

# Estimating a density, a hazard rate, and a transition intensity via the $\rho$ -estimation method

Mathieu Sart

► **To cite this version:**

Mathieu Sart. Estimating a density, a hazard rate, and a transition intensity via the  $\rho$ -estimation method. 2017. <hal-01557973v2>

**HAL Id: hal-01557973**

**<https://hal.archives-ouvertes.fr/hal-01557973v2>**

Submitted on 21 Dec 2017 (v2), last revised 31 Jan 2018 (v3)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# ESTIMATING A DENSITY, A HAZARD RATE, AND A TRANSITION INTENSITY VIA THE $\rho$ -ESTIMATION METHOD

MATHIEU SART

ABSTRACT. We propose a unified study of three statistical settings by widening the  $\rho$ -estimation method developed in [BBS17]. More specifically, we aim at estimating a density, a hazard rate (from censored data), and a transition intensity of a time inhomogeneous Markov process. We relate the performance of  $\rho$ -estimators to deviations of an empirical process. We deduce non-asymptotic risk bounds for an Hellinger-type loss when the models consist, for instance, of piecewise polynomial functions, multimodal functions, or piecewise convex-concave functions. Under convex-type assumptions on the models, maximum likelihood estimators may coincide with  $\rho$ -estimators, and satisfy therefore our risk bounds. However, our results also apply to some models where the maximum likelihood method does not work. Besides, the robustness properties of  $\rho$ -estimators are not, in general, shared by maximum likelihood estimators. Subsequently, we present an alternative way, based on estimator selection, to define a piecewise polynomial estimator. We control the risk of the estimator and carry out some numerical simulations to compare our approach with a more classical one based on maximum likelihood only.

## 1. INTRODUCTION

1.1. **Statistical settings.** In the present paper, we are interesting in estimating a unknown function  $s$  that appears in one of the following frameworks.

**Framework 1** (Density Estimation). *Let  $X$  be a real-valued random variable with density function  $s$  with respect to the Lebesgue measure  $\mu$ . Our aim is to estimate the density  $s$  from the observation of  $n$  independent copies  $X_1, \dots, X_n$  of  $X$ .*

**Framework 2** (Hazard rate estimation for right censored data). *Let  $(T_1, C_1), \dots, (T_n, C_n)$  be  $n$  independent copies of a pair  $(T, C)$  of non-negative random variables. The variable  $C$  may take the value  $+\infty$ . We suppose that  $T$  is independent of  $C$  and that  $T$  admits a density  $f$  with respect to the Lebesgue measure  $\mu$ . The target function is the hazard rate  $s$  defined for  $t \geq 0$  by*

$$s(t) = \frac{f(t)}{\mathbb{P}(T \geq t)}.$$

*The observations are  $(X_i, D_i)_{1 \leq i \leq n}$  where  $X_i = \min\{T_i, C_i\}$  and  $D_i = \begin{cases} 1 & \text{if } T_i \leq C_i, \\ 0 & \text{otherwise.} \end{cases}$*

**Framework 3** (Estimation of the transition intensity of a Markov process). *We consider a (possibly inhomogeneous) Markov process  $\{X_t, t \geq 0\}$  with the following properties:*

---

*Date:* December, 2017.

*2010 Mathematics Subject Classification.* 62G07, 62G35, 62N02, 62M05.

*Key words and phrases.*  $\rho$ -estimator,  $T$ -estimator, robust tests, maximum likelihood, qualitative assumptions, piecewise polynomial estimation.

- The process is cadlag with finite state space, says  $\{0, 1, \dots, m\}$ .
- The state 0 is absorbing.
- Let, for each interval  $I \subset [0, +\infty)$ ,  $A_I$  be the event: “the process jumps at least two times on  $I$ ”. Then,  $\mathbb{P}(A_I) = o(\mu(I))$  when the length  $\mu(I)$  of  $I$  tends to 0.
- The transition time

$$T_{1,0} = \inf \{t > 0, X_{t-} = 1, X_t = 0\},$$

which has values in  $[0, +\infty]$ , is absolutely continuous with respect to the Lebesgue measure  $\mu$  on  $\mathbb{R}$  and satisfies therefore for all Borel set  $A$  of  $\mathbb{R}$ ,

$$\mathbb{P}(T_{1,0} \in A) = \int_A f(u) du,$$

for a suitable non-negative measurable function  $f$ .

We consider an observation interval  $I_{obs} \subset [0, +\infty)$  either of the form  $I_{obs} = [0, T]$  with  $T \in (0, +\infty)$  or  $I_{obs} = [0, +\infty)$ . Our aim is to estimate the transition rate  $s$  from state 1 to 0 defined for  $t > 0$  by

$$s(t) = \frac{f(t)}{\mathbb{P}(X_{t-} = 1)},$$

from the observation of  $n$  independent copies  $\{X_t^{(i)}, t \in I_{obs}\}$  of  $\{X_t, t \in I_{obs}\}$ .

In all these frameworks, we will always suppose that  $n \geq 3$ . Although numerous estimation strategies can be considered, we will rather focus in this paper on a particular method developed in [BBS17] and named “ $\rho$ -estimation”.

**1.2. On  $\rho$ -estimation in framework 1.** We begin by outlining some of the underlying ideas of this estimation procedure. First, the method fits into the scheme of a series of papers using tests to construct estimators. Given two densities  $f$  and  $g$ , a test is, intuitively, a decision rule that decides which one is the best for estimating  $s$ . In order to measure the quality of estimation, we consider the Hellinger distance  $h$ , which is defined for two non-negative integrable functions  $f$  and  $g$  by

$$h^2(f, g) = \frac{1}{2} \int_{\mathbb{R}} \left( \sqrt{f(t)} - \sqrt{g(t)} \right)^2 dt.$$

We try, in  $\rho$ -estimation, to estimate  $h^2(s, f) - h^2(s, g)$ . The smaller this difference, the better  $f$ . Conversely, the larger this difference, the better  $g$ . Unfortunately, it seems difficult to build consistent estimators of  $h^2(s, f) - h^2(s, g)$  with good properties (without additional assumptions on  $s$ ). A way to circumvent the issue, and which follows from [Bar11], is to estimate a good approximation  $T_E(f, g)$  of  $h^2(s, f) - h^2(s, g)$ . We will denote such estimators by  $T(f, g)$ .

We then consider models  $S$ , that is collections of densities, which translate, in mathematical terms, the knowledge we have on the target  $s$ . A model may correspond to different assumptions on  $s$ , such as parametric, regularity, or qualitative ones. When  $S$  only consists of two densities,  $S = \{f, g\}$ , an estimator  $\hat{s}$  on  $S$  may be deduced from the test: we may set  $\hat{s} = g$  when  $T(f, g) > 0$  and  $\hat{s} = f$  when  $T(f, g) < 0$ . When the model  $S$  is more general, several methods exist in the literature to define an estimator. A key theoretical reference is [Bir06]. In  $\rho$ -estimation,  $T(f, g)$  is viewed as an approximate value of  $h^2(s, f) - h^2(s, g)$ . We may then form the criterion  $\gamma(f) = \sup_{g \in S} T(f, g)$  and interpret it as an approximation of  $h^2(s, f) - \inf_{g \in S} h^2(s, g)$ . It then remains to minimize  $\gamma$  to define the  $\rho$ -estimator  $\hat{s}$  (if such a minimizer does not exist, take an approximate minimizer).

Before going any further, we need to mention that we may construct several variants of the  $\rho$ -estimation procedure that may lead to similar theoretical results, at least in density estimation. For instance, there exist several ways to define the estimator  $T(f, g)$  as explained in [BB17]. It turns out that  $\rho$ -estimators (including some variants) satisfy interesting statistical properties. We briefly present below four of them: generality, optimality, robustness, and superminimaxity.

First of all,  $\rho$ -estimators exist on some models  $S$  for which the maximum likelihood method does not work. Numerous examples are known in the literature. A very simple one is

$$(1) \quad S = \{f \mathbb{1}_I, I \text{ is a closed interval of } \mathbb{R} \text{ and } f \text{ a non-increasing density on } I\},$$

where  $\mathbb{1}_I$  denotes the indicator function of  $I$ . In this model, the log likelihood can be made arbitrarily large, and the maximum likelihood estimator does not exist. By contrast, we may define, and study,  $\rho$ -estimators on  $S$ .

The quality of a  $\rho$ -estimator  $\hat{s}$  relies on the behaviour of the error  $T(f, g) - T_E(f, g)$ . This error can be controlled according to different notions that aim at measuring the “complexity”, or “massiveness” of the model  $S$  (entropy with bracketing, metric dimension, covering numbers...). We may deduce an upper-bound  $R_S(n)$  of the maximal risk  $\sup_{s \in S} \mathbb{E}[h^2(s, \hat{s})]$ . This result can then be compared to the minimax bound  $\inf_{\tilde{s}} \sup_{s \in S} \mathbb{E}[h^2(s, \tilde{s})]$ , where the infimum is taken over all estimators  $\tilde{s}$  with values in  $S$ . The rate of convergence of  $\inf_{\tilde{s}} \sup_{s \in S} \mathbb{E}[h^2(s, \tilde{s})]$  to 0 is usually called the optimal minimax rate of convergence. Yet,  $R_S(n)$  achieves this rate, up to possible logarithmic factors, in all cases we know.

This minimax point of view supposes that  $s$  does belong to  $S$ . Such an assumption corresponds to a perfect modelling of the statistical problem, which is scarcely the case in practice. It makes therefore more sense to study the risk of the estimator  $\hat{s}$  not only when  $s$  lies in  $S$  but more generally when  $s$  is close to the model  $S$ . It turns out that the Hellinger quadratic risk of a  $\rho$ -estimator  $\hat{s}$  can be bounded above by

$$\mathbb{E}[h^2(s, \hat{s})] \leq C \inf_{f \in S} h^2(s, f) + R_S(n) \quad \text{whatever the density } s,$$

where  $C$  is a universal constant (that is a number). This inequality asserts that a small error in the choice of the model  $S$  induces a small error in the estimation of  $s$ . This is a robustness property. Such a property is not shared in general by the maximum likelihood estimator. It may indeed perform very poorly when  $s \notin S$  but is close to  $S$  (when this closeness is measured by the Hellinger distance, see Section 2.3 of [Bir06] for an example).

The rate given by  $R_S(n)$  stands for the worst-case rate over all densities  $s$  of  $S$ . This rate may therefore be very pessimistic in the sense that the estimation may be much faster for some densities  $s \in S$ . One may actually refine the preceding risk bound to take into account this phenomenon (named superminimaxity in [BB16]). More precisely, it is shown in [BB16] a non-asymptotic risk bound of the form

$$(2) \quad \mathbb{E}[h^2(s, \hat{s})] \leq C' \inf_{f \in \bar{S}} \{h^2(s, f) + R_S(f, n)\} \quad \text{whatever the density } s,$$

where  $C'$  is a universal constant,  $\bar{S}$  a suitable subset of  $S$ , and where  $R_S(f, n)$  tends to 0 at a (usually) faster rate than  $R_S(n)$ . For illustration purposes, consider the model  $S$  defined by (1). Then,  $\bar{S}$  consists of piecewise constant densities belonging to  $S$ ,

$$R_S(f, n) = \frac{d(f)}{n} \log_+^3 \left( \frac{n}{d(f)} \right),$$

$d(f)$  is the number of pieces of  $f$ , and  $\log_+ x = \max\{\log x, 1\}$ . In particular, when the support of  $s$  is an interval on which  $s$  is non-increasing and piecewise constant, that is when  $s \in \bar{S}$ , the rate of estimation is parametric (up to some power of  $\log_+(n/d(f))$ ). If we now suppose that  $s$  belongs to  $S \setminus \bar{S}$ , some computations allows to bound (2) from above by  $C_s(\log^2 n)n^{-2/3}$ , where  $C_s$  only depends on  $s$ . This corresponds, up to a logarithmic factor, to the expected rate of convergence. The previous reasoning is not only valid for this particular model  $S$  but is more general and holds true for other models  $S$  corresponding to qualitative assumptions on the density (for instance,  $s$  may be piecewise monotone or  $\sqrt{s}$  may be piecewise convex-concave).

There are moreover two additional properties of  $\rho$ -estimators we now briefly mention. First,  $\rho$ -estimators can be related to maximum likelihood ones. Second, it is possible to introduce penalties into the criterion  $\gamma$ , leading to penalized  $\rho$ -estimators and allowing to cope with model selection.

**1.3. On hazard rate and transition intensity estimation.** In this paper, we propose to extend the scope of  $\rho$ -estimation to these two statistical settings. The first one, namely hazard rate estimation, frequently appears in different domains such as reliability or survival analysis. Typically, in medical studies,  $T$  may represent the lifetime of a patient, and the hazard rate  $s$  at time  $t$ ,

$$\begin{aligned} s(t) &= \frac{f(t)}{\mathbb{P}(T \geq t)}, \\ &= \lim_{h \rightarrow 0} \frac{\mathbb{P}(t \leq T \leq t+h \mid T \geq t)}{h}, \end{aligned}$$

measures the tendency of dying just after  $t$ , given survival to time  $t$ . In practice, some patients may leave the study before dying, which makes the data censored. The random variable  $C$  then gives the time of leaving and  $D = \mathbb{1}_{T \leq C}$  indicates whether the patient dies ( $D = 1$ ) or leaves the study ( $D = 0$ ).

The problem of transition intensity estimation of a Markov process may also be encountered in various domains. For example, in medical trials, a Markov process  $\{X_t, t > 0\}$  may be used to model the evolution of a disease, the state 0 representing (for instance) the death of the patient. The transition rate  $s$  at time  $t$ ,

$$\begin{aligned} s(t) &= \frac{f(t)}{\mathbb{P}(X_{t-} = 1)}, \\ &= \lim_{h \rightarrow 0} \frac{\mathbb{P}(X_{t+h} = 0 \mid X_{t-} = 1)}{h}, \end{aligned}$$

has similar interpretation than the hazard rate: it measures the risk of dying just after  $t$ , given the disease is in state 1 at time  $t-$ . This framework is actually more general than the one of hazard rate estimation (when the data are uncensored) as  $s$  coincides with the hazard rate of  $T$  when the Markov process is defined by  $X_t = \mathbb{1}_{T \geq t}$ .

In the literature, numerous estimators have been proposed to deal with (at least) one of these two frameworks. We may cite wavelet estimators, Kernel estimators, maximum likelihood estimators, procedures based on  $\mathbb{L}^2$  contrasts. . . However, non-asymptotic studies seem be rather scarce. We refer to [BC05, RB06, BC08, Pla09, AD10] for results concerning procedures based on (penalized)  $\mathbb{L}^2$  contrasts. We may cite [vdG95, DR02] for a study of non-asymptotic properties of maximum likelihood estimators. We refer to [BB09] for results concerning a selection rule based on pairwise comparisons of histogram type estimators.

**1.4. A generalized procedure.** As in [BB09, AD10, Bar11], we consider in the present paper the problem of estimating the intensity of a random measure. This is a very convenient statistical setting to study the three frameworks in a unified way.

We measure the risks of our estimators by means of a (possibly random) Hellinger-type distance  $h$  adapted to the framework. In framework 1,  $h$  is the usual Hellinger distance, in framework 2,

$$h^2(f, g) = \frac{1}{2} \int_0^\infty \left( \sqrt{f(t)} - \sqrt{g(t)} \right)^2 \left( \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{X_i \geq t} \right) dt,$$

and in framework 3,

$$h^2(f, g) = \frac{1}{2} \int_{I_{\text{obs}}} \left( \sqrt{f(t)} - \sqrt{g(t)} \right)^2 \left( \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{X_{t^-}^{(i)} = 1} \right) dt.$$

The quality of an estimator  $\hat{s}$  is therefore assessed by  $h^2(s, \hat{s})$ : the smaller  $h^2(s, \hat{s})$ , the better the estimator.

We define an approximation  $T_E(f, g)$  of  $h^2(s, f) - h^2(s, g)$  and an estimator  $T(f, g)$  of  $T_E(f, g)$ . Given a model  $S$ , that is a collection of possible candidates for estimating  $s$ , we define our estimator  $\hat{s}$  on  $S$  as a minimizer, or more precisely as an approximate minimizer of  $\gamma(f) = \sup_{g \in S} T(f, g)$ . Similarly to framework 1, bounding above the Hellinger-type risk  $h(s, \hat{s})$  of a  $\rho$ -estimator  $\hat{s}$  requires to control the deviations of the error  $T(f, g) - T_E(f, g)$ . We carry out a uniform exponential inequality to control these deviations. We deduce risk bounds for  $\rho$ -estimators over a particular class of models  $S$ . This class includes, for instance, models consisting of piecewise polynomial functions, multimodal functions, piecewise convex-concave functions or piecewise log-concave functions. We establish in the three frameworks a risk bound akin to the one obtained in density estimation by [BB16]. Actually, in density estimation, the control of the empirical process is a bit more accurate than in this last paper, leading to a slightly sharper estimation term  $R_S(f, n)$  (although probably not optimal). Moreover, the closeness of the results in the three frameworks allows to transfer the rates of convergence obtained in [BB16] in framework 1 when  $s$  is multimodal or when  $\sqrt{s}$  is piecewise convex-concave to frameworks 2 and 3 (with a minor improvement).

**1.5. On maximum likelihood estimation.** The  $\rho$ -estimation procedure differs from that of maximum likelihood. Nevertheless, the two approaches are very close in some situations.

This phenomenon may be understood by looking at the local behaviour of the estimator  $T(f, g)$ . Indeed, under some assumptions,  $T(f, g)$  roughly behaves as the difference of log likelihoods  $L(g) - L(f)$  when  $f$  and  $g$  are very close densities in framework 1. If we could replace  $T(f, g)$  by  $L(g) - L(f)$ ,  $\gamma(f) = \sup_{g \in S} T(f, g)$  would become  $\sup_{g \in S} L(g) - L(f)$ , and a minimizer of  $\gamma(\cdot)$  would be a maximum likelihood estimator.

We show that  $\rho$ -estimation may coincide with maximum likelihood estimation when the model is either convex or of the form  $S = \{f^2, f \in \mathcal{F}\}$  where  $\mathcal{F}$  is convex and consists of non-negative functions. In particular, a maximum likelihood estimator on such a model  $S$  is a  $\rho$ -estimator. A similar result for convex sets of densities was obtained by Su Weijie in the context of framework 1 and was recently included in [BB17]. Consequently, our theoretical risk bounds, which apply to  $\rho$ -estimators, may also apply to maximum likelihood ones for these models (in the three frameworks).

**1.6. Estimator selection.** The practical computation of  $\rho$ -estimators seems unfortunately to be numerically out of reach in numerous models. In some cases, a model  $S$  may be written as a union  $S = \bigcup_{m \in \mathcal{M}} S_m$  of models  $S_m$  satisfying the convex-type assumptions described in the preceding section. Thereby, one may maximize the likelihood to define a  $\rho$ -estimator  $\hat{s}_m$  on  $S_m$ . All that remains then is to select an estimator among  $\{\hat{s}_m, m \in \mathcal{M}\}$ . This solution may be more numerically friendly (when  $\mathcal{M}$  is not too large), but the quality of selected estimator needs to be shown.

In the present paper, we will not address this problem in general but rather focus on the particular model  $S_{\ell,r}$  consisting of non-negative piecewise polynomial functions of degree at most  $r$  and based on  $\ell$  consecutive intervals. Although maximum likelihood estimators do not exist on  $S_{\ell,r}$ ,  $\rho$ -estimators do exist, and we even control their Hellinger-type risks. Unfortunately, we do not know how to build these  $\rho$ -estimators in practice. Alternatively, we may consider collections  $m$  of intervals, and define models  $S_{\ell,r,m} \subset S_{\ell,r}$  consisting of functions which are polynomial on each interval  $I$  of  $m$ . Then,  $S_{\ell,r} = \bigcup_{m \in \mathcal{M}_\ell} S_{\ell,r,m}$  where the union is taken over a suitable (infinite) family  $\mathcal{M}_\ell$  of collections  $m$ . For each  $m \in \mathcal{M}_\ell$ , we may define a  $\rho$ -estimator  $\hat{s}_m$  on the convex model  $S_{\ell,r,m}$  by maximizing the likelihood. Selecting among all the estimators  $\hat{s}_m$  is theoretically feasible but does not yield a practical procedure as  $\mathcal{M}_\ell$  is infinite. This is the reason why we will replace  $\mathcal{M}_\ell$  by a finite, but usually very large, family  $\widehat{\mathcal{M}}_\ell \subset \mathcal{M}_\ell$ . We will carry out a new procedure, inspired from [Sar14], to select among the estimators  $\hat{s}_m, m \in \widehat{\mathcal{M}}_\ell$ . Although the large cardinal of  $\widehat{\mathcal{M}}_\ell$ , dynamic programming makes it possible the computation of the selected estimator in polynomial time (at least when one knows how to maximize a likelihood on a convex model). We prove an oracle inequality for the selected estimator from which, we deduce, when  $r = 0$ , a risk bound very similar to the one we would obtain for the  $\rho$ -estimator. Besides, we carry out a numerical study in which we compare, in the context of density estimation, our procedure with a selection rule based on maximum likelihood only.

We finally explain how we can modify this procedure to select adaptively the number  $\ell$  of pieces from the data. In particular, we show that we can build an estimator that performs well when  $s$  belongs, or is close to, the model  $S_r = \bigcup_{\ell=1}^{\infty} S_{\ell,r}$ . More precisely, the risk bound we get corresponds to the one we would obtain for the best estimator of the family  $\{\hat{s}_{\ell,r}, \ell \geq 1\}$  where  $\hat{s}_{\ell,r}$  denotes the  $\rho$ -estimator built on  $S_{\ell,r}$  (up to mild modifications).

**1.7. Organization of the paper.** We carry out in Section 2 the general statistical setting that encompasses the three frameworks. We then explain the estimation procedure and relate it to the maximum likelihood one. In Section 3, we present the probabilistic tool that enables us to control the risk of  $\rho$ -estimators. We then present the required assumptions on the models as well as our main result on the theoretical performances of  $\rho$ -estimators. In Section 4, we deal with estimator selection to define a piecewise polynomial estimator as explained in Section 1.6. Section 5 is devoted to numerical simulations. The proofs are deferred to Section 6.

## 2. THE $\rho$ -ESTIMATION METHOD

**2.1. Statistical setting and notations.** We consider an abstract probability space  $(\Omega, \mathcal{E}, \mathbb{P})$  on which are defined the random variables appearing in the different frameworks. We associate to each framework, and each borel set  $A \in \mathcal{B}(\mathbb{R})$  two random variables  $N(A)$  and  $M(A)$ . More precisely,

we set in density estimation,

$$N(A) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_A(X_i), \quad M(A) = \mu(A),$$

and in hazard rate estimation,

$$N(A) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_A(X_i) \mathbb{1}_{D_i=1}, \quad M(A) = \frac{1}{n} \sum_{i=1}^n \int_A \mathbb{1}_{X_i \geq t} \mathbb{1}_{[0, +\infty)}(t) dt.$$

In framework 3, we define the jump time of the  $i^{\text{th}}$  process

$$T_{1,0}^{(i)} = \inf \left\{ t > 0, X_{t-}^{(i)} = 1, X_t^{(i)} = 0 \right\},$$

and consider

$$N(A) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{T_{1,0}^{(i)} \in A} \mathbb{1}_{I_{\text{obs}}}(T_{1,0}^{(i)}), \quad M(A) = \frac{1}{n} \sum_{i=1}^n \int_A \mathbb{1}_{X_{t-}^{(i)}=1} \mathbb{1}_{I_{\text{obs}}}(t) dt.$$

These formulas define two random measures  $N$  and  $M$  on  $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$  such that

$$\mathbb{E}[N(A)] = \mathbb{E} \left[ \int_A s(t) dM(t) \right] \quad \text{for all } A \in \mathcal{B}(\mathbb{R}).$$

In each of the frameworks, the statistical problem may be reduced to that of estimating  $s$  from the observation of the random measures  $N$  and  $M$ .

As explained in the introduction, we will evaluate the quality of the estimators by using an Hellinger-type loss. This Hellinger-type distance  $h$  can be written simultaneously in the three statistical settings as

$$h^2(f, g) = \frac{1}{2} \int_{\mathbb{R}} \left( \sqrt{f(t)} - \sqrt{g(t)} \right)^2 dM(t),$$

for all non-negative measurable functions  $f$ , and  $g$  which are integrable with respect to the measure  $M$ .

We now introduce some notations that will be used all along the paper. We define  $\mathbb{R}_+ = [0, +\infty)$ , and set for  $x, y \in \mathbb{R}$ ,  $x \wedge y = \min(x, y)$ ,  $x \vee y = \max(x, y)$ . The positive part of a real valued function  $f$  is denoted by  $f_+$  and its negative part by  $f_-$ . The distance between a point  $x$  and a set  $A$  in a metric space  $(E, d)$  is denoted by  $d(x, A) = \inf_{y \in A} d(x, y)$ . We denote the cardinal of a set  $A$  by  $|A|$ , and its complement by  $A^c$ . We set  $\log_+ x = \max\{\log x, 1\}$  for all  $x > 0$ . The notations  $c, c', C, C', \dots$  are for the constants. These constants may change from line to line.

**2.2. Heuristics.** Let  $\mathcal{S} = \mathbb{L}_+^1(\mathbb{R}, \mu)$  be the cone of non-negative Lebesgue integrable functions in frameworks 1 and 3, and  $\mathcal{S}$  be the cone of measurable non-negative functions which are locally integrable with respect to  $\mu$  in framework 2. Let now  $S$  be a subset of  $\mathcal{S}$ . Such set will be named model. Our aim is to define an estimator  $\hat{s}$  with values in  $S$  such that  $h(s, \hat{s})$  is as small as possible.

Consider two arbitrary functions  $f, g$  of  $\mathcal{S}$ . As explained in the introduction, we begin by defining an approximation  $T_E(f, g)$  of  $h^2(s, f) - h^2(s, g)$ .



Let  $\psi$  be the real-valued function defined for  $x \geq 0$  by  $\psi(x) = \frac{\sqrt{x}-1}{\sqrt{x+1}}$ , and  $\psi(+\infty) = 1$ . For  $f, g \in \mathcal{S}$ , we set

$$T_E(f, g) = \int_{\mathbb{R}} \psi \left( \frac{g(x)}{f(x)} \right) s(x) \, dM(x) - \frac{1}{4} \int_{\mathbb{R}} (g(x) - f(x)) \, dM(x).$$

In this definition, and throughout the paper, we use the conventions  $0/0 = 1$  and  $x/0 = +\infty$  for all  $x > 0$ . Some computations show:

**Lemma 1.** *For all  $f, g \in \mathcal{S}$ ,*

$$(3) \quad \frac{1}{3}h^2(s, f) - 3h^2(s, g) \leq T_E(f, g) \leq 3h^2(s, f) - \frac{1}{3}h^2(s, g).$$

In particular, if  $T_E(f, g)$  is non-negative, then  $h^2(s, g) \leq 9h^2(s, f)$ . Conversely, if  $T_E(f, g)$  is non-positive, then  $h^2(s, f) \leq 9h^2(s, g)$ . In other words, the sign of  $T_E(f, g)$  allows us to know which function among  $f, g$  is the closest of  $s$  (up to a multiplicative constant).

Let  $S$  be a model and  $f \in S$ . We are interested in evaluating  $h^2(s, f) - h^2(s, S)$ . The smaller this number, the better  $f$ . As  $T_E(f, g)$  is roughly of the order of  $h^2(s, f) - h^2(s, g)$ , it is natural to approximate  $h^2(s, f) - h^2(s, S)$  by  $\gamma_E(f) = \sup_{g \in S} T_E(f, g)$  and to study the properties of the minimizers of  $\gamma_E$ .

We deduce from the above lemma that for all  $f \in S$ ,

$$\frac{1}{3}h^2(s, f) - 3h^2(s, S) \leq \gamma_E(f) \leq 3h^2(s, f) - \frac{1}{3}h^2(s, S).$$

Minimizing  $\gamma_E$  over  $S$  yields a function  $\bar{f} \in S$  (assuming such a function exists) such that,

$$\frac{1}{3}h^2(s, \bar{f}) - 3h^2(s, S) \leq \gamma_E(\bar{f}) \leq \inf_{f \in S} \gamma_E(f) \leq 3 \inf_{f \in S} h^2(s, f) - \frac{1}{3}h^2(s, S) = \frac{8}{3}h^2(s, S).$$

Therefore,  $h^2(s, \bar{f}) \leq 17h^2(s, S)$ , which means that  $\bar{f}$  is, up to a multiplicative constant, the closest function of  $s$  among the ones of  $S$ .

The approximation  $T_E(f, g)$  is certainly unknown in practice as it involves  $s$ . It can however be suitably estimated by

$$(4) \quad T(f, g) = \int_{\mathbb{R}} \psi \left( \frac{g(x)}{f(x)} \right) \, dN(x) - \frac{1}{4} \int_{\mathbb{R}} (g(x) - f(x)) \, dM(x).$$

Similarly,  $\gamma_E(f)$  is unknown but can be estimated by  $\gamma(f) = \sup_{g \in S} T(f, g)$ . It then remains to minimize this criterion to define the estimator as described below.

**2.3. The procedure.** Let for  $f, g \in \mathcal{S}$ ,  $T(f, g)$  be given by (4). Let  $S$  be a model and  $\gamma(f) = \sup_{g \in S} T(f, g)$ . Any estimator  $\hat{s} \in S$  satisfying

$$(5) \quad \gamma(\hat{s}) \leq \inf_{f \in S} \gamma(f) + 1/n$$

is called  $\rho$ -estimator.

**Remark 1.** We do not assume that  $S$  consists of densities in framework 1 for more flexibility in the choice of models. Likewise, the functions of  $S$  may not be hazard rates in framework 2, or transition intensities in framework 3.

