

Universités de Paris 6 & Paris 7 - CNRS (UMR 7599)

**PRÉPUBLICATIONS DU LABORATOIRE  
DE PROBABILITÉS & MODÈLES ALÉATOIRES**

4, place Jussieu - Case 188 - 75 252 Paris cedex 05

<http://www.proba.jussieu.fr>

**Classification under polynomial entropy and  
margin assumptions and randomized estimators**

**J.-Y. AUDIBERT**

**AVRIL 2004**

Prépublication n° 908

Laboratoire de Probabilités et Modèles Aléatoires, CNRS-UMR 7599,  
Université Paris VI & Université Paris VII,  
4, place Jussieu, Case 188, F-75252 Paris Cedex 05.

&

CREST, ENSAE, Laboratoire de Finance et Assurance,  
15 Bd Gabriel Péri, F-92245 Malakoff Cedex.

# CLASSIFICATION UNDER POLYNOMIAL ENTROPY AND MARGIN ASSUMPTIONS AND RANDOMIZED ESTIMATORS

J.-Y. AUDIBERT

*Université Paris VI and CREST*

ABSTRACT. The aim of this paper is two-fold. First we want to develop the PAC-Bayesian point of view [13, 3, 4, 1] and show how the efficiency of a Gibbs estimator relies on the weights given by the prior distribution to the balls centered at the best function in the model and associated with the pseudo-distance  $(f_1, f_2) \mapsto \mathbb{P}[f_1(X) \neq f_2(X)]$ .

Secondly, we show how to recover and improve results under empirical and non empirical polynomial entropy assumptions and Tsybakov's margin assumption. We also study the links between empirical and non empirical nets and give an observable version of the integral entropy [6, 9, 14].

## CONTENTS

1. Setup and notations	2
1.1. Measurability	4
1.2. Covering, packing and bracketing nets and entropies	4
2. Known PAC-Bayesian bounds	5
3. Convergence rate of classifiers under complexity and margin assumptions	6
3.1. Complexity and margin assumptions	6
3.1.1. Complexity assumptions	6
3.1.2. Margin assumptions	6
3.2. Gibbs classifier	7
3.2.1. Under Assumptions $(MA4)$ and $(CA3)$ for $q > 0$	7
3.2.2. Under Assumptions $(MA2)$ and $(CA3)$ for $q = 0$	8
3.2.3. Under Assumption $(MA2)$ and a local complexity assumption	9
3.2.4. Adaptive choice of the temperature	10
3.3. Empirical risk minimization on nets	10
3.3.1. Under Assumptions $(MA3)$ and $(CA1)$ for $q > 0$	10
3.3.2. Under Assumptions $(MA2)$ and $(CA1)$ for $q = 0$	11
3.4. Chaining	11
3.5. Bracketing entropy	14
4. Classification under empirical complexity assumptions	16

---

2000 *Mathematics Subject Classification*. Primary 62H30, Secondary 68Q32.

*Key words and phrases*. Gibbs classifiers, entropy assumptions, margin assumptions, PAC-Bayesian bounds, chaining, oracle inequalities, VC theory.

I would like to thank very deeply Professor Olivier Catoni, my PhD advisor, for the numerous hours he spent with me. I am also very grateful to Olivier Bousquet and Professor Alexandre Tsybakov for the numerous discussions which have fed this work. I would also like to thank Professor Bernhard Schölkopf who invited me twice at the Max Planck Institute of Tuebingen. Meeting people there has been a key point for widening my scope of interest during my PhD.

4.1.	Concentration of the empirical entropies	16
4.2.	Chaining empirical quantities...	17
4.2.1.	...in the transductive learning	17
4.2.2.	...in the inductive learning	18
4.3.	Application to VC-classes	19
5.	Assouad's lemma	20
6.	Proofs	21
6.1.	Proof of Lemma 3.1	21
6.2.	Proof of Theorem 3.3	23
6.3.	Proof of Theorem 3.5	24
6.3.1.	First case: $\log \pi^{-1}(\Delta R \leq x) = -C' \log x + C''' + o(x^s)$	24
6.3.2.	Second case: $\log \pi^{-1}(\Delta R \leq x) = C' x^{-\frac{q}{\kappa}} + C''' + o(1)$	25
6.4.	Proof of Theorem 3.6	26
6.5.	Proof of Theorem 3.7	26
6.6.	Proof of Theorem 3.9	27
6.7.	Proof of Theorem 3.10	28
6.7.1.	First step: upper bounds due to the chaining technique	28
6.7.2.	Second step: determining the radius of the nets	30
6.8.	Proof of Theorem 3.15	31
6.9.	Proof of Theorem 3.16	32
6.10.	Proof of Theorem 4.1	32
6.11.	Proof of Theorem 4.3	33
6.12.	Proof of Theorem 4.4	34
6.13.	Proof of Corollary 4.5	35
6.14.	Proof of Corollary 4.6	36
6.15.	Proof of Lemma 5.1	36
Appendix A.	Proof of inequality (6.2)	38
Appendix B.	Proof of inequality (6.4)	39
Appendix C.	Proof of inequality (6.6)	39
Appendix D.	Proof of inequality (6.7)	41
Appendix E.	Another way of getting the right order	42
Appendix F.	Proof of Theorem 5.2	43
References		44

## 1. SETUP AND NOTATIONS

We assume that we observe an i.i.d. sample  $Z_1^N \triangleq (X_i, Y_i)_{i=1}^N$  of random variables distributed according to a product probability measure  $\mathbb{P}^{\otimes N}$ , where  $\mathbb{P}$  is a probability distribution on  $(\mathcal{Z}, \mathcal{B}_{\mathcal{Z}}) \triangleq (\mathcal{X} \otimes \mathcal{Y}, \mathcal{B}_{\mathcal{X}} \otimes \mathcal{B}_{\mathcal{Y}})$ ,  $(\mathcal{X}, \mathcal{B}_{\mathcal{X}})$  is a measurable space called the pattern space,  $\mathcal{Y} = \{1, \dots, |\mathcal{Y}|\}$  is the (finite) label space and  $\mathcal{B}_{\mathcal{Y}}$  is the sigma algebra of all subsets of  $\mathcal{Y}$ . Let  $\mathbb{P}(dY|X)$  denote a regular version of the conditional probabilities (which we will use in the following without further mention).

