, \overline{\mathbf{P}}) \in \mathcal{P}^2.$$

In particular, since $|\psi| \leq 1$, the expectation in (14) is never larger than $2n$ so that $w(\mathcal{R}, \mathcal{Q}, \mathbf{P}, \bar{\mathbf{P}}, y) \leq a_1 y^2 / 8$ for $y \geq 4\sqrt{(n/a_1)}$ and (33) always holds with

$$\sqrt{D} = 4\beta\sqrt{(n/a_1)} = \sqrt{na_1}/a_2 \quad \text{or equivalently} \quad D = na_1/a_2^2 \leq n/6.$$

Finally, whatever the choices of \mathcal{Q} and ψ ,

$$(34) \quad D^{\mathcal{Q}}(\mathbf{P}, \bar{\mathbf{P}}) \leq n/6 \quad \text{for all } (\mathbf{P}, \bar{\mathbf{P}}) \in \mathcal{P}^2.$$

More precise bounds will now be given that depend on some specific features of \mathcal{Q} .

7.1. *The finite case.* Given a finite subset $\mathcal{Q} \subset \mathcal{P}$, let us set

$$(35) \quad \mathcal{H}(\mathcal{Q}, y) = \sup_{\mathbf{P} \in \mathcal{P}} \log_+(2|\mathcal{Q} \cap \mathcal{B}(\mathbf{P}, y)|) \quad \text{for all } y > 0$$

and, for $x_0 = \sqrt{2}[\sqrt{1 + (\beta/a_2)} + 1]$,

$$(36) \quad \bar{\eta} = \sup \left\{ z > 0, \sqrt{\mathcal{H}(\mathcal{Q}, z/\beta)} > z/x_0 \right\}.$$

Since \mathcal{Q} is finite, the function $y \mapsto \mathcal{H}(\mathcal{Q}, y)$ is bounded by $\log(2|\mathcal{Q}|)$ and since $\beta/a_2 = a_1/(4a_2^2) \leq 1/24$,

$$(37) \quad \bar{\eta} \leq x_0 \sqrt{\log(2|\mathcal{Q}|)} < 3\sqrt{\log(2|\mathcal{Q}|)}.$$

PROPOSITION 6. *If \mathcal{Q} is a finite subset of \mathcal{P} and $\bar{\eta}$ is defined by (36),*

$$D^{\mathcal{Q}}(\mathbf{P}, \bar{\mathbf{P}}) \leq D_n(\mathcal{Q}) = \bar{\eta}^2 \vee 1 < 9 \log(2|\mathcal{Q}|) \quad \text{for all } (\mathbf{P}, \bar{\mathbf{P}}) \in \mathcal{P}^2.$$

The proof of this result is given in Section D.6 of the of the Supplementary Material [Baraud and Birgé (2018)]. The first upper bound $\bar{\eta}^2 \vee 1$ for $D^{\mathcal{Q}}(\mathbf{P}, \bar{\mathbf{P}})$ neither depends on \mathbf{P} nor on $\bar{\mathbf{P}}$ but might depend on β . The second bound only depends on the cardinality of \mathcal{Q} and, therefore, holds whatever ψ .

If a model $\bar{\mathcal{Q}}$ is a totally bounded subset of the metric space $(\mathcal{P}, \mathbf{h})$ and $\eta > 0$, one can cover $\bar{\mathcal{Q}}$ by a finite number of closed balls of radius η and the set $\mathcal{Q}[\eta]$ of their centers is an η -net for $\bar{\mathcal{Q}}$ (see Definition 6), which means that $\mathbf{h}(\mathbf{Q}, \mathcal{Q}[\eta]) \leq \eta$ for all $\mathbf{Q} \in \bar{\mathcal{Q}}$. The function $y \mapsto \mathcal{H}(\mathcal{Q}[\eta], y)$ measures in a sense the massiveness of $\mathcal{Q}[\eta]$ and turns out to be a useful tool to measure that of $\bar{\mathcal{Q}}$. We shall in particular use the following classical notions of dimension based on the metric structure of $\bar{\mathcal{Q}}$.

DEFINITION 7. A set $\bar{\mathcal{Q}} \subset \mathcal{P}$ is said to have a metric dimension bounded by \tilde{D} , where \tilde{D} is a right-continuous function from $(0, +\infty)$ to $[1/2, +\infty]$, if, for any positive η , there exists an η -net $\mathcal{Q}[\eta]$ for $\bar{\mathcal{Q}}$ which satisfies

$$(38) \quad \mathcal{H}(\mathcal{Q}[\eta], y) \leq (y/\eta)^2 \tilde{D}_n(\eta) \quad \text{for all } y \geq 2\eta.$$

We shall say that $\overline{\mathcal{Q}}$ has an entropy dimension bounded by $V \geq 0$ if, for any $\eta > 0$, there exists an η -net $\mathcal{Q}[\eta]$ of $\overline{\mathcal{Q}}$ such that

$$(39) \quad \mathcal{H}(\mathcal{Q}[\eta], y) \leq V \log(y/\eta) \quad \text{for all } y \geq 2\eta.$$

For the sake of convenience, we have slightly modified the original definition of the metric dimension due to Birgé [(2006), Definition 6, page 293] which is actually obtained by replacing the left-hand side of (38) by $\mathcal{H}(\mathcal{Q}[\eta], y) - \log 2$. Since in both definitions the metric dimension is not smaller than $1/2$, it is easy to check that, if $\overline{\mathcal{Q}}$ has a metric dimension bounded by D_M in Birgé’s sense, it has a metric dimension bounded by $\tilde{D} = (1 + (\log 2)/2)D_M$ in our sense and, conversely, if $\overline{\mathcal{Q}}$ has a metric dimension bounded by \tilde{D} in our sense, it also has a metric dimension bounded by \tilde{D} in Birgé’s sense. Hence, changing \tilde{D} into D_M only changes the numerical constants.

The logarithm being a slowly varying function, it is not difficult to see that the notion of metric dimension is more general than the entropy one in the sense that if $\overline{\mathcal{Q}}$ has an entropy dimension bounded by V , then it also has a metric dimension bounded by $\tilde{D}_n(\cdot)$ with

$$(40) \quad \tilde{D}_n(\eta) \leq (1/2) \vee [V(\log 2)/4] \quad \text{for all } \eta > 0.$$

If $\overline{\mathcal{Q}}$ has a metric dimension bounded by \tilde{D} and if η is a positive number satisfying

$$(41) \quad \tilde{D}_n(\eta) \leq (\beta\eta/x_0)^2,$$

with x_0 given by (36), we deduce from (38) that there exists an η -net $\mathcal{Q}[\eta]$ for $\overline{\mathcal{Q}}$ for which

$$\sqrt{\mathcal{H}(\mathcal{Q}[\eta], z/\beta)} \leq z/x_0 \quad \text{for all } z \geq 2\beta\eta.$$

It then follows that $\bar{\eta}$, as defined in (36), satisfies $\bar{\eta} \leq 2\beta\eta$ and we deduce from Proposition 6 that the ρ -dimension function $D^{\mathcal{Q}}$ of $\mathcal{Q} = \mathcal{Q}[\eta] \subset \overline{\mathcal{Q}}$ satisfies

$$(42) \quad D^{\mathcal{Q}}(\mathbf{P}, \bar{\mathbf{P}}) \leq D_n(\mathcal{Q}) = (2\beta\eta)^2 \vee 1 \quad \text{for all } (\mathbf{P}, \bar{\mathbf{P}}) \in \mathcal{P}^2.$$

If, in particular, $\overline{\mathcal{Q}}$ has an entropy dimension bounded by $V \geq 0$ we deduce from (40) that (41) holds for

$$(43) \quad \eta^2 = \frac{x_0^2}{2\beta^2} \left(1 \vee \frac{V \log 2}{2} \right) < \frac{9}{2\beta^2} \left(1 \vee \frac{V \log 2}{2} \right)$$

and we derive from (42) that

$$(44) \quad D^{\mathcal{Q}}(\mathbf{P}, \bar{\mathbf{P}}) \leq D_n(\mathcal{Q}) = 18 \left(1 \vee \frac{V \log 2}{2} \right) \quad \text{for all } (\mathbf{P}, \bar{\mathbf{P}}) \in \mathcal{P}^2.$$

Since in both cases $\mathbf{h}(\mathbf{P}, \mathcal{Q}) \leq \mathbf{h}(\mathbf{P}, \overline{\mathcal{Q}}) + \eta$ for all $\mathbf{P} \in \mathcal{P}$ because \mathcal{Q} is a η -net for $\overline{\mathcal{Q}}$, we obtain from (42), (44) and (25) the following result.