The procedure may also be used to estimate the restriction of  $s$  on an interval  $I$ . Indeed, let  $N'$  be defined by  $N'(A) = N(A \cap I)$  for all  $A \in \mathcal{B}(\mathbb{R})$ . Then,  $\mathbb{E}[N'(A)] = \mathbb{E}[\int_A s \mathbb{1}_I dM]$  and the target function becomes  $s \mathbb{1}_I$ . Let now  $\mathcal{F}$  be a collection of functions and  $S$  be a model of the form  $S = \{f \mathbb{1}_I, f \in \mathcal{F}\}$ . Since the functions of  $S$  vanish outside  $I$ , we may replace  $N$  in the procedure by  $N'$  without changing the estimator. Thereby, when all functions of  $S$  vanish outside  $I$ , the estimator  $\hat{s}$  actually estimates  $s \mathbb{1}_I$ .

The procedures are quite close in the three frameworks. Indeed, let us consider framework 2 and suppose that the data are not censored (that is  $C = +\infty$ ). Define  $G_n(t) = n^{-1} \sum_{i=1}^n \mathbb{1}_{X_i \geq t}$  and consider a collection  $S$  consisting of non-negative integrable functions. We may use the procedure in density estimation with the (random) model  $S' = \{f G_n, f \in S\}$  to estimate the density of  $X$ . This leads to an estimator of the form  $\hat{s} G_n$  with  $\hat{s} \in S$ . Then,  $\hat{s}$  is also a  $\rho$ -estimator on  $S$  in framework 2. A similar reasoning applies to framework 3.

Remark 2. Two ingredients are required to define the procedure. First, we need to approximate  $h^2(s, f) - h^2(s, g)$  by a quantity  $T_E(f, g)$  that satisfies an inequality akin to (3). Second, we need a random variable  $T(f, g)$  that can be computed in practice and that is close enough to  $T_E(f, g)$  (for more details about the meaning of “close enough”, we refer to Section 3.1). When  $f$  and  $g$  are supposed to be densities in framework 1,  $T_E(f, g)$  becomes  $\int_{\mathbb{R}} \psi(g/f) s d\mu$ . We then recover an approximation of  $h^2(s, f) - h^2(s, g)$  that appears in [BB17].

**2.4. Connection with maximum likelihood estimation.** The  $\rho$ -estimation method may be related to that of maximum likelihood. Since a model  $S$  may not only consist of densities in framework 1, a known solution to define the maximum likelihood estimator is to add a Lagrange term in the log likelihood. More precisely, we define in the three frameworks

$$(6) \quad L(f) = \int_{\mathbb{R}} \log f dN - \int_{\mathbb{R}} f dM \quad \text{for all } f \in \mathcal{S},$$

and call maximum likelihood estimator any estimator maximizing  $L(\cdot)$  on  $S$ . In the above formula, and throughout the paper, the convention  $\log 0 = -\infty$  is used. In framework 2,  $L(f)$  may also be interpreted as a log likelihood (it is equal to the usual log likelihood when  $f$  is a hazard rate, up to some terms constant in  $f$ ). The same is true for framework 3, see the literature of counting processes *e.g.* equation (3.2) of [Ant89] (using that  $s$  is an Aalen’s multiplicative intensity).

We may write

$$T(f, g) = \int_{\mathbb{R}} \tanh\left(\frac{\log g - \log f}{4}\right) dN - \frac{1}{4} \int_{\mathbb{R}} (g - f) dM \quad \text{for all } f, g \in \mathcal{S}.$$

As  $\tanh(x) \simeq x$  when  $x \simeq 0$ , we deduce that if  $\tilde{s}$  maximizes  $L$  and  $g \simeq \tilde{s}$ ,

$$\begin{aligned} T(\tilde{s}, g) &\simeq \frac{1}{4} \left( \int_{\mathbb{R}} \log g dN - \int_{\mathbb{R}} \log \tilde{s} dN \right) - \frac{1}{4} \left( \int_{\mathbb{R}} g dM - \int_{\mathbb{R}} \tilde{s} dM \right) \\ &\simeq \frac{1}{4} (L(g) - L(\tilde{s})). \end{aligned}$$

Thereby,  $T(\tilde{s}, g)$  is likely non-positive. Under suitable properties of  $S$ , this result does not only occur when  $g \simeq \tilde{s}$ , but also for all  $g \in S$ , which implies that  $\gamma(\tilde{s}) = 0$ . In particular,  $\tilde{s}$  is a  $\rho$ -estimator.

**Theorem 1.** *Suppose that  $S$  is a convex subset of  $\mathcal{S}$ . Let  $\mathcal{X}$  be a subset of  $\mathbb{R}$  such that  $\{x \in \mathbb{R}, f(x) \neq 0\} \subset \mathcal{X}$  for all  $f \in S$ . Define*

$$L_{\mathcal{X}}(f) = \int_{\mathcal{X}} \log f \, dN - \int_{\mathcal{X}} f \, dM \quad \text{for all } f \in S,$$

and suppose that  $\sup_{g \in S} L_{\mathcal{X}}(g) \notin \{-\infty, +\infty\}$ .

*If there exists an estimator  $\tilde{s} \in S$  such that  $L_{\mathcal{X}}(g) \leq L_{\mathcal{X}}(\tilde{s})$  for all  $g \in S$ , then  $\gamma(\tilde{s}) = 0$  and  $\tilde{s}$  is a  $\rho$ -estimator. Conversely, assume that there exists a  $\rho$ -estimator  $\hat{s} \in S$  such that  $\gamma(\hat{s}) = 0$ . Then, for all  $g \in S$ ,  $L_{\mathcal{X}}(g) \leq L_{\mathcal{X}}(\hat{s})$ , and  $\hat{s}$  maximizes  $L_{\mathcal{X}}(\cdot)$  over  $S$ .*

When  $\mathcal{X} = \mathbb{R}$ ,  $L_{\mathcal{X}} = L$ , which means that results on maximum likelihood estimators may be derived from that of  $\rho$ -estimators and vice versa. We recover the result of Su Weijie when  $S$  consists of densities in framework 1. Using sets  $\mathcal{X}$  not equal to  $\mathbb{R}$  may be of interest to remove some observations that would make the log likelihood identically equal to  $-\infty$ . In that case, we rather estimate the restriction of  $s$  to  $\mathcal{X}$  as illustrated in the example below.

We consider the convex model  $S$  in framework 1 defined by

$$(7) \quad S = \{f \mathbb{1}_{(0,+\infty)}, f \text{ is a non-increasing function of } \mathcal{S} \text{ on } \mathbb{R}\}.$$

When the random variables  $X_i$  are positive, which in particular holds true almost surely if  $s$  does belong to  $S$ , the maximum likelihood estimator exists on  $S$  and is known as the Grenander estimator, see [Gre56]. We deduce from the above theorem with  $\mathcal{X} = \mathbb{R}$  that this estimator is, in this case, a  $\rho$ -estimator. When some of the random variables  $X_i$  are non-positive,  $L(g) = -\infty$  for all  $g \in S$ , and we cannot maximize  $L(\cdot)$  over  $S$  to design an estimator. However, the  $\rho$ -estimation approach works and still coincides with the maximum likelihood one, up to minor modifications. Indeed, in this case, the preceding theorem can be used with  $\mathcal{X} = (0, +\infty)$ . Then,  $L_{\mathcal{X}}(f)$  takes the form

$$L_{\mathcal{X}}(f) = \frac{1}{n} \sum_{\substack{i \in \{1, \dots, n\} \\ X_i > 0}} \log f(X_i) - \int_0^{\infty} f(t) \, dt \quad \text{for all } f \in S.$$

Let  $\tilde{s}$  be the Grenander estimator based on the random variables  $X_1, \dots, X_n$  that are positive. This estimator is a density and maximizes the map

$$f \mapsto \frac{1}{n_0} \sum_{\substack{i \in \{1, \dots, n\} \\ X_i > 0}} \log f(X_i)$$

over the densities  $f$  of  $S$ , where  $n_0$  is the number of positive random variables among  $X_1, \dots, X_n$ . One can verify that the estimator that maximizes  $L_{\mathcal{X}}(\cdot)$  over  $S$ , and which is thus the  $\rho$ -estimator on  $S$ , is  $\hat{s} = (n_0/n)\tilde{s}$ . Note that  $\int_{\mathbb{R}} \hat{s} \, d\mu = n_0/n$ , which means that the  $\rho$ -estimator is not a density unless all the observations  $X_i$  are positive. This is due to the fact that  $\hat{s}$  estimates in this case the restriction of  $s$  to  $(0, +\infty)$  (which is not a density when some observations  $X_i$  are negative).

It is sometimes convenient to consider models  $S$  of the form  $S = \{f^2, f \in \mathcal{F}\}$  where  $\mathcal{F}$  consists of non-negative functions. The set  $\mathcal{F}$  can then be interpreted as a translation of the knowledge one has on  $\sqrt{s}$ . For instance, if  $\mathcal{F}$  denotes the set of non-negative concave functions on  $[0, +\infty)$  vanishing on  $(-\infty, 0)$ , the assumption  $s \in S$  means that  $\sqrt{s}$  is concave on  $[0, +\infty)$  with support in  $[0, +\infty)$ . It turns out that the connection between  $\rho$ - and maximum likelihood estimators established by Theorem 1 remains valid when the convexity assumption is put on  $\mathcal{F}$  instead of  $S$ .

**Theorem 2.** *Let  $\mathcal{F}$  be a convex set of non-negative functions such that  $S = \{f^2, f \in \mathcal{F}\}$  is included in  $\mathcal{S}$ . Let  $\mathcal{X}$  be a subset of  $\mathbb{R}$  such that  $\{x \in \mathbb{R}, f(x) \neq 0\} \subset \mathcal{X}$  for all  $f \in \mathcal{F}$ . Then, if  $\sup_{g \in S} L_{\mathcal{X}}(g) \notin \{-\infty, +\infty\}$ , the conclusions of Theorem 1 apply to  $S$ : any maximizer  $\hat{s} \in S$  of  $L_{\mathcal{X}}(\cdot)$  on  $S$  vanishes  $\gamma(\cdot)$ , and any  $\hat{s} \in S$  vanishing  $\gamma(\cdot)$  maximizes  $L_{\mathcal{X}}(\cdot)$  over  $S$ .*

### 3. RISK BOUNDS OF $\rho$ -ESTIMATORS

**3.1. A uniform exponential inequality.** We recall that the definition of  $\rho$ -estimators is based on the minimization of a criterion  $\gamma$  on  $S$ . This criterion  $\gamma$  uses the approximation  $T(f, g) \simeq T_E(f, g)$  where  $f, g \in S$  as explained in Section 2.2. Bounding above the risk of the  $\rho$ -estimator requires to bound above the error due to the approximation of  $T_E$  by  $T$ .

We introduce for any bounded function  $\varphi \in \mathcal{S}$ , the random variable

$$Z(\varphi) = \int_{\mathbb{R}} \varphi(x) dN(x) - \int_{\mathbb{R}} \varphi(x) s(x) dM(x).$$

This variable is centered in each of the statistical settings. Note that  $Z(\varphi)$  measures the approximation error of  $T_E(f, g)$  by  $T(f, g)$  when  $\varphi = \psi(g/f)$ . The theorem below allows to control the deviations of  $Z(\varphi)$  under very simple (but somewhat restrictive) assumptions.

**Theorem 3.** *Let  $\mathcal{F} \subset \mathcal{S}$  be a set of functions  $\varphi$  such that  $|\varphi(x)| \leq 1$  for all  $\varphi \in \mathcal{F}$ ,  $x \in \mathbb{R}$ . Assume that for all  $t \in (0, 1)$ , the sets  $\{x \in \mathbb{R}, \varphi_+(x) > t\}$  and  $\{x \in \mathbb{R}, \varphi_-(x) > t\}$  are unions of at most  $d$  intervals ( $d \geq 1$ ). Let, for  $\varphi \in \mathcal{F}$ ,*

$$v(\varphi) = \int_{\mathbb{R}} \varphi^2(x) s(x) dM(x).$$

*Then, there exists for all  $\xi > 0$  an event which holds true with probability larger than  $1 - e^{-n\xi}$  and on which: for all  $\varphi \in \mathcal{F}$ ,*

$$(8) \quad |Z(\varphi)| \leq C \left\{ \sqrt{v(\varphi) \log_+(1/v(\varphi)) \left( \frac{d \log_+(n/d)}{n} + \xi \right)} + \frac{d \log_+(n/d)}{n} + \xi \right\}.$$

*Moreover, for all  $\varepsilon \in (0, 1]$ ,*

$$(9) \quad |Z(\varphi)| \leq \varepsilon v(\varphi) + C'_\varepsilon \left\{ \frac{d \log_+^2(n/d)}{n} + \xi \log_+(1/\xi) \right\}.$$

*In the above inequalities,  $C$  is a universal constant while  $C'_\varepsilon$  only depends on  $\varepsilon$ .*

In this theorem, we relate the complexity of  $\mathcal{F}$  to that of a collection of sets of the form  $A = \{x \in \mathbb{R}, \varphi_+(x) > t\}$  or  $A = \{x \in \mathbb{R}, \varphi_-(x) > t\}$ . We then measure the complexity of this family of sets by using the notion of “union of intervals”. This condition may seem stringent but is general enough to control the risks of  $\rho$ -estimators in several models  $S$  of interest (see the next section). Moreover, this theorem will allow us to perform (polynomial) estimator selection in Section 4. It also allows to refine some results already proved in density estimation.

As a by-product of the proof of the theorem, we get the following proposition which may be of independent interest:

**Proposition 4.** *Consider framework 1 and an at most countable set  $\mathcal{F} \subset \mathcal{S}$  of functions  $\varphi$  such that  $|\varphi(x)| \leq 1$  for all  $x \in \mathbb{R}$ ,  $\varphi \in \mathcal{F}$ . Let for  $t \in (0, 1)$ ,  $\mathcal{A}_t$  be the collection of sets defined by*

$$\mathcal{A}_t = \{\{x \in \mathbb{R}, \varphi_+(x) > t\}, \varphi \in \mathcal{F}\} \cup \{\{x \in \mathbb{R}, \varphi_-(x) > t\}, \varphi \in \mathcal{F}\},$$

and  $S_{\mathcal{A}_t}(2n)$  be the Vapnik-Chervonenkis shatter coefficient defined by

$$S_{\mathcal{A}_t}(2n) = \max_{x_1, \dots, x_{2n} \in \mathbb{R}} |\{\{x_1, \dots, x_{2n}\} \cap A, A \in \mathcal{A}_t\}|.$$

Let  $\sigma^2 = \sup_{\varphi \in \mathcal{F}} \mathbb{E}[\varphi^2(X)]$  and  $p_t = \sup_{\varphi \in \mathcal{F}} \max\{\mathbb{P}(\varphi_+ > t), \mathbb{P}(\varphi_- > t)\}$ . Then, there exist universal constants  $C, C'$  such that

$$\begin{aligned} \mathbb{E} \left[ \sup_{\varphi \in \mathcal{F}} |Z(\varphi)| \right] &\leq \frac{C}{\sqrt{n}} \inf_{\eta \in (0,1)} \left\{ \sigma \sqrt{\int_{\eta}^1 \frac{\log_+ |S_{\mathcal{A}_t}(2n)|}{t} dt} + \int_0^{\eta} \sqrt{p_t \log_+ |S_{\mathcal{A}_t}(2n)|} dt \right\} \\ &\quad + \frac{C}{n} \int_0^1 \log_+ |S_{\mathcal{A}_t}(2n)| dt \\ &\leq C' \left[ \sigma \sqrt{\frac{\sup_{t \in (0,1)} \log_+ |S_{\mathcal{A}_t}(2n)| \log_+(1/\sigma)}{n}} + \frac{\sup_{t \in (0,1)} \log_+ |S_{\mathcal{A}_t}(2n)|}{n} \right]. \end{aligned}$$

In the literature, numerous bounds on this expectation were showed under different assumptions on  $\mathcal{F}$ . We refer to [GG01], [GK06], Chapter 13 of [BLM13], [Bar16] and the references therein. In the last paper, the complexity of  $\mathcal{F}$  is also measured through the one of  $\mathcal{A}_t$ , and this corollary is similar to his result (when the  $X_i$  are identically distributed). Actually, our bound is sharper when  $\mathcal{A}_t$  is Vapnik-Chervonenkis with dimension  $d$  (apart from constants). Such an assumption corresponds to a notion of (weak) VC-major class (see [Bar16]). In that case, Sauer lemma [Sau72] implies

$$(10) \quad \mathbb{E} \left[ \sup_{\varphi \in \mathcal{F}} |Z(\varphi)| \right] \leq C'' \left[ \sigma \sqrt{\frac{d \log_+(n/d) \log_+(1/\sigma)}{n}} + \frac{d \log_+(n/d)}{n} \right],$$

where  $C''$  is a number. If we put aside the constant  $C''$ , the main difference between this bound and Inequality (2.8) of [Bar16] lies in the position of the logarithmic term  $\log_+(1/\sigma)$ : it is here involved inside the square root while it is outside in [Bar16].

Theorem 3 is well tailored for bounding the risk of a  $\rho$ -estimator. Indeed, when  $\varphi = \psi(g/f)$ , the random variable  $v(\varphi)$  can be bounded above as follows.

**Lemma 2.** *For all  $f, g \in \mathcal{S}$ ,*

$$\int_{\mathbb{R}} \psi^2(g/f) s dM \leq 4 (h^2(s, f) + h^2(s, g)).$$

Now, under suitable assumptions on the collection  $\mathcal{F} = \{\psi(g/f), f, g \in \mathcal{S}\}$ , Inequality (9) roughly says that with high probability (and  $\varepsilon = 1/24$ ):

$$(11) \quad |T(f, g) - T_E(f, g)| \leq \frac{1}{6} (h^2(s, f) + h^2(s, g)) + D_S(n) \quad \text{for all } f, g \in \mathcal{S}.$$

The term  $D_S(n)$  depends on the probability of the event on which (11) holds true and the complexity of  $S$ . The approximation  $T(f, g) \simeq T_E(f, g)$  is then accurate enough to control the risk of a  $\rho$ -estimator  $\hat{s}$ . We may do as in Section 2.2: we deduce from (3), that for all  $f, g \in S$ ,

$$\frac{1}{6}h^2(s, f) - \frac{19}{6}h^2(s, g) - D_S(n) \leq T(f, g) \leq \frac{19}{6}h^2(s, f) - \frac{1}{6}h^2(s, g) + D_S(n).$$

Therefore,

$$\frac{1}{6}h^2(s, f) - \frac{19}{6}h^2(s, S) - D_S(n) \leq \gamma(f) \leq \frac{19}{6}h^2(s, f) - \frac{1}{6}h^2(s, S) + D_S(n),$$

and hence,

$$\begin{aligned} \frac{1}{6}h^2(s, \hat{s}) - \frac{19}{6}h^2(s, S) - D_S(n) &\leq \gamma(\hat{s}) \\ &\leq \inf_{f \in S} \gamma(f) + 1/n \\ &\leq \frac{18}{6}h^2(s, S) + D_S(n) + 1/n. \end{aligned}$$

Finally, the risk of a  $\rho$ -estimator  $\hat{s}$  is bounded above by

$$h^2(s, \hat{s}) \leq 37h^2(s, S) + 12D_S(n) + 6/n.$$

It remains to explain the assumptions to put on the model  $S$  to make inequality (11) more precise and rigorous.

**3.2. Assumptions on models.** A convenient assumption to control the risks of  $\rho$ -estimators is the following.

**Assumption 1.** *There exists a collection  $\bar{S}$  of functions such that for all  $t \geq 0$ ,  $f \in \bar{S}$ ,  $g \in S$  the set  $\{x \in \mathbb{R}, g(x) > tf(x)\}$  is a union of at most  $d_S(f) \geq 1$  intervals.*

Not all models  $S$  satisfy this assumption. However, it applies for several models of interest, including some which are well suited for estimating functions under smooth or shape constraints. We carry out below some examples.

Let  $\mathcal{P}_{\ell, r}$  be the model defined for  $\ell \geq 1$ ,  $r \geq 0$  by

$$(12) \quad \mathcal{P}_{\ell, r} = \left\{ \sum_{j=1}^{\ell} P_j \mathbb{1}_{K_j}, P_j \text{ is a polynomial function of degree at most } r, \text{ which is} \right. \\ \left. \text{non-negative on an interval } K_j \text{ of } \mathbb{R} \text{ of finite length} \right\}.$$

**Proposition 5.** *Assumption 1 is fulfilled with  $S = \mathcal{P}_{\ell, r}$ ,  $\bar{S} \subset \mathcal{P}_{\ell, r}$  and for all  $f \in \bar{S}$ ,  $d_S(f) = (r+2)(2\ell+1)$ .*

We may also consider the model consisting of piecewise monotone functions defined for  $k \geq 1$  by

$$(13) \quad \mathcal{F}_k = S \cap \left\{ \sum_{j=1}^k f_j \mathbb{1}_{K_j}, K_j \text{ is an interval, and } f_j \text{ is monotone on } K_j \right\}.$$

Note that multimodal functions with  $k-1$  modes belong to  $\mathcal{F}_k$ .

We introduce for all interval  $K$  the set

$$\mathcal{G}(K) = \{f, f \text{ is either convex or concave on the interior of } K\},$$

and define

$$\mathcal{G}_k = \mathcal{S} \cap \left\{ \sum_{j=1}^k f_j \mathbb{1}_{K_j}, K_j \text{ is an interval, and } f_j \in \mathcal{G}(K_j) \right\},$$

$$\mathcal{G}'_k = \mathcal{S} \cap \left\{ \sum_{j=1}^k e^{f_j} \mathbb{1}_{K_j}, K_j \text{ is an interval, and } f_j \in \mathcal{G}(K_j) \right\}.$$

It is shown in [BB16] that  $\mathcal{P}_{\ell,0}$ ,  $\mathcal{F}_k$ ,  $\mathcal{G}_k$  and  $\mathcal{G}'_k$  satisfy Assumption 1 and the proposition below ensues from their results.

**Proposition 6.** *Let  $\mathcal{P}'_{\ell,1} = \{f, f \in \mathcal{S}, e^f \in \mathcal{P}_{\ell,1}\}$ . Assumption 1 is fulfilled with:*

- $S = \mathcal{F}_k$ ,  $\bar{S} \subset \cup_{\ell=1}^{\infty} \mathcal{P}_{\ell,0}$  and for all  $f \in \mathcal{P}_{\ell,0}$ ,  $d_S(f) = 2(k + \ell + 1)$ .
- $S = \mathcal{G}_k$ ,  $\bar{S} \subset \cup_{\ell=1}^{\infty} \mathcal{P}_{\ell,1}$  and for all  $f \in \mathcal{P}_{\ell,1}$ ,  $d_S(f) = 4(k + \ell + 1)$ .
- $S = \mathcal{G}'_k$ ,  $\bar{S} \subset \cup_{\ell=1}^{\infty} \mathcal{P}'_{\ell,1}$  and for all  $f \in \mathcal{P}'_{\ell,1}$ ,  $d_S(f) = 4(k + \ell + 1)$ .

### 3.3. A risk bound.

**Theorem 7.** *Let  $S$  be a model such that Assumption 1 is satisfied with  $\bar{S} \subset S$ . Then, for all  $\xi > 0$ , there exists an event which holds true with probability larger than  $1 - e^{-n\xi}$  and on which any  $\rho$ -estimator  $\hat{s}$  on  $S$  satisfies*

$$(14) \quad h^2(s, \hat{s}) \leq \inf_{f \in \bar{S}} \left\{ c_1 h^2(s, f) + c_2 \frac{d_S(f)}{n} \log_+^2 \left( \frac{n}{d_S(f)} \right) + c_3 \xi \log_+(1/\xi) \right\}.$$

In particular,

$$(15) \quad \mathbb{E} [h^2(s, \hat{s})] \leq \inf_{f \in \bar{S}} \left\{ c_1 \mathbb{E} [h^2(s, f)] + c'_2 \frac{d_S(f)}{n} \log_+^2 \left( \frac{n}{d_S(f)} \right) \right\}.$$

In the above inequalities,  $c_1, c_2, c'_2, c_3$  are universal positive constants.

Remark 1. When  $s \in \bar{S}$ , the risk of a  $\rho$ -estimator on  $S$  is bounded by

$$\mathbb{E} [h^2(s, \hat{s})] \leq C \frac{d_S(s)}{n} \log_+^2 \left( \frac{n}{d_S(s)} \right).$$

The rate of estimation of  $s$  then becomes parametric (up to the logarithmic term  $\log_+^2(n/d_S(s))$ ). When  $s \notin \bar{S}$ , the estimator automatically achieves the best trade-off between the bias (approximation) term  $h^2(s, f)$  and the variance (complexity) term  $(d_S(f)/n) \log_+^2(n/d_S(f))$ :

$$(16) \quad \mathbb{E} [h^2(s, \hat{s})] \leq CR(s) \quad \text{with } R(s) = \inf_{f \in \bar{S}} \left\{ \mathbb{E} [h^2(s, f)] + \frac{d_S(f)}{n} \log_+^2 \left( \frac{n}{d_S(f)} \right) \right\}.$$

It remains to compute  $R(s)$  to deduce (an upper bound of) the rate of convergence of the  $\rho$ -estimator when  $s \in S$ . This rate may be much slower than the rate we would obtain if  $s$  does belong to  $\bar{S}$  (see Section 3.5 for an example).

This phenomenon (faster rate of convergence when  $s \in \bar{S}$ ) has been put forward in [BB16] for  $\rho$ -estimators in density estimation and has been named superminimaxity. We show here that this

phenomenon also occurs for  $\rho$ -estimators in frameworks 2 and 3. If we now restrict to density estimation, we observe that our risk bound slightly improves the one of [BB16] in the sense that our approximation term involves a smaller exponent on the logarithm. We need, however, a more stringent assumption on the models.

Remark 2. When the framework varies, the Hellinger loss  $h$  in (15) changes, but the variance term remains the same. Moreover, the bias term in (15) in frameworks 2 and 3 is not larger than the bias term in density estimation (when one extends the definition of  $s$  in frameworks 2 and 3 to  $(-\infty, 0)$  by setting  $s(x) = 0$  when  $x < 0$ ). We may therefore only focus on density estimation to bound above  $R(s)$  in (16). The risk of the  $\rho$ -estimator is then bounded by this upper bound in the three frameworks (up to a multiplicative constant).

Remark 3. It follows from a crude application of the triangular inequality and from (15) that

$$\mathbb{E} [h^2(s, \hat{s})] \leq C \inf_{g \in S} \{ \mathbb{E} [h^2(s, g)] + R(g) \}.$$

If we know how to bound  $R(g)$  for  $g \in S$ , this inequality says that the risk of the  $\rho$ -estimator is not only controlled when  $s$  does belong to  $S$  but also when there exists  $g \in S$  such that  $s \simeq g$ , that is when  $s$  is close to  $S$ . In other words, the Hellinger risk of  $\hat{s}$  cannot substantially increase when  $s$  does not belong to  $S$  but is close to  $S$ . Such a result may be interpreted as a robustness property.

In particular, this robustness property applies to maximum likelihood estimators when assumptions of Theorem 1 or 2 are met with  $\mathcal{X} = \mathbb{R}$ . It is worth mentioning that this result is not true in general for maximum likelihood estimators. We refer to Section 2.3 of [Bir06] for an example in density estimation. We carry out below another example in framework 2. Let  $\alpha > 0$  and  $s_\alpha(t) = t^{-1} \mathbb{1}_{t \geq \alpha}$  be the hazard rate of a Pareto distribution. Suppose that the density  $f$  of  $T$  is a mixture of two Pareto distributions,

$$f(t) = \varepsilon \frac{\eta}{t^2} \mathbb{1}_{t \geq \eta} + (1 - \varepsilon) \frac{1}{t^2} \mathbb{1}_{t \geq 1},$$

where  $\varepsilon = 1/n$  and where  $\eta \in (0, 1/2)$ . Suppose moreover that the data are not censored:  $C = +\infty$  almost surely. The survival function  $G$  and the hazard rate  $s$  of  $T$  are given by the formulas

$$G(t) = \left(1 - \varepsilon + \varepsilon \frac{\eta}{t}\right) \mathbb{1}_{t \in [\eta, 1)} + \frac{1 - \varepsilon + \varepsilon \eta}{t} \mathbb{1}_{t \geq 1}$$

$$s(t) = \frac{\varepsilon \eta}{(1 - \varepsilon)t + \varepsilon \eta} \frac{\mathbb{1}_{t \in [\eta, 1)}}{t} + \frac{\mathbb{1}_{t \geq 1}}{t}.$$

The maximum likelihood estimator on  $S = \{s_\alpha, \alpha > 0\}$  is  $\tilde{s} = s_{\tilde{\alpha}}$  where  $\tilde{\alpha} = \min_{1 \leq i \leq n} X_i$ . We are now interested in its risk. Let  $\mathcal{A}_n$  be the event on which there exists at least one observation  $X_i$  smaller than  $2\eta$  and  $\mathcal{B}_n$  be the event on which  $n/2$  observations are larger than  $t_0$  where  $G(t_0) = 3/4$ . Then,  $\mathbb{P}(\mathcal{A}_n) = 1 - (1 - \varepsilon/2)^n = 1 - (1 - 1/(2n))^n > 1 - e^{-1/2}$  and  $\mathbb{P}(\mathcal{B}_n) = \mathbb{P}(Y \geq n/2)$  where  $Y$  is binomially distributed with parameters  $(n, 3/4)$ . Therefore, for all  $c \in (0, 1)$ , and  $n$  large enough,  $\mathbb{P}(\mathcal{B}_n) \geq c$ . Since  $t_0 \geq 2\eta$ , we deduce that on  $\mathcal{A}_n \cap \mathcal{B}_n$ :

$$h^2(s, \tilde{s}) \geq \frac{1}{4} \int_{2\eta}^1 \left( \sqrt{1/t} - \sqrt{s(t)} \right)^2 dt \geq \frac{1}{4} \int_{2\eta}^1 \frac{1}{t} \left( 1 - \sqrt{\frac{\varepsilon \eta}{(1 - \varepsilon)t + \varepsilon \eta}} \right)^2 dt.$$

Some elementary computations show (using that  $\varepsilon \leq 1/2$ ),

$$\inf_{t \in [2\eta, 1]} \left( 1 - \sqrt{\frac{\varepsilon \eta}{(1 - \varepsilon)t + \varepsilon \eta}} \right)^2 \geq \left( 1 - \sqrt{\frac{\varepsilon \eta}{2(1 - \varepsilon)\eta + \varepsilon \eta}} \right)^2 > 0.17.$$



Therefore,  $h^2(s, \tilde{s}) > 0.04 \log(1/(2\eta))$  and  $\mathbb{E}[h^2(s, \tilde{s})] > 0.04 \log(1/(2\eta))\mathbb{P}(\mathcal{A}_n \cap \mathcal{B}_n)$ . Finally, there exist universal constants  $c', n_0$  such that for all  $n \geq n_0$ :

$$\mathbb{E}[h^2(s, \tilde{s})] > c' \log(1/(2\eta)).$$

In particular, the risk of a maximum likelihood estimator  $\tilde{s}$  can be made arbitrarily large by playing with  $\eta$ . We now turn to  $\rho$ -estimation. The model  $S$  fulfils Assumption 1 with  $d_S(f) = 1$ ,  $\bar{S} = S$ . A  $\rho$ -estimator  $\hat{s}$  on  $S$  satisfies therefore  $\mathbb{E}[h^2(s, \hat{s})] \leq C(\mathbb{E}[h^2(s, S)] + \log^2 n/n)$ . The variance term  $\log^2 n/n$  is likely not optimal. However, the bias term can be bounded above by  $\mathbb{E}[h^2(s, S)] \leq \mathbb{E}[h^2(s, s_1)] \leq \mathbb{P}(T \in [\eta, 1]) = \varepsilon(1 - \eta) \leq 1/n$ . Consequently,

$$\sup_{\eta \in (0, 1/2)} \mathbb{E}[h^2(s, \hat{s})] \leq C' \frac{\log^2 n}{n},$$

where  $C'$  is universal. This example illustrates the lack of robustness of the maximum likelihood estimator in this model.