Let  $\mathcal{F}(\mathcal{X}, \mathcal{Y})$  denote the set of all measurable functions mapping  $\mathcal{X}$  into  $\mathcal{Y}$ . The aim of a classification procedure is to build a function  $f \in \mathcal{F}(\mathcal{X}, \mathcal{Y})$  from the learning sample such that  $f(X)$  predicts the label  $Y$  associated with  $X$ . The quality of the

prediction is measured by the expected risk

$$R(f) \triangleq \mathbb{P}[Y \neq f(X)].$$

A function  $f_{\mathbb{P}}^*$  such that for any  $x \in \mathcal{X}$ ,

$$f_{\mathbb{P}}^*(x) \in \operatorname{argmax}_{y \in \mathcal{Y}} \mathbb{P}(Y = y | X = x)$$

minimizes the expected risk. This function is not necessarily unique. We assume that there exists a measurable one. We will once for all fix it, refer to it as the Bayes classifier and often denote it  $f^*$  to shorten. Since we have no prior information about the distribution  $\mathbb{P}$  of  $(X, Y)$ , this classifier is unknown.

Since there is generally no measurable estimator  $\hat{f} : \mathcal{Z}^N \rightarrow \mathcal{F}(\mathcal{X}, \mathcal{Y})$  such that

$$\lim_{N \rightarrow +\infty} \sup_{\mathbb{P} \in \mathcal{M}_+^1(\mathcal{Z})} \left\{ \mathbb{P}^{\otimes(N+1)}[Y_{N+1} \neq \hat{f}(Z_1^N)(X_{N+1})] - \inf_{f \in \mathcal{F}(\mathcal{X}, \mathcal{Y})} \mathbb{P}[Y \neq f(X)] \right\} = 0,$$

we have to work with a prescribed set of classification functions  $\mathcal{F}$ , called the model. This set is just some subset of the set of all measurable functions  $\mathcal{F}(\mathcal{X}, \mathcal{Y})$ . Let us denote  $\tilde{f}$  the best function in the model, i.e. a function minimizing the expected risk:

$$\tilde{f} \in \operatorname{argmin}_{\mathcal{F}} R.$$

For sake of simplicity, we assume that it exists<sup>1</sup>. Let

$$\bar{\mathbb{P}} \triangleq \frac{1}{N} \sum_{i=1}^N \delta_{(X_i, Y_i)}$$

be the empirical distribution. The empirical risk

$$r(f) \triangleq \bar{\mathbb{P}}[Y \neq f(X)]$$

gives an estimate of the expected risk : from the law of large numbers, for any measurable function, it tends to the expected risk almost surely. An estimator which minimizes the empirical risk

$$\hat{f}_{\text{ERM}} \in \operatorname{argmin}_{\mathcal{F}} r$$

is called an ERM<sup>2</sup>-classifier. The regression function will be denoted

$$\eta^*(x) \triangleq \mathbb{P}(Y | X = x).$$

In the binary classification setting ( $\mathcal{Y} = \{0; 1\}$ ), we have  $\eta^*(x) = \mathbb{P}(Y = 1 | X = x)$ .

Since we will study randomized estimators, we assume that we have a  $\sigma$ -algebra  $\mathcal{T}$  such that  $(\mathcal{F}, \mathcal{T})$  is a measurable space containing the sets  $\{f\}$  for any  $f \in \mathcal{F}$  and such that the function

$$\begin{aligned} \mathcal{F} \times \mathcal{X} &\rightarrow \mathcal{Y} \\ (f, x) &\mapsto f(x) \end{aligned}$$

is measurable. A randomized estimator consists in drawing a function in  $\mathcal{F}$  according to some random distribution  $\hat{\rho} : \mathcal{Z}^N \rightarrow \mathcal{M}_+^1(\mathcal{F})$ , where  $\mathcal{M}_+^1(\mathcal{F})$  is the set of probability distributions on the measurable space  $(\mathcal{F}, \mathcal{T})$ .

<sup>1</sup>Otherwise we would have to introduce some small positive real  $\beta$  and consider  $\tilde{f}$  as an estimator minimizing the expected risk up to  $\beta$ . This real  $\beta$  would then appear in all the equations related to this function and make things needlessly messy.

<sup>2</sup>ERM = Empirical Risk Minimization

To shorten, we will use  $\mu h$  to denote the expectation of the random variable  $h$  under the probability distribution  $\mu$ :  $\mu h \triangleq \int h(x)d\mu(x)$ . The Kullback-Leibler divergence between two probability distributions is defined as  $K(\mu, \nu) = \mu \log \frac{d\mu}{d\nu}$  when  $\mu$  is absolutely continuous with respect to  $\nu$  and  $K(\mu, \nu) = +\infty$  otherwise.

The symbol  $C$  will denote a positive universal constant whose value may differ from line to line whereas the symbol  $\check{C}$  will denote a positive constant whose value depends on other constants and may also differ from line to line.

We define

$$\pi_h \triangleq \frac{\exp(h)}{\pi \exp(h)} \cdot \pi$$

for any measurable real function  $h$  such that  $\exp(h)$  is  $\pi$ -integrable. The randomized estimators associated with the posterior distributions  $\pi_{-C_T}$  will be called the standard Gibbs estimators with temperature  $\frac{1}{C}$ .

**1.1. Measurability.** Finally, to circumvent some measurability problems, we will consider inner and outer expectations. Let  $(A, \mathcal{A}, \mu)$  be a measure space and  $\mathcal{C}(A; \mathbb{R})$  be the class of real measurable functions. For any (measurable or not) function  $f$ , its inner and outer expectation wrt  $\mu$  are respectively  $\mu_*(h) \triangleq \sup \{\mu(g) : g \in \mathcal{C}(A; \mathbb{R}), g \leq h\}$  and  $\mu^*(h) \triangleq \inf \{\mu(g) : g \in \mathcal{C}(A; \mathbb{R}), g \geq h\}$ . Naturally, for any set  $B \subset A$ ,  $\mu_*(B)$  and  $\mu^*(B)$  are defined by  $\mu_*(B) = \mu_*(\mathbb{1}_B)$  and  $\mu^*(B) = \mu^*(\mathbb{1}_B)$ . Note that  $\mu_*$  and  $\mu^*$  are not measures but satisfy  $\mu^*(B) + \mu_*(B^c) = 1$  and  $\mu^*(B_1 \cup B_2) \leq \mu^*(B_1) + \mu^*(B_2)$ . Besides, if  $\mu^*(h) < +\infty$ , then there exists a random variable  $h^*$  such that  $\mu^*(h) = \mu(h^*)$ . For more details on properties of inner and outer expectations, see [17].