COROLLARY 2. *Let ψ be a function satisfying Assumption 2.*

(i) *If $\overline{\mathcal{Q}}$ has a metric dimension bounded by \tilde{D} and η satisfies (41), any ρ -estimator $\widehat{\mathbf{P}}$ based on a suitable η -net \mathcal{Q} for $\overline{\mathcal{Q}}$ satisfies for all $\mathbf{P} \in \mathcal{P}$ and $\xi > 0$:*

$$(45) \quad \mathbb{P}\left[\text{Ch}^2(\mathbf{P}, \widehat{\mathbf{P}}) \leq \mathbf{h}^2(\mathbf{P}, \overline{\mathcal{Q}}) + (\eta^2 \vee 1) + \xi\right] \geq 1 - e^{-\xi}.$$

(ii) *If $\overline{\mathcal{Q}}$ has an entropy dimension bounded by V and η satisfies (43), any ρ -estimator $\widehat{\mathbf{P}}$ based on a suitable η -net \mathcal{Q} for $\overline{\mathcal{Q}}$ satisfies for all $\mathbf{P} \in \mathcal{P}$ and $\xi > 0$:*

$$(46) \quad \mathbb{P}\left[\text{Ch}^2(\mathbf{P}, \widehat{\mathbf{P}}) \leq \mathbf{h}^2(\mathbf{P}, \overline{\mathcal{Q}}) + (V \vee 1) + \xi\right] \geq 1 - e^{-\xi}.$$

In both cases, C is a constant depending only on the choice of ψ .

7.2. *Bounds based on the VC-index.* In this section, we investigate the case of a model $\overline{\mathcal{Q}}$ given by a specific representation $(\boldsymbol{\mu}, \overline{\mathcal{Q}})$ where the density set $\overline{\mathcal{Q}}$ is possibly uncountable but satisfies some property to be described below.

A density $\mathbf{q} = (q_1, \dots, q_n) \in \mathcal{L}(\boldsymbol{\mu})$ can be viewed as a function on $\overline{\mathcal{X}} = \bigcup_{i=1}^n (\{i\} \times \mathcal{X}_i)$ defined, for $\bar{x} = (i, x)$ with $x \in \mathcal{X}_i$, by $\mathbf{q}(i, x) = q_i(x)$ so that a subset $\overline{\mathcal{Q}} \subset \mathcal{L}(\boldsymbol{\mu})$ is now viewed as a class of real-valued functions on $\overline{\mathcal{X}}$. A common notion of dimension for the class $\overline{\mathcal{Q}}$ is the following one.

DEFINITION 8. A class \mathcal{F} of functions from a set \mathcal{X} with values in $(-\infty, +\infty]$ is VC-subgraph with index \overline{V} (or equivalently with dimension $\overline{V} - 1 \geq 0$) if the class of subgraphs $\{(x, u) \in \mathcal{X} \times \mathbb{R} \mid f(x) > u\}$ as f varies in \mathcal{F} is a VC-class of sets in $\mathcal{X} \times \mathbb{R}$ with index \overline{V} (or dimension $\overline{V} - 1$).

We recall that, by definition, the index \overline{V} of a VC-class is a positive integer, hence its dimension $\overline{V} - 1 \in \mathbb{N}$. For additional information about VC-classes and related notions, we refer to van der Vaart and Wellner (1996) and Baraud, Birgé and Sart (2017), Section 8.

PROPOSITION 7. *Let ψ satisfy Assumption 2 and $\overline{\mathcal{Q}} \subset \mathcal{L}(\boldsymbol{\mu})$ be a VC-subgraph class of densities on $\overline{\mathcal{X}}$ with index not larger than \overline{V} . For any ρ -model $\mathcal{Q} \subset \overline{\mathcal{Q}} = \{\mathbf{q} \cdot \boldsymbol{\mu}, \mathbf{q} \in \overline{\mathcal{Q}}\}$, for all $(\mathbf{P}, \overline{\mathbf{P}}) \in \mathcal{P} \times \mathcal{P}^\mu$*

$$D^{\mathcal{Q}}(\mathbf{P}, \overline{\mathbf{P}}) \leq D_n(\overline{\mathcal{Q}}) = C_1(\overline{V} \wedge n)[1 + \log_+(n/\overline{V})],$$

where C_1 is a universal constant.

The proof is given in Section D.7 of the Supplementary Material [Baraud and Birgé (2018)]. A nice feature of this bound lies in the fact that it neither depends on the choices of ψ nor on the cardinality of \mathcal{Q} which can therefore be arbitrarily large. In particular, when \mathcal{Q} is \mathcal{V} -dense in $\overline{\mathcal{Q}}$ we deduce the following result from Proposition 7, (25) (with $\eta = 0$) and Theorem 3.

COROLLARY 3. *Let ψ be a function satisfying Assumption 2 and $\overline{\mathcal{Q}} \subset \mathcal{L}(\mu)$ a VC-subgraph class of densities on $\overline{\mathcal{X}}$ with index \overline{V} . Any ρ -estimator $\widehat{\mathbf{P}}$ built on a countable and \mathcal{V} -dense subset \mathcal{Q} of $\overline{\mathcal{Q}} = \{\mathbf{q} \cdot \mu, \mathbf{q} \in \mathcal{Q}\}$ satisfies, for all $\mathbf{P} \in \mathcal{P}$ and $\xi > 0$,*

$$(47) \quad \mathbb{P}\left[\mathbf{Ch}^2(\mathbf{P}, \widehat{\mathbf{P}}) \leq \mathbf{h}^2(\mathbf{P}, \overline{\mathcal{Q}}) + (\overline{V} \wedge n)[1 + \log_+(n/\overline{V})] + \xi\right] \geq 1 - e^{-\xi},$$

where the constant C only depends on the choice of ψ . If, moreover, $\overline{\mathcal{Q}}$ is universally separable, then (47) holds for any ρ -estimator relative to $((\mu, \overline{\mathcal{Q}}), \mathbf{0})$.

In the particular case of density estimation, the following result is useful in view of applying Proposition 7.

PROPOSITION 8. *If $\overline{\mathcal{Q}} \subset \mathcal{L}(\mu)$ is VC-subgraph on $\overline{\mathcal{X}}$ with index \overline{V} , then the set $\overline{\mathcal{Q}} = \{\mathbf{q} = (q, \dots, q), q \in \overline{\mathcal{Q}}\} \subset \mathcal{L}(\mu)$ is VC-subgraph on $\overline{\mathcal{X}}$ with an index not larger than \overline{V} .*

PROOF. If the class of subgraphs $\{(\overline{x}, u) \in \overline{\mathcal{X}} \times \mathbb{R} \mid \mathbf{q}(\overline{x}) > u\}$, with \mathbf{q} running in $\overline{\mathcal{Q}}$, shatters the subset $\{(\overline{x}_1, u_1), \dots, (\overline{x}_k, u_k)\}$ of $\overline{\mathcal{X}} \times \mathbb{R}$, then whatever $J \subset \{1, \dots, k\}$, there exists $\mathbf{q} \in \overline{\mathcal{Q}}$ such that $j \in J$ is equivalent to $\mathbf{q}(\overline{x}_j) = q(x_j) > u_j$. Hence, the class of subgraphs $\{(x, u) \in \mathcal{X} \times \mathbb{R} \mid q(x) > u\}$ with q running in $\overline{\mathcal{Q}}$ shatters the subset $\{(x_1, u_1), \dots, (x_k, u_k)\}$ of $\mathcal{X} \times \mathbb{R}$ and, therefore, $k + 1 \leq \overline{V}$. \square

7.3. *More bounds.* As we have observed in the previous sections [see (45), (46) and (47)], there are various situations for which, given a model $\overline{\mathcal{Q}} \subset \mathcal{P}$, it is possible to build a ρ -estimator $\widehat{\mathbf{P}}$ with values in $\overline{\mathcal{Q}}$ satisfying for all $\mathbf{P} \in \mathcal{P}$ and $\xi > 0$,

$$\mathbb{P}\left[\mathbf{Ch}^2(\mathbf{P}, \widehat{\mathbf{P}}) \leq \mathbf{h}^2(\mathbf{P}, \overline{\mathcal{Q}}) + D_n(\overline{\mathcal{Q}}) + \xi\right] \geq 1 - e^{-\xi}$$

for some quantity $D_n(\overline{\mathcal{Q}}) \geq 1$ only depending on the specific features of $\overline{\mathcal{Q}}$ and some constant $C > 0$ depending on the choice of ψ . Such an inequality leads to a risk bound of the following form (with $C' > 0$):

$$(48) \quad \mathbb{E}\left[\mathbf{h}^2(\mathbf{P}, \widehat{\mathbf{P}})\right] \leq C'\left[\mathbf{h}^2(\mathbf{P}, \overline{\mathcal{Q}}) + D_n(\overline{\mathcal{Q}})\right] \quad \text{for all } \mathbf{P} \in \mathcal{P}$$

and allows us to bound from above the minimax risk over $\overline{\mathcal{Q}}$ by $C'D_n(\overline{\mathcal{Q}})$.