When a  $\rho$ -estimator  $\hat{s}$  vanishes the criterion  $\gamma$ , which typically happens when  $\hat{s}$  maximizes  $L_{\mathcal{X}}$  over some models, the constant  $c_1$  appearing in front of the bias term  $h^2(s, f)$  in Theorem 7 can be improved:

**Proposition 8.** *Let  $S$  be a model such that Assumption 1 is satisfied with  $\bar{S} \subset S$ . Suppose that there exists  $\hat{s} \in S$  such that  $\gamma(\hat{s}) = 0$ .*

*Then, (14) and (15) hold for all  $\varepsilon > 0$  with  $c_1 = c_{1,\varepsilon}$ ,  $c_2 = c_{2,\varepsilon}$ ,  $c'_2 = c'_{2,\varepsilon}$  such that  $c_{1,\varepsilon} \geq 9$  and  $\lim_{\varepsilon \rightarrow 0} c_{1,\varepsilon} = 9$ .*

The constant  $c_1 = c_{1,\varepsilon}$  may therefore be made as close as 9 as wished. We do not know to what extent this result can be improved for estimators  $\hat{s}$  maximizing  $L_{\mathcal{X}}$  under convex type assumptions on the models. However,  $c_1 = c_{1,\varepsilon}$  cannot be smaller than 2 in general as shown by the following elementary example.

Consider framework 1, a finite collection of disjoint intervals  $m$ , and the model  $S$  consisting of piecewise constant densities based on  $m$ . Since  $S$  is convex, the usual histogram estimator  $\hat{s}$  is a  $\rho$ -estimator that vanishes  $\gamma$ . Moreover, it is known that the usual histogram estimator may not converge to the closest function of  $s$  in the model (for the Hellinger distance), see [BR06]. In particular, let  $p \in (0, 1)$  to be specified later,  $s = \mathbb{1}_{[0,p]} + (1-p)\mathbb{1}_{(1,2]}$  and  $S$  be the model of piecewise constant densities based on the partition  $\{[0, 1], (1, 2]\}$ . Then  $\mathbb{E}[h^2(s, \hat{s})]$  converges to  $h^2(s, \bar{s}_1)$  where  $\bar{s}_1 = p\mathbb{1}_{[0,1]} + (1-p)\mathbb{1}_{(1,2]}$ . Define

$$\bar{s}_2 = \frac{p^2}{p^2 + 1 - p} \mathbb{1}_{[0,1]} + \frac{1-p}{p^2 + 1 - p} \mathbb{1}_{(1,2]} \in S.$$

The fraction

$$\frac{h^2(s, \bar{s}_1)}{h^2(s, \bar{s}_2)} = \frac{p(1 - \sqrt{p})}{1 - \sqrt{1 - p(1 - p)}}$$

can be made arbitrarily close to 2 by choosing  $p$  small enough. Therefore, for all  $\eta \in (1, 2)$ , and  $p$  small enough

$$\lim_{n \rightarrow +\infty} \mathbb{E}[h^2(s, \hat{s})] \geq \eta h^2(s, \bar{s}_2) \geq \eta h^2(s, S),$$

and thus  $c_1 \geq \eta$ .

**3.4. Risks of  $\rho$ -estimators for models consisting of step functions.** For illustration purposes, consider the collection  $\mathcal{P}_{\ell,0}$  of step functions defined by (12). Then, any  $\rho$ -estimator  $\hat{s}$  on  $\mathcal{P}_{\ell,0}$  satisfies

$$(17) \quad \mathbb{E} [h^2(s, \hat{s})] \leq C \left\{ \mathbb{E} [h^2(s, \mathcal{P}_{\ell,0})] + \frac{\ell}{n} \log_+^2 \left( \frac{n}{\ell} \right) \right\}.$$

The first term  $\mathbb{E} [h^2(s, \mathcal{P}_{\ell,0})]$  can be interpreted as an approximation term that is small if  $s$  is close to a step function. When  $s$  does belong to  $\mathcal{P}_{\ell,0}$ , the bound becomes

$$\mathbb{E} [h^2(s, \hat{s})] \leq C \frac{\ell}{n} \log_+^2 \left( \frac{n}{\ell} \right).$$

This result is, in general, slightly suboptimal as it is possible to do better in density estimation. Indeed, we may relate in framework 1 the performance of a  $\rho$ -estimator to the Vapnik-Chervonenkis dimension of the subgraphs of the model. Yet, this dimension is of the order of  $\ell$  for  $\mathcal{P}_{\ell,0}$  (see [BB17]). We then derive from the proof of Theorem 12 of [BBS17],

$$\mathbb{E} [h^2(s, \hat{s})] \leq C' \left\{ h^2(s, \mathcal{P}_{\ell,0}) + \frac{\ell}{n} \log_+ \left( \frac{n}{\ell} \right) \right\},$$

where  $C'$  is universal. This last inequality is also shown and discussed in Section 7.4 of [BB17]. In particular, the logarithm in this inequality cannot be avoided since the variance term corresponds to the optimal minimax rate of convergence on  $\mathcal{P}_{\ell,0}$  (Proposition 2 of [BM98]).

**3.5. Risks of  $\rho$ -estimators for models consisting of piecewise monotone functions.** Let  $(\mathbb{L}^2(\mathbb{R}, M_E), d_2)$  be the metric space of square integrable functions on  $\mathbb{R}$  with respect to the measure  $M_E$  defined by  $M_E(A) = \mathbb{E}[M(A)]$  for all  $A \in \mathcal{B}(\mathbb{R})$ . Let  $S$  be a model satisfying the assumptions of Theorem 7. We know from (16) that a  $\rho$ -estimator  $\hat{s}$  satisfies

$$(18) \quad \mathbb{E} [h^2(s, \hat{s})] \leq CR(s) \quad \text{with } R(s) = \inf_{f \in \bar{S}} \left\{ \frac{1}{2} d_2^2(\sqrt{s}, f) + \frac{d_S(f)}{n} \log_+^2 \left( \frac{n}{d_S(f)} \right) \right\},$$

where  $C$  is universal. It then remains to bound  $R(s)$  to control the risk of  $\hat{s}$ .

Two bounds on  $R(s)$  in density estimation can be deduced from the results of [BB16]: when  $s$  is piecewise monotone and when  $\sqrt{s}$  is piecewise convex-concave. Dealing with the two other statistical settings requires little supplementary work (see Remark 2 in Section 3.3). We obtain bounds that are roughly of the order of  $(\log^2 n/n)^{2/3}$  when  $s$  is piecewise monotone, and  $(\log^2 n/n)^{4/5}$  when  $\sqrt{s}$  is piecewise convex-concave. In the sequel, we propose to make explicit the first bound only.

We define for  $K \subset \mathbb{R}$  and  $f \in \mathcal{S}$ ,  $V_K(f) = \sup_{x \in K} f(x) - \inf_{x \in K} f(x)$ . Let  $\mathcal{M}_k$  be the family gathering all the collections  $m$  of at most  $k$  disjoint intervals of  $\mathbb{R}$ , and for  $m \in \mathcal{M}_k$ ,

$$\mathcal{F}(m) = \left\{ \sum_{K \in m} f_K \mathbb{1}_K, \text{ where } f_K \in \mathcal{S} \text{ is monotone on } K \right\}.$$

We define for  $m \in \mathcal{M}_k$  and  $f \in \mathcal{F}(m)$  of the form  $f = \sum_{K \in m} f_K \mathbb{1}_K$ ,

$$L_m(f) = \sum_{K \in m} [M_E(K) V_K^2(f_K)]^{1/3}.$$

In this equality, we use the convention  $+\infty \times 0 = 0$  when  $M_E(K) = +\infty$ . For  $f \in \mathcal{F}_k$ , we set

$$L(f) = \inf_{m \in \mathcal{M}_k} L_m(f).$$

The result is the following.

**Corollary 1.** *Any  $\rho$ -estimator  $\hat{s}$  on  $\mathcal{F}_k$  ( $k \geq 1$ ) satisfies*

$$(19) \quad \mathbb{E} [h^2(s, \hat{s})] \leq C \inf_{f \in \mathcal{F}_k} \left\{ d_2^2(\sqrt{s}, f) + L(f) \left( \frac{\log^2 n}{n} \right)^{2/3} + \frac{k \log^2 n}{n} \right\}.$$

*In particular, if  $s$  does belong to  $\mathcal{F}_k$ , then  $f = \sqrt{s}$  also belongs to  $\mathcal{F}_k$  and hence,*

$$(20) \quad \mathbb{E} [h^2(s, \hat{s})] \leq C \left[ L(\sqrt{s}) \left( \frac{\log^2 n}{n} \right)^{2/3} + \frac{k \log^2 n}{n} \right].$$

*In the preceding inequalities,  $C$  is a universal constant.*

We may also make inequality (20) more explicit when  $k = 2$  and  $s$  is unimodal by bounding  $L(\sqrt{s})$  from above. We distinguish the cases according to the different frameworks.

Consider framework 1. Then,  $M_E(K) = \mu(K)$ . Therefore, if the support of  $s$  is of finite length  $L_{supp}$ ,

$$L(\sqrt{s}) \leq 2L_{supp}^{1/3} (\sup_{x \in \mathbb{R}} s(x))^{1/3}.$$

Consider now framework 2 and suppose that  $X$  has finite expectation. Then, for all interval  $K \subset [0, +\infty)$ ,  $M_E(K)$  is not larger than  $\mathbb{E}(X)$  and hence

$$L(\sqrt{s}) \leq 2(\mathbb{E}(X))^{1/3} \left( \sup_{x \geq 0} s(x) \right)^{1/3}.$$

As to framework 3, define the time  $T_1 = \int_{I_{\text{obs}}} \mathbb{1}_{X_{t-}=1} dt$  during which the (left limit) of the observed Markov process is in state 1 and suppose that  $\mathbb{E}(T_1) < \infty$ . Then,  $M_E(K) \leq \mathbb{E}(T_1)$  and thus

$$L(\sqrt{s}) \leq 2(\mathbb{E}(T_1))^{1/3} \left( \sup_{x \in I_{\text{obs}}} s(x) \right)^{1/3}.$$

#### 4. SELECTING AMONG ESTIMATORS

It is often difficult in practice to find a global minimum of  $\gamma$  and thus to build  $\rho$ -estimators. In particular, we do not know how to construct a  $\rho$ -estimator on the model  $S = \mathcal{P}_{\ell, r}$ . In this section, we propose an alternative way, more numerically friendly, to define a piecewise polynomial estimator and study its properties.

**4.1. A uniform risk bound.** The quality of a  $\rho$ -estimator is, in the present paper, assessed by means of Theorem 7. We have not yet mentioned the following result: the event on which (14) is valid depends neither on  $S$  nor on  $\bar{S}$ . In other words, the theorem can be rewritten as:

**Theorem 9.** *For all  $\xi > 0$ , there exists an event which holds true with probability larger than  $1 - e^{-n\xi}$  and on which: for all model  $S$  satisfying Assumption 1 with  $\bar{S} \subset S$ , and all  $\rho$ -estimator  $\hat{s}$  on  $S$ ,*

$$h^2(s, \hat{s}) \leq \inf_{f \in \bar{S}} \left\{ c_1 h^2(s, f) + c_2 \frac{d_S(f)}{n} \log_+^2 \left( \frac{n}{d_S(f)} \right) + c_3 \xi \log_+(1/\xi) \right\}.$$

*In the above inequality,  $c_1, c_2, c_3$  are universal positive constants.*

Thereby, this risk bound simultaneously holds true for all  $\rho$ -estimators and models  $S$  satisfying Assumption 1. In particular, this allows the model  $S$  to vary randomly among the class of models for which this assumption is fulfilled. It may then be chosen according to the data.

An application of this result is polynomial estimator selection: consider  $\ell \geq 1$ ,  $r \geq 0$  and an at most countable collection  $\{\hat{s}_\lambda, \lambda \in \Lambda\} \subset \mathcal{P}_{\ell,r}$  of non-negative piecewise polynomial estimators of degree at most  $r$  based on at most  $\ell$  pieces. Building a  $\rho$ -estimator on the (random) model  $S = \{\hat{s}_\lambda, \lambda \in \Lambda\}$  amounts to selecting an estimator among  $\{\hat{s}_\lambda, \lambda \in \Lambda\}$ . Moreover, we deduce from Proposition 5 that we may set  $d_S(\hat{s}_\lambda) = (r+2)(2\ell+1)$ . Thereby, any  $\rho$ -estimator  $\hat{s}$  is of the form  $\hat{s} = \hat{s}_{\hat{\lambda}}$  and satisfies

$$(21) \quad \mathbb{E} [h^2(s, \hat{s}_{\hat{\lambda}})] \leq C \left\{ \inf_{\lambda \in \Lambda} \mathbb{E} [h^2(s, \hat{s}_\lambda)] + \frac{(r+1)\ell \log_+^2(n/(\ell(r+1)))}{n} \right\},$$

where  $C$  is a universal constant.

This risk bound is always worse than the one we would obtain for a  $\rho$ -estimator  $\hat{s}$  on  $S = \mathcal{P}_{\ell,r}$ :

$$(22) \quad \mathbb{E} [h^2(s, \hat{s})] \leq C' \left\{ \mathbb{E} [h^2(s, \mathcal{P}_{\ell,r})] + \frac{(r+1)\ell \log_+^2(n/(\ell(r+1)))}{n} \right\},$$

where  $C'$  is universal. The interest of  $\hat{s}_{\hat{\lambda}}$  is practical: the construction of  $\hat{s}$  seems to be numerically difficult whereas the selected estimator  $\hat{s}_{\hat{\lambda}}$  can be computed in a reasonable amount of time as soon as  $\Lambda$  is finite and not too large (and when the computation of the estimators  $\hat{s}_\lambda$  is fast enough).

**4.2. Selecting among a special collection of piecewise polynomial estimators.** As we see in (21), we should take  $\Lambda$  as large as possible to improve on the theoretical performances of the selected estimator. In this section, our aim is to select an estimator among a special, but possibly very large, collection of piecewise polynomial  $\rho$ -estimators  $\{\hat{s}_\lambda, \lambda \in \Lambda\}$ .

The preceding procedure could be a solution. However, its numerical complexity depends heavily on  $|\Lambda|$ , which is, in this section, usually very large. This procedure may therefore be numerically intractable in practice. We propose in this section an alternative way inspired from [Sar14] that improves on its numerical cost.

Let  $\mathcal{M}$  be the class of finite (non-empty) collections  $m$  of disjoint intervals  $K$  that are right-closed, not reduced to a singleton, and of finite length. Let  $r \geq 0$ ,  $m \in \mathcal{M}$  and

$$\mathcal{P}_r(m) = \left\{ \sum_{K \in m} f_K \mathbb{1}_K, \text{ for all } K \in m, f_K \text{ is a polynomial function of degree at most } r, \right. \\ \left. \text{non-negative on } K \right\}.$$

We may compute a piecewise polynomial  $\rho$ -estimator  $\hat{s}_m$  on the convex model  $\mathcal{P}_r(m)$ :

**Lemma 3.** *Let  $m \in \mathcal{M}$  and for  $K \in m$ ,*

$$\mathcal{P}_r(K) = \{f \mathbb{1}_K, f \text{ is a polynomial function of degree at most } r \text{ and non-negative on } K\}.$$

*Then,  $\sup_{f \in \mathcal{P}_r(K)} L_K(f)$  is finite and achieved at a point  $\hat{s}_K$ . Moreover,  $\hat{s}_m = \sum_{K \in m} \hat{s}_K$  maximizes  $L_{\mathcal{X}}$  over  $\mathcal{P}_r(m)$  where  $\mathcal{X} = \bigcup_{K \in m} K$ . It is a  $\rho$ -estimator on the model  $S = \mathcal{P}_r(m)$  that vanishes  $\gamma$ .*

We now consider some (possibly random) collection of distinct random variables  $\{Y_i, i \in \hat{I}\}$  where  $\hat{I}$  is a (possibly random) set such that  $\hat{n} = |\hat{I}| \geq 2$ . Since the random variables  $(Y_i)_{i \in \hat{I}}$  are distinct almost surely, we may order them:  $Y_{(1)} < Y_{(2)} < \dots < Y_{(\hat{n})}$ . We define the collection  $\widehat{\mathcal{M}}$  that gathers all the partitions  $m$  of  $[Y_{(1)}, Y_{(\hat{n})}]$  of the form

$$m = \{[Y_{(1)}, Y_{(n_1)}], (Y_{(n_1)}, Y_{(n_2)}], (Y_{(n_2)}, Y_{(n_3)}], \dots, (Y_{(n_k)}, Y_{(\hat{n})})\},$$

where  $k \geq 0$  and  $1 < n_1 < n_2 < \dots < n_k < \hat{n}$  with the convention that  $m = \{[Y_{(1)}, Y_{(\hat{n})}]\}$  when  $k = 0$ . We set for  $\ell \in \{1, \dots, \hat{n} - 1\}$ ,

$$\widehat{\mathcal{M}}_\ell = \{m \in \widehat{\mathcal{M}}, |m| = \ell\}.$$

We consider a random variable  $\hat{\ell}$  with values in  $\{1, \dots, \hat{n} - 1\}$ . For each  $m \in \widehat{\mathcal{M}}_{\hat{\ell}}$ , we define the  $\rho$ -estimator  $\hat{s}_m$  on the model  $\mathcal{P}_r(m)$  as in Lemma 3. The aim of this section is to explain how we can select an estimator among the family  $\{\hat{s}_m, m \in \widehat{\mathcal{M}}_{\hat{\ell}}\}$ .

We define for  $m \in \widehat{\mathcal{M}}$ ,  $K \in m$  and  $m_K \in \widehat{\mathcal{M}}$ , the partition  $m_K \vee K$  of  $K$  by

$$(23) \quad m_K \vee K = \{K' \cap K, K' \in m_K, K' \cap K \neq \emptyset\}.$$

We now consider a positive number  $L$  and define the criterion  $\gamma_2$  for  $m \in \widehat{\mathcal{M}}_{\hat{\ell}}$  by

$$(24) \quad \gamma_2(\hat{s}_m) = \sum_{K \in m} \sup_{m' \in \widehat{\mathcal{M}}_{\hat{\ell}}} \left\{ T(\hat{s}_m \mathbb{1}_K, \hat{s}_{m'} \mathbb{1}_K) - L(r+1) \frac{|m' \vee K| \log_+^2(n/(r+1))}{n} \right\}.$$

The selected estimator is then any estimator  $\hat{s}_{\hat{m}}$  of the collection  $\{\hat{s}_m, m \in \widehat{\mathcal{M}}_{\hat{\ell}}\}$  minimizing  $\gamma_2$ :

$$(25) \quad \gamma_2(\hat{s}_{\hat{m}}) = \min_{m \in \widehat{\mathcal{M}}_{\hat{\ell}}} \gamma_2(\hat{s}_m).$$

Note that the above minimum is achieved since  $\widehat{\mathcal{M}}_{\hat{\ell}}$  is finite.

**Theorem 10.** *There exists a universal constant  $L_0$  such that if  $L \geq L_0$ , any estimator  $\hat{s}_{\hat{m}}$  minimizing (25) satisfies for all  $\xi > 0$ , and probability larger than  $1 - e^{-n\xi}$ ,*

$$(26) \quad h^2(s, \hat{s}_{\hat{m}}) \leq C \left\{ \inf_{m \in \widehat{\mathcal{M}}_{\hat{\ell}}} h^2(s, \mathcal{P}_r(m)) + L \frac{(r+1)\hat{\ell} \log_+^2(n/(r+1))}{n} + \xi \log_+(1/\xi) \right\}.$$

In particular,

$$(27) \quad \mathbb{E} [h^2(s, \hat{s}_{\hat{m}})] \leq C' \mathbb{E} \left[ \inf_{m \in \widehat{\mathcal{M}}_{\hat{\ell}}} h^2(s, \mathcal{P}_r(m)) + L \frac{(r+1)\hat{\ell} \log_+^2(n/(r+1))}{n} \right].$$

In the above inequalities,  $C$  and  $C'$  are universal constants.

This last inequality also says

$$\mathbb{E} [h^2(s, \hat{s}_{\hat{m}})] \leq C' \mathbb{E} \left[ \inf_{m \in \widehat{\mathcal{M}}_{\hat{\ell}}} h^2(s, \hat{s}_m) + L \frac{(r+1)\hat{\ell} \log_+^2(n/(r+1))}{n} \right].$$

This risk bound is very similar to the one (21) obtained for the first selection rule (we only slightly loose on the variance term).

The numerical complexity of this procedure may be significantly reduced by using an algorithm of dynamic programming. The estimator  $\hat{s}_{\hat{m}}$  may then be computed in practice when the  $\hat{s}_K$  can be constructed sufficiently fast and when  $n, \hat{n}, \hat{\ell}$  are small enough (the numerical complexity of the algorithm increases polynomially with  $\hat{n}$  and  $\hat{\ell}$ ). Unfortunately, the practical computation of the estimator may take too long in other situations. For more informations on this algorithm, we refer to [Kan92] and Section 4.2 of [CR04].

We now consider framework 1. When  $Y_{(1)} \leq \min_{1 \leq i \leq n} X_i \leq \max_{1 \leq i \leq n} X_i \leq Y_{(\hat{n})}$ ,  $\hat{s}_m$  maximizes  $L(\cdot)$  and is therefore a maximum likelihood estimator. It is then natural to compare our estimator  $\hat{s}_{\hat{m}}$  to the one  $\hat{s}_{\hat{m}}$  that maximizes the log likelihood  $L(\hat{s}_m)$  over  $m \in \widehat{\mathcal{M}}_{\hat{\ell}}$ . We refer to Section 5 for numerical simulations when  $\{Y_i, i \in \hat{I}\} = \{X_i, i \in \{1, \dots, n\}\}$  and  $r = 0$  (and when the parameter  $L$  is chosen as in the next section). We do not know theoretical results for  $\hat{s}_{\hat{m}}$ . However, when  $\hat{I}$  is a deterministic subset of  $\{1, \dots, n\}$ , when the  $Y_i$  are deterministic numbers, and when  $\hat{\ell}$  is deterministic, then results concerning the maximizer of  $m \mapsto L(\hat{s}_m)$  over  $m \in \widehat{\mathcal{M}}_{\hat{\ell}}$  may be found in the literature. We refer to Theorem 3.2 of [Cas99] (when  $r = 0$ ) and Theorem 2 of [BBM99] (when  $r \geq 0$ ) for upper-bounds of the Hellinger risk in density estimation. Note that they put restrictions either on  $s$ , or on the minimal length of the intervals  $K$  of the partitions  $m \in \widehat{\mathcal{M}}_{\hat{\ell}}$ . Besides, contrary to ours, their upper-bounds involve the Kullback Leibler divergence.

The bias term in (26) depends on the collection  $\widehat{\mathcal{M}}_{\hat{\ell}}$  and thus on the choice of  $\{Y_i, i \in \hat{I}\}$ . In general, this bias term may be larger than the one we would obtain for a  $\rho$ -estimator on  $\mathcal{P}_{\hat{\ell}, r}$ . Nevertheless, it may be controlled in favourable situations as explained below.

Suppose that  $\{Y_i, i \in \hat{I}\}$  is rich enough to satisfy

$$(28) \quad N(A) = N(\{Y_i, i \in \hat{I}\} \cap A) \quad \text{for all } A \in \mathcal{B}(\mathbb{R}).$$

For instance, we may define  $\{Y_i, i \in \hat{I}\}$  as follows:

- in framework 1, we may set  $\hat{I} = \{1, \dots, n\}$ , and for all  $i \in \hat{I}$ ,  $Y_i = X_i$ ,
- in framework 2, we may consider a set  $\hat{I} \subset \{1, \dots, n\}$ , such that  $|\hat{I}| \geq 2$  and  $\hat{I} \supset \{i \in \{1, \dots, n\}, D_i = 1\}$  and define for all  $i \in \hat{I}$ ,  $Y_i = X_i$ ,
- in framework 3, we may consider a set  $\hat{I} \subset \{1, \dots, n\}$ , such that  $|\hat{I}| \geq 2$  and  $\hat{I} \supset \{i \in \{1, \dots, n\}, T_{1,0}^{(i)} \in I_{\text{obs}}\}$ , and define for all  $i \in \hat{I}$ ,  $Y_i = T_{1,0}^{(i)}$ .

The lemma below allows to bound the bias term when  $r = 0$ .

**Lemma 4.** *Suppose that condition (28) is met. Let  $r = 0$ , and for  $\ell \geq 1$ ,  $\mathcal{M}'_{\ell}$  be the family that gathers all the collections  $m$  of disjoint intervals of the form*

$$(29) \quad m = \{[x_1, x_2], (x_2, x_3], (x_3, x_4], \dots, (x_{\ell}, x_{\ell+1}]\},$$

where  $x_1 < x_2 < \dots < x_{\ell+1}$  (with the convention that  $m = \{[x_1, x_2]\}$  when  $\ell = 1$ ). Then, for all  $\xi > 0$ , the following holds with probability larger than  $1 - e^{-n\xi}$ : for all  $\ell \in \{1, \dots, \hat{n} - 1\}$ ,  $m \in \mathcal{M}'_{\ell}$  written as in (29) and such that  $Y_{(1)}$  and  $Y_{(\hat{n})}$  belong to  $[x_1, x_{\ell+1}]$ , there exists  $m' \in \widehat{\mathcal{M}}_{\ell}$  such that the  $\rho$ -estimator  $\hat{s}_{m'}$  satisfies

$$(30) \quad h^2(s, \hat{s}_{m'}) \leq C \left\{ h^2(s, \mathcal{P}_0(m)) + \frac{\ell \log_+^2(n/\ell)}{n} + \xi \log_+(1/\xi) \right\}.$$

In particular, suppose that  $s$  vanishes outside  $[0, 1]$ , and  $r = 0$ . Let  $\mathcal{P}'_{\hat{\ell}, 0}$  be the collection of step functions based on partitions  $m$  of  $[0, 1]$  belonging to  $\mathcal{M}'_{\hat{\ell}}$ . Then, we deduce from (27) and the preceding lemma when (28) holds,  $\{Y_i, i \in \hat{I}\} \subset [0, 1]$  and  $L \geq \max\{L_0, 1\}$ ,

$$\mathbb{E} [h^2(s, \hat{s}_{\hat{m}})] \leq C'' \mathbb{E} \left[ h^2(s, \mathcal{P}'_{\hat{\ell}, 0}) + L \frac{\hat{\ell} \log^2 n}{n} \right],$$

where  $C''$  is a universal constant. This risk bound, corresponds, up to slight modifications (in the constant and the variance term), to the one we would obtain for the  $\rho$ -estimator on  $\mathcal{P}'_{\hat{\ell}, 0}$ .

**4.3. Selecting among a special collection of piecewise polynomial estimators: a calibration free approach.** The preceding procedure suffers from a major drawback: the choice of  $L$ . This parameter, is, indeed, involved in the construction of  $\hat{s}_{\hat{m}}$  and a bad choice of  $L$  may deteriorate the estimator.

A simple solution to avoid this pitfall, is to proceed as follows. We consider a (non-empty, but at most countable) collection  $\mathcal{L}$  of positive numbers. For each  $L \in \mathcal{L}$ , we may use the procedure described in the preceding section with the parameter  $L$  in (24) to select an estimator among the collection  $\{\hat{s}_m, m \in \widehat{\mathcal{M}}_{\hat{\ell}}\}$ . The selected estimator is written  $\hat{s}_{\hat{m}_L}$  to emphasize that it depends on  $L$ . We then select an estimator among  $\{\hat{s}_{\hat{m}_L}, L \in \mathcal{L}\}$  as explained in Section 4.1.