**1.2. Covering, packing and bracketing nets and entropies.** Let  $\mathbb{Q}$  denote a probability distribution on the measurable space  $(\mathcal{X}, \mathcal{B}_{\mathcal{X}})$ . The mapping  $\mathbb{Q}_{\cdot, \cdot}$  from  $\mathcal{F} \times \mathcal{F}$  into  $\mathbb{R}_+$  defined as

$$\mathbb{Q}_{f_1, f_2} \triangleq \mathbb{Q}[f_1(X) \neq f_2(X)] \quad \text{for any } f_1, f_2 \in \mathcal{F}$$

is a pseudo-distance. For any  $u \geq 0$ , a set of measurable functions  $\mathcal{G} \subset \mathcal{F}(\mathcal{X}, \mathcal{Y})$  such that

$$\sup_{f \in \mathcal{F}} \inf_{g \in \mathcal{G}} \mathbb{Q}_{f, g} \leq u$$

is called a  $u$ -covering net of the set  $\mathcal{F}$  wrt the pseudo-distance  $\mathbb{Q}$ .

The log-cardinal  $H(u, \mathcal{F}, \mathbb{Q}_{\cdot, \cdot})$  of the smallest  $u$ -covering net (possibly infinite) is called the  $u$ -covering entropy. A  $u$ -covering net with log-cardinal equal to  $H(u, \mathcal{F}, \mathbb{Q}_{\cdot, \cdot})$  is called a minimal  $u$ -covering net<sup>3</sup>.

In bracketing nets, we require in addition that any function in  $\mathcal{F}$  can be encapsulated by two functions of the net. Specifically, for any  $u \geq 0$ , a set of measurable functions  $\mathcal{G} \subset \mathcal{F}(\mathcal{X}, \mathcal{Y})$  such that for any function  $f \in \mathcal{F}$ , there exist  $f_L, f_U \in \mathcal{G}$  satisfying  $f_L \leq f \leq f_U$  and  $\mathbb{Q}_{f_L, f_U} \leq u$ , is called a  $u$ -bracketing net of the set  $\mathcal{F}$  wrt the pseudo-distance  $\mathbb{Q}$ . The log-cardinal  $H^{[\cdot]}(u, \mathcal{F}, \mathbb{Q}_{\cdot, \cdot})$  of the smallest  $u$ -bracketing net (possibly infinite) is called the  $u$ -bracketing entropy. A  $u$ -bracketing net with log-cardinal equal to  $H^{[\cdot]}(u, \mathcal{F}, \mathbb{Q}_{\cdot, \cdot})$  is called a minimal  $u$ -bracketing net.

<sup>3</sup>Here the functions in the net can be taken outside  $\mathcal{F}$ . This is not so important since it is well-known that a  $2u$ -covering net with functions in  $\mathcal{F}$  can be constructed from any  $u$ -covering net.

Packing nets are covering nets such that for any functions  $f_1, f_2$  in the net, we have  $\mathbb{Q}_{f_1, f_2} > u$ . The packing entropy  $H_p(u, \mathcal{F}, \mathbb{Q}, \cdot)$  is the log-cardinal of a minimal packing net.

We have  $H(u, \mathcal{F}, \mathbb{Q}, \cdot) \leq H_p(u, \mathcal{F}, \mathbb{Q}, \cdot) \leq H(\frac{u}{2}, \mathcal{F}, \mathbb{Q}, \cdot)$ . Any  $u$ -bracketing net is a  $u$ -covering net. The converse is false since it is easy to find a set  $\mathcal{F}$  with finite  $u$ -covering entropy and infinite  $u$ -bracketing entropy.

Finally, we will say that a family of  $u_N$ -nets,  $N \in \mathbb{N}$ , is almost minimal when the log-cardinal of the size of the  $u_N$ -net has the same order as the  $(u_N, \mathcal{F}, \mathbb{P}, \cdot)$ -entropy.

The paper is organized as follows. Section 2 recalls some PAC-Bayesian concentration inequalities which are extracted from [1]. In Section 3, we assume that we have Tsybakov's margin assumption and that the  $\mathbb{P}, \cdot$ -entropies are polynomial. In this setting, we study the convergence rate of standard Gibbs estimators and classifiers minimizing the empirical risk on  $\mathbb{P}, \cdot$ -covering nets. In particular, we stresses on the influence of the chaining trick and the differences between bracketing and covering entropy assumptions. Section 4 tries to answer the questions: what happens when we relieve the polynomial  $\mathbb{P}, \cdot$ -entropy assumption? Can we give an empirical equivalent (i.e. with  $\bar{\mathbb{P}}, \cdot$ -entropies) of the previous results? Section 5 gives a version of Assouad's lemma dedicated to classification. The proofs are gathered in Section 6.

## 2. KNOWN PAC-BAYESIAN BOUNDS

In this section, we recall some results of [1] which will be useful in this paper.

**Theorem 2.1.** *Let  $g(u) \triangleq \frac{\exp(u)-1-u}{u^2}$  for any  $u > 0$ . For any  $\lambda > 0$ ,  $\epsilon > 0$  and  $\pi_1, \pi_2 \in \mathcal{M}_+^1(\mathcal{F})$ , with  $(\mathbb{P}^{\otimes N})_*$ -probability at least  $1 - \epsilon$ , for any  $\rho_1, \rho_2 \in \mathcal{M}_+^1(\mathcal{F})$ , we have*

$$(2.1) \quad \rho_2 R - \rho_1 R + \rho_1 r - \rho_2 r \leq \frac{\lambda}{N} g\left(\frac{\lambda}{N}\right) (\rho_1 \otimes \rho_2) \mathbb{P}, \cdot + \frac{K(\rho_1, \pi_1) + K(\rho_2, \pi_2) + \log(\epsilon^{-1})}{\lambda}$$

As a consequence, with  $(\mathbb{P}^{\otimes N})_*$ -probability at least  $1 - \epsilon$ , for any  $\rho_1, \rho_2 \in \mathcal{M}_+^1(\mathcal{F})$ ,

$$(2.2) \quad \rho_2 R - \rho_1 R + \rho_1 r - \rho_2 r \leq \min_{\lambda \in [\sqrt{N}, N]} \left\{ 0.8 \frac{\lambda}{N} (\rho_1 \otimes \rho_2) \mathbb{P}, \cdot + 1.7 \frac{K(\rho_1, \pi_1) + K(\rho_2, \pi_2) + \log[\log(eN)\epsilon^{-1}]}{\lambda} \right\}.$$

Besides, let  $\mathcal{S}_1$  and  $\mathcal{S}_2$  be finite subsets of  $\mathcal{F}$ , with  $(\mathbb{P}^{\otimes N})_*$ -probability at least  $1 - \epsilon$ , for any  $(f_1, f_2) \in \mathcal{S}_1 \times \mathcal{S}_2$ , we have

$$(2.3) \quad R(f_2) - R(f_1) + r(f_1) - r(f_2) \leq \sqrt{\frac{2 \log(|\mathcal{S}_1| |\mathcal{S}_2| \epsilon^{-1}) \mathbb{P}_{f_1, f_2}}{N}} + \frac{\log(|\mathcal{S}_1| |\mathcal{S}_2| \epsilon^{-1})}{3N}.$$