However, not all statistical models admit a finite minimax risk and for such models there is consequently no hope to bound from above the ρ -dimension function uniformly as we did in the previous sections. One such example is the set of probabilities on $(0, +\infty)$ with nonincreasing densities with respect to the Lebesgue measure. More examples can also be found in Baraud and Birgé (2016). For some of these models, it is possible to build a ρ -estimator the risk of which does not

degenerate, a typical example being the Grenander estimator which is, as already seen, a ρ -estimator.

Following Baraud (2016), we introduce this definition.

DEFINITION 9. A class of functions \mathcal{F} defined on a set \mathcal{X} and with values in $[-\infty, +\infty]$ is said to be weak VC-major with dimension not larger than $k \in \mathbb{N}$ if, for all $u \in \mathbb{R}$, the class of subsets

$$\mathcal{C}_u(\mathcal{F}) = \{ \{x \in \mathcal{X} \mid f(x) > u\}, f \in \mathcal{F} \}$$

is a VC-class with dimension not larger than k (index not larger than $k + 1$). The weak dimension of \mathcal{F} is the smallest of such integers k .

DEFINITION 10. Let \mathcal{F} be a class of real-valued functions on \mathcal{X} . We shall say that an element $\bar{f} \in \mathcal{F}$ is extremal in \mathcal{F} with degree $d \in \mathbb{N}$ if the class of functions

$$(\mathcal{F}/\bar{f}) = \{ f/\bar{f}, f \in \mathcal{F} \}$$

is weak VC-major with dimension d .

For $\mu \in \mathcal{M}$, we consider a density set $\overline{\mathcal{Q}} \subset \mathcal{L}(\mu)$ which is viewed as a class of real-valued functions on $\overline{\mathcal{X}} = \bigcup_{i=1}^n (\{i\} \times \mathcal{X}_i)$ as we did in Section 7.2. The corresponding model $\overline{\mathcal{Q}}$ is $\{ \mathbf{q} \cdot \mu, \mathbf{q} \in \overline{\mathcal{Q}} \}$ and, for $d \in \{1, \dots, n\}$, we denote by $\overline{\mathcal{Q}}_d$ the subset of $\overline{\mathcal{Q}}$ of those densities \mathbf{q} which are extremal in $\overline{\mathcal{Q}}$ with degree d . We set $\overline{\mathcal{Q}}_d = \{ \mathbf{q} \cdot \mu, \mathbf{q} \in \overline{\mathcal{Q}}_d \}$ and let \mathcal{D} be the subset of $\{1, \dots, n\}$ consisting of those d such that $\overline{\mathcal{Q}}_d \neq \emptyset$.

PROPOSITION 9. Let ψ satisfy Assumption 2. For all ρ -models $\mathcal{Q} \subset \overline{\mathcal{Q}}$ and all $d \in \mathcal{D}$,

$$D^{\mathcal{Q}}(\mathbf{P}, \overline{\mathbf{P}}) \leq 33d[\log(e^2 n/d)]^3 \quad \text{for all } (\mathbf{P}, \overline{\mathbf{P}}) \in \mathcal{P} \times \overline{\mathcal{Q}}_d.$$

The proof is given in Section D.8 of the Supplementary Material [Baraud and Birgé (2018)]. This upper bound, although depending on the specific features of $\overline{\mathcal{Q}}$, is free from the choices of ψ . We immediately derive from Proposition 9 and Theorem 1 with a suitable choice of $\overline{\mathbf{P}}$ the following result.

COROLLARY 4. Let ψ satisfy Assumption 2 and assume that $\mathcal{D} \neq \emptyset$ and that $\overline{\mathcal{Q}}_d$ contains a countable and \mathcal{V} -dense subset \mathcal{Q}_d for all $d \in \mathcal{D}$. Any ρ -estimator $\overline{\mathbf{P}}$ on a ρ -model $\mathcal{Q} \subset \overline{\mathcal{Q}}$ containing $\bigcup_{d \in \mathcal{D}} \mathcal{Q}_d$ satisfies, for all $\mathbf{P} \in \mathcal{P}$ and $\xi > 0$,

$$(49) \quad \mathbb{P} \left[C\mathbf{h}^2(\mathbf{P}, \overline{\mathbf{P}}) \leq \inf_{d \in \mathcal{D}} \left[\mathbf{h}^2(\mathbf{P}, \overline{\mathcal{Q}}_d) + d[\log(e^2 n/d)]^3 \right] + \xi \right] \geq 1 - e^{-\xi},$$

for some constant C depending on ψ only. If, moreover, $\overline{\mathcal{Q}}$ is universally separable, any ρ -estimator relative to $((\mu, \overline{\mathcal{Q}}), \mathbf{0})$ also satisfies (49).

PROOF. Proposition 9 and Theorem 1 lead to (49). When $\overline{\mathcal{Q}}$ is universally separable there exists a countable and \mathcal{F} -dense subset $\mathcal{Q}' \subset \overline{\mathcal{Q}}$. The countable set $\mathcal{Q} = \mathcal{Q}' \cup (\bigcup_{d \in \mathcal{D}} \mathcal{Q}_d)$ is still countable and \mathcal{F} -dense in $\overline{\mathcal{Q}}$ and the corresponding ρ -model \mathcal{Q} is \mathcal{V} -dense in $\overline{\mathcal{Q}}$ and contains $\bigcup_{d \in \mathcal{D}} \mathcal{Q}_d$. By Theorem 3, any ρ -estimator relative to $((\mu, \overline{\mathcal{Q}}), \mathbf{0})$ is also a ρ -estimator relative to $((\mu, \mathcal{Q}), \mathbf{0})$ and, therefore, satisfies (49). \square

Note that the bound depends on the initial representation $(\mu, \overline{\mathcal{Q}})$ because the sets $\overline{\mathcal{Q}}_d$, hence the sets $\overline{\mathcal{Q}}_d$, depend on $\overline{\mathcal{Q}}$. This result looks like a model selection result among the sets $\{\overline{\mathcal{Q}}_d, d \in \mathcal{D}\}$ although the ρ -estimator $\widehat{\mathbf{P}}$ is built on a single ρ -model $\mathcal{Q} \subset \overline{\mathcal{Q}}$ with a null penalty function. It implies that the minimax risk over each set $\overline{\mathcal{Q}}_d$ is necessarily finite, while that on $\overline{\mathcal{Q}}$ might not be.

In the particular case of density estimation, the following result turns to be useful in view of applying Proposition 9.

PROPOSITION 10. *Let $\overline{\mathcal{Q}}$ be a subset of $\mathcal{L}(\mu)$ viewed as a class of functions on \mathcal{X} . If \overline{p} is extremal in $\overline{\mathcal{Q}}$ with degree d , $\overline{\mathbf{p}} = (\overline{p}, \dots, \overline{p})$ is extremal in $\overline{\mathcal{Q}} = \{\mathbf{q} = (q, \dots, q), q \in \overline{\mathcal{Q}}\}$, viewed as a class of functions on \mathcal{X} , with degree not larger than d .*

PROOF. Let $u \in \mathbb{R}$. If $\mathcal{C}_u((\overline{\mathcal{Q}}/\overline{\mathbf{p}}))$ shatters $\{\overline{x}_1, \dots, \overline{x}_k\} \subset \overline{\mathcal{X}}$, for all $J \subset \{1, \dots, k\}$ there exists $\mathbf{q} \in \overline{\mathcal{Q}}$ such that $j \in J$ if and only if $(\mathbf{q}/\overline{\mathbf{p}})(\overline{x}_j) = (q/\overline{p})(x_j) > u$. Hence, $\mathcal{C}_u((\overline{\mathcal{Q}}/\overline{\mathbf{p}}))$ shatters $\{x_1, \dots, x_k\} \subset \mathcal{X}$ which is only possible for $k \leq d$. \square

7.4. *Some examples of statistical models.* Let us restrict ourselves here to the density framework where X_1, \dots, X_n are assumed to be i.i.d. with common distribution P on \mathcal{X} and we have at hand a set of candidate probabilities $\mathcal{Q} = q \cdot \underline{\mu}$ with $q \in \overline{\mathcal{Q}} \subset \mathcal{L}(\mu)$ for P . We shall provide here some examples of density sets $\overline{\mathcal{Q}}$ to which Proposition 7 or 9 applies.