The  $\rho$ -estimator on the model  $\{\hat{s}_{\hat{m}_L}, L \in \mathcal{L}\}$  is of the form  $\hat{s} = \hat{s}_{\hat{m}_L}$  and satisfies for all  $\xi > 0$  and probability larger than  $1 - e^{-n\xi}$ ,

$$h^2(s, \hat{s}) \leq C \left[ \inf_{L \in \mathcal{L}} \{h^2(s, \hat{s}_{\hat{m}_L})\} + \frac{(r+1)\hat{\ell} \log_+^2(n/(r+1))}{n} + \xi \log_+(1/\xi) \right],$$

where  $C$  is a universal constant. If  $\mathcal{L}$  contains at least one number  $L$  larger than  $L_0$ , we derive from (26),

$$h^2(s, \hat{s}) \leq C' \left[ \inf_{m \in \widehat{\mathcal{M}}_{\hat{\ell}}} \{h^2(s, \mathcal{P}_r(m))\} + \max \left( 1, \inf_{\substack{L \in \mathcal{L}, \\ L \geq L_0}} L \right) \frac{\hat{\ell}(r+1) \log_+^2(n/(r+1))}{n} + \xi \log_+(1/\xi) \right],$$

where  $C'$  is a universal constant.

This estimator  $\hat{s}$  does not depend on the particular choice of a calibration parameter  $L$  but rather on a collection  $\mathcal{L}$ . The larger  $\mathcal{L}$ , the better the risk bound. However, the numerical complexity of the whole procedure increases with the size of  $\mathcal{L}$ .

**4.4. Adaptive piecewise polynomial estimation.** In this section, we modify the previous procedure in order to choose the number  $\hat{\ell}$  of pieces from the data.

We define for  $k \in \{1, \dots, \hat{n} - 1\}$  the collection  $\widehat{\mathcal{M}}_{k, lower}$  of partitions  $m \in \widehat{\mathcal{M}}$  whose cardinal is at most  $k$ ,

$$\widehat{\mathcal{M}}_{k, lower} = \{m \in \widehat{\mathcal{M}}, |m| \leq k\} = \bigcup_{\ell=1}^k \widehat{\mathcal{M}}_{\ell}.$$

We consider a random variable  $\hat{k}$  with values in  $\{1, \dots, \hat{n} - 1\}$  and aim at selecting an estimator among  $\{\hat{s}_m, m \in \widehat{\mathcal{M}}_{\hat{k}, lower}\}$ .

We consider some  $L > 0$  and set for  $m \in \widehat{\mathcal{M}}_{\hat{k}, \text{lower}}$ ,

$$\gamma_3(\hat{s}_m) = \sum_{K \in m} \sup_{m' \in \widehat{\mathcal{M}}_{\hat{k}, \text{lower}}} \left\{ T(\hat{s}_m \mathbb{1}_K, \hat{s}_{m'} \mathbb{1}_K) - L(r+1) \frac{|m' \vee K| \log_+^2(n/(r+1))}{n} \right\}.$$

The selected estimator  $\hat{s}_{\hat{m}}$  is any estimator of the family satisfying

$$(31) \quad \begin{aligned} \gamma_3(\hat{s}_{\hat{m}}) + 2L(r+1) \frac{|\hat{m}| \log_+^2(n/(r+1))}{n} \\ = \inf_{m \in \widehat{\mathcal{M}}_{\hat{k}, \text{lower}}} \left\{ \gamma_3(\hat{s}_m) + 2L(r+1) \frac{|m| \log_+^2(n/(r+1))}{n} \right\}. \end{aligned}$$

**Theorem 11.** *There exists a universal constant  $L_0$  such that if  $L \geq L_0$ , any estimator  $\hat{s}_{\hat{m}}$  satisfying (31) satisfies for all  $\xi > 0$ , and probability larger than  $1 - e^{-n\xi}$ ,*

$$h^2(s, \hat{s}_{\hat{m}}) \leq C \inf_{m \in \widehat{\mathcal{M}}_{\hat{k}, \text{lower}}} \left\{ h^2(s, \mathcal{P}_r(m)) + L \frac{(r+1)|m| \log_+^2(n/(r+1))}{n} + \xi \log_+(1/\xi) \right\}.$$

In particular,

$$(32) \quad \mathbb{E} [h^2(s, \hat{s}_{\hat{m}})] \leq C' \mathbb{E} \left[ \inf_{m \in \widehat{\mathcal{M}}_{\hat{k}, \text{lower}}} \left\{ h^2(s, \mathcal{P}_r(m)) + L \frac{(r+1)|m| \log_+^2(n/(r+1))}{n} \right\} \right].$$

In the above inequalities,  $C$  and  $C'$  are universal constants.

This risk bound improves when  $\hat{k}$  grows up. Moreover, (32) implies

$$\mathbb{E} [h^2(s, \hat{s}_{\hat{m}})] \leq C' \mathbb{E} \left[ \inf_{1 \leq \hat{\ell} \leq \hat{k}} \left\{ \inf_{m \in \widehat{\mathcal{M}}_{\hat{\ell}}} \{h^2(s, \mathcal{P}_r(m))\} + L \frac{(r+1)\hat{\ell} \log_+^2(n/(r+1))}{n} \right\} \right].$$

The right-hand side of this inequality corresponds to the bound (27) achieved by the estimator of Section 4.2 when the choice of  $\hat{\ell}$  is the best possible among  $\{1, \dots, \hat{k}\}$ .

The quality and the construction of the estimator  $\hat{s}_{\hat{m}}$  still depends on  $\{Y_i, i \in \hat{I}\}$ . However, when  $\hat{k} = \hat{n}$ , and when  $\{Y_i, i \in \hat{I}\}$  is rich enough, the infimum in (32) can be taken over the infinite collection  $\mathcal{M}$  (up to a modification of  $C'$ ), as shown below.

**Lemma 5.** *Suppose that  $\{Y_i, i \in \hat{I}\}$  is chosen in such a way that  $N$  satisfies (28). There exists a universal constant  $C$  such that for all  $\xi > 0$  and probability larger than  $1 - e^{-n\xi}$ : for all  $m \in \mathcal{M}$ , there exists  $m' \in \widehat{\mathcal{M}}$  such that*

$$h^2(s, \hat{s}_{m'}) \leq C \left\{ h^2(s, \mathcal{P}_r(m)) + \frac{(r+1)|m| \log_+^2(n/(|m|(r+1)))}{n} + \xi \log_+(1/\xi) \right\}.$$

Moreover,  $|m'| \leq 2|m| + 3$ .

Therefore, when  $\hat{k} = \hat{n}$ , when (28) is satisfied, and when  $L \geq \max\{1, L_0\}$ , we deduce from (32) that  $\hat{s}_{\hat{m}}$  satisfies

$$\mathbb{E} [h^2(s, \hat{s}_{\hat{m}})] \leq C'' \mathbb{E} \left[ \inf_{m \in \mathcal{M}} \left\{ h^2(s, \mathcal{P}_r(m)) + L \frac{|m|(r+1) \log_+^2(n/(r+1))}{n} \right\} \right],$$



where  $C''$  is a universal constant. Now,  $\cup_{\substack{m \in \mathcal{M}, \\ |m|=\ell}} \mathcal{P}_r(m)$  is dense in  $\mathcal{P}_{\ell,r}$  in the metric space  $(\mathcal{S}, h)$ , and hence

$$\begin{aligned} \mathbb{E} [h^2(s, \hat{s}_{\hat{m}})] &\leq C'' \mathbb{E} \left[ \inf_{\ell \geq 1} \left\{ h^2(s, \mathcal{P}_{\ell,r}) + L \frac{\ell(r+1) \log_+^2(n/(r+1))}{n} \right\} \right], \\ &\leq C'' \inf_{\ell \geq 1} \mathcal{R}(\ell), \end{aligned}$$

where

$$\mathcal{R}(\ell) = \mathbb{E} [h^2(s, \mathcal{P}_{\ell,r})] + L \frac{\ell(r+1) \log_+^2(n/(r+1))}{n}.$$

This term  $\mathcal{R}(\ell)$  can be interpreted as an upper-bound of the risk of the  $\rho$ -estimator on  $\mathcal{P}_{\ell,r}$  (up to constants), barely worse than the one given by Theorem 7 and that is written in (22).

Remark. As in Section 4.2, we could compare this estimator with the one that maximizes a penalized log-likelihood criterion of the form  $m \mapsto L(\hat{s}_m) - \text{pen}(m)$  over  $m \in \widehat{\mathcal{M}}_{\hat{k}, \text{lower}}$ . We do not know theoretical results for this estimator when  $\{Y_i, i \in \hat{I}\}$  is random, but refer to [RMG10] for a numerical study in framework 1.

## 5. NUMERICAL SIMULATIONS

We consider framework 1,  $r = 0$ ,  $\ell \in \{1, \dots, n\}$ ,  $\{Y_i, i \in \hat{I}\} = \{X_1, \dots, X_n\}$  and the (random) collection  $\widehat{\mathcal{M}}_\ell$  consisting of partitions of  $[X_{(1)}, X_{(n)}]$  of size  $\ell$  defined in Section 4.2. For each  $m \in \widehat{\mathcal{M}}_\ell$ , we consider the  $\rho$ - and maximum likelihood estimator  $\hat{s}_m$  on  $\mathcal{P}_0(m)$  defined by

$$\hat{s}_m = \sum_{K \in m} \frac{N(K)}{\mu(K)} \quad \text{with } N(K) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_K(X_i).$$

We carry out in this section a numerical study to compare two selection rules described in Sections 4.2 and 4.3.

- The first procedure is based on the likelihood. We select the partition  $\hat{m}^{(1,\ell)} \in \widehat{\mathcal{M}}_\ell$  by maximizing the map

$$m \mapsto L(\hat{s}_m) = \frac{1}{n} \sum_{i=1}^n \log \hat{s}_m(X_i) \quad \text{over } m \in \widehat{\mathcal{M}}_\ell.$$

- The second procedure is based on the  $\rho$ -estimation method. We consider a set  $A$  consisting of 300 equally spaced points over  $[0, 3]$ , and define

$$\mathcal{L} = \left\{ \frac{a}{\log^2 n}, a \in A \right\}.$$

For each  $L \in \mathcal{L}$ , we use the procedure of Section 4.2 specified in (24) and (25) to get a partition  $\hat{m}_L \in \widehat{\mathcal{M}}_\ell$ . We then use the procedure of Section 4.1 to pick out an estimator among  $\{\hat{s}_{\hat{m}_L}, L \in \mathcal{L}\}$  as explained in Section 4.3. This leads to a selected partition of the form  $\hat{m}_{\hat{L}} \in \widehat{\mathcal{M}}_\ell$  that will be denoted in the sequel by  $\hat{m}^{(2,\ell)}$ .

We consider four densities  $s$ :

**Example 1.**  $s$  is the density of a Normal distribution

$$s(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2} \quad \text{for all } x \in \mathbb{R}.$$

**Example 2.**  $s$  is the density of a log Normal distribution

$$s(x) = \frac{1}{x\sqrt{2\pi}} e^{-\frac{1}{2}\log^2 x} \mathbb{1}_{(0,+\infty)}(x) \quad \text{for all } x \in \mathbb{R}.$$

**Example 3.**  $s$  is the density of an exponential distribution

$$s(x) = e^{-x} \mathbb{1}_{[0,+\infty)} \quad \text{for all } x \in \mathbb{R}.$$

**Example 4.**  $s$  is the density of a mixture of uniform distributions

$$s(x) = \frac{1}{2} \times 3\mathbb{1}_{[0,1/3]} + \frac{1}{8} \times 3\mathbb{1}_{[1/3,2/3]} + \frac{3}{8} \times 3\mathbb{1}_{[2/3,1]} \quad \text{for all } x \in \mathbb{R}.$$

We simulate  $N_{\text{rep}}$  samples  $(X_1, \dots, X_n)$  according to a density  $s$  defined above, and compute, in each of these samples the two selected estimators. Let, for  $k \in \{1, 2\}$  and  $i \in \{1, \dots, N_{\text{rep}}\}$ ,  $\hat{s}_{\hat{m}^{(k,\ell,i)}}$  be the value of the estimator corresponding to the  $k^{\text{th}}$  procedure and the  $i^{\text{th}}$  sample. We evaluate the quality of the estimators by

$$\hat{R}(k, \ell) = \frac{1}{N_{\text{rep}}} \sum_{i=1}^{N_{\text{rep}}} h^2(s, \hat{s}_{\hat{m}^{(k,\ell,i)}}).$$

We estimate the probability that the two estimators coincide by

$$\hat{P}_{\text{equal}}(\ell) = \frac{1}{N_{\text{rep}}} \sum_{i=1}^{N_{\text{rep}}} \mathbb{1}_{\hat{m}^{(2,\ell,i)} = \hat{m}^{(1,\ell,i)}}$$

Results are summarized in Figures 1 (when  $n = 50$ ) and 2 (when  $n = 100$ ).

	Ex 1	Ex 2	Ex 3	Ex 4		Ex 1	Ex 2	Ex 3	Ex 4
$\hat{R}(1, 2)$	0.057	0.078	0.064	0.052	$\hat{R}(1, 5)$	0.062	0.063	0.061	0.060
$\hat{R}(2, 2)$	0.057	0.080	0.065	0.051	$\hat{R}(2, 5)$	0.059	0.062	0.059	0.060
$\frac{\hat{R}(2,2)}{\hat{R}(1,2)}$	1.00	1.02	1.02	0.99	$\frac{\hat{R}(2,5)}{\hat{R}(1,5)}$	0.95	0.98	0.98	1.00
$\hat{P}_{\text{equal}}(2)$	0.76	0.75	0.80	0.78	$\hat{P}_{\text{equal}}(5)$	0.27	0.33	0.32	0.39
$\hat{R}(1, 3)$	0.052	0.056	0.053	0.048	$\hat{R}(1, 6)$	0.067	0.068	0.066	0.065
$\hat{R}(2, 3)$	0.047	0.055	0.052	0.047	$\hat{R}(2, 6)$	0.065	0.067	0.065	0.065
$\frac{\hat{R}(2,3)}{\hat{R}(1,3)}$	0.91	0.98	0.97	0.99	$\frac{\hat{R}(2,6)}{\hat{R}(1,6)}$	0.97	0.99	0.99	1.00
$\hat{P}_{\text{equal}}(3)$	0.63	0.64	0.66	0.57	$\hat{P}_{\text{equal}}(6)$	0.28	0.33	0.33	0.37
$\hat{R}(1, 4)$	0.057	0.058	0.056	0.054	$\hat{R}(1, 7)$	0.071	0.072	0.071	0.070
$\hat{R}(2, 4)$	0.052	0.055	0.053	0.053	$\hat{R}(2, 7)$	0.070	0.072	0.070	0.071
$\frac{\hat{R}(2,4)}{\hat{R}(1,4)}$	0.92	0.94	0.95	0.98	$\frac{\hat{R}(2,7)}{\hat{R}(1,7)}$	0.99	1.00	1.00	1.00
$\hat{P}_{\text{equal}}(4)$	0.32	0.40	0.40	0.43	$\hat{P}_{\text{equal}}(7)$	0.32	0.36	0.35	0.41

FIGURE 1. Results for simulated data with  $n = 50$ ,  $N_{\text{rep}} = 10000$ .

	Ex 1	Ex 2	Ex 3	Ex 4		Ex 1	Ex 2	Ex 3	Ex 4
$\widehat{R}(1, 2)$	0.055	0.074	0.056	0.035	$\widehat{R}(1, 5)$	0.038	0.038	0.037	0.033
$\widehat{R}(2, 2)$	0.056	0.076	0.057	0.034	$\widehat{R}(2, 5)$	0.035	0.034	0.035	0.033
$\frac{\widehat{R}(2,2)}{\widehat{R}(1,2)}$	1.03	1.02	1.02	0.98	$\frac{\widehat{R}(2,5)}{\widehat{R}(1,5)}$	0.92	0.94	0.95	1.00
$\widehat{P}_{equal}(2)$	0.63	0.60	0.70	0.80	$\widehat{P}_{equal}(5)$	0.15	0.18	0.17	0.23
$\widehat{R}(1, 3)$	0.034	0.042	0.037	0.023	$\widehat{R}(1, 6)$	0.041	0.040	0.039	0.037
$\widehat{R}(2, 3)$	0.033	0.042	0.036	0.024	$\widehat{R}(2, 6)$	0.039	0.040	0.038	0.037
$\frac{\widehat{R}(2,3)}{\widehat{R}(1,3)}$	0.96	1.00	0.98	1.01	$\frac{\widehat{R}(2,6)}{\widehat{R}(1,6)}$	0.95	0.97	0.98	1.00
$\widehat{P}_{equal}(3)$	0.71	0.63	0.63	0.57	$\widehat{P}_{equal}(6)$	0.10	0.15	0.11	0.19
$\widehat{R}(1, 4)$	0.036	0.035	0.034	0.028	$\widehat{R}(1, 7)$	0.044	0.043	0.043	0.40
$\widehat{R}(2, 4)$	0.032	0.034	0.032	0.028	$\widehat{R}(2, 7)$	0.043	0.043	0.042	0.40
$\frac{\widehat{R}(2,4)}{\widehat{R}(1,4)}$	0.90	0.96	0.94	0.98	$\frac{\widehat{R}(2,7)}{\widehat{R}(1,7)}$	0.96	0.99	0.98	1.00
$\widehat{P}_{equal}(4)$	0.29	0.39	0.35	0.33	$\widehat{P}_{equal}(7)$	0.09	0.11	0.11	0.16

FIGURE 2. Results for simulated data with  $n = 100$ ,  $N_{\text{rep}} = 1000$ .

Numerically, we observe in these examples that the two estimators  $\hat{s}_{\widehat{m}(1,\ell)}$  and  $\hat{s}_{\widehat{m}(2,\ell)}$  perform similarly. Their risks are close and the estimators may even coincide. In Example 4,  $s$  does belong to  $\mathcal{P}_{3,0}$  and the fractions  $\widehat{R}(2,\ell)/\widehat{R}(1,\ell)$  are very close to 1. In the other examples,  $s$  is not piecewise constant, and the robustness properties of the second procedure may be useful. The fractions  $\widehat{R}(2,\ell)/\widehat{R}(1,\ell)$  suggest indeed that the second procedure improves the risk of the first one by a few percent, at least when the size  $\ell$  of the partitions is well adapted to the underlying density, that is when  $\ell$  corresponds to the smallest values of  $\widehat{R}(1,\ell)$  and  $\widehat{R}(2,\ell)$ .

Remark. The fractions  $\widehat{R}(2,\ell)/\widehat{R}(1,\ell)$  are computed with all significant digits and are then rounded.

## 6. PROOFS

6.1. **Proof of Lemma 1.** Let  $\sqrt{q} = (\sqrt{f} + \sqrt{g})/2$  and

$$\mathcal{X} = \{x \in \mathbb{R}, g(x) \neq 0 \text{ or } f(x) \neq 0\}.$$

Then,

$$\begin{aligned} \frac{1}{2} \int_{\mathcal{X}} \frac{\sqrt{g} - \sqrt{f}}{\sqrt{q}} (\sqrt{s} - \sqrt{q})^2 dM &= \frac{1}{2} \int_{\mathcal{X}} \frac{\sqrt{g} - \sqrt{f}}{\sqrt{q}} s dM + \frac{1}{2} \int_{\mathcal{X}} (\sqrt{g} - \sqrt{f}) \sqrt{q} dM \\ &\quad - \int_{\mathcal{X}} (\sqrt{g} - \sqrt{f}) \sqrt{s} dM. \end{aligned}$$

Note that

$$h^2(s, g) - h^2(s, f) = \frac{1}{2} \int_{\mathcal{X}} (g - f) dM + \int_{\mathcal{X}} \sqrt{s} (\sqrt{f} - \sqrt{g}) dM.$$

Therefore,

$$\begin{aligned}
\frac{1}{2} \int_{\mathcal{X}} \frac{\sqrt{g} - \sqrt{f}}{\sqrt{q}} (\sqrt{s} - \sqrt{q})^2 \, dM &= \frac{1}{2} \int_{\mathcal{X}} \frac{\sqrt{g} - \sqrt{f}}{\sqrt{q}} s \, dM + \frac{1}{2} \int_{\mathcal{X}} (\sqrt{g} - \sqrt{f}) \sqrt{q} \, dM \\
&\quad - \frac{1}{2} \int_{\mathcal{X}} (g - f) \, dM + h^2(s, g) - h^2(s, f) \\
(33) \qquad \qquad \qquad &= T_E(f, g) + h^2(s, g) - h^2(s, f).
\end{aligned}$$

Now,

$$\begin{aligned}
\frac{1}{2} \int_{\mathcal{X}} \frac{\sqrt{g} - \sqrt{f}}{\sqrt{q}} (\sqrt{s} - \sqrt{q})^2 \, dM &= \int_{\mathcal{X}} \frac{\sqrt{g} - \sqrt{f}}{\sqrt{f} + \sqrt{g}} \left( \sqrt{s} - \frac{\sqrt{f} + \sqrt{g}}{2} \right)^2 \, dM \\
&\leq \int_{\mathcal{X}} \left( \sqrt{s} - \frac{\sqrt{f} + \sqrt{g}}{2} \right)^2 \, dM \\
&\leq \frac{1}{4} \int_{\mathcal{X}} \left( (\sqrt{s} - \sqrt{f}) + (\sqrt{s} - \sqrt{g}) \right)^2 \, dM.
\end{aligned}$$

By using the inequality  $(x + y)^2 \leq (1 + \alpha)x^2 + (1 + \alpha^{-1})y^2$ ,

$$\begin{aligned}
\frac{1}{2} \int_{\mathcal{X}} \frac{\sqrt{g} - \sqrt{f}}{\sqrt{q}} (\sqrt{s} - \sqrt{q})^2 \, dM &\leq \frac{1 + \alpha}{4} \int_{\mathcal{X}} (\sqrt{s} - \sqrt{f})^2 \, dM + \frac{1 + \alpha^{-1}}{4} \int_{\mathcal{X}} (\sqrt{s} - \sqrt{g})^2 \, dM \\
&\leq \frac{1 + \alpha}{2} h^2(s, f) + \frac{1 + \alpha^{-1}}{2} h^2(s, g).
\end{aligned}$$

We now plug this inequality into (33) to get

$$T_E(f, g) \leq \frac{3 + \alpha}{2} h^2(s, f) - \frac{1 - \alpha^{-1}}{2} h^2(s, g).$$

The right-hand side of (3) follows from this inequality with  $\alpha = 3$ . As to the left-hand side, note that we also have (setting  $\alpha = 3$ , and exchanging the role of  $f$  and  $g$ ),

$$T_E(g, f) \leq 3h^2(s, g) - \frac{1}{3}h^2(s, f).$$

Yet,  $T_E(f, g) = -T_E(g, f)$  and hence  $T_E(f, g) \geq \frac{1}{3}h^2(s, f) - 3h^2(s, g)$  as wished.  $\square$

**6.2. Proof of Theorem 1.** In each framework, the measure  $N$  can be put of the form  $N(A) = n^{-1} \sum_{i \in \hat{I}} \mathbb{1}_A(Y_i)$  where  $\hat{I} \subset \{1, \dots, n\}$ , and where the  $Y_i$  are suitable real-valued random variables. For instance, in framework 1,  $\hat{I} = \{1, \dots, n\}$ ,  $Y_i = X_i$ , in framework 2,  $\hat{I} = \{i \in \{1, \dots, n\}, D_i = 1\}$ ,  $Y_i = X_i$ , and in framework 3,  $\hat{I} = \{i \in \{1, \dots, n\}, T_{1,0}^{(i)} \in I_{\text{obs}}\}$ ,  $Y_i = T_{1,0}^{(i)}$ .

Set  $\hat{J} = \{i \in \hat{I}, Y_i \in \mathcal{X}\}$ . Then, for  $f, g \in \mathcal{S}$ ,  $T(f, g)$  and  $L_{\mathcal{X}}(f)$  take the form

$$\begin{aligned}
T(f, g) &= \frac{1}{n} \sum_{j \in \hat{J}} \psi \left( \frac{g(Y_j)}{f(Y_j)} \right) - \frac{1}{4} \int_{\mathcal{X}} (g(x) - f(x)) \, dM(x) \\
L_{\mathcal{X}}(f) &= \frac{1}{n} \sum_{j \in \hat{J}} \log f(Y_j) - \int_{\mathcal{X}} f(x) \, dM(x).
\end{aligned}$$

The proof is straightforward if  $\hat{J} = \emptyset$  since then  $4T(f, g) = L_{\mathcal{X}}(g) - L_{\mathcal{X}}(f)$  and  $4\gamma(f) = \sup_{g \in S} L_{\mathcal{X}}(g) - L_{\mathcal{X}}(f)$ . We suppose from now on that  $\hat{J} \neq \emptyset$ .

**Claim 1.** *Let  $\bar{S} = \{f \in S, L_{\mathcal{X}}(f) \neq -\infty\}$  and  $\bar{f} \in \bar{S}$ . Then,  $\sup_{g \in \bar{S}} T(\bar{f}, g) = 0$  if and only if  $\sup_{g \in \bar{S}} L_{\mathcal{X}}(g) - L_{\mathcal{X}}(\bar{f}) = 0$ .*

*Proof.* Suppose that  $\sup_{g \in \bar{S}} L_{\mathcal{X}}(g) - L_{\mathcal{X}}(\bar{f}) = 0$ .

Let  $\bar{S}_1 = \{g \in \bar{S}, g = \bar{f}, N \text{ a.s.}\}$  and  $\bar{S}_2 = \bar{S} \setminus \bar{S}_1$ . When  $g \in \bar{S}_1$ ,

$$\begin{aligned} T(\bar{f}, g) &= -\frac{1}{4} \int_{\mathcal{X}} (g(x) - \bar{f}(x)) \, dM(x) \\ &= \frac{1}{4} (L_{\mathcal{X}}(g) - L_{\mathcal{X}}(\bar{f})). \end{aligned}$$

Therefore,  $T(\bar{f}, g) \leq 0$ .

Let now  $g \in \bar{S}_2$ ,  $u \in [0, 1]$  and  $\zeta = g - \bar{f}$ . Note that  $\bar{f} + u\zeta = (1-u)\bar{f} + ug \in \bar{S}$  and thus  $L_{\mathcal{X}}(\bar{f} + u\zeta) - L_{\mathcal{X}}(\bar{f}) \leq 0$ . We introduce the real-valued map  $\wp_1$  for  $u \in [0, 1]$  by

$$\begin{aligned} \wp_1(u) &= L_{\mathcal{X}}(\bar{f} + u\zeta) - L_{\mathcal{X}}(\bar{f}) \\ &= \frac{1}{n} \sum_{j \in \hat{J}} \log \left( \frac{\bar{f}(Y_j) + u\zeta(Y_j)}{\bar{f}(Y_j)} \right) - u \int_{\mathcal{X}} \zeta(x) \, dM(x). \end{aligned}$$

We now define  $\wp_2$  for  $u \in [0, 1]$  by

$$\begin{aligned} \wp_2(u) &= 4T(\bar{f}, \bar{f} + u\zeta) \\ &= \frac{4}{n} \sum_{j \in \hat{J}} \psi \left( \frac{\bar{f}(Y_j) + u\zeta(Y_j)}{\bar{f}(Y_j)} \right) - u \int_{\mathcal{X}} \zeta(x) \, dM(x). \end{aligned}$$

Some computations show that  $\wp_1$  and  $\wp_2$  are twice differentiable on  $[0, 1]$  and

$$\begin{aligned} \wp_1(0) &= \wp_2(0) = 0 \\ \wp_1'(0) &= \wp_2'(0) = \frac{1}{n} \sum_{j \in \hat{J}} \frac{\zeta(Y_j)}{\bar{f}(Y_j)} - \int_{\mathcal{X}} \zeta(x) \, dM(x) \\ \wp_1''(0) &= \wp_2''(0) = -\frac{1}{n} \sum_{j \in \hat{J}} \left( \frac{\zeta(Y_j)}{\bar{f}(Y_j)} \right)^2. \end{aligned}$$

Therefore,  $\wp_1''(0)$  and  $\wp_2''(0)$  are always negative.

Since  $\wp_1(u)$  is non-positive for all  $u \in [0, 1]$ ,  $\wp_1'(0) \leq 0$ . The above computations show the existence of  $u_1 \in (0, 1]$  such that  $\wp_2(u) \leq 0$  for all  $u \in [0, u_1]$ . Now,  $\wp_2$  is concave, and hence non-positive on  $[0, 1]$ . In particular,  $\wp_2(1) = T(\bar{f}, g) \leq 0$ .

Likewise,  $\sup_{g \in \bar{S}} T(\bar{f}, g) = 0$  implies  $\sup_{g \in \bar{S}} L_{\mathcal{X}}(g) - L_{\mathcal{X}}(\bar{f}) = 0$ .

□

Let  $\tilde{s} \in S$  such that  $L_{\mathcal{X}}(\tilde{s}) \geq L_{\mathcal{X}}(g)$  for all  $g \in S$  and  $L_{\mathcal{X}}(\tilde{s}) \neq -\infty$ . The above claim then shows that  $T(\tilde{s}, g) \leq 0$  for all  $g \in S$  such that  $L_{\mathcal{X}}(g) \neq -\infty$ . Choose now  $g \in S$  such that  $L_{\mathcal{X}}(g) = -\infty$ .

Define for  $u \in [0, 1]$ ,  $f_u = (1 - u)\tilde{s} + ug \in S$  and note that  $f_1 = g$ . If  $u \in [0, 1)$ ,  $L_{\mathcal{X}}(f_u) \neq -\infty$  and thus  $T(\tilde{s}, f_u) \leq 0$ . The continuity of the map  $u \in [0, 1] \mapsto T(\tilde{s}, f_u)$  ensures that  $T(\tilde{s}, g) \leq 0$ . Finally,  $\gamma(\tilde{s}) = 0$ .