*Proof.* The first part comes from Theorem 4.8 in [1]. Then the second part is obtained by a union bound on the set of parameters  $\Lambda \triangleq \{\sqrt{N} e^{k/2}; 0 \leq k \leq \log N\}$  (see Section 4.2 in [1] for details). The third part comes from Theorem 8.1 in [1] applied to  $\mathcal{W}[(f_1, f_2), Z] = \mathbb{1}_{Y \neq f_2(X)} - \mathbb{1}_{Y \neq f_1(X)}$  and  $\nu$  equal to the uniform measure on  $\mathcal{S}_1 \times \mathcal{S}_2$ .  $\square$

The following theorem ([1, Theorem 6.4]) brackets the efficiency of a standard Gibbs classifier

**Theorem 2.2.** *For any  $\lambda > 0$  and  $0 < \chi \leq 1$ , we have*

$$\pi_{-(1+\chi)\lambda R} R - \frac{K(\pi_{-\lambda r}, \pi_{-\lambda R})}{\chi\lambda} \leq \pi_{-\lambda r} R \leq \pi_{-(1-\chi)\lambda R} R + \frac{K(\pi_{-\lambda r}, \pi_{-\lambda R})}{\chi\lambda},$$

and for any  $\epsilon > 0$ ,  $0 < \gamma < \frac{1}{2}$  and  $0 < \lambda \leq 0.39\gamma N$ , with  $\mathbb{P}^{\otimes N}$ -probability at least  $1 - \epsilon$ , we have

$$(2.4) \quad \begin{aligned} K(\pi_{-\lambda r}, \pi_{-\lambda R}) &\leq \frac{4}{1-\gamma} \log \pi_{-\lambda R} \exp\left(\frac{4.1\lambda^2}{\gamma N} \pi_{-\lambda R} \mathbb{P}_{\cdot, \cdot}\right) + \frac{5\gamma}{1-\gamma} \log(4\epsilon^{-1}) \\ &\leq 16 \log \pi_{-\lambda R} \exp\left(\frac{4.1\lambda^2}{\gamma N} \mathbb{P}_{\cdot, \tilde{f}}\right) + 10\gamma \log(4\epsilon^{-1}). \end{aligned}$$

### 3. CONVERGENCE RATE OF CLASSIFIERS UNDER COMPLEXITY AND MARGIN ASSUMPTIONS

**3.1. Complexity and margin assumptions.** The following assumptions have the same form as the one used in the pioneering work of Mammen and Tsybakov ([11]). The margin assumption appears to be the key assumption to obtain fast rates of convergence (i.e.  $N^{-\beta}$  with  $\beta > \frac{1}{2}$ ).

**3.1.1. Complexity assumptions.** Let  $q \geq 0$ . Define

$$h_q(u) \triangleq \begin{cases} \log(eu^{-1}) & \text{when } q = 0 \\ u^{-q} & \text{when } q > 0 \end{cases}.$$

We will alternatively use the following complexity assumptions.

(CA1) : there exists  $C' > 0$  such that the covering entropy of the model  $\mathcal{F}$  for the distance  $\mathbb{P}_{\cdot, \cdot}$  satisfies for any  $u > 0$ ,  $H(u, \mathcal{F}, \mathbb{P}_{\cdot, \cdot}) \leq C' h_q(u)$ .

(CA2) : there exists  $C' > 0$  such that the bracketing entropy of the model  $\mathcal{F}$  for the distance  $\mathbb{P}_{\cdot, \cdot}$  satisfies for any  $u > 0$ ,  $H^{[1]}(u, \mathcal{F}, \mathbb{P}_{\cdot, \cdot}) \leq C' h_q(u)$ .

(CA3) : there exist  $C' > 0$  and  $\pi \in \mathcal{M}_+^1(\mathcal{F})$  such that for any  $t > 0$ , for any  $f' \in \mathcal{F}$ , we have  $\pi(\mathbb{P}_{\cdot, f'} \leq t) \geq \exp[-C' h_q(t)]$ .

We have<sup>4</sup>: (CA2)  $\Rightarrow$  (CA1)  $\Leftrightarrow$  (CA3). Let  $t$  and  $C'$  be positive reals. We will say that a probability distribution  $\pi$  satisfies  $(t, C')$ -(CA3) when we have

$$\pi(\mathbb{P}_{\cdot, \tilde{f}} \leq t) \geq \exp[-C' h_q(t)].$$

Note that this last assumption is, unlike the others, a local complexity assumption.

**3.1.2. Margin assumptions.** We will consider variants of Tsybakov's margin assumption ([11, 15]). Let  $\alpha \in \mathbb{R}_+ \cup \{+\infty\}$  and  $\kappa \in [1; +\infty]$ . We define

$$\Delta R(f) \triangleq R(f) - R(\tilde{f}).$$

(MA1) :  $\mathcal{Y} = \{0; 1\}$  and there exists  $C'' > 0$  such that for any  $t > 0$ ,

$$\mathbb{P}(0 < |\eta^*(X) - 1/2| \leq t) \leq C'' t^\alpha.$$

(MA2) : there exists  $C'' > 0$  such that for any function  $f \in \mathcal{F}$ ,

$$\mathbb{P}_{f, \tilde{f}} \leq C'' [\Delta R(f)]^{\frac{1}{\kappa}}.$$

(MA3) : there exist  $c'', C'' > 0$  such that for any function  $f \in \mathcal{F}$ ,

$$(3.1) \quad c'' [\Delta R(f)]^{\frac{1}{\kappa}} \leq \mathbb{P}_{f, \tilde{f}} \leq C'' [\Delta R(f)]^{\frac{1}{\kappa}}.$$

<sup>4</sup>To prove (CA1)  $\Rightarrow$  (CA3): for any  $k \in \mathbb{N}^*$ , introduce  $\pi_k$  the uniform distribution on a  $(2^{-k}, \mathcal{F}, \mathbb{P}_{\cdot, \cdot})$ -minimal covering net. The prior distribution  $\pi \triangleq \sum_{k \geq 1} \frac{\pi_k}{k(k+1)}$  satisfies the claim.



(MA4) : there exist  $c'', C'' > 0$  such that  $\mathbb{P}_{\cdot, \tilde{f}} \leq C''[\Delta R]^{\frac{1}{\kappa}}$ , and for any  $t > 0$ ,  $\pi(\Delta R \leq t) \geq c''\pi(\mathbb{P}_{\cdot, \tilde{f}} \leq C''t^{\frac{1}{\kappa}})^5$ .