Piecewise constant functions. Let k be a positive integer and \mathcal{X} an arbitrary interval of \mathbb{R} (possibly $\mathcal{X} = \mathbb{R}$). We define \mathcal{F}_k as the class of functions f on \mathcal{X} such that there exists a partition $\mathcal{I}(f)$ of \mathcal{X} into at most k intervals (of positive lengths) with f constant on each of these intervals. Note that $\mathcal{I}(f)$ depends on f . The following result is to be proved in Section C.1 of the Supplementary Material [Baraud and Birgé (2018)].

PROPOSITION 11. *The set \mathcal{F}_k is VC-subgraph with dimension bounded by $2k$.*

Let us apply this to histogram estimation on $\mathcal{X} = \mathbb{R}$. For a positive integer D , we denote by $\overline{\mathcal{Q}}_D$ the subset of \mathcal{F}_{D+2} of right-continuous densities with respect to the Lebesgue measure μ , that is, the set of right-continuous and piecewise

constant densities on \mathbb{R} with at most D pieces and by $\overline{\mathcal{Q}}_D = \{q \cdot \mu, q \in \overline{\mathcal{Q}}_D\}$ the corresponding model for P . We derive from Propositions 7 and 11 that, for some universal constant $C > 0$ and all ρ -models $\mathcal{Q} \subset \overline{\mathcal{Q}}_D$,

$$D^{\mathcal{Q}}(\mathbf{P}, \overline{\mathbf{P}}) \leq C(D \wedge n)[1 + \log_+(n/D)] \quad \text{for all } (\mathbf{P}, \overline{\mathbf{P}}) \in \mathcal{P} \times \mathcal{P}^\mu.$$

Hence, by Corollary 3, for all ρ -estimators $\widehat{\mathbf{P}}$ on some countable and \mathcal{V} -dense subset \mathcal{Q}_D of $\overline{\mathcal{Q}}_D$,

$$(50) \quad C\mathbb{E}[\mathbf{h}^2(\mathbf{P}, \widehat{\mathbf{P}})] \leq \mathbf{h}^2(\mathbf{P}, \overline{\mathcal{Q}}_D) + (D \wedge n)[1 + \log_+(n/D)].$$

Since $\overline{\mathcal{Q}}_D$ is universally separable (see Section B of the Supplementary Material [Baraud and Birgé (2018)]), we deduce from Theorem 3 that (50) remains true for any ρ -estimator $\widehat{\mathbf{P}}$ on the noncountable model $\overline{\mathcal{Q}}_D$ relative to the representation $(\mu, \overline{\mathcal{Q}}_D)$ (with a null penalty function).

The logarithmic factor in this bound turns out to be necessary. The argument is as follows. When $\mathbf{P} \in \overline{\mathcal{Q}}_D$, it follows from (50) that $\mathbb{E}[\mathbf{h}^2(\mathbf{P}, \widehat{\mathbf{P}})] \leq C'(D \wedge n)[1 + \log_+(n/D)]$ for some universal constant $C' > 0$. This inequality appears to be optimal (up to the numerical constant C') in view of the lower bound established in Proposition 2 of Birgé and Massart (1998). This also shows that the logarithmic factor involved in the bound of the ρ -dimension function established in Proposition 7 is necessary, at least for some VC-subgraph classes.

Piecewise exponential families. Using similar arguments based on Corollaries 3 and 4 as we did above, we may establish risk bounds of the same flavour as (50) with the following density sets.

DEFINITION 11. Let g_1, \dots, g_J be $J \geq 1$ real-valued functions on a set \mathcal{X} . We shall say that a class \mathcal{F} of positive functions on \mathcal{X} is an exponential family based on g_1, \dots, g_J if the elements f of \mathcal{F} are of the form

$$(51) \quad f = \exp\left[\sum_{j=1}^J \beta_j g_j\right] \quad \text{for } \beta_1, \dots, \beta_J \in \mathbb{R}.$$

If \mathcal{X} is a nontrivial interval of \mathbb{R} and k a positive integer, we shall say that \mathcal{F} is a k -piecewise exponential family based on g_1, \dots, g_J if for any $f \in \mathcal{F}$ there exists a partition $\mathcal{I}(f)$ of \mathcal{X} into at most k intervals such that for all $I \in \mathcal{I}(f)$, the restriction f_I of f to I is of the form (51) with coefficients β_j depending on I .

The properties of exponential and piecewise exponential families are described by the following proposition to be proven in Section C of the Supplementary Material [Baraud and Birgé (2018)].

PROPOSITION 12. Let $\overline{\mathcal{Q}}$ be a class of functions on \mathcal{X} .

(i) If $\overline{\mathcal{Q}}$ is an exponential family based on $J \geq 1$ functions, $\overline{\mathcal{Q}}$ is VC-subgraph with index not larger than $J + 2$.

(ii) Let \mathcal{I} be a partition of \mathcal{X} with cardinality not larger than $k \geq 1$. If for all $I \in \mathcal{I}$, the family $\overline{\mathcal{Q}}_I$ consisting of the restrictions of the functions q in $\overline{\mathcal{Q}}$ to the set I is an exponential family on I based on $J \geq 1$ functions, $\overline{\mathcal{Q}}$ is VC-subgraph with index not larger than $k(J + 2)$.

(iii) If \mathcal{X} is a nontrivial interval of \mathbb{R} and $\overline{\mathcal{Q}}$ is a k -piecewise exponential family based on J functions, all densities $\overline{p} \in \overline{\mathcal{Q}}$ are extremal in $\overline{\mathcal{Q}}$ with degree d not larger than $\lceil 9.4k(J + 2) \rceil = \inf\{j \in \mathbb{N}, j \geq 9.4k(J + 2)\}$.

8. Estimating a conditional distribution.

8.1. *Description of the framework.* Let us now apply our result to the estimation of a conditional distribution. We consider i.i.d. pairs $X_i = (W_i, Y_i)$, $i = 1, \dots, n$ of random variables with values in the product space $(\mathcal{W} \times \mathcal{Y}, \mathcal{B}(\mathcal{W}) \otimes \mathcal{B}(\mathcal{Y}))$ and common distribution P , assuming that truly $\mathbf{P} = P^{\otimes n}$. We denote by P_W the marginal distribution of W and assume the existence of a conditional distribution P_w of Y when $W = w$, which means that for all bounded measurable functions f on \mathcal{Y} ,

$$\mathbb{E}[f(Y)|W = w] = \int_{\mathcal{Y}} f(y) dP_w(y) \quad P_W\text{-a.s.}$$

and for all bounded measurable functions g on $\mathcal{W} \times \mathcal{Y}$,

$$\mathbb{E}[g(W, Y)] = \int_{\mathcal{W}} \left[\int_{\mathcal{Y}} g(w, y) dP_w(y) \right] dP_W(w).$$