Conversely, let  $\hat{s}$  be a  $\rho$ -estimator satisfying  $\gamma(\hat{s}) = 0$ . We begin by proving that  $L_{\mathcal{X}}(\hat{s}) \neq -\infty$ . Consider  $g \in S$  such that  $L_{\mathcal{X}}(g) \neq -\infty$  and define for  $u \in [0, 1]$ ,  $f_u = (1 - u)\hat{s} + ug \in S$ ,

$$\begin{aligned} \wp_3(u) &= T(\hat{s}, f_u) \\ &= \frac{1}{n} \sum_{j \in \hat{J}} \psi \left( \frac{(1 - u)\hat{s}(Y_j) + ug(Y_j)}{\hat{s}(Y_j)} \right) - \frac{1}{4} \int_{\mathcal{X}} (f_u(x) - \hat{s}(x)) \, dM(x). \end{aligned}$$

When  $j \in \hat{J}$ ,  $g(Y_j) > 0$ . Therefore, if  $\hat{J}' = \{j \in \hat{J}, \hat{s}(Y_j) = 0\}$  and  $u \in (0, 1]$ ,

$$\wp_3(u) = \frac{|\hat{J}'|}{n} + \frac{1}{n} \sum_{j \in \hat{J} \setminus \hat{J}'} \psi \left( \frac{(1 - u)\hat{s}(Y_j) + ug(Y_j)}{\hat{s}(Y_j)} \right) - \frac{1}{4} \int_{\mathcal{X}} (f_u(x) - \hat{s}(x)) \, dM(x).$$

Therefore, if  $\hat{J}' \neq \emptyset$  choosing  $u > 0$  small enough leads to  $\wp_3(u) \geq |\hat{J}'|/(2n) > 0$ , which is impossible as  $\gamma(\hat{s}) = 0$ . Therefore,  $\hat{J}' = \emptyset$  and  $L_{\mathcal{X}}(\hat{s}) \neq -\infty$ . The claim then asserts that for all  $g \in S$  such that  $L_{\mathcal{X}}(g) \neq -\infty$ ,  $L_{\mathcal{X}}(g) \leq L_{\mathcal{X}}(\hat{s})$ . This inequality being true if  $L_{\mathcal{X}}(g) = -\infty$ , the proof is complete.  $\square$

**6.3. Sketch of the proof of Theorem 2.** We define the elements  $Y_i$ ,  $\hat{I}$ ,  $\hat{J}$  as in the proof of Theorem 1. Let for  $x \in [0, +\infty]$ ,  $\psi_2(x) = \psi(x^2)$  and for  $f, g \in \mathcal{F}$ ,

$$\begin{aligned} T_2(f, g) &= T(f^2, g^2) = \frac{1}{n} \sum_{j \in \hat{J}} \psi_2 \left( \frac{g(Y_j)}{f(Y_j)} \right) - \frac{1}{4} \int_{\mathcal{X}} (g^2(x) - f^2(x)) \, dM(x), \\ L_{\mathcal{X},2}(f) &= L_{\mathcal{X}}(f^2) = \frac{2}{n} \sum_{j \in \hat{J}} \log f(Y_j) - \int_{\mathcal{X}} f^2(x) \, dM(x). \end{aligned}$$

The proof is very similar to the one of Theorem 1. The main change lies in the replacement of the symbols  $S$ ,  $T$ ,  $L_{\mathcal{X}}$  by  $\mathcal{F}$ ,  $T_2$ ,  $L_{\mathcal{X},2}$ . We will only give some insight into why Claim 1 remain valid under these modifications.

As in the proof of Theorem 1, we may suppose that  $\hat{J} \neq \emptyset$ .

**Claim 2.** Let  $\bar{\mathcal{F}} = \{g \in \mathcal{F}, L_{\mathcal{X},2}(g) \neq -\infty\}$  and  $\bar{f} \in \bar{\mathcal{F}}$ . Then,  $\sup_{g \in \bar{\mathcal{F}}} T_2(\bar{f}, g) = 0$  if and only if  $\sup_{g \in \bar{\mathcal{F}}} L_{\mathcal{X},2}(g) - L_{\mathcal{X},2}(\bar{f}) = 0$ .

*Sketch of the proof.* We prove that  $\sup_{g \in \bar{\mathcal{F}}} L_{\mathcal{X},2}(g) - L_{\mathcal{X},2}(\bar{f}) = 0$  implies  $\sup_{g \in \bar{\mathcal{F}}} T_2(\bar{f}, g) = 0$ . The proof of the converse is similar. Let  $\bar{\mathcal{F}}_1 = \{g \in \bar{\mathcal{F}}, g = \bar{f}, N \text{ a.s}\}$  and  $\bar{\mathcal{F}}_2 = \bar{\mathcal{F}} \setminus \bar{\mathcal{F}}_1$ . As in the proof of Claim 1,  $T_2(\bar{f}, g) = (L_{\mathcal{X},2}(g) - L_{\mathcal{X},2}(\bar{f}))/4$  when  $g \in \bar{\mathcal{F}}_1$  and is thus non-positive. Let now  $g \in \bar{\mathcal{F}}_2$ ,  $u \in [0, 1]$  and  $\zeta = g - \bar{f}$ . Note that  $\bar{f} + u\zeta = (1 - u)\bar{f} + ug \in \bar{\mathcal{F}}$  and thus  $L_{\mathcal{X},2}(\bar{f} + u\zeta) - L_{\mathcal{X},2}(\bar{f}) \leq 0$ .

We introduce the real-valued map  $\wp_1$  for  $u \in [0, 1]$  by

$$\begin{aligned}\wp_1(u) &= L_{\mathcal{X},2}(\bar{f} + u\zeta) - L_{\mathcal{X},2}(\bar{f}) \\ &= \frac{2}{n} \sum_{j \in \mathcal{J}} \log \left( \frac{\bar{f}(Y_j) + u\zeta(Y_j)}{f(Y_j)} \right) - u^2 \int_{\mathcal{X}} \zeta^2(x) dM(x) - 2u \int_{\mathcal{X}} \zeta(x) \bar{f}(x) dM(x).\end{aligned}$$

We now define  $\wp_2$  for  $u \in [0, 1]$  by

$$\begin{aligned}\wp_2(u) &= 4T_2(\bar{f}, \bar{f} + u\zeta) \\ &= \frac{4}{n} \sum_{j \in \mathcal{J}} \psi_2 \left( \frac{\bar{f}(Y_j) + u\zeta(Y_j)}{f(Y_j)} \right) - u^2 \int_{\mathcal{X}} \zeta^2(x) dM(x) - 2u \int_{\mathcal{X}} \zeta(x) \bar{f}(x) dM(x).\end{aligned}$$

Some computations show that  $\wp_1(0) = \wp_2(0) = 0$ ,  $\wp_1'(0) = \wp_2'(0)$ ,  $\wp_1''(0) = \wp_2''(0) < 0$ .

As  $\wp_1(u)$  is non-positive for all  $u \in [0, 1]$ ,  $\wp_1'(0) \leq 0$ . There exists therefore  $u_1 \in (0, 1]$  such that  $\wp_2(u) \leq 0$  for all  $u \in [0, u_1]$ . Since  $\psi_2$  is concave,  $\wp_2$  is also concave, and  $\wp_2$  is non-positive on  $[0, 1]$ . In particular,  $\wp_2(1) = T_2(\bar{f}, g) \leq 0$ .

□

**6.4. Proof of Theorem 3.** We introduce the random measure  $M_s$  defined by

$$M_s(A) = \int_A s(t) dM(t) \quad \text{for all } A \in \mathcal{B}(\mathbb{R}).$$

Note that  $\mathbb{E}[N(A)] = \mathbb{E}[M_s(A)]$  for all  $A \in \mathcal{B}(\mathbb{R})$ .

We begin by showing that the problem boils down to a suitable control of the deviations of  $N(A) - \mathbb{E}[N(A)]$  and  $M_s(A) - \mathbb{E}[M_s(A)]$ .

**Lemma 6.** *Let  $\mathcal{F}$  be a collection of functions of  $\mathcal{S}$  such that  $|\varphi| \leq 1$  for all  $\varphi \in \mathcal{F}$ . Consider a collection  $\mathcal{A} \subset \mathcal{B}(\mathbb{R})$  such that*

$$\mathcal{A} \supset \bigcup_{t \in (0,1)} \{ \{x \in \mathbb{R}, \varphi_+(x) > t\}, \varphi \in \mathcal{F} \} \cup \{ \{x \in \mathbb{R}, \varphi_-(x) > t\}, \varphi \in \mathcal{F} \}.$$

*Suppose that there exist  $\alpha, \beta$  and an event on which: for all  $A \in \mathcal{A}$ ,*

$$(34) \quad |N(A) - \mathbb{E}[N(A)]| \leq \sqrt{\frac{\alpha}{n}} \left( \sqrt{N(A)} + \sqrt{\mathbb{E}[N(A)]} \right) + \frac{\beta}{n}.$$

$$(35) \quad |M_s(A) - \mathbb{E}[M_s(A)]| \leq \sqrt{\frac{\alpha}{n}} \left( \sqrt{M_s(A)} + \sqrt{\mathbb{E}[M_s(A)]} \right) + \frac{\beta}{n}.$$

*Then, on this event, for all  $\varphi \in \mathcal{F}$ ,*

$$|Z(\varphi)| \leq C \left\{ \sqrt{\frac{\alpha}{n} v(\varphi) \log_+(1/v(\varphi))} + \frac{\alpha + \beta}{n} \right\}.$$

*The constant  $C$  above is universal.*

The proof of this result is delayed to Section 6.5 below. It remains to verify that (34) and (35) hold true in our different statistical settings. As  $N$  is an empirical measure, we can prove that (34) is valid on an event of high probability when  $\mathcal{A}$  is a Vapnik-Chervonenkis class of finite dimension (this result follows, for example, from Vapnik-Chervonenkis inequalities for relative deviations).

The proof of (35) in frameworks 2 and 3 is more involved. It is based on the notion of entropy with bracketing. This is the reason why in the lemma below, we only prove (35) for a class  $\mathcal{A}$  of unions of at most  $d$  intervals. Its proof is deferred to Section 6.6.

**Lemma 7.** *Let  $\mathcal{A} \subset \mathcal{B}(\mathbb{R})$  be an at most countable collection of Borel sets and  $S_{\mathcal{A}}(2n)$  be the Vapnik-Chervonenkis shatter coefficient defined by*

$$S_{\mathcal{A}}(2n) = \max_{x_1, \dots, x_{2n} \in \mathbb{R}} |\{\{x_1, \dots, x_{2n}\} \cap A, A \in \mathcal{A}\}|.$$

Let  $\xi > 0$ . There exist a universal constant  $c$  and an event  $\Omega_{\xi,1}$  such that  $\mathbb{P}[\Omega_{\xi,1}] \geq 1 - e^{-n\xi}$  and on which (34) holds for all  $A \in \mathcal{A}$  with  $\alpha = c[\log_+ |S_{\mathcal{A}}(2n)| + n\xi]$  and  $\beta = 0$ .

Let for  $d \geq 1$ ,  $\overline{\mathcal{I}}_d$  be the class of unions of at most  $d$  intervals with endpoints in  $\mathbb{Q} \cup \{-\infty, +\infty\}$ . Then,  $\overline{\mathcal{I}}_d$  is at most countable, and Sauer lemma implies

$$(36) \quad \log_+ |S_{\overline{\mathcal{I}}_d}(2n)| \leq 4d \log_+(n/d).$$

Moreover, there exist a universal constant  $c'$  and an event  $\Omega_{\xi,2}$  such that  $\mathbb{P}[\Omega_{\xi,2}] \geq 1 - e^{-n\xi}$  and on which (35) holds true for all  $A \in \overline{\mathcal{I}}_d$  with  $\alpha = \beta = c'[d \log_+(n/d) + n\xi]$ .

We may now finish the proof of Theorem 3. Lemma 7 implies that (34) and (35) hold true for the collection  $\overline{\mathcal{I}}_d$  with  $\alpha = \beta = c'[d \log_+(n/d) + n\xi]$  on an event  $\Omega_{\xi}$  such that  $\mathbb{P}[\Omega_{\xi}] \geq 1 - e^{-n\xi}$ . Let now  $\mathcal{I}_d$  be the class of unions of at most  $d$  intervals. Then, for all  $\epsilon > 0$ ,  $A \in \mathcal{I}_d$ , there exists  $\bar{A} \in \overline{\mathcal{I}}_d$  such that

$$|N(A) - N(\bar{A})| \leq \epsilon \quad \text{and} \quad |M_s(A) - M_s(\bar{A})| \leq \epsilon.$$

Thereby, (34) and (35) hold with  $\mathcal{A} = \mathcal{I}_d$  (up to a modification of  $\beta$ ). Lemma 6 finally implies (8).

Elementary computations then show (9). As  $x \mapsto x \log_+(1/x)$  is non-decreasing, for all  $x \leq y$ ,  $xy \log_+(1/x) \leq y^2 \log_+(1/y)$ . Moreover, when  $x \geq y$ ,  $\log_+(1/x) \leq \log_+(1/y)$  and hence  $xy \log_+(1/x) \leq xy \log_+(1/y)$ . Therefore, for all  $x, y > 0$ ,

$$\begin{aligned} xy \log_+(1/x) &\leq \max\{x, y\} y \log_+(1/y) \\ &\leq (x + y) y \log_+(1/y). \end{aligned}$$

We thus obtain for all  $\varepsilon > 0$ ,

$$\begin{aligned} 2\sqrt{xy \log_+(1/x)} &\leq \varepsilon(x + y) + \varepsilon^{-1} y \log_+(1/y) \\ &\leq \varepsilon x + C_{\varepsilon} y \log_+(1/y), \end{aligned}$$

where  $C_{\varepsilon}$  depends on  $\varepsilon$ . By using this inequality with  $x = v(\varphi)$  and  $y = (d/n) \log_+(n/d)$ ,

$$\begin{aligned} 2\sqrt{v(\varphi) \log_+(1/v(\varphi)) \left(\frac{d \log_+(n/d)}{n}\right)} &\leq \varepsilon v(\varphi) + C_{\varepsilon} \left(\frac{d \log_+(n/d)}{n}\right) \log_+ \left(\frac{1}{\frac{d \log_+(n/d)}{n}}\right) \\ &\leq \varepsilon v(\varphi) + C_{\varepsilon} \frac{d \log_+^2(n/d)}{n}. \end{aligned}$$

Similarly,

$$2\sqrt{v(\varphi) \log_+(1/v(\varphi)) \xi} \leq \varepsilon v(\varphi) + C_{\varepsilon} \xi \log_+(1/\xi).$$



Therefore,

$$2\sqrt{v(\varphi) \log_+(1/v(\varphi)) \left( \frac{d \log_+(n/d)}{n} + \xi \right)} \leq 2\varepsilon v(\varphi) + C_\varepsilon \left[ \frac{d \log_+^2(n/d)}{n} + \xi \log_+(1/\xi) \right],$$

and (9) follows from (8).  $\square$

**6.5. Proof of Lemma 6.** For convenience, and to make the proof more readable, we introduce a new notation. Given  $x, y \in \mathbb{R}$ , the assertion: there exists a universal constant  $C$  such that  $x \leq Cy$  is written in the sequel as  $x \preceq y$ . The claim below follows from elementary computations.

**Claim 3.** *When (34) and (35) hold true,*

$$(37) \quad |N(A) - \mathbb{E}[N(A)]| \preceq \sqrt{\frac{\alpha}{n} \min\{N(A), \mathbb{E}[N(A)]\}} + \frac{\alpha + \beta}{n}$$

$$(38) \quad |M_s(A) - \mathbb{E}[M_s(A)]| \preceq \sqrt{\frac{\alpha}{n} \min\{M_s(A), \mathbb{E}[M_s(A)]\}} + \frac{\alpha + \beta}{n}$$

$$(39) \quad |N(A) - M_s(A)| \preceq \sqrt{\frac{\alpha}{n} \min\{N(A), \mathbb{E}[N(A)], M_s(A)\}} + \frac{\alpha + \beta}{n}.$$

*Proof of Claim 3.* For reasons of symmetry, we may suppose that  $N(A) \geq \mathbb{E}[N(A)]$  to prove (37). We derive from (34),

$$\begin{aligned} |N(A) - \mathbb{E}[N(A)]| &\preceq \sqrt{\frac{\alpha}{n} N(A)} + \frac{\alpha + \beta}{n} \\ &\preceq \sqrt{\frac{\alpha}{n} (N(A) - \mathbb{E}[N(A)])} + \sqrt{\frac{\alpha}{n} \mathbb{E}[N(A)]} + \frac{\alpha + \beta}{n}. \end{aligned}$$

For all  $\varepsilon > 0$ , we deduce from the inequality  $2\sqrt{xy} \leq \varepsilon x + \varepsilon^{-1}y$ , that

$$|N(A) - \mathbb{E}[N(A)]| \leq \frac{1}{2}|N(A) - \mathbb{E}[N(A)]| + C \left[ \sqrt{\frac{\alpha}{n} \mathbb{E}[N(A)]} + \frac{\alpha + \beta}{n} \right],$$

where  $C$  is universal. This proves (37). The proof of (38) is identical.

As to (39), we use  $\mathbb{E}[N(A)] = \mathbb{E}[M_s(A)]$ , the triangular inequality

$$|N(A) - M_s(A)| \leq |N(A) - \mathbb{E}[N(A)]| + |M_s(A) - \mathbb{E}[M_s(A)]|,$$

and (37), (38):

$$|N(A) - M_s(A)| \preceq \sqrt{\frac{\alpha}{n} \mathbb{E}[N(A)]} + \frac{\alpha + \beta}{n}.$$

Moreover, when  $N(A) \geq M_s(A)$ ,

$$\begin{aligned} |N(A) - M_s(A)| &\preceq \sqrt{\frac{\alpha}{n} N(A)} + \frac{\alpha + \beta}{n} \\ &\preceq \sqrt{\frac{\alpha}{n} (N(A) - M_s(A))} + \sqrt{\frac{\alpha}{n} M_s(A)} + \frac{\alpha + \beta}{n}. \end{aligned}$$

As in the end of the proof of (37), we may deduce

$$|N(A) - M_s(A)| \preceq \sqrt{\frac{\alpha}{n} M_s(A)} + \frac{\alpha + \beta}{n},$$

which ends the proof.  $\square$

Without loss of generality, we prove the lemma when the functions  $\varphi$  of  $\mathcal{F}$  are non-negative. We suppose moreover that we are on an event on which (34) and (35) hold true. Let for  $t \in (0, 1)$ ,  $A_{\varphi,t} = \{x \in \mathbb{R}, \varphi(x) > t\}$ . As in [Bar16], the notion of integral is a great help: for all  $x \in \mathbb{R}$ ,

$$\varphi(x) = \int_0^1 \mathbb{1}_{A_{\varphi,t}}(x) dt.$$

Let  $\varepsilon > 0$  and  $\eta \in (0, 1]$  to be specified later. Since

$$\varphi^2(x) = 2 \int_0^1 t \mathbb{1}_{A_{\varphi,t}}(x) dt,$$

we get

$$\begin{aligned} |Z(\varphi)| - \varepsilon v(\varphi) &= \left| \int_0^1 (N(A_{\varphi,t}) - M_s(A_{\varphi,t})) dt \right| - 2\varepsilon \int_0^1 t M_s(A_{\varphi,t}) dt \\ &\leq \int_0^1 \{|N(A_{\varphi,t}) - M_s(A_{\varphi,t})| - 2\varepsilon t M_s(A_{\varphi,t})\} dt \\ (40) \quad &\leq \int_0^\eta |N(A_{\varphi,t}) - M_s(A_{\varphi,t})| dt + \int_\eta^1 \{|N(A_{\varphi,t}) - M_s(A_{\varphi,t})| - 2\varepsilon t M_s(A_{\varphi,t})\} dt. \end{aligned}$$

It follows from (39) and the inequality  $2\sqrt{xy} \leq \varepsilon x + \varepsilon^{-1}y$ ,

$$|N(A_{\varphi,t}) - M_s(A_{\varphi,t})| - 2\varepsilon t M_s(A_{\varphi,t}) \preceq \frac{\alpha}{n\varepsilon t} + \frac{\alpha + \beta}{n}.$$

We deduce,

$$(41) \quad |Z(\varphi)| - \varepsilon v(\varphi) \preceq \sqrt{\frac{\alpha}{n}} \int_0^\eta \sqrt{\mathbb{E}[N(A_{\varphi,t})]} dt + \frac{\alpha}{n\varepsilon} \log(1/\eta) + \frac{\alpha + \beta}{n}.$$

We now optimize this result with respect to  $\varepsilon$  and  $\eta$ :

$$|Z(\varphi)| \preceq \sqrt{\frac{\alpha}{n}} \inf_{\eta \in (0,1]} \left\{ \sqrt{v(\varphi) \log(1/\eta)} + \int_0^\eta \sqrt{\mathbb{E}[N(A_{\varphi,t})]} dt \right\} + \frac{\alpha + \beta}{n}.$$

It remains to use  $\mathbb{E}[N(A_{\varphi,t})] \leq 1$  to conclude.  $\square$

## 6.6. Proof of Lemma 7.

6.6.1. *Proof of (34).* A possible way to prove (34) is to use the celebrated Vapnik-Chervonenkis inequalities for relative deviation (see for instance page 24 of [DL12]). We recall them below:

**Theorem 12** (Vapnik-Chervonenkis inequalities for relative deviation). *Let  $Z_1, \dots, Z_n$  be  $n$  independent and identically distributed random variables with values in a space  $\mathcal{X}$ . Let  $\mathcal{A}'$  be an at most*

countable collection of measurable sets. Define the empirical measure  $\nu_n(A') = n^{-1} \sum_{i=1}^n \mathbb{1}_{A'}(Z_i)$ ,  $\nu(A') = \mathbb{E}[\mu_n(A')]$  and the Vapnik-Chervonenkis shatter coefficient

$$S_{\mathcal{A}'}(2n) = \max_{z_1, \dots, z_{2n} \in \mathcal{X}} |\{\{z_1, \dots, z_{2n}\} \cap A', A' \in \mathcal{A}'\}|.$$

Then, for all  $t > 0$ ,

$$\begin{aligned} \mathbb{P} \left( \sup_{A' \in \mathcal{A}'} \frac{\nu(A') - \nu_n(A')}{\sqrt{\nu(A')}} \geq t \right) &\leq 4S_{\mathcal{A}'}(2n)e^{-nt^2/4} \\ \mathbb{P} \left( \sup_{A' \in \mathcal{A}'} \frac{\nu_n(A') - \nu(A')}{\sqrt{\nu_n(A')}} \geq t \right) &\leq 4S_{\mathcal{A}'}(2n)e^{-nt^2/4}. \end{aligned}$$

In particular,

$$(42) \quad \mathbb{P} \left( \sup_{A' \in \mathcal{A}'} \left| \sqrt{\nu_n(A')} - \sqrt{\nu(A')} \right| \geq t \right) \leq 8S_{\mathcal{A}'}(2n)e^{-nt^2/4}.$$

Assume that we are within framework 1. Then, the random measure  $N$  is the empirical measure of  $X_1, \dots, X_n$ . Now (42) with  $\mathcal{A}' = \mathcal{A}$ ,

$$t^2 = \frac{4}{n} (\log 8 + \log_+ |S_{\mathcal{A}}(2n)| + n\xi)$$

shows that (34) holds true with probability larger than  $1 - e^{-n\xi}$ ,  $\alpha = nt^2$ ,  $\beta = 0$ .

The proof in frameworks 2 and 3 is very similar since  $N$  is an empirical measure for suitable random variables with values in  $\mathcal{X} = \mathbb{R} \times \{0, 1\}$ :  $Z_i = (X_i, \mathbb{1}_{D_i=1})$  in framework 2 and  $Z_i = (T_{1,0}^{(i)} \mathbb{1}_{T_{1,0}^{(i)} \in I_{\text{obs}}}, \mathbb{1}_{T_{1,0}^{(i)} \in I_{\text{obs}}})$  in framework 3. We apply (42) with  $\mathcal{A}' = \{A \times \{1\}, A \in \mathcal{A}\}$ . Moreover,

$$\begin{aligned} |S_{\mathcal{A}'}(2n)| &\leq \max_{x_1, \dots, x_{2n} \in \mathbb{R}} |\{\{x_1, \dots, x_{2n}\} \cap A, A \in \mathcal{A}\}| \\ &\leq |S_{\mathcal{A}}(2n)|. \end{aligned}$$

We end the proof as in framework 1.  $\square$

6.6.2. *Proof of (36).* It is known that the Vapnik-Chervonenkis dimension of  $\overline{\mathcal{I}}_d$  is at most  $2d$ . By using Sauer lemma, we deduce

$$\left| S_{\overline{\mathcal{I}}_d}(2n) \right| \leq \sum_{j=0}^{2d} C_{2n}^j.$$

By using a classical inequality (see, for instance, exercise 2.14 of [BLM13]), we deduce when  $d \leq n$ ,  $|S_{\overline{\mathcal{I}}_d}(2n)| \leq (en/d)^{2d}$ , and when  $d \geq n$ ,  $|S_{\overline{\mathcal{I}}_d}(2n)| \leq e^{2d}$ . In particular, (36) holds true.  $\square$

6.6.3. *Proof of (35).* There is nothing to prove in framework 1 as  $M_s$  is deterministic and we therefore focus on frameworks 2 and 3. We define  $V_i(t) = \mathbb{1}_{X_i \geq t} \mathbb{1}_{[0, +\infty)}(t)$  in framework 2 and  $V_i(t) = \mathbb{1}_{X_{t^-}^{(i)} = 1} \mathbb{1}_{I_{\text{obs}}}(t)$  in framework 3. Then,  $M_s$  is of the form

$$M_s(A) = \frac{1}{n} \sum_{i=1}^n \int_A s(t) V_i(t) dt \quad \text{for all } A \in \mathcal{B}(\mathbb{R}).$$

There exist independent random variables  $Z_1, \dots, Z_n$  such that

$$M_s(A) = \frac{1}{n} \sum_{i=1}^n f_A(Z_i) \quad \text{with } f_A(Z_i) = \int_A s(t) V_i(t) dt.$$

We show that the assumptions of Bernstein's deviation inequality are satisfied:

**Lemma 8.** *For all  $k \geq 1$ ,  $i \in \{1, \dots, n\}$  and  $A \in \mathcal{B}(\mathbb{R})$ ,*

$$\mathbb{E} \left[ (f_A(Z_i))^k \right] \leq k! \mathbb{E} [M_s(A)].$$

We then measure the complexity of the family  $\{f_A, A \in \overline{\mathcal{I}}_d\}$  by means of the notion of entropy with bracketing:

**Lemma 9.** *For all  $\delta > 0$ , there exists a collection  $\mathcal{C}_\delta$  of functions of the form  $f_A$  with  $A \in \overline{\mathcal{I}}_d$ . The cardinal of this set can be bounded by  $\log |\mathcal{C}_\delta| \leq cd \log_+(1/\delta^2)$ , where  $c$  is a universal constant. Moreover, for all  $A \in \overline{\mathcal{I}}_d$ , there exist  $f_{A_1}, f_{A_2} \in \mathcal{C}_\delta$  such that  $f_{A_1} \leq f_A \leq f_{A_2}$  and such that for all  $k \geq 1$ ,*

$$\mathbb{E} \left[ (f_{A_2}(Z_1) - f_{A_1}(Z_1))^k \right] \leq \frac{k!}{2} \delta^2.$$

The proofs of the two lemmas are given in Sections 6.6.4, 6.6.5 and 6.6.6.

Finally, it remains to use several times an exponential inequality of [Mas07] to control the deviations of  $M_s(A) - \mathbb{E}[M_s(A)]$ . We keep the notation  $\preceq$  introduced at the beginning of Section 6.5.