This last assumption makes sense only in the bayesian context where a prior distribution  $\pi$  is put on the model. It is easy to check the following implications: (MA3)  $\Rightarrow$  (MA4)  $\Rightarrow$  (MA2). Besides when  $f^* \in \mathcal{F}$  (= no bias assumption), we have: (MA1)  $\Rightarrow$  (MA2) for  $\kappa = \frac{1+\alpha}{\alpha}$ . When  $\kappa = +\infty$ , Assumption (MA2) is empty<sup>6</sup> and Assumptions (MA3) and (MA4) are not satisfied by non-trivial models. The margin Assumptions (MA1) and (MA2) are all the stronger as  $\kappa$  is small. When  $\kappa = 1$ , the lower bound in inequality (3.1) holds trivially for  $c'' = 1$  and we have: (MA3)  $\Leftrightarrow$  (MA4)  $\Leftrightarrow$  (MA2).

*Remark 3.1.* For sake of simplicity, we have assumed that there exists a function  $\tilde{f} \in \mathcal{F}$  such that  $R(\tilde{f}) = \inf_{\mathcal{F}} R$ . Then, under Assumption (MA2), this function needs to be unique. In fact this is not more necessary than the existence of the minimum. To be more specific, the results in this paper under Assumption (MA2) will still hold when this assumption is replaced with: there exists  $k \in \mathbb{N}^*$  such that for any  $\beta > 0$ , there exists  $f_1, \dots, f_k \in \mathcal{F}$

$$\forall f \in \mathcal{F}, \exists i \in \{1, \dots, k\}, \mathbb{P}_{f, f_i} \leq C''[R(f) - \inf_{\mathcal{F}} R]^{\frac{1}{\kappa}} + \beta.$$

Note that this implies that for any  $i \in \{1, \dots, k\}$ ,  $R(f_i) - \inf_{\mathcal{F}} R \leq \beta$ . Similarly, we can give weakened versions of Assumptions (MA3) and (MA4). Naturally, the value of  $k$  will influence the value of the constants in the results under Assumption (MA2).

## 3.2. Gibbs classifier.

3.2.1. *Under Assumptions (MA4) and (CA3) for  $q > 0$ .* In this paper, we will often consider prior distributions  $\pi^{(N)}$  which may depend on  $N$ . To shorten, we will simply write it  $\pi$ . The following lemma guarantees the efficiency of the standard Gibbs estimator for a temperature appropriately chosen.

**Lemma 3.1.** *Let  $\pi$  be a probability distribution such that*

$$(3.2) \quad \pi[\Delta R \leq \check{C}_1 N^{-\frac{\kappa}{2\kappa-1+q}}] \geq e^{-\check{C}_2 N^{\frac{q}{2\kappa-1+q}}}$$

and  $\lambda_N$  have the same order as  $N^{\frac{\kappa+q}{2\kappa-1+q}}$ , i.e. such that

$$(3.3) \quad \check{C}_3 N^{\frac{\kappa+q}{2\kappa-1+q}} \leq \lambda_N \leq \check{C}_4 N^{\frac{\kappa+q}{2\kappa-1+q}}$$

for some positive constants  $\check{C}_i, i = 1, \dots, 4$ . Then, under Assumption (MA2), the standard Gibbs classifier in which the prediction function is drawn according to the posterior distribution  $\pi_{-\lambda_N r}$  has the convergence rate  $N^{-\frac{\kappa}{2\kappa-1+q}}$  to the extent that

$$\mathbb{P}^{\otimes N} \pi_{-\lambda_N r} R - R(\tilde{f}) \leq \check{C} N^{-\frac{\kappa}{2\kappa-1+q}}$$

for some constant  $\check{C} > 0$  (depending only on  $c'', C'', \kappa$  and  $\check{C}_i, i = 1, \dots, 4$ ).

More precisely, with  $\mathbb{P}^{\otimes N}$ -probability at least  $1 - \epsilon$ , with  $\pi_{-\lambda_N r}$ -probability at least  $1 - \epsilon$ , we have

$$(3.4) \quad R - R(\tilde{f}) \leq \check{C} \log(e\epsilon^{-1}) N^{-\frac{\kappa}{2\kappa-1+q}},$$

for some constant  $\check{C} > 0$  (depending on  $C'', \kappa$  and  $\check{C}_i, i = 1, \dots, 4$ ).

<sup>5</sup>As a consequence,  $\pi(\Delta R \leq t)$  has the same order as  $\pi(\mathbb{P}_{\cdot, \tilde{f}} \leq C''t^{\frac{1}{\kappa}})$ .

<sup>6</sup>since the inequality trivially holds for  $C'' = 1$

*Proof.* See Section 6.1.  $\square$

**Theorem 3.2.** *Let  $\pi$  be a distribution satisfying Assumptions (MA4) and  $(\check{C}_1 N^{-\frac{1}{2\kappa-1+q}}, \check{C}_2)$ -(CA3) for  $q > 0$  and  $\lambda_N$  be a real satisfying inequality (3.3) for given positive constants  $\check{C}_i, i = 1, \dots, 4$ . Then we have*

$$\mathbb{P}^{\otimes N} \pi_{-\lambda_N r} R - R(\tilde{f}) \leq \check{C} N^{-\frac{\kappa}{2\kappa-1+q}}$$

for some constant  $\check{C} > 0$  (depending only on  $c', C'', \kappa$  and  $\check{C}_i, i = 1, \dots, 4$ ).

*Proof.* It suffices to check that, under these assumptions, we can apply Lemma 3.1.  $\square$

*Remark 3.2.* To make the link with previous works about non randomized sieve estimators, one can choose  $\pi$  as the uniform distribution on an almost minimal  $(\check{C}_1 N^{-\frac{1}{2\kappa-1+q}}, \mathcal{F}, \mathbb{P}, \cdot)$ -covering net. Then Assumption (CA1) implies that the distribution  $\pi$  satisfies Assumption  $(\check{C}_1 N^{-\frac{1}{2\kappa-1+q}}, \check{C}_2)$ -(CA3) for some constant  $\check{C}_2 > 0$  (depending on  $\check{C}_1$ , on the almost minimality constant and on the constant  $C'$  involved in (CA1)). Note that, as in Mammen and Tsybakov's work ([11, 15]), the computation of the estimator requires that, without knowing  $\mathbb{P}(dX)$  exactly, one can construct a  $(t, \mathcal{F}, \mathbb{P}, \cdot)$ -net with log-cardinality of order  $H(t, \mathcal{F}, \mathbb{P}, \cdot)$ .

The convergence rate of the standard Gibbs estimator in Theorem 3.2 is optimal since the following lower bound holds.