Our purpose is to estimate the conditional distribution P_w without the knowledge of P_W which may therefore be completely arbitrary. To do so, we consider a reference measure λ on $(\mathcal{Y}, \mathcal{B}(\mathcal{Y}))$ and the set $\mathcal{L}_c(\mathcal{W}, \lambda)$ of conditional densities with respect to λ , that is, the set of measurable functions t from $(\mathcal{W} \times \mathcal{Y}, \mathcal{B}(\mathcal{W}) \otimes \mathcal{B}(\mathcal{Y}))$ to \mathbb{R}_+ such that for all $w \in \mathcal{W}$, the function $t_w : y \mapsto t(w, y) \in \mathcal{L}(\lambda)$. Then to each element $t \in \mathcal{L}_c(\mathcal{W}, \lambda)$ is associated a conditional distribution $t_w \cdot \lambda$ for Y . In order to build our estimators, we first introduce a countable family $\{S_m, m \in \mathcal{M}\}$ of countable subsets S_m of $\mathcal{L}_c(\mathcal{W}, \lambda)$, and a non-negative weight function Δ on \mathcal{M} satisfying (21). To each S_m , we associate the ρ -density model $\mathcal{Q}_m = \{Q_t, t \in S_m\}$ for P , where the probability Q_t on $\mathcal{W} \times \mathcal{Y}$ is given by

$$Q_t(A \times B) = \int_A \left[\int_B t_w(y) d\lambda(y) \right] dP_W(w),$$

or equivalently

$$\frac{dQ_t}{dP_W \otimes d\lambda}(w, y) = t(w, y).$$

This means that Q_t has a marginal distribution P_W on \mathscr{W} and a conditional distribution given $W = w$ with density t_w with respect to λ . Note that the ρ -models \mathscr{Q}_m depend on the unknown distribution P_W but the densities with respect to the dominating measure $P_W \otimes \lambda$ do not. This leads to a family of ρ -models \mathscr{Q}_m for \mathbf{P} and a reference ρ -model $\mathscr{Q} = \bigcup_{m \in \mathcal{M}} \mathscr{Q}_m$. If we introduce a suitable penalty pen on \mathscr{Q} , we may build a ρ -estimator of \mathbf{P} from our sample X_1, \dots, X_n according to the recipe of Section 2.3 since its values only depend on the family of densities in $\mathscr{Q} = \bigcup_{m \in \mathcal{M}} S_m$. As a consequence, our estimation strategy neither needs to know P_W nor to estimate it. Such a ρ -estimator will be of the form $Q_{\hat{s}}^{\otimes n}$ with $Q_{\hat{s}} = \hat{s} \cdot (P_W \otimes \lambda)$ and will provide an estimator $\hat{s}_w \cdot \lambda$ of the conditional probability P_w .

Within this framework, the Hellinger distance between the probabilities at hand writes, for any measure ν that dominates both P and $P_W \otimes \lambda$,

$$\begin{aligned} \frac{1}{2} \int_{\mathscr{W} \times \mathscr{Y}} \left(\sqrt{\frac{dP_W(w) dP_w(y)}{d\nu}} - \sqrt{t(w, y) \frac{dP_W(w) d\lambda(y)}{d\nu}} \right)^2 d\nu(w, y) \\ = h^2(P, Q_t) = \int_{\mathscr{W}} h^2(P_w, t_w \cdot \lambda) dP_W(w). \end{aligned}$$

Therefore,

$$(52) \quad h^2(P, \mathscr{Q}_m) = \inf_{t \in S_m} \int_{\mathscr{W}} h^2(P_w, t_w \cdot \lambda) dP_W(w).$$

Note that $h^2(P, Q_t)$ can actually be viewed as a loss function for the conditional distributions, of the form $\ell(P_w, t_w \cdot \lambda)$ since it actually only depends on P_w and t_w .

8.2. *Assumptions and results.* Let us assume the following.

ASSUMPTION 5. For all $m \in \mathcal{M}$, S_m is VC-subgraph with index not larger \bar{V}_m .

We may then deduce from Theorem A.1 of the Supplementary Material [Baraud and Birgé (2018)] the following result.

COROLLARY 5. Let $\{S_m, m \in \mathcal{M}\}$ be a family of countable subsets of $\mathcal{L}_c(\mathscr{W}, \lambda)$ satisfying Assumption 5, Δ be a weight function on \mathcal{M} which satisfies (21), ψ a function satisfying Assumption 2, $\mathscr{Q} = \bigcup_{m \in \mathcal{M}} \{Q_t, t \in S_m\}$ and $\text{pen} : \mathscr{Q} \rightarrow \mathbb{R}_+$ given, for all $Q \in \mathscr{Q}$, by

$$\text{pen}(Q) = \kappa \inf_{\{m \in \mathcal{M} \mid Q = Q_t \text{ with } t \in S_m\}} \left[\frac{C_1}{4.7} (\bar{V}_m \wedge n) \left[1 + \log_+ \left(\frac{n}{\bar{V}_m} \right) \right] + \Delta(m) \right],$$

where κ is given by (19) and C_1 is the constant appearing in Proposition 7. Then any density ρ -estimator $Q_{\hat{s}}$ relative to $[(\mu, \mathscr{Q}), \text{pen}]$ satisfies, for some constant $C' > 0$ depending on the choice of ψ only,

$$\mathbb{E}[h^2(P, Q_{\hat{s}})] \leq C' \inf_{m \in \mathcal{M}} \left[h^2(P, \mathscr{Q}_m) + \frac{\bar{V}_m \wedge n}{n} \left(1 + \log_+ \left(\frac{n}{\bar{V}_m} \right) \right) + \frac{\Delta(m)}{n} \right].$$

Note that this result does not require any information or assumption on the distribution of W . If, in particular, the conditional probability P_w is absolutely continuous with respect to λ for almost all w with density $dP_w/d\lambda = s_w$, P_W -a.s., one can write

$$\begin{cases} h^2(P, Q_{\hat{s}}) = \int_{\mathcal{W}} h^2(s_w \cdot \lambda, \hat{s}_w \cdot \lambda) dP_W(w), \\ h^2(P, \mathcal{Q}_m) = \inf_{t \in S_m} \int_{\mathcal{W}} h^2(s_w \cdot \lambda, t_w \cdot \lambda) dP_W(w). \end{cases}$$

PROOF OF COROLLARY 5. If we apply Propositions 7 and 8 to each ρ -model \mathcal{Q}_m with $m \in \mathcal{M}$, we obtain under Assumption 5 the existence of a universal constant $C_1 > 0$ such that for all $(\mathbf{P}, \bar{\mathbf{P}}) \in \mathcal{P} \times \mathcal{P}^\mu$:

$$D^{\mathcal{Q}_m}(\mathbf{P}, \bar{\mathbf{P}}) \leq D_n(m) = C_1(\bar{V}_m \wedge n)[1 + \log_+(n/\bar{V}_m)].$$

Inequality (22) is fulfilled with $K = 0$ and the penalty function therefore satisfies (23) with $\kappa_1 = 0$ for all $m \in \mathcal{M}$. The result follows from Theorem 2, then an integration of (24), with respect to $\xi > 0$. \square

9. Regression with a random design. In this section, we assume that the observations $X_i = (W_i, Y_i)$, $1 \leq i \leq n$ are i.i.d. copies of a random pair

$$(53) \quad X = (W, Y) \quad \text{with } Y = f(W) + \epsilon,$$

where W is a random variable with distribution P_W on a measurable space $(\mathcal{W}, \mathcal{B}(\mathcal{W}))$, f is an unknown regression function mapping \mathcal{W} into \mathbb{R} and ϵ is a real-valued random variable with distribution P_ϵ , which is independent of W . Both distributions P_W and P_ϵ are assumed to be unknown. We shall use the specific notation introduced in Section 3.3 when the data are i.i.d. and denote by μ the product measure $P_W \otimes \lambda$ where λ is the Lebesgue measure on \mathbb{R} . Note that μ is unknown since it depends on the distribution P_W of the design W .

If ϵ had a density s with respect to λ , the distribution P of $X = (W, Y)$ would be absolutely continuous with respect to μ with density p given by

$$(54) \quad p(w, y) = s(y - f(w)) \quad \text{for } (w, y) \in \mathcal{X},$$

depending thus on two parameters: the density s of the errors and the regression function f .

Denoting by \mathcal{D} the set of all densities on $(\mathbb{R}, \mathcal{B}(\mathbb{R}), \lambda)$ and \mathcal{F} the set of all measurable functions mapping \mathcal{W} into \mathbb{R} , our aim is to estimate P assuming that it is close to some distribution of the form $p \cdot \mu$ with p given by (54) for some $s \in \mathcal{D}$ and $f \in \mathcal{F}$. Besides, when P is truly of this form we shall also derive estimators for both s and f .