Set for  $\delta > 0$ ,  $\mathcal{B}_\delta = \mathcal{C}_\delta \cup \{-f, f \in \mathcal{C}_\delta\}$ . Note that

$$\log |\mathcal{B}_\delta| \leq \log 2 + \log |\mathcal{C}_\delta| \leq c_1 d \log_+(1/\delta^2),$$

where  $c_1$  is a universal constant. We set  $H(\delta) = c_1 d \log_+(1/\delta^2)$  and for  $\sigma \in (0, 1]$ ,

$$E = \sqrt{n} \int_0^\sigma \sqrt{H(u) \wedge n} du + 2(1 + \sigma)H(\sigma).$$

Simple arguments allow to bound  $E$  from above, see for instance page 190 of [GN15]: the fundamental theorem of calculus shows

$$\begin{aligned} \sigma \sqrt{\log(e/\sigma)} &= \int_0^\sigma \left( \sqrt{\log(e/u)} - \frac{1}{2\sqrt{\log(e/u)}} \right) du \\ &\geq \int_0^\sigma \sqrt{\log_+(1/u)} du - \sigma/2. \end{aligned}$$

Consequently,

$$(43) \quad E \preceq \sigma \sqrt{nd \log_+(1/\sigma^2)} + d \log_+(1/\sigma^2).$$

Consider  $\xi > 0$  and define  $J$  as the (possibly empty) set of non-negative integers  $j$  such that  $2^{-j} \geq d/(2n)$ . Let, for  $j \in J$ ,  $x_j = 2 \log(j+1) + 1 + n\xi$ ,  $\overline{\mathcal{A}}_j = \{A \in \overline{\mathcal{I}}_d, 2^{-j-1} \leq \mathbb{E}[M_s(A)] \leq 2^{-j}\}$ . The assumptions of Corollary 6.9 of [Mas07] are satisfied with  $\mathcal{F} = \{f_A, -f_A, A \in \overline{\mathcal{A}}_j\}$ ,  $\sigma^2 = 2^{-j+1}$ ,  $b = 1$ , and  $H(\delta) = c_1 d \log_+(1/\delta^2)$ . There exists an event  $\Omega_j$  such that  $\mathbb{P}(\Omega_j) \geq 1 - e^{-x_j}$  and on which: for all  $A \in \overline{\mathcal{A}}_j$ ,

$$n |M_s(A) - \mathbb{E}[M_s(A)]| \preceq E + \sigma \sqrt{nx_j} + x_j.$$

Therefore,

$$|M_s(A) - \mathbb{E}[M_s(A)]| \preceq \sigma \sqrt{\frac{d \log_+(1/\sigma^2) + x_j}{n}} + \frac{d \log_+(1/\sigma^2) + x_j}{n}.$$

As  $\sigma^2 \leq 4\mathbb{E}[M_s(A)]$ ,  $\sigma^2 \geq d/n$ , and  $x_j \preceq \log_+(n/d) + n\xi$ ,

$$|M_s(A) - \mathbb{E}[M_s(A)]| \preceq \sqrt{\mathbb{E}[M_s(A)]} \sqrt{\frac{d \log_+(n/d) + n\xi}{n}} + \frac{d \log_+(n/d) + n\xi}{n}.$$

Let now  $\bar{\mathcal{A}} = \{A \in \bar{\mathcal{I}}_d, \mathbb{E}[M_s(A)] \leq d/(2n)\}$ . We apply Corollary 6.9 of [Mas07] with  $\mathcal{F} = \{f_A, -f_A, A \in \bar{\mathcal{A}}\}$ ,  $b = 1$ ,  $\sigma^2 = \min\{d/n, 2\}$ . We deduce that there exists an event  $\Omega'$  such that  $\mathbb{P}(\Omega') \geq 1 - (1/2)e^{-n\xi}$  and on which: for all  $A \in \bar{\mathcal{A}}$ ,

$$|M_s(A) - \mathbb{E}[M_s(A)]| \preceq \sigma \sqrt{\frac{d \log_+(1/\sigma^2) + n\xi + \log 2}{n}} + \frac{d \log_+(1/\sigma^2) + n\xi + \log 2}{n}.$$

Since  $\sigma \leq \sqrt{d/n} \leq \sqrt{(d \log_+(n/d) + n\xi)/n}$ ,

$$|M_s(A) - \mathbb{E}[M_s(A)]| \preceq \frac{d \log_+(n/d) + n\xi}{n}.$$

We deduce from these computations that (35) holds true on the event  $\Omega' \cap (\cap_{j \in J} \Omega_j)$  for all  $A \in \cup_{j \in J} \bar{\mathcal{A}}_j \cup \bar{\mathcal{A}}$  with  $\alpha, \beta$  of the form  $c_2(d \log_+(n/d) + n\xi)/n$ . Now,  $\bar{\mathcal{I}}_d = \cup_{j \in J} \bar{\mathcal{A}}_j \cup \bar{\mathcal{A}}$ , and

$$\begin{aligned} \mathbb{P} \left[ \left( \Omega' \cap \left( \bigcap_{j \in J} \Omega_j \right) \right)^c \right] &\leq \mathbb{P}[\Omega'^c] + \sum_{j \in J} \mathbb{P}[\Omega_j^c] \\ &\leq \frac{e^{-n\xi}}{2} + \sum_{j=1}^{\infty} \frac{e^{-n\xi}}{j^2 e} \\ &\leq e^{-n\xi}. \end{aligned}$$

□

6.6.4. *Proof of Lemma 8 in framework 2.* Without loss of generality, we suppose from now on that  $A \subset [0, +\infty)$ . We define for  $k \geq 1$ ,

$$J_k = \int_{\substack{u_1, \dots, u_k \in A \\ u_1 < u_2 < \dots < u_k}} \left( \prod_{j=1}^k s(u_j) \right) \mathbb{P}(X \geq u_k) du_1 du_2 \dots du_k.$$

We have,

$$\begin{aligned}
\mathbb{E} \left[ (f_A(Z_i))^k \right] &= \mathbb{E} \left[ \left( \int_A s(u) \mathbb{1}_{X \geq u} du \right)^k \right] \\
&= \mathbb{E} \left[ \int_{A^k} \prod_{j=1}^k s(u_j) \mathbb{1}_{X \geq u_j} du_1 du_2 \dots du_k \right] \\
&= \int_{A^k} \left( \prod_{j=1}^k s(u_j) \right) \mathbb{P}(X \geq \max\{u_1, \dots, u_k\}) du_1 du_2 \dots du_k \\
&= k! J_k.
\end{aligned}$$

Now,

$$J_k \leq \int_{\substack{u_1, \dots, u_{k-1} \in A \\ u_1 < u_2 < \dots < u_{k-1}}} \left( \prod_{j=1}^{k-1} s(u_j) \right) \left( \int_{u_{k-1}}^{\infty} s(u_k) \mathbb{P}(X \geq u_k) du_k \right) du_1 du_2 \dots du_{k-1},$$

and

$$\begin{aligned}
\int_{u_{k-1}}^{\infty} s(u_k) \mathbb{P}(X \geq u_k) du_k &= \int_{u_{k-1}}^{\infty} f(u_k) \mathbb{P}(C \geq u_k) du_k \\
&\leq \left( \int_{u_{k-1}}^{\infty} f(u_k) du_k \right) \mathbb{P}(C \geq u_{k-1}) \\
&\leq \mathbb{P}(T \geq u_{k-1}) \mathbb{P}(C \geq u_{k-1}) \\
&\leq \mathbb{P}(X \geq u_{k-1}).
\end{aligned}$$

Therefore,

$$\begin{aligned}
J_k &\leq \int_{\substack{u_1, \dots, u_{k-1} \in A \\ u_1 < u_2 < \dots < u_{k-1}}} \left( \prod_{j=1}^{k-1} s(u_j) \right) \mathbb{P}(X \geq u_{k-1}) du_1 du_2 \dots du_{k-1} \\
&\leq J_{k-1}.
\end{aligned}$$

By induction,  $J_k \leq J_1 = \mathbb{E}[M_s(A)]$ . □

#### 6.6.5. Proof of Lemma 8 in framework 3.

**Claim 4.** Let  $t > 0$ ,  $\mathcal{F}_t = \sigma(X_v, v \leq t)$  be the  $\sigma$ -algebra generated by the family of random variables  $X_v$ ,  $v \in [0, t]$ . Let  $B$  be an event  $\mathcal{F}_t$ -measurable. Let  $\mu_B$  be the measure defined for all  $A \in \mathcal{B}(\mathbb{R})$  by

$$\mu_B(A) = \mathbb{P}(B \text{ and } T_{1,0} \in A).$$

Then, for  $\mu$ -almost all  $u > t$ ,

$$(44) \quad \frac{d\mu_B}{du}(u) = \mathbb{P}(B \text{ and } X_{u-} = 1) s(u).$$

*Proof.* First of all,  $\mu_B$  is absolutely continuous with respect to the Lebesgue measure  $\mu$  and admits therefore a Radon-Nikodym derivative. We now aim to show that this derivative is given by (44) for almost all  $u > t$ .

Let  $Z_h(u)$  be the random variable giving the number of jumps of the Markov process in  $[u - h, u + h]$ . Then,  $\mathbb{P}(Z_h(u) \geq 2) = o(h)$  when  $h \rightarrow 0$ . We deduce,

$$\mu_B([u, u + h]) = \mathbb{P}(B, Z_h(u) = 1, T_{1,0} \in [u, u + h]) + o(h).$$

When  $Z_h(u) = 1$ ,  $T_{1,0} \in [u, u + h]$  is equivalent to  $X_{u-} = 1$  and  $X_{u+h} = 0$ . This yields

$$\begin{aligned} \mu_B([u, u + h]) &= \mathbb{P}(B, Z_h(u) = 1, X_{u-} = 1, X_{u+h} = 0) + o(h) \\ &= \mathbb{P}(B, X_{u-} = 1, X_{u+h} = 0) + o(h) \\ &= \mathbb{P}(B, X_{u-} = 1) \mathbb{P}(X_{u+h} = 0 \mid B, X_{u-} = 1) + o(h). \end{aligned}$$

As  $B$  is  $\mathcal{F}_t$ -measurable and  $u > t$ ,

$$\begin{aligned} \mu_B([u, u + h]) &= \mathbb{P}(B, X_{u-} = 1) \mathbb{P}(X_{u+h} = 0 \mid X_{u-} = 1) + o(h) \\ (45) \quad &= \mathbb{P}(B, X_{u-} = 1) \frac{\mathbb{P}(X_{u-} = 1, X_{u+h} = 0)}{\mathbb{P}(X_{u-} = 1)} + o(h). \end{aligned}$$

Now,

$$\begin{aligned} \mathbb{P}(X_{u-} = 1, X_{u+h} = 0) &= \mathbb{P}(X_{u-} = 1, X_{u+h} = 0, Z_h(u) = 1) + o(h) \\ &= \mathbb{P}(T_{1,0} \in [u, u + h], Z_h(u) = 1) + o(h) \\ &= \mathbb{P}(T_{1,0} \in [u, u + h]) + o(h). \end{aligned}$$

Finally, by plugging this inequality into (45),

$$\begin{aligned} \mu_B([u, u + h]) &= \mathbb{P}(B, X_{u-} = 1) \frac{\mathbb{P}(T_{1,0} \in [u, u + h])}{\mathbb{P}(X_{u-} = 1)} + o(h) \\ &= \mathbb{P}(B, X_{u-} = 1) \frac{hf(u)}{\mathbb{P}(X_{u-} = 1)} + o(h) \\ &= h\mathbb{P}(B, X_{u-} = 1)s(u) + o(h), \end{aligned}$$

which proves (44).  $\square$

We now return to the proof of Lemma 8. Without loss of generality, we suppose that  $A \subset [0, +\infty)$ . Define for  $k \geq 1$ ,

$$J_k = \int_{\substack{u_1, \dots, u_k \in A \\ u_1 < u_2 < \dots < u_k}} \left( \prod_{j=1}^k s(u_j) \right) \mathbb{P}(X_{u_1-} = 1, \dots, X_{u_k-} = 1) du_1 du_2 \dots du_k.$$

We have,

$$\begin{aligned} \mathbb{E} \left[ (f_A(Z_i))^k \right] &= \mathbb{E} \left[ \left( \int_A s(u) \mathbb{1}_{X_{u-} = 1} du \right)^k \right] \\ &= \mathbb{E} \left[ \int_{A^k} \prod_{j=1}^k s(u_j) \mathbb{1}_{X_{u_j-} = 1} du_1 du_2 \dots du_k \right] \\ &= \int_{A^k} \left( \prod_{j=1}^k s(u_j) \right) \mathbb{P}(X_{u_1-} = 1, \dots, X_{u_k-} = 1) du_1 du_2 \dots du_k \\ &= k! J_k. \end{aligned}$$

Yet,

$$J_k \leq \int_{\substack{u_1, \dots, u_{k-1} \in A \\ u_1 < u_2 < \dots < u_{k-1}}} \left( \prod_{j=1}^{k-1} s(u_j) \right) \left( \int_{u_{k-1}}^{\infty} s(u_k) \mathbb{P}(X_{u_1-} = 1, \dots, X_{u_k-} = 1) du_k \right) du_1 du_2 \dots du_{k-1}.$$

Let  $B = [X_{u_1-} = 1, \dots, X_{u_{k-1}-} = 1] \in \mathcal{F}_{u_{k-1}}$ . Then,

$$\begin{aligned} \int_{u_{k-1}}^{\infty} s(u_k) \mathbb{P}(X_{u_1-} = 1, \dots, X_{u_k-} = 1) du_k &= \int_{u_{k-1}}^{\infty} \frac{d\mu_B}{du_k}(u_k) du_k \\ &= \mu_B([u_{k-1}, +\infty)) \\ &= \mathbb{P}(X_{u_1-} = 1, \dots, X_{u_{k-1}-} = 1 \text{ and } T_{1,0} \geq u_{k-1}) \\ &\leq \mathbb{P}(X_{u_1-} = 1, \dots, X_{u_{k-1}-} = 1). \end{aligned}$$

Therefore,  $J_k \leq J_{k-1}$  and by induction  $J_k \leq J_1 = \mathbb{E}[M_s(A)]$ .  $\square$

6.6.6. *Proof of Lemma 9.* First of all, we only need to prove the lemma when  $\delta$  is smaller than 1, what we will do in the sequel.

We endow  $\overline{\mathcal{I}}_d$  with the distance  $dist$  defined for  $A_1, A_2 \in \overline{\mathcal{I}}_d$  by

$$dist(A_1, A_2) = \mathbb{E}[M_s(A_1 \Delta A_2)] \quad \text{where } A_1 \Delta A_2 = (A_1 \setminus A_2) \cup (A_2 \setminus A_1).$$

We may write  $dist(A_1, A_2)$  as

$$dist(A_1, A_2) = \int_{\mathbb{R}} |\mathbb{1}_{A_1}(t) - \mathbb{1}_{A_2}(t)| f(t) dt,$$

where  $f(t) = s(t)\mathbb{E}[V_1(t)]\mathbb{1}_{t \geq 0}$  is a non-negative function satisfying  $\int_{\mathbb{R}} f(t) dt \leq 1$ .

We introduce the real valued function  $F$  defined by

$$F(x) = \int_{-\infty}^x f(t) dt \quad \text{for all } x \in \mathbb{R}.$$

Since  $F$  is a continuous non-decreasing function such that  $F(\mathbb{R}) \subset [0, 1]$ , there exist an even integer  $\ell \in [2, 4d/\delta^2 + 2]$ , and  $\ell$  numbers  $(x_1, x_2, \dots, x_{\ell-1}, x_{\ell}) \in \{-\infty\} \times \mathbb{Q}^{\ell-2} \times \{+\infty\}$  such that

$$\max_{1 \leq i \leq \ell-1} \{F(x_{i+1}) - F(x_i)\} \leq \delta^2/(4d).$$

We may suppose that  $\ell \geq d$ . Let  $\mathcal{X} = \{x_1, x_2, \dots, x_{\ell}\}$ , and  $\overline{\mathcal{I}}_{dis}$  be the collection of unions of at most  $d$  closed intervals whose endpoints belong to  $\mathcal{X}$ .

When  $k \leq \ell/2$ , choosing  $k$  disjoint closed intervals whose endpoints belong to  $\mathcal{X}$  amounts to choosing  $2k$  numbers among  $\mathcal{X}$ . When  $k > \ell/2$ , we cannot find  $k$  disjoint closed intervals with endpoints in  $\mathcal{X}$ . The cardinality of  $\overline{\mathcal{I}}_{dis}$  is therefore bounded by

$$|\overline{\mathcal{I}}_{dis}| \leq \sum_{k=0}^d C_{\ell}^{2k}.$$

Standard arguments (see, for instance, exercise 2.14 of [BLM13]) show that  $|\overline{\mathcal{I}}_{dis}| \leq (\ell e/d)^d$ . Using now that  $\ell \leq 4d/\delta^2 + 2$ , we derive that

$$\log |\overline{\mathcal{I}}_{dis}| \leq cd \log_+(1/\delta^2)$$

for a suitable universal constant  $c$ .



For each set  $A \in \overline{\mathcal{I}}_d$ , we now show that there exist  $A_1, A_2 \in \overline{\mathcal{I}}_{dis}$  such that  $f_{A_1} \leq f_A \leq f_{A_2}$  and  $dist(A_1, A_2) \leq \delta^2/2$ . Let  $A \in \overline{\mathcal{I}}_d$  be written as  $A = \bigcup_{k=1}^d A_k$  where  $A_k$  is an interval whose endpoints are  $a_k \leq b_k$ . For each  $k \in \{1, \dots, d\}$ , there exist  $a_k^{(1)} \leq a_k^{(2)} \leq b_k^{(1)} \leq b_k^{(2)} \in \mathcal{X}$  such that

$$a_k^{(1)} \leq a_k \leq a_k^{(2)}, \quad b_k^{(1)} \leq b_k \leq b_k^{(2)},$$

and

$$F(a_k^{(2)}) - F(a_k^{(1)}) \leq \delta^2/(4d), \quad F(b_k^{(2)}) - F(b_k^{(1)}) \leq \delta^2/(4d).$$

Define the closed intervals

$$A_k^{(1)} = \left\{ x \in \mathbb{R}, a_k^{(2)} \leq x \leq b_k^{(1)} \right\}, \quad A_k^{(2)} = \left\{ x \in \mathbb{R}, a_k^{(1)} \leq x \leq b_k^{(2)} \right\}.$$

Then,  $A_1 = \bigcup_{k=1}^d A_k^{(1)}$  and  $A_2 = \bigcup_{k=1}^d A_k^{(2)}$  belong to  $\overline{\mathcal{I}}_{dis}$  and satisfy  $f_{A_1} \leq f_A \leq f_{A_2}$ . Moreover,

$$A_2 \Delta A_1 \subset \bigcup_{k=1}^d [a_k^{(1)}, a_k^{(2)}) \cup (b_k^{(1)}, b_k^{(2)}],$$

and hence,

$$\begin{aligned} dist(A_1, A_2) &\leq \sum_{k=1}^d \int_{[a_k^{(1)}, a_k^{(2)}) \cup (b_k^{(1)}, b_k^{(2)}]} f(t) dt \\ &\leq \sum_{k=1}^d \left( F(a_k^{(2)}) - F(a_k^{(1)}) + F(b_k^{(2)}) - F(b_k^{(1)}) \right) \\ &\leq \sum_{k=1}^d (\delta^2/(4d) + \delta^2/(4d)) \\ &\leq \delta^2/2. \end{aligned}$$

Now,

$$\begin{aligned} \mathbb{E} \left[ (f_{A_2}(Z_1) - f_{A_1}(Z_1))^k \right] &= \mathbb{E} \left[ (f_{A_2 \setminus A_1}(Z_1))^k \right] \\ &\leq k! \mathbb{E} [M_s(A_2 \setminus A_1)] \quad \text{thanks to Lemma 8} \\ &\leq k! dist(A_1, A_2) \\ &\leq k! \delta^2/2, \end{aligned}$$

which completes the proof with  $\mathcal{C}_\delta = \{f_A, A \in \overline{\mathcal{I}}_{dis}\}$ .  $\square$

**6.7. Proof of Proposition 4.** The proof follows closely the one of Theorem 3. Suppose without loss of generality that the functions  $\varphi$  are non-negative. Consider  $\varepsilon > 0$  and  $\eta \in (0, 1)$ . According to (40), for all  $\varphi \in \mathcal{F}$ , and  $t \in (0, 1)$ , there exists  $A_t \in \mathcal{A}_t$  (we omit the subscript  $\varphi$ ) such that

$$\begin{aligned} |Z(\varphi)| &\leq \varepsilon \sigma^2 + \int_0^\eta |N(A_t) - \mathbb{E}[N(A_t)]| dt \\ &\quad + \int_\eta^1 \{|N(A_t) - \mathbb{E}[N(A_t)]| - 2\varepsilon t \mathbb{E}[N(A_t)]\} dt. \end{aligned}$$

Therefore,

$$(46) \quad \mathbb{E} \left[ \sup_{\varphi \in \mathcal{F}} |Z(\varphi)| \right] \leq \varepsilon \sigma^2 + \int_0^\eta \mathbb{E} \left[ \sup_{A_t \in \mathcal{A}_t} |N(A_t) - \mathbb{E}[N(A_t)]| \right] dt \\ + \int_\eta^1 \mathbb{E} \left[ \sup_{A_t \in \mathcal{A}_t} \{ |N(A_t) - \mathbb{E}[N(A_t)]| - 2\varepsilon t \mathbb{E}[N(A_t)] \} \right] dt.$$

Let now  $\xi > 0$ . As (34) holds true for all  $A \in \mathcal{A}_t$ , on an event  $\Omega_{\xi,t}$  such that  $\mathbb{P}[\Omega_{\xi,t}] \geq 1 - e^{-n\xi}$ , with  $\alpha = c[\log_+ |S_{\mathcal{A}_t}(2n)| + n\xi]$ ,  $\beta = 0$ , we deduce from Claim 3 that on this event: for all  $A_t \in \mathcal{A}_t$ .

$$|N(A_t) - \mathbb{E}[N(A_t)]| - 2\varepsilon t \mathbb{E}[N(A_t)] \leq C \left\{ \frac{\log_+ |S_{\mathcal{A}_t}(2n)| + n\xi}{n} + \frac{\log_+ |S_{\mathcal{A}_t}(2n)| + n\xi}{n\varepsilon t} \right\}.$$

and

$$|N(A_t) - \mathbb{E}[N(A_t)]| \leq C \left[ \sqrt{\frac{\log_+ |S_{\mathcal{A}_t}(2n)| + n\xi}{n}} \sqrt{\mathbb{E}[N(A_t)]} + \frac{\log_+ |S_{\mathcal{A}_t}(2n)| + n\xi}{n} \right]$$

In these two inequalities,  $C$  denotes a universal constant.

We integrate these inequalities with respect to  $\xi$  to get

$$\mathbb{E} \left[ \sup_{A_t \in \mathcal{A}_t} \{ |N(A_t) - \mathbb{E}[N(A_t)]| - 2\varepsilon t \mathbb{E}[N(A_t)] \} \right] \leq C' \left[ \frac{\log_+ |S_{\mathcal{A}_t}(2n)|}{n} + \frac{\log_+ |S_{\mathcal{A}_t}(2n)|}{n\varepsilon t} \right]$$

and

$$\mathbb{E} \left[ \sup_{A_t \in \mathcal{A}_t} |N(A_t) - \mathbb{E}[N(A_t)]| \right] \leq C' \left[ \sqrt{\frac{\log_+ |S_{\mathcal{A}_t}(2n)|}{n}} \sup_{A_t \in \mathcal{A}_t} \sqrt{\mathbb{E}[N(A_t)]} + \frac{\log_+ |S_{\mathcal{A}_t}(2n)|}{n} \right],$$

where  $C'$  is universal. We now plug the two latter inequalities into (46)

$$\mathbb{E} \left[ \sup_{\varphi \in \mathcal{F}} |Z(\varphi)| \right] \leq \varepsilon \sigma^2 + \frac{C'}{n\varepsilon} \int_\eta^1 \frac{\log_+ |S_{\mathcal{A}_t}(2n)|}{t} dt \\ + \frac{C'}{\sqrt{n}} \int_0^\eta \sqrt{\left( \sup_{A_t \in \mathcal{A}_t} \mathbb{E}[N(A_t)] \right) (\log_+ |S_{\mathcal{A}_t}(2n)|)} dt + \frac{C'}{n} \int_0^1 \log_+ |S_{\mathcal{A}_t}(2n)| dt.$$

As  $\varepsilon > 0$  and  $\eta \in (0, 1)$  are arbitrary,

$$\mathbb{E} \left[ \sup_{\varphi \in \mathcal{F}} |Z(\varphi)| \right] \leq \frac{C''}{\sqrt{n}} \inf_{\eta \in (0,1)} \left\{ \sigma \sqrt{\int_\eta^1 \frac{\log_+ |S_{\mathcal{A}_t}(2n)|}{t} dt} + \int_0^\eta \sqrt{\left( \sup_{A_t \in \mathcal{A}_t} \mathbb{E}[N(A_t)] \right) (\log_+ |S_{\mathcal{A}_t}(2n)|)} dt \right\} \\ + \frac{C''}{n} \int_0^1 \log_+ |S_{\mathcal{A}_t}(2n)| dt,$$

where  $C''$  is a universal constant. □

**6.8. Proof of Lemma 2.** The proof of this lemma follows from some computations as in Section 8.4 of [Bar11] (see also Proposition 3 of [BB17]). Let  $\sqrt{q} = (\sqrt{f} + \sqrt{g})/2$  and

$$\mathcal{X} = \{x \in \mathbb{R}, g(x) \neq 0 \text{ or } f(x) \neq 0\}.$$

Then,

$$\begin{aligned} \int_{\mathbb{R}} \psi^2 \left( \frac{g}{f} \right) s \, dM &= \int_{\mathcal{X}} \psi^2 \left( \frac{g}{f} \right) s \, dM \\ &= \frac{1}{4} \int_{\mathcal{X}} (\sqrt{g} - \sqrt{f})^2 \frac{s}{q} \, dM \\ &= \frac{1}{4} \int_{\mathcal{X}} (\sqrt{g} - \sqrt{f})^2 \left( \sqrt{\frac{s}{q}} - 1 + 1 \right)^2 \, dM \\ &\leq \frac{1}{2} \int_{\mathcal{X}} (\sqrt{g} - \sqrt{f})^2 \left( \sqrt{\frac{s}{q}} - 1 \right)^2 \, dM + \frac{1}{2} \int_{\mathcal{X}} (\sqrt{g} - \sqrt{f})^2 \, dM \\ &\leq \frac{1}{2} \int_{\mathcal{X}} \frac{(\sqrt{g} - \sqrt{f})^2}{q} (\sqrt{s} - \sqrt{q})^2 \, dM + h^2(f, g) \\ &\leq 2 \int_{\mathcal{X}} (\sqrt{s} - \sqrt{q})^2 \, dM + h^2(f, g) \\ &\leq \frac{1}{2} \int_{\mathcal{X}} \left( (\sqrt{s} - \sqrt{f}) + (\sqrt{s} - \sqrt{g}) \right)^2 \, dM + h^2(f, g) \\ &\leq \int_{\mathcal{X}} (\sqrt{s} - \sqrt{f})^2 \, dM + \int_{\mathcal{X}} (\sqrt{s} - \sqrt{g})^2 \, dM + h^2(f, g) \\ &\leq 2h^2(s, f) + 2h^2(s, g) + h^2(f, g). \end{aligned}$$

We complete the proof by using  $h^2(f, g) \leq 2h^2(s, f) + 2h^2(s, g)$ .  $\square$

**6.9. Proof of Proposition 5 for  $S = \mathcal{P}_{\ell, r}$ .** Let  $f, g \in \mathcal{P}_{\ell, r}$ . There exist two partitions  $m_1, m_2$  of  $\mathbb{R}$  into intervals such that  $|m_1| \leq 2\ell + 1$  and  $|m_2| \leq 2\ell + 1$  and such that  $f$  (respectively  $g$ ) is polynomial on each element  $K_1 \in m_1$  (respectively  $K_2 \in m_2$ ). Let

$$m = \{K_1 \cap K_2, (K_1, K_2) \in m_1 \times m_2, K_1 \cap K_2 \neq \emptyset\}.$$

Then,  $m$  is a partition of  $\mathbb{R}$  into intervals such that  $|m| \leq |m_1| + |m_2| \leq 4\ell + 2$ . Moreover, we may write  $f$  and  $g$  as

$$f = \sum_{K \in m} P_K \mathbb{1}_K \quad \text{and} \quad g = \sum_{K \in m} Q_K \mathbb{1}_K,$$

where  $P_K$  and  $Q_K$  are non-negative polynomial functions on  $K$  of degree at most  $r$ . Let  $R_K = P_K - tQ_K$ . Now,

$$\{x \in \mathbb{R}, g(x) > tf(x)\} = \bigcup_{K \in m} \{x \in K, R_K(x) > 0\}.$$

Let  $\mathcal{X}$  be the set gathering the zeros of  $R_K$ . If  $\mathcal{X} = \emptyset$ , then  $R_K$  is either positive, or negative on  $\mathbb{R}$  and the set  $\{x \in K, R_K(x) > 0\}$  is either empty or the interval  $K$ . If  $\mathcal{X} = \mathbb{R}$ , then  $R_K = 0$  and  $\{x \in K, R_K(x) > 0\} = \emptyset$ . Suppose now that  $\mathcal{X} \neq \emptyset$  and  $\mathcal{X} \neq \mathbb{R}$ . We may write  $\mathcal{X} = \{b_1, \dots, b_k\}$  with  $b_1 < b_2 < \dots < b_k$  and  $k \leq r$ . We set  $b_0 = -\infty$  and  $b_{k+1} = +\infty$ . For all  $j \in \{0, \dots, k\}$ ,  $R_K$  is either

positive or negative on  $(b_j, b_{j+1})$ , and its sign changes with  $j$ . Therefore, the set  $\{x \in K, R_K(x) > 0\}$  is a union of at most  $k/2 + 1$  intervals.

Finally, for all  $K \in m$ ,  $\{x \in K, R_K(x) > 0\}$  is always a union of at most  $r/2 + 1$  intervals, which implies that  $\{x \in \mathbb{R}, g(x) > tf(x)\}$  is a union of at most  $(r/2 + 1)(4\ell + 2)$  intervals.  $\square$

**6.10. Proof of Theorem 7.** Let for  $d \geq 1$ ,

$$\vartheta(d) = \frac{d}{n} \log_+^2 \left( \frac{n}{d} \right).$$

We need to prove that there exists an event  $\Omega_\xi$  such that  $\mathbb{P}(\Omega_\xi) \geq 1 - e^{-n\xi}$  and on which any  $\rho$ -estimator  $\hat{s}$  on  $S$  satisfies

$$(47) \quad h^2(s, \hat{s}) \leq \inf_{f \in \bar{S}} \{c_1 h^2(s, f) + c_2 \vartheta(d_S(f)) + c_3 \xi \log_+(1/\xi)\}.$$

We introduce the following notations. Let for  $d \geq 1$ ,  $\mathcal{I}_d$  be the class of unions of at most  $d$  intervals. Let  $f, g \in \mathcal{S}$ . Suppose that there exists  $d \geq 1$  such that for all  $t > 0$ , the set  $\{x \in \mathbb{R}, g(x) > tf(x)\}$  belongs to  $\mathcal{I}_d$ . Then,  $d_g(f)$  stands for any number  $d$  such that

$$\{x \in \mathbb{R}, g(x) > tf(x)\}$$

belongs to  $\mathcal{I}_d$  (for all  $t > 0$ ). If the preceding assumption does not hold, we set  $d_g(f) = +\infty$ .

We define for  $d \geq 2$ ,

$$\mathcal{G}_d = \{\psi(g/f), g \in \mathcal{S}, f \in \mathcal{S}, d_g(f) = d - 1\}.$$

We will apply Theorem 3 to the class  $\mathcal{F} = \mathcal{G}_d$ . In order to verify its assumptions, we begin with the following elementary claim:

**Claim 5.** *We have,*

- For all  $J \in \mathcal{I}_d$ ,  $\mathbb{R} \setminus J \in \mathcal{I}_{d+1}$ .
- For all non-increasing sequence  $(J_n)_{n \geq 1}$  of  $\mathcal{I}_d$ ,  $\bigcap_{n \geq 1} J_n$  belongs to  $\mathcal{I}_d$ .