**Theorem 3.3.** *Let  $q \geq 0$  and  $\kappa \in [1; +\infty]$ . There exist an input space  $(\mathcal{X}, \mathcal{B}_{\mathcal{X}})$ , a model  $\mathcal{F}$  and a set  $\mathcal{P}$  be the set of probability distributions satisfying*

- for any  $\mathbb{P} \in \mathcal{P}$ ,  $f_{\mathbb{P}}^* \in \mathcal{F}$
- Assumptions (CA2), (MA3) and (MA1) with  $\alpha = \frac{1}{\kappa-1} \in [0; +\infty]$

such that for any measurable estimator  $\hat{f} : \mathcal{Z}^N \rightarrow \mathcal{F}(\mathcal{X}, \mathcal{Y})$ ,

$$\sup_{\mathbb{P} \in \mathcal{P}} \{ \mathbb{P}^{\otimes N} R(\hat{f}) - R(\tilde{f}) \} \geq C N^{-\frac{\kappa}{2\kappa-1+q}}.$$

*Proof.* See Section 6.2.  $\square$

*Remark 3.3.* In [15], the same result is proved (using the classes of boundary fragments) for the set of probability distributions such that the Bayes classifier is in the model and Assumptions (CA2) and (MA1) with  $\alpha = \frac{1}{\kappa-1}$  hold.

*Remark 3.4.* The previous theorem is stronger than what is required to prove that the convergence rate obtained in Theorem 3.2 is optimal since the set  $\mathcal{P}$  in Theorem 3.3 is smaller than the set of probability distributions  $\mathbb{P}$  such that there exists a distribution  $\pi$  satisfying Assumptions (MA4) and  $(\check{C}_1 N^{-\frac{1}{2\kappa-1+q}}, \check{C}_2)$ -(CA3).

3.2.2. *Under Assumptions (MA2) and (CA3) for  $q = 0$ .* Using the same tools as in the previous section, we can prove

**Theorem 3.4.** *Let  $\pi$  be a distribution satisfying Assumption  $(\check{C}_1 N^{-\frac{\kappa}{2\kappa-1}}, \check{C}_2)$ -(CA3) for  $q = 0$  and  $\lambda_N$  be a real satisfying*

$$(3.5) \quad \check{C}_3 N^{\frac{\kappa}{2\kappa-1}} \leq \lambda_N \leq \check{C}_4 N^{\frac{\kappa}{2\kappa-1}}$$

for given positive constants  $\check{C}_i, i = 1, \dots, 4$ . Under Assumption (MA2), we have

$$\mathbb{P}^{\otimes N} \pi_{-\lambda_N r} R - R(\tilde{f}) \leq \check{C} (\log N) N^{-\frac{\kappa}{2\kappa-1}}$$

for some constant  $\check{C} > 0$  (depending only on  $C'', \kappa$  and  $\check{C}_i, i = 1, \dots, 4$ ).

*Proof.* We use Lemma 6.1 and the inequalities

$$\begin{cases} T_2(\pi) & \leq \check{C}\lambda\left(\frac{\lambda}{N}\right)^{\frac{\kappa}{\kappa-1}} \\ T_1(\pi) & \leq -\log\left[\pi\left(\Delta R \leq \check{C}_1 N^{-\frac{\kappa}{2\kappa-1}}\right)\right] + \lambda\check{C}_1 N^{-\frac{\kappa}{2\kappa-1}} \\ & \leq -\log\left[\pi\left(\mathbb{P}_{\cdot, \tilde{f}} \leq \check{C}_1 N^{-\frac{\kappa}{2\kappa-1}}\right)\right] + \lambda\check{C}_1 N^{-\frac{\kappa}{2\kappa-1}} \end{cases} .$$

□

From Theorem 3.3, since we have  $(CA2) \Rightarrow (CA3)$  and  $(MA3) \Rightarrow (MA2)$ , this convergence rate is optimal up to the  $\log N$  factor.

3.2.3. *Under Assumption (MA2) and a local complexity assumption.* The following theorem considers a local complexity assumption and its first and second parts respectively complete Theorem 3.4 and Lemma 3.1.

**Theorem 3.5.** *Let  $\epsilon > 0$ ,  $s \geq 0$ ,  $C' > 0$ ,  $C''' \in \mathbb{R}$  and  $1 \leq \kappa \leq +\infty$ . Consider  $\lambda$  depending on  $N$  such that  $\lambda \xrightarrow{N \rightarrow +\infty} +\infty$  and that Assumption (MA2) holds.*

*First, assume that  $\log \pi^{-1}\{R - R(\tilde{f}) \leq x\} = -C' \log x + C''' + \underset{x \rightarrow 0}{o}(x^s)$ . Then we have*

- for  $\lambda = \underset{N \rightarrow +\infty}{o}\left(N^{\frac{\kappa}{2\kappa-1}}\right)$ , with  $\mathbb{P}^{\otimes N}$ -probability at least  $1 - \epsilon$ ,

$$\pi_{-\lambda r} R = \frac{C'}{\lambda} + \frac{1}{\lambda} \left\{ \underset{N \rightarrow +\infty}{o}\left(\lambda^{-\frac{s}{2}}\right) + \underset{N \rightarrow +\infty}{O}\left(\lambda^{\frac{(2\kappa-1)C'}{2(2\kappa C' + \kappa - 1)}} N^{-\frac{\kappa C'}{2(2\kappa C' + \kappa - 1)}}\right) \log(\epsilon \epsilon^{-1}) \right\}$$

- when  $\lambda = cN^{\frac{\kappa}{2\kappa-1}}$ : for any  $\beta > 0$ , there exist  $c > 0$  and  $N_0 > 0$  such that for any  $N > N_0$ , with  $\mathbb{P}^{\otimes N}$ -probability at least  $1 - \epsilon$ , :

$$\frac{C' - \beta}{\lambda} \leq \pi_{-\lambda r} R \leq \frac{C' + \beta}{\lambda}.$$

*Secondly, assume that  $\log \pi^{-1}\{R - R(\tilde{f}) \leq x\} = C' x^{-\frac{q}{\kappa}} + C''' + \underset{x \rightarrow 0}{o}(1)$  with  $q > 0$  and  $\kappa \neq +\infty$ . Then we have*

- for  $\lambda = \underset{N \rightarrow +\infty}{o}\left(N^{-\frac{\kappa+q}{2\kappa-1+q}}\right)$ , with  $\mathbb{P}^{\otimes N}$ -probability at least  $1 - \epsilon$ ,

$$\pi_{-\lambda r} R = \left(\frac{qC' + o(1)}{\kappa\lambda}\right)^{\frac{\kappa}{\kappa+q}} + \underset{N \rightarrow +\infty}{O}\left(1\right) \frac{\log(\epsilon^{-1})}{\lambda}$$