9.1. *The main result.* For $r \in \mathcal{D}$ and $g \in \mathcal{F}$, we set

$$Q_{r,g} = q_{r,g} \cdot \mu \quad \text{with } q_{r,g}(w, y) = r(y - g(w)),$$

which means that $Q_{r,g}$ is the distribution of X in (53) when $f = g$ and ϵ is distributed according to $R = r \cdot \lambda$. Given a density $r \in \mathcal{D}$ and a countable subset F of \mathcal{F} , we define the ρ -model:

$$\mathcal{Q}_m = \{Q_{r,g}, g \in F\} \quad \text{for } m = (r, F).$$

Given a countable subset \mathcal{D} of \mathcal{D} and a countable family \mathbb{F} of countable subsets F of \mathcal{F} , we estimate P on the basis of the collection of ρ -models $\{\mathcal{Q}_m, m \in \mathcal{M}\}$ with $\mathcal{M} \subset \mathcal{D} \times \mathbb{F}$. We endow the family $\{\mathcal{Q}_m, m \in \mathcal{M}\}$ with a weight function Δ satisfying (21) and assume the following.

ASSUMPTION 6.

- (i) The densities $r \in \mathcal{D}$ are unimodal.
- (ii) Each F in \mathbb{F} is VC-subgraph with index $\bar{V}(F)$.
- (iii) The function ψ satisfies Assumption 2 with $\mathcal{Q} = \{Q^{\otimes n}, Q \in \mathcal{Q} = \bigcup_{m \in \mathcal{M}} \mathcal{Q}_m\}$.

Under Assumptions 6-(i) and 6-(ii), the family of densities Q_m is a VC-subgraph on \mathcal{X} with index not larger than

$$(55) \quad \bar{V}_m = 9.41 \bar{V}(F) \quad \text{for all } m = (r, F) \in \mathcal{M}.$$

This result derives from Baraud, Birgé and Sart (2017), Proposition 42. Besides, under Assumption 6-(iii), Proposition 7 applies and implies that, for some universal constant $C_1 > 0$, all $m \in \mathcal{M}$, $\mathbf{P} \in \mathcal{P}$ and $\bar{\mathbf{P}} \in \mathcal{Q}$, $D^{\mathcal{Q}_m}(\mathbf{P}, \bar{\mathbf{P}}) \leq D_n(m)$ with

$$(56) \quad D_n(m) = C_1 (\bar{V}_m \wedge n) [1 + \log_+(n/\bar{V}_m)] \quad \text{for all } m = (r, F) \in \mathcal{M},$$

so that (22) holds with $K = 0$. Setting

$$\text{pen}(\mathbf{Q}) = \kappa \inf_{\{m \in \mathcal{M} \mid \mathcal{Q}_m \ni Q\}} \left\{ \frac{D_n(m)}{4.7} + \Delta(m) \right\},$$

we may apply Theorem 2 with $\kappa_1 = 0$, which leads, in this particular case, to the following analogue of (24).

THEOREM 4. *Assume that Assumption 6 holds. For any distribution $P \in \mathcal{P}$ and $\mathbf{P} = P^{\otimes n}$, any ρ -estimator $\hat{\mathbf{P}} = (\hat{P}, \dots, \hat{P})$ satisfies, for all $\xi > 0$, with probability at least $1 - e^{-\xi}$,*

$$(57) \quad \begin{aligned} &Ch^2(P, \hat{P}) \\ &\leq \inf_{m \in \mathcal{M}} \left[h^2(P, \mathcal{Q}_m) + \frac{\bar{V}_m \wedge n}{n} \left[1 + \log_+ \left(\frac{n}{\bar{V}_m} \right) \right] + \frac{\Delta(m)}{n} \right] + \frac{\xi}{n}, \end{aligned}$$

for some constant $C > 0$ only depending on the choice of ψ .

At this stage, some comments are in order:

(a) This result holds without any assumption on the distribution P_W of the design.

(b) The result is true even if the regression framework (53) is not exact as long as the X_i are i.i.d. In particular, the distribution P needs not have a density with respect to $\mu = P_W \otimes \lambda$.

(c) If r admits k modes with $k > 1$ and F is a VC-subgraph with index not larger than \bar{V} , $\mathcal{Q}_{(r,F)}$ remains a VC-subgraph and its index is still bounded by $C(k)\bar{V}$ for some constant $C(k)$ that now depends on k . Consequently, the above result generalizes to families \mathcal{D} of densities admitting more than a single mode in which case $\bar{V}(F)$ should be replaced by $c(r)\bar{V}(F)$ where $c(r)$ is a positive number depending on the number of modes of the density r .

(d) With Theorem 4 at hand, we could obtain in the present random design context an analogue of Corollary 39 in Baraud, Birgé and Sart (2017) which was established when the W_i were deterministic (fixed design regression).

9.2. *Estimation of s and f .* Let us now consider the situation where the regression framework (53) is exact and ϵ has an unknown density s with respect to the Lebesgue measure λ . Then $P = Q_{s,f}$ admits a density $q_{s,f}$ with respect to μ which is given by (54) with s belonging to \mathcal{D} and f to \mathcal{F} but not necessarily to our ρ -models \mathcal{D} and $\mathcal{F} = \bigcup_{F \in \mathbb{F}} F$, respectively. Since we may choose our ρ -estimator of the form

$$\hat{P} = q_{\hat{s}, \hat{f}} \cdot \mu \quad \text{with } (\hat{s}, \hat{f}) \in \mathcal{D} \times \mathcal{F},$$

our procedure results in estimators \hat{s} and \hat{f} for s and f , respectively, and our aim in this section is to establish risk bounds for these two estimators.

Since the map $(r, g) \mapsto Q_{r,g}$ is not necessarily one-to-one from $\mathcal{D} \times \mathcal{F}$ to \mathcal{P} , an identifiability condition is required on our ρ -model \mathcal{Q} so that the equality $Q_{r,g} = Q_{r',g'}$ with $r, r' \in \mathcal{D}$ and $g, g' \in \mathcal{F}$ implies that $r = r'$ λ -a.e. and $g = g'$ P_W -a.s. In order to state this identifiability condition, let us introduce the following notation. For $r \in \mathcal{D}$ and $a \in \mathbb{R}$, we shall denote by R_a the probability on $(\mathbb{R}, \mathcal{B}(\mathbb{R}), \lambda)$ with density $r_a(\cdot) = r(\cdot - a)$. When ϵ has density r and $a = g(w)$ for some $w \in \mathcal{W}$, R_a can be viewed as the conditional distribution of $Y = g(W) + \epsilon$ given $W = w$. Given $r, r' \in \mathcal{D}$, $g, g' \in \mathcal{F}$ and $w \in \mathcal{W}$, the Hellinger distance between the probabilities $R_{g(w)}$, and $R'_{g'(w)}$ is given by

$$h^2(R_{g(w)}, R'_{g'(w)}) = \frac{1}{2} \int_{\mathbb{R}} \left[\sqrt{r(y - g(w))} - \sqrt{r'(y - g'(w))} \right]^2 d\lambda(y)$$

and the Hellinger distance between the corresponding probabilities $Q_{r,g}$ and $Q_{r',g'}$ on $(\mathcal{X}, \mathcal{B})$ writes

$$(58) \quad h^2(Q_{r,g}, Q_{r',g'}) = \int_{\mathcal{W}} h^2(R_{g(w)}, R'_{g'(w)}) dP_W(w).$$

We recall that the Hellinger distance is translation invariant which means that for all densities $r, r' \in \mathcal{D}$, $a, a' \in \mathbb{R}$,

$$(59) \quad h^2(R_a, R_{a'}) = h^2(R_{a-a'}, R').$$

In particular, taking $a = g(w)$ and $a' = g'(w)$ for $g, g' \in \mathcal{F}$ and $w \in \mathcal{W}$ and integrating (59) with respect to P_W we get for all $(g, g') \in \mathcal{F}^2$ and $(r, r') \in \mathcal{D}^2$

$$(60) \quad h^2(Q_{r,g}, Q_{r',g'}) = h^2(Q_{r,g-g'}, Q_{r',0}).$$

In order to warrant identifiability, we assume the following.

ASSUMPTION 7. There exists a positive constant A such that, for all $r, r' \in \mathcal{D}$, $R = r \cdot \lambda$ and $R' = r' \cdot \lambda$,

$$h(R, R') \leq A \inf_{a \in \mathbb{R}} h(R_a, R').$$

When $Q_{r,g} = Q_{r',g'}$, (58) asserts that $h(R_{g(w)}, R'_{g'(w)}) = 0$ for P_W -almost all $w \in \mathcal{W}$ and (59) implies that $h(R_{g(w)-g'(w)}, R') = 0$ for all such $w \in \mathcal{W}$. Applying Assumption 7 with $a = g(w) - g'(w)$ leads to $R = R'$ and $g(w) = g'(w)$ which solves our identifiability problem.