The set  $\mathcal{G}_d$  enjoys the following properties:

**Claim 6.** *The functions  $\varphi \in \mathcal{G}_d$  satisfy  $|\varphi| \leq 1$ . Moreover, for all  $t \in (0, 1)$ ,  $\varphi \in \mathcal{G}_d$ ,  $\{x \in \mathbb{R}, \varphi_+(x) > t\} \in \mathcal{I}_{d-1}$  and  $\{x \in \mathbb{R}, \varphi_-(x) > t\} \in \mathcal{I}_d$ .*

*Proof.* Let  $\varphi \in \mathcal{G}_d$  written as  $\varphi = \psi(g/f)$ . Then,

$$\begin{aligned} \{x \in \mathbb{R}, \varphi_+(x) > t\} &= \{x \in \mathbb{R}, \psi_+(g(x)/f(x)) > t\} \\ &= \{x \in \mathbb{R}, f(x) \neq 0, \psi_+(g(x)/f(x)) > t\} \cup \{x \in \mathbb{R}, f(x) = 0, g(x) > 0\} \\ &= \{x \in \mathbb{R}, f(x) \neq 0, g(x) > uf(x)\} \cup \{x \in \mathbb{R}, f(x) = 0, g(x) > 0\}, \end{aligned}$$

where  $u = \psi^{-1}(t)$ . Therefore,

$$\{x \in \mathbb{R}, \varphi_+(x) > t\} = \{x \in \mathbb{R}, g(x) > uf(x)\} \in \mathcal{I}_{d-1},$$

as  $d_g(f) = d - 1$ . Now, note that  $\psi_-(x) = \psi_+(1/x)$ . Hence,

$$\{x \in \mathbb{R}, \varphi_-(x) > t\} = \{x \in \mathbb{R}, \psi_+(f(x)/g(x)) > t\}.$$

By exchanging the role of  $f$  and  $g$  in the above computations, we derive

$$\begin{aligned} \{x \in \mathbb{R}, \varphi_-(x) > t\} &= \{x \in \mathbb{R}, f(x) > ug(x)\} \\ &= \{x \in \mathbb{R}, g(x) < (1/u)f(x)\}. \end{aligned}$$

By using the first point of Claim 5, for all  $n \geq 1$ ,

$$\{x \in \mathbb{R}, g(x) \leq (1/u + 1/n)f(x)\} \in \mathcal{I}_d.$$

Yet,

$$\{x \in \mathbb{R}, g(x) < (1/u)f(x)\} = \bigcap_{n=1}^{\infty} \{x \in \mathbb{R}, g(x) \leq (1/u + 1/n)f(x)\}.$$

The second point of Claim 5 ensures that  $\{x \in \mathbb{R}, g(x) < (1/u)f(x)\}$  belongs to  $\mathcal{I}_d$ , which completes the proof.  $\square$

The lemma below is at the core of the proof of Theorem 7.

**Lemma 10.** *For all  $\xi > 0$ , there exists an event  $\Omega_\xi$  such that  $\mathbb{P}(\Omega_\xi) \geq 1 - e^{-n\xi}$  and on which: for all  $\varepsilon \in (0, 1/12)$ ,  $f, g \in \mathcal{S}$ ,*

$$(48) \quad T(f, g) \leq (3 + 4\varepsilon)h^2(s, f) - \frac{4 - 3\varepsilon}{12}h^2(s, g) + c_1\vartheta(d_g(f)) + c_2\xi \log_+(1/\xi).$$

*In the above inequality,  $c_1$  and  $c_2$  only depend on  $\varepsilon$ . Besides, we use the convention  $\vartheta(+\infty) = +\infty$  when  $d_g(f) = \infty$ .*

*Proof.* Let  $d \geq 2$ . Theorem 3 shows the existence of an event  $\Omega_\xi(d)$  such that  $\mathbb{P}[\Omega_\xi(d)] \geq 1 - e^{-n\xi}$  and on which: for all  $\varepsilon > 0$ ,  $\varphi \in \mathcal{G}_d$  of the form  $\varphi = \psi(g/f)$ , with  $f, g \in \mathcal{S}$ ,

$$|Z(\varphi)| \leq \varepsilon v(\varphi) + c[\vartheta(d_g(f) + 1) + \xi \log_+(1/\xi)].$$

In this inequality,  $c$  only depends on  $\varepsilon$ . Since

$$\vartheta(d_g(f) + 1) \leq 2\vartheta(d_g(f)),$$

we get

$$|Z(\varphi)| \leq \varepsilon v(\varphi) + 2c\vartheta(d_g(f)) + c\xi \log_+(1/\xi).$$

Let  $\Omega_\xi = \bigcap_{d=2}^{\infty} \Omega_{\xi+(2\log(1+d))/n}(d)$ . Then,

$$\mathbb{P}[(\Omega_\xi)^c] \leq \sum_{d=2}^{\infty} \mathbb{P}[(\Omega_{\xi+(2\log(1+d))/n}(d))^c] \leq \sum_{d=2}^{\infty} \frac{e^{-n\xi}}{(1+d)^2} \leq e^{-n\xi}.$$

Moreover, on  $\Omega_\xi$ : for all  $f, g \in \mathcal{S}$ ,  $\varphi = \psi(g/f)$  such that  $d_g(f) < \infty$ ,

$$\begin{aligned} |Z(\varphi)| &\leq \varepsilon v(\varphi) + 2c\vartheta(d_g(f)) + c \left[ \left( \xi + \frac{2\log(1+d_g(f))}{n} \right) \log_+ \left( \frac{1}{\xi + \frac{2\log(1+d_g(f))}{n}} \right) \right] \\ &\leq \varepsilon v(\varphi) + 2c\vartheta(d_g(f)) + \frac{2c\log(1+d_g(f))}{n} \log_+ \left( \frac{n}{2\log(1+d_g(f))} \right) + c\xi \log_+(1/\xi) \\ (49) \quad &\leq \varepsilon v(\varphi) + c'\vartheta(d_g(f)) + c\xi \log_+(1/\xi), \end{aligned}$$

where  $c'$  only depends on  $\varepsilon$ . This last inequality remains true when  $d_g(f) = \infty$  using the convention  $\vartheta(+\infty) = +\infty$ .

Now, it follows from (3) that for all  $f, g \in S$ ,

$$(50) \quad T(f, g) \leq 3h^2(s, f) - \frac{1}{3}h^2(s, g) + Z(\psi(g/f)).$$

Therefore, we deduce from Lemma 2 and from (49) that on  $\Omega_\xi$ : for all  $f, g \in S$ ,

$$T(f, g) \leq (3 + 4\varepsilon)h^2(s, f) - \frac{4 - 3\varepsilon}{12}h^2(s, g) + c'\vartheta_1(d_g(f)) + c\xi \log_+(1/\xi),$$

which proves (48) with  $c_1 = c'$  and  $c_2 = c$ .  $\square$

We now finish the proof of Theorem 7. Lemma 10 implies that on  $\Omega_\xi$ : for all  $f, g \in S$ ,

$$(51) \quad T(f, g) \leq (3 + 4\varepsilon)h^2(s, f) - \frac{4 - 3\varepsilon}{12}h^2(s, g) + c_1\vartheta(d_g(f)) + c_2\xi \log_+(1/\xi).$$

Thus, for all  $f \in S$ ,

$$(52) \quad \gamma(f) \leq (3 + 4\varepsilon)h^2(s, f) - \frac{4 - 3\varepsilon}{12}h^2(s, S) + c_1 \sup_{g \in S} \vartheta(d_g(f)) + c_2\xi \log_+(1/\xi).$$

By using  $T(f, g) = -T(g, f)$ , we deduce from (51) that for all  $f, g \in S$ ,

$$\frac{4 - 3\varepsilon}{12}h^2(s, g) - (3 + 4\varepsilon)h^2(s, f) - c_1\vartheta(d_g(f)) - c_2\xi \log_+(1/\xi) \leq T(g, f).$$

Any  $\rho$ -estimator  $\hat{s}$  satisfies on  $\Omega_\xi$ : for all  $f \in S$ ,

$$(53) \quad \begin{aligned} \frac{4 - 3\varepsilon}{12}h^2(s, \hat{s}) - (3 + 4\varepsilon)h^2(s, f) - c_1 \sup_{g \in S} \vartheta(d_g(f)) - c_2\xi \log_+(1/\xi) &\leq T(\hat{s}, f) \\ &\leq \gamma(\hat{s}) \\ &\leq \gamma(f) + 1/n. \end{aligned}$$

Using now (52) and  $1/n \leq \vartheta(d_g(f))$ , we deduce when  $\varepsilon \in (0, 1/12)$ ,

$$(54) \quad h^2(s, \hat{s}) \leq \inf_{f \in S} \left\{ c_{1,\varepsilon}h^2(s, f) - h^2(s, S) + c_{2,\varepsilon} \sup_{g \in S} \vartheta(d_g(f)) + c_{2,\varepsilon}\xi \log_+(1/\xi) \right\},$$

with  $c_{1,\varepsilon} = 24(3 + 4\varepsilon)/(4 - 3\varepsilon)$ , and with  $c_{2,\varepsilon}$  depending only on  $\varepsilon$ .

When  $f \in \bar{S}$  and  $g \in S$ , Assumption 1 says that  $d_g(f)$  may be defined by  $d_g(f) = d_S(f)$ . Therefore, (54) becomes

$$h^2(s, \hat{s}) \leq \inf_{f \in \bar{S}} \left\{ c_{1,\varepsilon}h^2(s, f) - h^2(s, S) + c_{2,\varepsilon}\vartheta(d_S(f)) + c_{2,\varepsilon}\xi \log_+(1/\xi) \right\},$$

and it remains to choose  $\varepsilon$  arbitrarily in  $(0, 1/12)$  to prove the theorem.  $\square$

**6.11. Proof of Proposition 8.** We keep the notations introduced in the proof of Theorem 7. It follows from (53) that on  $\Omega_\xi$ : for all  $f \in S$ ,

$$\frac{4 - 3\varepsilon}{12}h^2(s, \hat{s}) - (3 + 4\varepsilon)h^2(s, f) - c_1 \sup_{g \in S} \vartheta(d_g(f)) - c_2\xi \log_+(1/\xi) \leq T(\hat{s}, f),$$

where  $c_1$  and  $c_2$  only depends on  $\varepsilon$ . Now  $\gamma(\hat{s}) = 0$  and hence  $T(\hat{s}, f) \leq 0$ . This leads to

$$h^2(s, \hat{s}) \leq \inf_{f \in S} \left\{ c_{1,\varepsilon}h^2(s, f) + c_{2,\varepsilon} \sup_{g \in S} \vartheta(d_g(f)) + c_{2,\varepsilon}\xi \log_+(1/\xi) \right\} \quad \text{with } c_{1,\varepsilon} = 12 \frac{3 + 4\varepsilon}{4 - 3\varepsilon}.$$

It then remains to use Assumption 1 to define  $d_g(f)$  by  $d_g(f) = d_S(f)$  when  $f \in \bar{S}$ ,  $g \in S$ .  $\square$

**6.12. Proofs of Theorems 10 and 11.** It is convenient for ease of demonstration to encompass the two procedures into a more general selection rule we now describe. Theorems 10 and 11 follow from Theorem 13 below. Their proofs are given in Sections 6.12.2 and 6.12.3.

We consider an arbitrary (possibly random) set  $\hat{\Lambda}$ . For each  $\lambda \in \hat{\Lambda}$ , we consider an estimator  $\hat{s}_\lambda$  with values in  $\mathcal{S}$ . Our aim is to select an estimator among the collection  $\{\hat{s}_\lambda, \lambda \in \hat{\Lambda}\}$ .

We consider for each  $\lambda \in \hat{\Lambda}$  a (possibly random) model  $\hat{S}_\lambda \subset \mathcal{S}$ . We associate to each  $\lambda \in \hat{\Lambda}$ ,  $\hat{g} \in \hat{S}_\lambda$ , two penalty terms  $\text{pen}_{1,\lambda}(\hat{g})$  and  $\text{pen}_2(\lambda)$ . We finally define the criterion  $\gamma_4$  by

$$\gamma_4(\hat{s}_\lambda) = \sup_{\hat{g} \in \hat{S}_\lambda} [T(\hat{s}_\lambda, \hat{g}) - \text{pen}_{1,\lambda}(\hat{g})].$$

The selected estimator  $\hat{s}_{\hat{\lambda}}$  is then any estimator among  $\{\hat{s}_\lambda, \lambda \in \hat{\Lambda}\}$  satisfying

$$\gamma_4(\hat{s}_{\hat{\lambda}}) + 2\text{pen}_2(\hat{\lambda}) \leq \inf_{\lambda \in \hat{\Lambda}} \{\gamma_4(\hat{s}_\lambda) + 2\text{pen}_2(\lambda)\} + 1/n.$$

The risk of this estimator is bounded above as follows.

**Theorem 13.** *Let for  $\xi > 0$ ,  $\Omega_\xi$  be the event defined by Lemma 10, which is such that  $\mathbb{P}(\Omega_\xi) \geq 1 - e^{-n\xi}$ . We assume that there exist two real valued maps,  $\Delta \geq 0$  on  $\hat{\Lambda}$ , and  $d \geq 1$  on  $\mathcal{S}$  such that*

$$(55) \quad d_{\hat{s}_\lambda}(\hat{g}) \leq d(\hat{g}) + \Delta(\lambda) \quad \text{for all } \lambda \in \hat{\Lambda}, \hat{g} \in \hat{S}_\lambda.$$

where  $d_{\hat{s}_\lambda}(\hat{g})$  is defined at the beginning of the proof of Theorem 7. We suppose that there exist a (possibly random) model  $\hat{S} \subset \bigcap_{\lambda \in \hat{\Lambda}} \hat{S}_\lambda$  and a map  $\text{pen}_1$  on  $\hat{S}$  such that

$$(56) \quad \text{pen}_{1,\lambda}(\hat{g}) \leq \text{pen}_1(\hat{g}) + \text{pen}_2(\lambda) \quad \text{for all } \hat{g} \in \hat{S}, \lambda \in \hat{\Lambda}.$$

There exists a universal constant  $L_1$  such that if for all  $\lambda \in \hat{\Lambda}$ ,  $\hat{g} \in \hat{S}_\lambda$ ,  $\hat{f} \in \hat{S}$ ,

$$(57) \quad \begin{aligned} \text{pen}_{1,\lambda}(\hat{g}) &\geq L_1 \frac{d(\hat{g})}{n} \log_+^2 \left( \frac{n}{d(\hat{g})} \right) \\ \text{pen}_1(\hat{f}) &\geq L_1 \frac{d(\hat{f})}{n} \log_+^2 \left( \frac{n}{d(\hat{f})} \right) \\ \text{pen}_2(\lambda) &\geq L_1 \frac{\Delta(\lambda)}{n} \log_+^2 \left( \frac{n}{\Delta(\lambda)} \right), \end{aligned}$$

then, for all  $\xi > 0$ , on  $\Omega_\xi$ :

$$h^2(s, \hat{s}_{\hat{\lambda}}) \leq c \left( \inf_{\lambda \in \hat{\Lambda}} \{h^2(s, \hat{s}_\lambda) + \text{pen}_2(\lambda)\} + \inf_{\hat{g} \in \hat{S}} \{h^2(s, \hat{g}) + \text{pen}_1(\hat{g})\} + \xi \log_+(1/\xi) \right).$$

In the above inequality,  $c$  is a universal constant and the convention  $0 \times \log_+^2(n/0) = 0$  is used when  $\Delta(\lambda) = 0$ .

6.12.1. *Proof of Theorem 13.* Let  $\varepsilon \in (0, 1/12)$ . We deduce from Lemma 10 that there exists an event  $\Omega_\xi$  such that  $\mathbb{P}(\Omega_\xi) \geq 1 - e^{-n\xi}$  and on which: for all  $f, g \in \mathcal{S}$  such that  $d_g(f) < +\infty$ ,

$$(58) \quad T(f, g) \leq (3 + 4\varepsilon)h^2(s, f) - \frac{4 - 3\varepsilon}{12}h^2(s, g) + c_1\vartheta(d_g(f)) + c_2\xi \log_+(1/\xi)$$

where  $c_1, c_2$  only depend on  $\varepsilon$ .

Let  $\lambda \in \hat{\Lambda}$  and  $\hat{g} \in \hat{S}_\lambda$ . It follows from Claim 5 and from (55) that we may set  $d_{\hat{g}}(\hat{s}_\lambda) = d_{\hat{s}_\lambda}(\hat{g}) + 1 = d(\hat{g}) + \Delta(\lambda) + 1$  (we may always increase  $d_{\hat{s}_\lambda}(\hat{g})$ ). We now use (58) with  $f = \hat{s}_\lambda, g = \hat{g}$ :

$$\begin{aligned} T(\hat{s}_\lambda, \hat{g}) &\leq (3 + 4\varepsilon)h^2(s, \hat{s}_\lambda) - \frac{4 - 3\varepsilon}{12}h^2(s, \hat{g}) + c_1\vartheta(d(\hat{g}) + \Delta(\lambda) + 1) + c_2\xi \log_+(1/\xi) \\ &\leq (3 + 4\varepsilon)h^2(s, \hat{s}_\lambda) + c_1\vartheta(d(\hat{g}) + \Delta(\lambda) + 1) + c_2\xi \log_+(1/\xi) \\ &\leq (3 + 4\varepsilon)h^2(s, \hat{s}_\lambda) + 2c_1\vartheta(d(\hat{g})) + c_1\vartheta(\Delta(\lambda)) + c_2\xi \log_+(1/\xi). \end{aligned}$$

If  $L_1$  is large enough,

$$T(\hat{s}_\lambda, \hat{g}) \leq (3 + 4\varepsilon)h^2(s, \hat{s}_\lambda) + \text{pen}_{1,\lambda}(\hat{g}) + \text{pen}_2(\lambda) + c_2\xi \log_+(1/\xi),$$

and hence

$$(59) \quad \gamma_4(\hat{s}_\lambda) \leq (3 + 4\varepsilon)h^2(s, \hat{s}_\lambda) + \text{pen}_2(\lambda) + c_2\xi \log_+(1/\xi).$$

We now derive from (58) that for all  $\hat{g} \in \hat{S}_{\hat{\lambda}}$ ,

$$T(\hat{g}, \hat{s}_{\hat{\lambda}}) \leq (3 + 4\varepsilon)h^2(s, \hat{g}) - \frac{4 - 3\varepsilon}{12}h^2(s, \hat{s}_{\hat{\lambda}}) + c_1\vartheta(d_{\hat{s}_{\hat{\lambda}}}(\hat{g})) + c_2\xi \log_+(1/\xi).$$

Using moreover that  $T(\hat{s}_{\hat{\lambda}}, \hat{g}) = -T(\hat{g}, \hat{s}_{\hat{\lambda}})$  we deduce,

$$\begin{aligned} \frac{4 - 3\varepsilon}{12}h^2(s, \hat{s}_{\hat{\lambda}}) &\leq T(\hat{s}_{\hat{\lambda}}, \hat{g}) + (3 + 4\varepsilon)h^2(s, \hat{g}) + c_1\vartheta(d_{\hat{s}_{\hat{\lambda}}}(\hat{g})) + c_2\xi \log_+(1/\xi) \\ &\leq T(\hat{s}_{\hat{\lambda}}, \hat{g}) + (3 + 4\varepsilon)h^2(s, \hat{g}) + c_1\vartheta(d(\hat{g})) + c_1\vartheta(\Delta(\hat{\lambda})) + c_2\xi \log_+(1/\xi). \end{aligned}$$

If  $L_1$  is large enough,

$$\begin{aligned} \frac{4 - 3\varepsilon}{12}h^2(s, \hat{s}_{\hat{\lambda}}) &\leq T(\hat{s}_{\hat{\lambda}}, \hat{g}) + (3 + 4\varepsilon)h^2(s, \hat{g}) + \frac{1}{2}\text{pen}_{1,\hat{\lambda}}(\hat{g}) + \frac{1}{2}\text{pen}_2(\hat{\lambda}) \\ &\quad + c_2\xi \log_+(1/\xi) - 1/n \\ &\leq \left[ T(\hat{s}_{\hat{\lambda}}, \hat{g}) - \text{pen}_{1,\hat{\lambda}}(\hat{g}) \right] + \frac{1}{2}\text{pen}_2(\hat{\lambda}) + \left[ (3 + 4\varepsilon)h^2(s, \hat{g}) + \frac{3}{2}\text{pen}_{1,\hat{\lambda}}(\hat{g}) \right] \\ &\quad + c_2\xi \log_+(1/\xi) - 1/n. \end{aligned}$$

Since this inequality is valid for all  $\hat{g} \in \hat{S}_{\hat{\lambda}}$  and  $\hat{S} \subset \hat{S}_{\hat{\lambda}}$ ,

$$\begin{aligned} \frac{4 - 3\varepsilon}{12}h^2(s, \hat{s}_{\hat{\lambda}}) &\leq \gamma_4(\hat{s}_{\hat{\lambda}}) + \frac{1}{2}\text{pen}_2(\hat{\lambda}) + \inf_{\hat{g} \in \hat{S}} \left\{ 3(1 + \varepsilon)h^2(s, \hat{g}) + \frac{3}{2}\text{pen}_{1,\hat{\lambda}}(\hat{g}) \right\} \\ &\quad + c_2\xi \log_+(1/\xi) - 1/n. \end{aligned}$$

We deduce from (56),

$$\begin{aligned} \frac{4 - 3\varepsilon}{12}h^2(s, \hat{s}_{\hat{\lambda}}) &\leq \gamma_4(\hat{s}_{\hat{\lambda}}) + 2\text{pen}_2(\hat{\lambda}) + \inf_{\hat{g} \in \hat{S}} \left\{ 3(1 + \varepsilon)h^2(s, \hat{g}) + \frac{3}{2}\text{pen}_1(\hat{g}) \right\} \\ &\quad + c_2\xi \log_+(1/\xi) - 1/n. \end{aligned}$$



By using the definition of  $\hat{\lambda}$  and (59), we get for all  $\lambda \in \hat{\Lambda}$ ,

$$\begin{aligned} \frac{4-3\varepsilon}{12}h^2(s, \hat{s}_\lambda) &\leq \gamma_4(\hat{s}_\lambda) + 2\text{pen}_2(\lambda) + \inf_{\hat{g} \in \hat{S}} \left\{ 3(1+\varepsilon)h^2(s, \hat{g}) + \frac{3}{2}\text{pen}_1(\hat{g}) \right\} + 2c_2\xi \log_+(1/\xi) \\ &\leq (3+4\varepsilon)h^2(s, \hat{s}_\lambda) + 3\text{pen}_2(\lambda) + \inf_{\hat{g} \in \hat{S}} \left\{ 3(1+\varepsilon)h^2(s, \hat{g}) + \frac{3}{2}\text{pen}_1(\hat{g}) \right\} \\ &\quad + 2c_2\xi \log_+(1/\xi). \end{aligned}$$

It remains to take the infimum over  $\lambda \in \hat{\Lambda}$  to finish the proof.  $\square$

6.12.2. *Proof of Theorem 10.* We will apply the selection rule developed in Section 6.12 to pick out an estimator among  $\{\hat{s}_\lambda, \lambda \in \hat{\Lambda}\} = \{\hat{s}_m, m \in \widehat{\mathcal{M}}_{\hat{\ell}}\}$ . For this purpose, we need to explain the values of the different parameters involved in the procedure. We set  $\hat{S} = \{\hat{s}_m, m \in \widehat{\mathcal{M}}_{\hat{\ell}}\}$ , and for  $m \in \widehat{\mathcal{M}}_{\hat{\ell}}$ ,

$$\hat{S}_m = \left\{ \sum_{K \in m} \hat{s}_{m_K} \mathbb{1}_K, m_K \in \widehat{\mathcal{M}}_{\hat{\ell}} \right\}.$$

Note that the assumption  $\hat{S} \subset \bigcap_{m \in \widehat{\mathcal{M}}_{\hat{\ell}}} \hat{S}_m$  of Theorem 13 is fulfilled. We define for  $m \in \widehat{\mathcal{M}}_{\hat{\ell}}$ ,  $K \in m$  and  $m_K \in \widehat{\mathcal{M}}_{\hat{\ell}}$ , the partition  $m_K \vee K$  of  $K$  by (23). A function  $\hat{g} \in \hat{S}_m$  of the form  $\hat{g} = \sum_{K \in m} \hat{s}_{m_K} \mathbb{1}_K$  is piecewise polynomial. In the sequel,  $m(\hat{g})$  designs a partition of  $\widehat{\mathcal{M}}$  of the form

$$m(\hat{g}) = \bigcup_{K \in m} m_K \vee K,$$

with minimal length that is such that

$$|m(\hat{g})| = \inf \left\{ \sum_{K \in m} |m_K \vee K|, \hat{g} = \sum_{K \in m} \hat{s}_{m_K} \mathbb{1}_K \right\}.$$

Let  $\bar{S} = \bigcup_{k=1}^{\infty} \mathcal{P}_{k,r}$  and note that  $\hat{S}_m \subset \bar{S}$  for all  $m \in \widehat{\mathcal{M}}_{\hat{\ell}}$ . Let  $f \in \bar{S}$  and  $k \geq 1$  be the smallest integer for which  $f \in \mathcal{P}_{k,r}$ . It follows from Proposition 5 that one may define  $d_{\mathcal{P}_{\hat{\ell} \vee k, r}}(f) = (r+2)(2(\hat{\ell} \vee k) + 1)$ . In particular, for all  $m \in \widehat{\mathcal{M}}_{\hat{\ell}}$  and  $f \in \bar{S}$ , we may set since  $\hat{s}_m \in \mathcal{P}_{\hat{\ell}, r}$ ,

$$d_{\hat{s}_m}(f) = (r+2)(2 \inf_{\substack{k \geq 1 \\ \mathcal{P}_{k,r} \ni f}} (\hat{\ell} \vee k) + 1).$$

We now define  $d$  for  $f \in \bar{S}$  and  $\Delta$  for  $m \in \widehat{\mathcal{M}}_{\hat{\ell}}$  by

$$d(f) = (r+2)(2 \inf_{\substack{k \geq 1 \\ \mathcal{P}_{k,r} \ni f}} k + 1), \quad \Delta(m) = 2\hat{\ell}(r+2).$$

We define  $d$  arbitrarily when  $f \notin \bar{S}$ . Note that (55) is satisfied. We now define  $L_0 = 6L_1$  and the penalties for  $L \geq L_0$ ,  $m \in \widehat{\mathcal{M}}_{\hat{\ell}}$  and  $\hat{g} \in \hat{S}_m$  by

$$\text{pen}_{1,m}(\hat{g}) = L \frac{(r+1)|m(\hat{g})| \log_+^2(n/(r+1))}{n}, \quad \text{pen}_2(m) = L \frac{(r+1)\hat{\ell} \log_+^2(n/(r+1))}{n}.$$

The first penalty satisfies the lower bound (57) since

$$d(\hat{g}) \leq (r+2)(2|m(\hat{g})| + 1) \leq 6(r+1)|m(\hat{g})| \quad \text{for all } \hat{g} \in \hat{S}_m.$$

It remains to define  $\text{pen}_1(\hat{g})$  for  $\hat{g} \in \hat{S} = \{\hat{s}_m, m \in \widehat{\mathcal{M}}_{\hat{\ell}}\}$ .

**Claim 7.** For all  $m, m' \in \widehat{\mathcal{M}}$ ,  $|m(\hat{s}_{m'})| \leq |m| + |m'|$ .

*Proof of Claim 7.* We have,

$$\begin{aligned} |m(\hat{s}_{m'})| &\leq \sum_{K \in m} |m'_K \vee K| \\ &\leq \sum_{K \in m} |\{K \cap K', K' \in m', K \cap K' \neq \emptyset\}| \\ &\leq |\{K \cap K', (K, K') \in m \times m', K \cap K' \neq \emptyset\}|. \end{aligned}$$

Since  $m$  and  $m'$  are partitions into intervals, we deduce that  $|m(\hat{s}_{m'})| \leq |m| + |m'|$ .  $\square$

It then follows that for all  $m, m' \in \widehat{\mathcal{M}}_{\hat{\ell}}$ ,

$$\text{pen}_{1,m}(\hat{s}_{m'}) \leq L \frac{(r+1)\hat{\ell} \log_+^2(n/(r+1))}{n} + \text{pen}_2(m).$$

The penalty defined by

$$\text{pen}_1(\hat{s}_{m'}) = L \frac{(r+1)\hat{\ell} \log_+^2(n/(r+1))}{n}$$

satisfies therefore (56).