- when  $\lambda = cN^{-\frac{\kappa+q}{2\kappa-1+q}}$ : for any  $0 < \beta \leq qC'$ , there exist  $c > 0$  and  $N_0 > 0$  such that for any  $N > N_0$ , with  $\mathbb{P}^{\otimes N}$ -probability at least  $1 - \epsilon$ , :

$$\left(\frac{qC' - \beta}{\kappa\lambda}\right)^{\frac{\kappa}{\kappa+q}} \leq \pi_{-\lambda r} R \leq \left(\frac{qC' + \beta}{\kappa\lambda}\right)^{\frac{\kappa}{\kappa+q}}.$$

*Proof.* See Section 6.3. □

It is interesting to note that this asymptotic behaviour only depends on the local complexity given by the weight of the sets  $\{f \in \mathcal{F} : R(f) - R(\tilde{f}) \leq x\}$  when  $x \rightarrow 0$ . Had we had  $\mathbb{P}_{f, \tilde{f}} \underset{x \rightarrow 0}{\sim} C'' [R - R(\tilde{f})]^{\frac{1}{\kappa}}$  on these sets, the complexity assumption would be similar to the ones introduced in Section 3.1.1 to the extent that we would have  $\log \pi^{-1}(\mathbb{P}_{\cdot, \tilde{f}} \leq x) \underset{x \rightarrow 0}{\sim} \check{C} h_q(x)$ .

In Theorems 3.2 and 3.4, we have seen how to choose the parameter  $\lambda$  depending on  $N$  such that the Gibbs classifier has the optimal convergence rate. The previous result shows that for  $\lambda$  smaller than these “optimal” parameters and a slightly modified complexity assumption, we can tightly bracket the efficiency of standard

Gibbs classifiers. For larger  $\lambda$ , the picture is not clear: it seems that the KL-divergence term in Theorem 2.2 becomes the leading term. This KL-divergence will in general explode for  $\lambda \gg N$ , and finally we just know that

$$\pi_{-\lambda r} R \xrightarrow{\lambda \rightarrow +\infty} \pi|_{r=\min_{\mathcal{F}} r} R \triangleq \frac{\int_{\mathcal{F}} \mathbb{1}_{r(f)=\min_{\mathcal{F}} r} R(f) d\pi(f)}{\int_{\mathcal{F}} \mathbb{1}_{r(f)=\min_{\mathcal{F}} r} d\pi(f)}.$$

*Remark 3.5.* The confidence level  $\epsilon$  does not appear in the main terms of the expansions of  $\pi_{-\lambda r} R$ , hence the asymptotic orders of  $\pi_{-\lambda r} R$  hold with exponential probability.

**3.2.4. Adaptive choice of the temperature.** Here we consider that Assumption (MA3) holds for an unknown margin parameter  $\kappa$  and we prove that under assumption (CA3) a standard Gibbs classifier with an appropriately chosen temperature is adaptive wrt this parameter, i.e. without prior knowledge of  $\kappa$ , the generalization error of the randomized estimator is upper bounded by  $\check{C} N^{\frac{\kappa}{2\kappa-1+q}}$  when  $q > 0$  and by  $\check{C}(\log N) N^{\frac{\kappa}{2\kappa-1}}$  when  $q = 0$ . The adaptation to the margin problem has also been studied in [15, 16]. In particular, in [16], Tsybakov and van de Geer proposed an adaptive penalized classifier using wavelets.

**Theorem 3.6.** *Under Assumptions (MA3) and (CA3), the algorithm given in Section 3.4.2 of [1] achieves an adaptive choice of the temperature of the standard Gibbs classifier wrt the margin parameter  $\kappa$ .*

*Proof.* See Section 6.4. □

### 3.3. Empirical risk minimization on nets.

**3.3.1. Under Assumptions (MA3) and (CA1) for  $q > 0$ .** This section shows that, by using inequality (2.1), we can recover results on sieve estimators given in [11, 15]. These results have to be compared with the ones in Section 3.2.1 (recall that (MA3)  $\Rightarrow$  (MA4) and (CA1)  $\Leftrightarrow$  (CA3)).

**Theorem 3.7.** *Under Assumptions (MA3) and (CA1) for  $q > 0$ , for any classifier  $\hat{f}$  minimizing the empirical risk among a  $u_N$ -covering net  $\mathcal{N}_{u_N}$  such that*

$$(3.6) \quad \check{C}_1 N^{-\frac{1}{2\kappa-1+q}} \leq u_N \leq \check{C}_2 N^{-\frac{1}{2\kappa-1+q}}$$

and

$$(3.7) \quad \log |\mathcal{N}_{u_N}| \leq \check{C}_3 h_q(u_N)$$

for some positive constants  $\check{C}_i, i = 1, \dots, 3$ , we have

$$\mathbb{P}^{\otimes N} [R(\hat{f}) - R(\tilde{f})] \leq \check{C} N^{-\frac{\kappa}{2\kappa-1+q}}$$

for some constant  $\check{C} > 0$  (depending only on  $C', c'', C''$  and  $\check{C}_i, i = 1, \dots, 3$ ).

*Proof.* See Section 6.5. □

*Remark 3.6.* Inequality (3.7) just says that the net  $\mathcal{N}_{u_N}$  is almost minimal.

3.3.2. *Under Assumptions (MA2) and (CA1) for  $q = 0$ .* Since  $(CA1) \Leftrightarrow (CA3)$ , this section gives results to be compared with the ones in Section 3.2.2.

**Theorem 3.8.** *Under Assumptions (MA2) and (CA1) for  $q = 0$ , for any classifier  $\hat{f}$  minimizing the empirical risk among a  $u_N$ -covering net  $\mathcal{N}_{u_N}$  such that*

$$(3.8) \quad N^{-\check{C}_1} \leq u_N \leq \check{C}_2(\log N)N^{-\frac{\kappa}{2\kappa-1}}$$

and

$$(3.9) \quad \log |\mathcal{N}_{u_N}| \leq \check{C}_3 h_0(u_N)$$

for some positive constants  $\check{C}_i, i = 1, \dots, 3$ , we have

$$\mathbb{P}^{\otimes N} [R(\hat{f}) - R(\tilde{f})] \leq \check{C}(\log N)N^{-\frac{\kappa}{2\kappa-1}}$$

for some constant  $\check{C} > 0$  (depending only on  $c'', C''$  and  $\check{C}_i, i = 1, \dots, 3$ ).