In order to evaluate the risk of our estimator \hat{f} of f , we endow \mathcal{F} with the loss function d_s defined on $\mathcal{F} \times \mathcal{F}$ by

$$d_s^2(g, g') = \frac{1}{2} \int_{\mathcal{W} \times \mathbb{R}} \left(\sqrt{s_{g(w)}(y)} - \sqrt{s_{g'(w)}(y)} \right)^2 dP_W(w) dy \quad \text{for } g, g' \in \mathcal{F}.$$

This loss function depends on the true density s of the errors ϵ and on the distribution P_W of the design, hence on P . We have seen in Section 6.3 of Baraud, Birgé and Sart (2017) that, if the density s is of order α with $\alpha \in (-1, 1]$ (see Definition 26 of that paper for the order of a function), the restriction of d_s to the $\mathbb{L}_\infty(P_W)$ -ball $\mathcal{B}_\infty(b)$ centred at 0 with radius b is equivalent (up to factors depending on b and s) to

$$\|g - g'\|_{1+\alpha, P_W}^{(1+\alpha)/2} \quad \text{with} \quad \|g - g'\|_{1+\alpha, P_W} = \left[\int_{\mathcal{W}} |g - g'|^{1+\alpha} dP_W \right]^{1/(1+\alpha)}.$$

In particular, if $\mathcal{F} \subset \mathcal{B}_\infty(b)$ and the true regression function f also belongs to $\mathcal{B}_\infty(b)$,

$$c(s, b) \|f - g\|_{1+\alpha, P_W}^{(1+\alpha)/2} \leq d_s(f, g) \leq C(s, b) \|f - g\|_{1+\alpha, P_W}^{(1+\alpha)/2} \quad \text{for all } g \in \mathcal{F}$$

and suitable positive numbers $c(s, b)$ and $C(s, b)$. Of special interest is the case of $\alpha = 1$ for which $d_s(f, g)$ is of the order of the $\mathbb{L}_2(P_W)$ -distance between f and g for all $g \in \mathcal{F}$. This situation is met when the translation ρ -model associated to s is regular which is the case when s is Cauchy, Gaussian, Laplace, etc. When s is uniform or exponential, $d_s^2(\cdot, \cdot)$ is then equivalent to the $\mathbb{L}_1(P_W)$ -norm.

Furthermore, when \mathcal{F} is a subset of a k -dimensional linear space $\overline{\mathcal{F}}$ generated by ϕ_1, \dots, ϕ_k , these norms can in turn be translated into a norm on \mathbb{R}^k between the coefficients relative to this basis. More precisely, if f belongs to $\overline{\mathcal{F}}$ and writes as $\sum_{j=1}^k \beta_j \phi_j$ and $\overline{f} = \sum_{j=1}^k \overline{\beta}_j \phi_j$ is an element of $\overline{\mathcal{F}}$, there exist two positive constants $c'(s, b, k)$ and $C'(s, b, k)$ such that

$$c'(s, b, k) \left[\max_{j=1, \dots, k} |\beta_j - \overline{\beta}_j| \right]^{(1+\alpha)/2} \leq d_s(f, \overline{f}) \leq C'(s, b, k) \left[\max_{j=1, \dots, k} |\beta_j - \overline{\beta}_j| \right]^{(1+\alpha)/2}.$$

In particular, if $\overline{f} = \overline{f}(n)$ converges towards f at a rate v_n with respect to the distance $d_s(\cdot, \cdot)$, the coefficients of \overline{f} converge in sup-norm towards those of f at rate $v_n^{2/(1+\alpha)}$; this latter rate being faster than v_n when $\alpha < 1$.

In view of evaluating the risk of our estimator \widehat{s} of the density s , we shall consider the loss between two densities $r, r' \in \mathcal{D}$ induced by the Hellinger distance between the two corresponding measures $r \cdot \lambda$ and $r' \cdot \lambda$ and we shall write this loss $h(r, r')$ so that

$$h(r, r') = h(r \cdot \lambda, r' \cdot \lambda) = h(Q_{r,0}, Q_{r',0}) \quad \text{for all } r, r' \in \mathcal{D}.$$

We deduce from Theorem 4 the following result.

COROLLARY 6. *Assume that the X_i are i.i.d. with density p given by (54) and that Assumptions 6 and 7 are satisfied. For all $\xi > 0$ and all ρ -estimators $Q_{\widehat{s}, \widehat{f}}$, with $\widehat{s} \in \mathcal{D}$ and $\widehat{f} \in \mathcal{F}$, based on the family of ρ -models \mathcal{Q}_m defined in Section 9.1, with probability at least $1 - e^{-\xi}$,*

$$C \max \left\{ d_s^2(f, \widehat{f}), h^2(s, \widehat{s}) \right\} \leq \inf_m \left\{ d_s^2(f, F) + h^2(s, r) + \frac{\overline{V}_m \wedge n}{n} \left[1 + \log_+ \left(\frac{n}{\overline{V}_m} \right) \right] + \frac{\Delta(m)}{n} \right\} + \frac{\xi}{n},$$

where the infimum runs among all $m = (r, F) \in \mathcal{M}$, C is a positive constant depending on A and the choice of ψ and \overline{V}_m is given by (55).

The risk bound is the same for the two estimators and depends on the approximation properties of \mathcal{F} and \mathcal{D} with respect to f and s , respectively. The proof of this corollary is given in Section D.9 of the Supplementary Material [Baraud and Birgé (2018)].

10. Estimator selection and aggregation. In the case of density estimation, ρ -estimators can also be used to perform selection or aggregation of preliminary estimators. In this case, we assume that we have at hand a set $\mathbf{X}_1 = (X_1, \dots, X_n)$

of n independent random variables with an unknown joint distribution \mathbf{P} to be estimated. We also have at hand a finite family $\mathcal{Q} = \{\mathbf{P}_j, j \in \mathcal{J}\}$ of probabilities that can be considered as candidate estimators for \mathbf{P} . These are completely arbitrary, but in a typical situation, it is assumed (although this may not be true) that the observations X_i are i.i.d. and the \mathbf{P}_j are preliminary estimators of the form $\mathbf{P}_j = P_j^{\otimes n}(X_2)$, where X_2 is a second sample independent from X_1 , and the $P_j = P_j(X_2)$ are estimators that derive from various procedures applied to the sample X_2 .

10.1. *Estimator selection.* Taking $\mathcal{M} = \mathcal{J}$, we view each probability \mathbf{P}_j as a ρ -model $\mathcal{Q}_j = \{\mathbf{P}_j\}$ with a single element. As a consequence, it follows from Proposition 6 that $D^{\mathcal{Q}_j}(\mathbf{P}, \mathbf{P}_j) \leq 9 \log 2 < 6.3$ so that (22) holds with $D_n(j) = 6.3$ for all j and $K = 0$. Then we choose the weights $\Delta(j)$ satisfying (21). We may choose $\Delta(j) = \log |\mathcal{J}|$ for all $j \in \mathcal{J}$ but other more Bayesian choices are possible, or choices based on the confidence we have in the various procedures used to build the preliminary estimators. To compute the penalized ρ -estimator $\hat{\mathbf{P}}$ of \mathbf{P} , we may use the penalty function $\text{pen}(\mathbf{P}_j) = \kappa \Delta(j)$ for all $j \in \mathcal{J}$ which is of the form (23) with $\kappa_1 = -\kappa(6.3/4.7)$. Finally, (24) shows that, for all $\xi > 0$ and $\mathbf{P} \in \mathcal{P}$,

$$\mathbb{P}\left[\mathbf{h}^2(\mathbf{P}, \hat{\mathbf{P}}) \leq C \inf_{j \in \mathcal{J}} \left(\mathbf{h}^2(\mathbf{P}, \mathbf{P}_j) + \Delta(j) + 1 + \xi\right)\right] \geq 1 - e^{-\xi},$$

where C denotes a suitable constant depending on ψ only.

10.2. *Convex estimator aggregation.* In this case, we set $\mathcal{J} = \{1, \dots, N\}$, $N \geq 2$ and \mathcal{C} for the N -dimensional simplex:

$$\mathcal{C} = \left\{ (\alpha_1, \dots, \alpha_N) \text{ such that } \alpha_j \geq 0 \text{ for } 1 \leq j \leq N \text{ and } \sum_{j=1}^N \alpha_j = 1 \right\}.$$

We select a dominating measure μ , densities $p_j = dP_j/d\mu$ and we then consider a single density model:

$$\overline{\mathcal{Q}} = \left\{ \sum_{j=1}^N \alpha_j p_j \text{ for } (\alpha_1, \dots, \alpha_N) \in \mathcal{C} \right\}.$$

The following result then holds.