Note now that the selection rules described in Sections 6.12 and 4.2 coincide. Theorem 13 controls the risk of the selected estimator: for all  $\xi > 0$ , there exists an event  $\Omega_\xi$  of probability larger than  $1 - e^{-n\xi}$ , and on which:

$$h^2(s, \hat{s}_{\hat{m}}) \leq C \left( \inf_{m \in \widehat{\mathcal{M}}_{\hat{\ell}}} \{h^2(s, \hat{s}_m) + \text{pen}_2(m)\} + \inf_{m \in \widehat{\mathcal{M}}_{\hat{\ell}}} \{h^2(s, \hat{s}_m) + \text{pen}_1(\hat{s}_m)\} + \xi \log_+(1/\xi) \right),$$

where  $C$  is a universal constant. By using the definition of the penalty terms,

$$h^2(s, \hat{s}_{\hat{m}}) \leq C' \left\{ \inf_{m \in \widehat{\mathcal{M}}_{\hat{\ell}}} \{h^2(s, \hat{s}_m)\} + L \frac{(r+1)\hat{\ell} \log_+^2(n/(r+1))}{n} + \xi \log_+(1/\xi) \right\},$$

where  $C'$  is a universal constant. It then remains to use the fact that  $\hat{s}_m$  is a  $\rho$ -estimator on  $\mathcal{P}_r(m)$  to get a bound on  $h^2(s, \hat{s}_m)$  on the same event  $\Omega_\xi$  (the event that appears in Theorem 7 to control the risk of a  $\rho$ -estimator is the same that the one that appears in Theorem 13. It is, in each case, defined by Lemma 10).  $\square$

6.12.3. *Proof of Theorem 11.* The proof is almost the same than the one of Theorem 10. The modifications are very mild, and this is the reason why we only specify the values of the different

parameters involved in the procedure of Section 6.12:

$$\begin{aligned}\hat{S} &= \{\hat{s}_m, m \in \widehat{\mathcal{M}}_{\hat{k}, \text{lower}}\} \\ \hat{S}_m &= \left\{ \sum_{K \in m} \hat{s}_{m_K} \mathbb{1}_K, m_K \in \widehat{\mathcal{M}}_{\hat{k}, \text{lower}} \right\} \quad \text{for all } m \in \widehat{\mathcal{M}}_{\hat{k}, \text{lower}} \\ \text{pen}_1(\hat{s}_m) = \text{pen}_2(m) &= L \frac{(r+1)|m| \log_+^2(n/(r+1))}{n} \quad \text{for all } m \in \widehat{\mathcal{M}}_{\hat{k}, \text{lower}}.\end{aligned}$$

□

**6.13. Proof of Lemma 3.** As in the proof of Theorem 1, the measure  $N$  can be put of the form  $N(A) = n^{-1} \sum_{i \in \hat{I}} \mathbb{1}_A(Y_i)$  where  $\hat{I} \subset \{1, \dots, n\}$ , and where the  $Y_i$  are suitable real-valued random variables.

Note that if  $K \cap \{Y_{(1)}, \dots, Y_{(\hat{n})}\} = \emptyset$  then,

$$L_K(f) = - \int_K f(t) dM(t),$$

and the supremum  $\sup_{f \in \mathcal{P}_r(K)} L_K(f)$  is achieved at  $\hat{s}_K = 0$  and equals 0. We now suppose that  $K \cap \{Y_{(1)}, \dots, Y_{(\hat{n})}\} \neq \emptyset$ .

Let  $V$  be the Radon–Nikodym derivative of  $M$  with respect to the Lebesgue measure  $\mu$ . Then,  $V = 1$  in framework 1,  $V(t) = n^{-1} \sum_{i=1}^n \mathbb{1}_{X_i \geq t} \mathbb{1}_{[0, +\infty)}(t)$  in framework 2 and  $V(t) = n^{-1} \sum_{i=1}^n \mathbb{1}_{X_{t-}^{(i)} = 1} \mathbb{1}_{I_{\text{obs}}}(t)$  in framework 3. Let  $k$  be the largest integer of  $\{1, \dots, \hat{n}\}$  such that  $Y_{(k)}$  belongs to  $K$  and  $K' = K \cap (-\infty, Y_{(k)}]$ . There exists some  $\alpha > 0$  such that  $(Y_{(k)} - \alpha, Y_{(k)}) \subset K'$ . Moreover, we can choose  $\alpha$  small enough to get  $V(t) \geq 1/n$  for all  $t \in (Y_{(k)} - \alpha, Y_{(k)})$ .

Let now  $f \in \mathcal{P}_r(K)$ . Then,  $L_K(f)$  takes the form

$$L_K(f) = \frac{1}{n} \sum_{i \in \hat{I}} (\log f(Y_i)) \mathbb{1}_K(Y_i) - \int_K f(t) V(t) dt,$$

and is bounded above by

$$L_K(f) \leq \log_+ \left( \sup_{t \in K'} f(t) \right) - \frac{1}{n} \int_{Y_{(k)} - \alpha}^{Y_{(k)}} f(t) dt.$$

We endow the linear space consisting of polynomial functions of degree at most  $r$  with the two following norms:

$$\|f\|_1 = \int_{Y_{(k)} - \alpha}^{Y_{(k)}} |f(t)| dt, \quad \|f\|_\infty = \sup_{t \in K'} |f(t)|.$$

This linear space being of finite dimension, there exists  $C$  such that  $\|f\|_\infty \leq C \|f\|_1$  for all  $f \in \mathcal{P}_r(K)$ . Now,

$$L_K(f) \leq \log_+ (C \|f\|_1) - \frac{\|f\|_1}{n}.$$

The continuous map  $L_K$  tends therefore to  $-\infty$  when  $\|f\|_1 \rightarrow +\infty$ . As there exists at least a function  $f \in \mathcal{P}_r(K)$  such that  $L_K(f) \neq -\infty$ ,  $\hat{s}_K$  does exist.

For the second part of the lemma, we use Theorem 1 to deduce that  $T(\hat{s}_K, f_K) \leq 0$  for all  $f_K \in \mathcal{P}_r(K)$ . If  $f \in \mathcal{P}_r(m)$  is of the form  $f = \sum_{K \in m} f_K$ ,

$$T(\hat{s}_m, f) = \sum_{K \in m} T(\hat{s}_K, f_K) \leq 0.$$

Thus,  $\gamma(\hat{s}_m) = 0$  and  $\hat{s}_m$  is a  $\rho$ -estimator on  $\mathcal{P}_r(m)$ .  $\square$

**6.14. Proof of Lemmas 4 and 5.** The following claim will be useful in the sequel.

**Claim 8.** *Let  $\xi > 0$  and  $\Omega_\xi$  be the event given by Lemma 10. Then,  $\mathbb{P}(\Omega_\xi) \geq 1 - e^{-n\xi}$ . Let  $\eta \geq 0$ ,  $r \geq 0$ , and  $m, m' \in \mathcal{M}$ . The following holds on  $\Omega_\xi$ : for all piecewise polynomial estimators  $\hat{s}_m \in \mathcal{P}_r(m)$ ,  $\hat{s}_{m'} \in \mathcal{P}_r(m')$  such that  $T(\hat{s}_m, \hat{s}_{m'}) \geq -\eta$ ,*

$$h^2(s, \hat{s}_{m'}) \leq C \left\{ h^2(s, \hat{s}_m) + \frac{(r+1)(|m| + |m'|)}{n} \log_+^2 \left( \frac{n}{(r+1)(|m| + |m'|)} \right) + \xi \log_+(1/\xi) + \eta \right\}.$$

Moreover, if  $\hat{s}_m$  is a  $\rho$ -estimator on  $\mathcal{P}_r(m)$ ,

$$h^2(s, \hat{s}_m) \leq C \left\{ h^2(s, \mathcal{P}_r(m)) + \frac{(r+1)|m|}{n} \log_+^2 \left( \frac{n}{(r+1)|m|} \right) + \xi \log_+(1/\xi) \right\}.$$

In the above inequalities,  $C$  is universal.

*Proof.* Let  $\varepsilon = 1/24$ . On  $\Omega_\xi$ :

$$(60) \quad T(\hat{s}_m, \hat{s}_{m'}) \leq (3 + 4\varepsilon)h^2(s, \hat{s}_m) - \frac{4 - 3\varepsilon}{12}h^2(s, \hat{s}_{m'}) + c_1\vartheta(d_{\hat{s}_{m'}}(\hat{s}_m)) + c_2\xi \log_+(1/\xi),$$

where  $c_1$  and  $c_2$  are universal constants. Now,  $\hat{s}_m$  and  $\hat{s}_{m'}$  belong to  $\mathcal{P}_r(m'')$  where

$$m'' = \{K \cap K', (K, K') \in m \times m', K \cap K' \neq \emptyset\}.$$

Yet,  $|m''| \leq |m| + |m'|$ . Thereby,  $\hat{s}_m$  and  $\hat{s}_{m'}$  belong to  $\mathcal{P}_{|m|+|m'|, r}$  and it follows from Proposition 5 that we may set

$$d_{\hat{s}_{m'}}(\hat{s}_m) = (r+2)(2(|m| + |m'|) + 1).$$

We now bound above  $\vartheta(d_{\hat{s}_{m'}}(\hat{s}_m))$  in (60), and then use  $T(\hat{s}_m, \hat{s}_{m'}) \geq -\eta$  to prove the first inequality of the claim. The second one follows from the proof of Theorem 7 (use (54) with  $S = \bar{S} = \mathcal{P}_r(m)$  and apply Proposition 5).  $\square$

*Proof of Lemma 4.* Let  $m \in \mathcal{M}'_\ell$  be a collection written as

$$m = \{[x_1, x_2], (x_2, x_3], (x_3, x_4], \dots, (x_\ell, x_{\ell+1}]\}$$

and such that  $x_1 \leq Y_{(1)}$ , and  $Y_{(\hat{n})} \leq x_{\ell+1}$ . We may define a partition  $\bar{m} \in \mathcal{M}'_\ell$  of the form

$$\bar{m} = \{[\bar{x}_1, \bar{x}_2], (\bar{x}_2, \bar{x}_3], (\bar{x}_3, \bar{x}_4], \dots, (\bar{x}_\ell, \bar{x}_{\ell+1}]\}$$

where  $\bar{x}_1 = Y_{(1)}$  and  $\bar{x}_{\ell+1} = Y_{(\hat{n})}$  and whose intervals are included into the ones of  $m$ .

Let  $\hat{s}_m$  and  $\hat{s}_{\bar{m}}$  be  $\rho$ -estimators on  $\mathcal{P}_0(m)$  and  $\mathcal{P}_0(\bar{m})$  respectively. We order the intervals of  $\bar{m}$  as follows. We define  $\ell$  intervals  $I_1, \dots, I_\ell$  such that  $\bar{m} = \{I_1, \dots, I_\ell\}$  and such that the value  $\hat{s}_{\bar{m}}$  on  $I_j$ , denoted by  $\hat{s}_{I_j}$ , is non-decreasing when  $j$  grows up. We denote the endpoints of  $I_j$  by  $a_j < b_j$ . We now define  $j_1$  as the largest integer of  $\{1, \dots, \hat{n}\}$  such that  $Y_{(j_1)} \leq a_j$  and  $j_2$  as the smallest integer such that  $Y_{(j_2)} \geq b_j$ . When  $j_1 = 1$ , we set  $K_j = [Y_{(j_1)}, Y_{(j_2)}]$  and when  $j_1 \neq 1$ , we set

$K_j = (Y_{(j_1)}, Y_{(j_2)})]$ . Note that  $K_j$  is the smallest interval containing  $I_j$  that is either of the form  $[Y_{(j_1)}, Y_{(j_2)}]$  or  $(Y_{(j_1)}, Y_{(j_2)})]$ .

Define  $J_1 = K_1$  and for  $j \in \{2, \dots, \ell\}$ ,  $J_j = K_j \setminus \bigcup_{i=1}^{j-1} K_i$ . Since  $K_i \not\subset K_j$  when  $i \neq j$ ,  $K_j \setminus K_i$  is an interval. Therefore,  $J_j = \bigcap_{i=1}^j (K_j \setminus K_i)$  is also an interval. When it is not empty, it is either of the form  $[Y_{(1)}, Y_{(i)}]$  with  $i > 1$  or  $(Y_{(i_1)}, Y_{(i_2)})]$  with  $i_1 < i_2$ . The collection  $\bar{m}' = \{J_j, j \in \{1, \dots, \ell\}\}$  defines therefore a partition of  $[Y_{(1)}, Y_{(\hat{n})}]$  that belongs to  $\widehat{\mathcal{M}}_{\ell'}$  with  $\ell' \leq \ell$  (we must remove the empty sets). Let  $f$  be the step function of  $\mathcal{P}_0(\bar{m}')$  defined by

$$f = \sum_{j=1}^{\ell} \hat{s}_{I_j} \mathbb{1}_{J_j}.$$

We now prove that  $f \leq \hat{s}_{\bar{m}}$ . When  $x \notin [Y_{(1)}, Y_{(\hat{n})}]$ ,  $f(x) = \hat{s}_{\bar{m}}(x) = 0$ . When  $x \in [Y_{(1)}, Y_{(\hat{n})}]$ , there exist  $j \in \{1, \dots, \ell\}$  such that  $x \in I_j$  and  $j' \leq j$  such that  $x \in J_{j'}$ . Therefore,  $f(x) = \hat{s}_{I_{j'}}$ . By using that  $\hat{s}_{I_{j'}} \leq \hat{s}_{I_j}$ , we finally deduce that  $f(x) \leq \hat{s}_{\bar{m}}(x)$ .

Consider an interval  $I_j \in \bar{m}$  and let us denote the cardinal of  $\{Y_{(i)}, Y_{(i)} \in I_j, i \in \{1, \dots, \hat{n}\}\}$  by  $k_j$ . When  $k_j \geq 3$ , there exists at least  $k_j - 2$  random variables  $Y_{(i)}$  that belong to  $I_j$  but not to  $\bigcup_{j' \in \{1, \dots, \hat{n}\}, j' \neq j} K_{j'}$ . Such  $Y_{(i)}$  belong therefore to  $J_j$  and satisfy  $f(Y_{(i)}) = \hat{s}_{\bar{m}}(Y_{(i)})$ . Therefore,

$$\begin{aligned} |\{Y_{(i)}, f(Y_{(i)}) \neq \hat{s}_{\bar{m}}(Y_{(i)}), i \in \{1, \dots, \hat{n}\}\}| &= \sum_{j=1}^{\ell} |\{Y_{(i)}, f(Y_{(i)}) \neq \hat{s}_{\bar{m}}(Y_{(i)}), Y_{(i)} \in I_j, i \in \{1, \dots, \hat{n}\}\}| \\ (61) \qquad \qquad \qquad &\leq 2\ell. \end{aligned}$$

It follows from  $f \leq \hat{s}_{\bar{m}}$ , (61) and (28) that  $T(f, \hat{s}_{\bar{m}}) \leq 2\ell/n$ . We now use Claim 8 to get on  $\Omega_{\xi}$

$$(62) \qquad h^2(s, f) \leq C \left\{ h^2(s, \hat{s}_{\bar{m}}) + \frac{\ell}{n} \log_+^2(n/\ell) + \xi \log_+(1/\xi) \right\},$$

where  $C$  is universal.

We may refine the partition  $\bar{m}' \in \widehat{\mathcal{M}}_{\ell'}$  to get  $m' \in \widehat{\mathcal{M}}_{\ell}$  such that  $\mathcal{P}_0(\bar{m}') \subset \mathcal{P}_0(m')$ . Let  $\hat{s}_{m'}$  and  $\hat{s}_{\bar{m}'}$  be  $\rho$ -estimators on  $\mathcal{P}_0(m')$  and  $\mathcal{P}_0(\bar{m}')$  respectively. There exists a universal constant  $C'$  such that on  $\Omega_{\xi}$ :

$$h^2(s, \hat{s}_{m'}) \leq C' \left\{ h^2(s, \mathcal{P}_0(m')) + \frac{\ell}{n} \log_+^2(n/\ell) + \xi \log_+(1/\xi) \right\}.$$

By using that  $f \in \mathcal{P}_0(\bar{m}') \subset \mathcal{P}_0(m')$  and (62),

$$\begin{aligned} h^2(s, \hat{s}_{m'}) &\leq C' \left\{ h^2(s, f) + \frac{\ell}{n} \log_+^2(n/\ell) + \xi \log_+(1/\xi) \right\}, \\ (63) \qquad \qquad \qquad &\leq C'' \left\{ h^2(s, \hat{s}_{\bar{m}}) + \frac{\ell}{n} \log_+^2(n/\ell) + \xi \log_+(1/\xi) \right\}, \end{aligned}$$

where  $C'''$  is universal. Note now that  $\hat{s}_m \mathbb{1}_{[Y_{(1)}, Y_{(\hat{n})}]} \in \mathcal{P}_0(\bar{m})$  and thus  $T(\hat{s}_{\bar{m}}, \hat{s}_m \mathbb{1}_{[Y_{(1)}, Y_{(\hat{n})}]}) \leq 0$  as  $\hat{s}_{\bar{m}}$  is a  $\rho$ -estimator on the convex model  $\mathcal{P}_0(\bar{m})$  (see Theorem 1 and Lemma 3). Now,

$$\begin{aligned} T(\hat{s}_{\bar{m}}, \hat{s}_m) &= T(\hat{s}_{\bar{m}}, \hat{s}_m \mathbb{1}_{[Y_{(1)}, Y_{(\hat{n})}]}) + \frac{1}{4} \left( \int_{\mathbb{R}} \hat{s}_m \mathbb{1}_{[Y_{(1)}, Y_{(\hat{n})}]} dM - \int_{\mathbb{R}} \hat{s}_m dM \right) \\ &\leq 0. \end{aligned}$$

Therefore, Claim 8 asserts that

$$(64) \quad h^2(s, \hat{s}_{\bar{m}}) \leq C''' \left\{ h^2(s, \hat{s}_m) + \frac{\ell}{n} \log_+^2(n/\ell) + \xi \log_+(1/\xi) \right\},$$

where  $C'''$  is universal. By using that  $\hat{s}_m$  is a  $\rho$ -estimator,

$$(65) \quad h^2(s, \hat{s}_m) \leq C''' \left\{ h^2(s, \mathcal{P}_0(m)) + \frac{\ell}{n} \log_+^2(n/\ell) + \xi \log_+(1/\xi) \right\}.$$

It remains to put inequalities (63), (64) and (65) together to finish the proof.  $\square$

*Proof of Lemma 5.* Note that we may always suppose that

$$\{K \cap \{Y_{(1)}, \dots, Y_{(\hat{n})}\}, K \in m\}$$

contains  $Y_{(1)}$  and  $Y_{(\hat{n})}$  (up to an increase of  $|m|$  by 2). Let

$$m_1 = \{K \in m, \{Y_{(1)}, \dots, Y_{(\hat{n})}\} \cap K \neq \emptyset\}.$$

Then,  $m_1 \neq \emptyset$  and we may write  $m_1 = \{K_j, j \in \{1, \dots, \ell\}\}$  where  $1 \leq \ell \leq |m|$  and where  $K_j$  is an interval with endpoints  $a_j, b_j$  satisfying  $a_1 < b_1 \leq a_2 < b_2 < \dots$ .

Let us recall that the  $\rho$ -estimator  $\hat{s}_m$  is of the form

$$\hat{s}_m = \sum_{K \in m} \hat{s}_K \quad \text{where } \hat{s}_K \text{ maximizes } L_K \text{ over } \mathcal{P}_r(K).$$

When  $K \in m$  does not belong to  $m_1$ ,  $\hat{s}_K = 0$  and hence

$$\hat{s}_m = \sum_{j=1}^{\ell} \hat{s}_{K_j}.$$

For each  $j \in \{1, \dots, \ell\}$ , we set  $\alpha_j = \min\{Y_{(i)}, Y_{(i)} \in K_j\}$ ,  $\beta_j = \max\{Y_{(i)}, Y_{(i)} \in K_j\}$ . We define for  $j \in \{2, \dots, \ell - 1\}$ ,  $J_{2j} = (\beta_j, \alpha_{j+1}]$  and for  $j \in \{2, \dots, \ell\}$ ,  $J_{2j-1} = (\alpha_j, \beta_j]$ . If  $\beta_1 = Y_{(1)}$ , we set  $J_1 = \emptyset$ ,  $J_2 = [\beta_1, \alpha_2]$  and if  $\beta_1 > Y_{(1)}$ ,  $J_1 = [Y_{(1)}, \beta_1]$ ,  $J_2 = (\beta_1, \alpha_2]$ . Note that  $J_{2j-1} \subset K_j$  for all  $j \in \{1, \dots, \ell\}$ . The collection  $m' = \{J_j, j \in \{1, \dots, 2\ell - 1\}\}$  defines a partition of  $\widehat{\mathcal{M}}$  such that  $|m'| \leq 2\ell - 1$ . We define the  $\rho$ -estimator

$$\hat{s}_{m'} = \sum_{j=1}^{\ell} \hat{s}_{J_{2j-1}} + \sum_{j=1}^{\ell-1} \hat{s}_{J_{2j}},$$

where  $\hat{s}_A$  maximizes  $L_A$  over  $\mathcal{P}_r(A)$  for all non-empty interval  $A$  with the convention that  $\hat{s}_{\emptyset} = 0$  when  $A = \emptyset$ . We now consider

$$\tilde{s}_{m'} = \sum_{j=1}^{\ell} \hat{s}_{J_{2j-1}}.$$

Note that  $\tilde{s}_{m'}$  also belongs to the random model  $\mathcal{P}_r(m')$  and hence  $T(\hat{s}_{m'}, \tilde{s}_{m'}) \leq 0$ . We deduce from Claim 8 that on  $\Omega_\xi$ :

$$(66) \quad h^2(s, \hat{s}_{m'}) \leq C \left\{ h^2(s, \tilde{s}_{m'}) + \frac{(r+1)|m|}{n} \log_+^2 \left( \frac{n}{(r+1)|m|} \right) + \xi \log_+(1/\xi) \right\},$$

where  $C$  is universal.

Now, for all  $j \in \{1, \dots, \ell\}$ , such that  $J_{2j-1} \neq \emptyset$ ,

$$(67) \quad T(\hat{s}_{J_{2j-1}}, \hat{s}_{K_j} \mathbb{1}_{J_{2j-1}}) \leq 0,$$

since  $\hat{s}_{J_{2j-1}}$  maximizes  $L_{J_{2j-1}}$  over  $\mathcal{P}_r(J_{2j-1})$  and that  $\hat{s}_{K_j} \mathbb{1}_{J_{2j-1}} \in \mathcal{P}_r(J_{2j-1})$ . When  $J_{2j-1} = \emptyset$ ,  $T(\hat{s}_{J_{2j-1}}, \hat{s}_{K_j} \mathbb{1}_{J_{2j-1}}) = 0$ , and thus (67) also holds.

We define

$$A = \bigcup_{j=1}^{\ell} J_{2j-1}.$$

We deduce from (67) that  $T(\tilde{s}_{m'} \mathbb{1}_A, \hat{s}_m \mathbb{1}_A) \leq 0$ . Therefore,

$$\begin{aligned} T(\tilde{s}_{m'}, \hat{s}_m) &= T(\tilde{s}_{m'} \mathbb{1}_A, \hat{s}_m \mathbb{1}_A) + T(0, \hat{s}_m \mathbb{1}_{A^c}) \\ &\leq 0 + T(0, \hat{s}_m \mathbb{1}_{A^c}) \\ &\leq \int_{A^c} \psi(\hat{s}_m/0) \, dN, \end{aligned}$$

where we recall the conventions  $\psi(0/0) = \psi(1) = 0$ ,  $\psi(x/0) = \psi(\infty) = 1$  for all  $x > 0$ . Let  $B = \bigcup_{j=1}^{\ell} K_j$ . Note that  $\hat{s}_m$  vanishes outside  $B$  and thus, as  $|\psi| \leq 1$ ,

$$(68) \quad T(\tilde{s}_{m'}, \hat{s}_m) \leq \int_{B \cap A^c} \psi(\hat{s}_m/0) \, dN \leq N(B \cap A^c).$$

Now,

$$N(B \cap A^c) = \sum_{j=1}^{\ell} \{N(K_j) - N(J_{2j-1})\}.$$

Since  $\alpha_j, \beta_j \in \{Y_i, i \in \hat{I}\}$ , we deduce from (28) that  $N(K_j) - N(J_{2j-1}) = N(\{\alpha_j\})$ . In each of the frameworks,  $N(\{\alpha_j\}) \leq 1/n$  and thus  $N(B \cap A^c) \leq \ell/n$ . By using (68), we get  $T(\tilde{s}_{m'}, \hat{s}_m) \leq \ell/n$ . Claim 8 with  $\eta = \ell/n \leq |m|/n$  ensures that on  $\Omega_\xi$ :

$$h^2(s, \tilde{s}_{m'}) \leq C' \left\{ h^2(s, \hat{s}_m) + \frac{(r+1)|m|}{n} \log_+^2 \left( \frac{n}{(r+1)|m|} \right) + \xi \log_+(1/\xi) \right\},$$

where  $C'$  is universal. Since  $\hat{s}_m$  is a  $\rho$ -estimator on  $\mathcal{P}_r(m)$ , we deduce that on the same event  $\Omega_\xi$ :

$$h^2(s, \hat{s}_m) \leq C'' \left\{ h^2(s, \mathcal{P}_r(m)) + \frac{(r+1)|m|}{n} \log_+^2 \left( \frac{n}{(r+1)|m|} \right) + \xi \log_+(1/\xi) \right\},$$

where  $C''$  is universal. It then remains to combine the two last inequalities with (66) to finish the proof.  $\square$

## REFERENCES

- [AD10] Nathalie Akakpo and Cécile Durot. Histogram selection for possibly censored data. *Mathematical Methods of Statistics*, 19(3):189–218, 2010.
- [Ant89] Anestis Antoniadis. A penalty method for nonparametric estimation of the intensity function of a counting process. *Annals of the Institute of Statistical Mathematics*, 41(4):781–807, 1989.
- [Bar11] Yannick Baraud. Estimator selection with respect to Hellinger-type risks. *Probability Theory and Related Fields*, 151(1-2):353–401, 2011.
- [Bar16] Yannick Baraud. Bounding the expectation of the supremum of an empirical process over a (weak) vc-major class. *Electronic journal of statistics*, 10(2):1709–1728, 2016.
- [BB09] Yannick Baraud and Lucien Birgé. Estimating the intensity of a random measure by histogram type estimators. *Probability Theory and Related Fields*, 143:239–284, 2009.
- [BB16] Yannick Baraud and Lucien Birgé.  $\rho$ -estimators for shape restricted density estimation. *Stochastic Processes and their Applications*, 126(12):3888–3912, 2016.
- [BB17] Yannick Baraud and Lucien Birgé.  $\rho$ -estimators revisited: general theory and applications. *arXiv preprint arXiv:1605.05051*, 2017.
- [BBM99] Andrew Barron, Lucien Birgé, and Pascal Massart. Risk bounds for model selection via penalization. *Probability theory and related fields*, 113(3):301–413, 1999.
- [BBS17] Yannick Baraud, Lucien Birgé, and Mathieu Sart. A new method for estimation and model selection:  $\rho$ -estimation. *Inventiones mathematicae*, 207(2):425–517, 2017.
- [BC05] Elodie Brunel and Fabienne Comte. Penalized contrast estimation of density and hazard rate with censored data. *Sankhyā: The Indian Journal of Statistics*, pages 441–475, 2005.
- [BC08] Elodie Brunel and Fabienne Comte. Adaptive estimation of hazard rate with censored data. *Communications in Statistics—Theory and Methods*, 37(8):1284–1305, 2008.
- [Bir06] Lucien Birgé. Model selection via testing: an alternative to (penalized) maximum likelihood estimators. *Annales de l’Institut Henri Poincaré. Probabilités et Statistique*, 42(3):273–325, 2006.
- [BLM13] Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. *Concentration inequalities: A nonasymptotic theory of independence*. Oxford university press, 2013.
- [BM98] Lucien Birgé and Pascal Massart. Minimum contrast estimators on sieves: exponential bounds and rates of convergence. *Bernoulli*, 4(3):329–375, 1998.
- [BR06] Lucien Birgé and Yves Rozenholc. How many bins should be put in a regular histogram. *ESAIM Probab. Stat.*, 10:24–45, 2006.
- [Cas99] Gwénaëlle Castellan. Modified akaike’s criterion for histogram density estimation. *Technical report*, 1999.
- [CR04] Fabienne Comte and Yves Rozenholc. A new algorithm for fixed design regression and denoising. *Annals of the Institute of Statistical Mathematics*, 56(3):449–473, 2004.
- [DL12] Luc Devroye and Gábor Lugosi. *Combinatorial methods in density estimation*. Springer Science & Business Media, 2012.
- [DR02] Sebastian Dohler and Ludger Ruschendorf. Adaptive estimation of hazard functions. *Probability and mathematical statistics - Wroclaw University*, 22(2):355–379, 2002.
- [GG01] Evarist Giné and Armelle Guillou. On consistency of kernel density estimators for randomly censored data: rates holding uniformly over adaptive intervals. *Annales de l’Institut Henri Poincaré. Probabilités et Statistique*, 37(4):503–522, 2001.
- [GK06] Evarist Giné and Vladimir Koltchinskii. Concentration inequalities and asymptotic results for ratio type empirical processes. *The Annals of Probability*, 34(3):1143–1216, 2006.



- [GN15] Evarist Giné and Richard Nickl. *Mathematical foundations of infinite-dimensional statistical models*, volume 40. Cambridge University Press, 2015.
- [Gre56] Ulf Grenander. On the theory of mortality measurement. *Scandinavian Actuarial Journal*, 1956(1):70–96, 1956.
- [Kan92] Yuichiro Kanazawa. An optimal variable cell histogram based on the sample spacings. *The Annals of Statistics*, pages 291–304, 1992.
- [Mas07] Pascal Massart. *Concentration inequalities and model selection*, volume 6. Springer, 2007.
- [Pla09] Sandra Placade. Non parametric estimation of hazard rate in presence of censoring. *hal preprint hal-00410799*, 2009.
- [RB06] Patricia Reynaud-Bouret. Penalized projection estimators of the aalen multiplicative intensity. *Bernoulli*, 12(4):633–661, 2006.
- [RMG10] Yves Rozenholc, Thoralf Mildenberger, and Ursula Gather. Combining regular and irregular histograms by penalized likelihood. *Computational Statistics & Data Analysis*, 54(12):3313–3323, 2010.
- [Sar14] Mathieu Sart. Estimation of the transition density of a Markov chain. *Annales de l’Institut Henri Poincaré. Probabilités et Statistique*, 50(3):1028–1068, 2014.
- [Sau72] Norbert Sauer. On the density of families of sets. *Journal of Combinatorial Theory, Series A*, 13(1):145–147, 1972.
- [vdG95] Sara van de Geer. Exponential inequalities for martingales, with application to maximum likelihood estimation for counting processes. *The Annals of Statistics*, pages 1779–1801, 1995.

UNIV LYON, UJM-SAINT-ÉTIENNE, CNRS UMR 5208, INSTITUT CAMILLE JORDAN, 10 RUE TRÉFILIERIE, CS 82301, F-42023 SAINT-ETIENNE CEDEX 2, FRANCE

*E-mail address:* mathieu.sart@univ-st-etienne.fr