PROPOSITION 13. *Let ψ satisfy Assumption 2. Any ρ -estimator $\hat{\mathbf{P}}$ on $\overline{\mathcal{Q}} = \{\mathbf{q} \cdot \boldsymbol{\mu}, \mathbf{q} \in \overline{\mathcal{Q}}\}$ relative to $((\boldsymbol{\mu}, \overline{\mathcal{Q}}), \mathbf{0})$ satisfies*

$$\mathbb{P}\left[Ch^2(\mathbf{P}, \hat{\mathbf{P}}) \leq h^2(\mathbf{P}, \overline{\mathcal{Q}}) + N \log n + \xi\right] \geq 1 - e^{-\xi} \quad \text{for all } \xi > 0$$

and some constant $C > 0$ depending on ψ only.

PROOF. The map $(\alpha_1, \dots, \alpha_n) \mapsto \sum_{j=1}^N \alpha_j p_j$ from \mathcal{C} to $\overline{\mathcal{Q}}$ is continuous if we equip \mathcal{C} with the usual Euclidean distance and $\overline{\mathcal{Q}}$ with the topology of pointwise convergence. Since \mathcal{C} is separable, $\overline{\mathcal{Q}}$ is universally separable and contains thus a countable subset \mathcal{Q} which is \mathcal{T} -dense in $\overline{\mathcal{Q}}$. The set $\overline{\mathcal{Q}}$ being furthermore a subset of an N -dimensional linear space, it is a VC-subgraph with index \overline{V} not larger than $N + 2$ and it follows from Proposition 7 (and Proposition 8) that, for all ρ -models $\mathcal{Q} \subset \overline{\mathcal{Q}}$,

$$D^{\mathcal{Q}}(\mathbf{P}, \overline{\mathbf{P}}) \leq D_n = CN \log n \quad \text{for all } (\mathbf{P}, \overline{\mathbf{P}}) \in \mathcal{P} \times \mathcal{P}^{\mu}.$$

We may apply Theorem 3 with $p \equiv 0$ to the single model $\overline{\mathcal{Q}}$. A ρ -estimator $\widehat{\mathbf{P}}$ on $\overline{\mathcal{Q}}$ relative to $((\mu, \overline{\mathcal{Q}}), \mathbf{0})$ is also a ρ -estimator on the ρ -model $\mathcal{Q} = \{\mathbf{q} \cdot \mu, \mathbf{q} \in \mathcal{Q}\}$ and we deduce from Theorem 1, more precisely from (20), that

$$\mathbb{P}\left[Ch^2(\mathbf{P}, \widehat{\mathbf{P}}) \leq h^2(\mathbf{P}, \mathcal{Q}) + N \log n + \xi\right] \geq 1 - e^{-\xi} \quad \text{for all } \xi > 0,$$

some constant C depending on ψ only and all $\mathbf{P} \in \mathcal{P}$. We conclude using the fact that \mathcal{Q} is \mathcal{V} -dense in $\overline{\mathcal{Q}}$ since \mathcal{Q} is \mathcal{T} -dense in $\overline{\mathcal{Q}}$. \square

It should be noted that there is no \mathbb{L}_2 -type argument here; the densities p_j can be absolutely anything and the true distribution \mathbf{P} should be a product measure but not necessarily of the form $P^{\otimes n}$.

Practical implementation. Since the set $\overline{\mathcal{Q}}$ is convex, Proposition 5 applies with $\psi = \psi_1$ or $\psi = \psi_2$ of Proposition 3 and the MLE on $\overline{\mathcal{Q}}$ is a ρ -estimator. It is obtained by maximizing over the convex set \mathcal{C} the concave map

$$(\alpha_1, \dots, \alpha_N) \mapsto \sum_{i=1}^n \log \left(\sum_{j=1}^N \alpha_j p_j(X_i) \right).$$

Acknowledgements. The authors are grateful to Weijie Su for letting them know about the nice connection between the MLE and ρ -estimators in the case of a convex parameter set and for allowing them to include his result (Proposition 5) in this paper.

SUPPLEMENTARY MATERIAL

Supplement to “Rho-estimators revisited: general theory and applications” (DOI: [10.1214/17-AOS1675SUPP](https://doi.org/10.1214/17-AOS1675SUPP); .pdf). This supplement provides the proofs of most results given in the paper and an additional section (D.10) devoted to robust tests.

REFERENCES

- AUDIBERT, J.-Y. and CATONI, O. (2011). Robust linear least squares regression. *Ann. Statist.* **39** 2766–2794. [MR2906886](#)
- BARAUD, Y. (2016). Bounding the expectation of the supremum of an empirical process over a (weak) VC-major class. *Electron. J. Stat.* **10** 1709–1728. [MR3522658](#)
- BARAUD, Y. and BIRGÉ, L. (2016). Rho-estimators for shape restricted density estimation. *Stochastic Process. Appl.* **126** 3888–3912.
- BARAUD, Y. and BIRGÉ, L. (2018). Supplement to “Rho-estimators revisited: General theory and applications.” DOI:10.1214/17-AOS1675SUPP.
- BARAUD, Y., BIRGÉ, L. and SART, M. (2017). A new method for estimation and model selection: ρ -estimation. *Invent. Math.* **207** 425–517.
- BIRGÉ, L. (1983). Approximation dans les espaces métriques et théorie de l’estimation. *Z. Wahrsch. Verw. Gebiete* **65** 181–237. [MR0722129](#)
- BIRGÉ, L. (2006). Model selection via testing: An alternative to (penalized) maximum likelihood estimators. *Ann. Inst. Henri Poincaré Probab. Stat.* **42** 273–325.
- BIRGÉ, L. and MASSART, P. (1998). Minimum contrast estimators on sieves: Exponential bounds and rates of convergence. *Bernoulli* **4** 329–375.
- GINÉ, E. and KOLTCHINSKII, V. (2006). Concentration inequalities and asymptotic results for ratio type empirical processes. *Ann. Probab.* **34** 1143–1216. [MR2243881](#)
- GYÖRFI, L., KOHLER, M., KRZYŻAK, A. and WALK, H. (2002). *A Distribution-Free Theory of Nonparametric Regression*. Springer, New York.
- KOLTCHINSKII, V. (2006). Local Rademacher complexities and oracle inequalities in risk minimization. *Ann. Statist.* **34** 2593–2656. [MR2329442](#)
- LE CAM, L. (1973). Convergence of estimates under dimensionality restrictions. *Ann. Statist.* **1** 38–53.
- LE CAM, L. (1975). On local and global properties in the theory of asymptotic normality of experiments. In *Stochastic Processes and Related Topics (Proc. Summer Res. Inst. Statist. Inference for Stochastic Processes, Indiana Univ., Bloomington, Ind., 1974, Vol. 1; Dedicated to Jerzy Neyman)* 13–54. Academic Press, New York.
- LE CAM, L. (1990). Maximum likelihood: An introduction. *Int. Stat. Rev.* **58** 153–171.
- POLLARD, D. (1984). *Convergence of Stochastic Processes*. Springer, New York.
- SART, M. (2017). Estimating the conditional density by histogram type estimators and model selection. *ESAIM, Probab. Stat.* **21** 34–55.
- VAN DE GEER, S. A. (2000). *Applications of Empirical Process Theory. Cambridge Series in Statistical and Probabilistic Mathematics* **6**. Cambridge Univ. Press, Cambridge. [MR1739079](#)
- VAN DER VAART, A. W. and WELLNER, J. A. (1996). *Weak Convergence and Empirical Processes. With Applications to Statistics*. Springer, New York.

UNIVERSITÉ CÔTE D’AZUR
 UMR CNRS 7351
 LABORATOIRE JEAN ALEXANDRE DIEUDONNÉ
 PARC VALROSE
 06108 NICE CEDEX 02
 FRANCE
 E-MAIL: baraud@unice.fr

SORBONNE UNIVERSITÉ ET CNRS
 LABORATOIRE DE PROBABILITÉS, STATISTIQUE
 ET MODÉLISATION (LPSM)
 CASE COURRIER 188
 F-75252 PARIS CEDEX 05
 FRANCE
 E-MAIL: lucien.birge@upmc.fr