*Proof.* It follows the lines of Section 6.5. This time, we take  $(\frac{\lambda}{N})^{\frac{\kappa}{\kappa-1}}$  and  $\frac{\log(eu^{-1})}{\lambda}$  of the same order and greater than  $u$ . This is realized when inequality (3.8) is satisfied and  $\lambda = N^{\frac{\kappa}{2\kappa-1}}$ .  $\square$

**3.4. Chaining.** When a class of functions has a polynomial entropy, there is a trick called the chaining ([6]) which allows us to improve the previous results. This technique is used to get tighter upper bounds of the difference  $R(f_1) - R(f_2)$  between the expected risk at two different functions  $f_1$  and  $f_2$ . It is based on finer and finer approximations of these functions. The advantage of considering rough approximation of these functions is that the set of all possible rough approximations is small (in other words, has a small complexity). On the contrary, the set of fine approximations is big, but the distance between the fine approximation and the function approximated is small. So there is a kind of bias/variance trade-off and for polynomial entropy classes of functions, it is interesting to have this trade-off on a sequence of links and not directly on the big link  $f_1 \cdots f_2$ .

Let us give some results due to this technique. Consider the context of Theorem 3.7. Let us see what happens if we replace the margin Assumption (MA3) with Assumption (MA2). Then we can no longer upper bound  $\Delta R$  with  $\text{Cst } \mathbb{P}_{\cdot, \tilde{f}}^{\kappa}$  (inequality which is used to obtain (6.8)). We only have  $\Delta R \leq \mathbb{P}_{\cdot, \tilde{f}}$ . This leads to the convergence rate  $N^{-\frac{\kappa}{2\kappa-1+q\kappa}}$  instead of  $N^{-\frac{\kappa}{2\kappa-1+q}}$ . Using the chaining trick, we will prove (see Theorem 3.10) that this rate is suboptimal and that, by minimizing the empirical risk on well chosen nets, we can still reach the rate  $N^{-\frac{\kappa}{2\kappa-1+q}}$  when  $0 < q < 1$  and the rate  $N^{-\frac{1}{1+q}}$  when  $q > 1$ .

*Remark 3.7.* The convergence rate  $N^{-\frac{\kappa}{2\kappa-1+q\kappa}}$  is optimal under Assumption (MA2) and the complexity assumption  $H(u_N) \leq C' h_q(u_N)$  for the radius  $u_N = N^{-\frac{\kappa}{2\kappa-1+q\kappa}}$ . The lower bound comes from Lemma 5.1 applied to a  $(N^{\frac{q\kappa}{2\kappa-1+q\kappa}}, N^{-\frac{1+q\kappa}{2\kappa-1+q\kappa}}, \frac{1}{2}N^{-\frac{\kappa-1}{2\kappa-1+q\kappa}})$ -constant hypercube. By slightly modifying the proof of Theorem 3.7, we can obtain that, under the previous margin and complexity assumptions, any classifier  $\hat{f}$  minimizing the empirical risk among a  $u_N$ -almost minimal net satisfies

$$\mathbb{P}^{\otimes N} [R(\hat{f}) - R(\tilde{f})] \leq \check{C} N^{-\frac{\kappa}{2\kappa-1+q\kappa}}.$$

The chaining technique appears to be the only tool which allows to take into account an entropy assumption which holds for any radius such as (CA1) and (CA2).

The chaining trick may also be used to prove that the empirical risk can be minimized on tighter nets (provided that they are still minimal, or at least almost minimal, under polynomial entropy assumptions).

Before giving results concerning nets, one can illustrate the chaining technique by considering randomized posteriors concentrated on small balls of fixed radius. For any  $u > 0$ , introduce  $B_{f,u} \triangleq \{f' \in \mathcal{F} : \mathbb{P}_{f,f'} \leq u\}$  and  $\pi_{f,u} \triangleq \frac{\mathbb{1}_{B_{f,u}}}{\pi(B_{f,u})} \cdot \pi$ . Define  $h(v) \triangleq \sup_{f \in \mathcal{F}} \log \pi^{-1}(B_{f,v})$  and  $h_+(v) \triangleq h(v) \vee 1$ .

**Theorem 3.9.** *Let  $u > 0$ ,  $L \triangleq \frac{\log(2/u)}{\log 2}$ ,  $C_1 \triangleq \sqrt{\frac{4h_+(u)}{3Nu}}$  and  $C_0 \triangleq 2\sqrt{3}[1 + 2g(C_1)]$ . For any  $\epsilon > 0$ , with  $\mathbb{P}^{\otimes N}$ -probability at least  $1 - \epsilon$ , for any  $f_1, f_2 \in \mathcal{F}$  such that  $\mathbb{P}_{f_1, f_2} > u$ , we have*

$$\begin{aligned} \pi_{f_2, u} R - \pi_{f_1, u} R + \pi_{f_1, u} r - \pi_{f_2, u} r \\ \leq \frac{C_0}{\sqrt{N}} \sum_{k \in \mathbb{N}: u2^k < \mathbb{P}_{f_1, f_2}} \sqrt{u2^k h_+(u2^k)} + 6\sqrt{\frac{\mathbb{P}_{f_1, f_2}}{N}} \log[L\epsilon^{-1}] \\ \leq \frac{2C_0}{\sqrt{N}} \int_{u/2}^{\mathbb{P}_{f_1, f_2}} \sqrt{\frac{h_+(v)}{v}} dv + 6\sqrt{\frac{\mathbb{P}_{f_1, f_2}}{N}} \log[L\epsilon^{-1}]. \end{aligned}$$

*Proof.* See Section 6.6. □

Had we not chained inequality (2.1), we would have obtained

$$\pi_{f_2, u} R - \pi_{f_1, u} R + \pi_{f_1, u} r - \pi_{f_2, u} r \leq \frac{\lambda}{N} g\left(\frac{\lambda}{N}\right) (\mathbb{P}_{f_1, f_2} + 2u) + \frac{2h(u) + \log(\epsilon^{-1})}{\lambda}$$

This upper bound is greater than  $\inf_{\lambda > 0} \left\{ \frac{\lambda}{2N} \mathbb{P}_{f_1, f_2} + \frac{2h(u)}{\lambda} \right\} = 2\sqrt{\frac{\mathbb{P}_{f_1, f_2} h(u)}{N}}$ , which is much bigger than the chained bound for polynomial entropies  $h(u) \approx u^{-q}$ ,  $q > 0$  when<sup>7</sup>  $N^{-\frac{1}{1+q}} \leq u \ll \mathbb{P}_{f_1, f_2}$ .

The following result is an extension of Theorems 3.7 and 3.8.

**Theorem 3.10.** *We assume that Assumptions (MA2) and (CA1) hold. When Assumption (MA3) also holds, we define*

$$(v_N, a_N) \triangleq \begin{cases} \left( \left[ \frac{\log N}{N} \right]^{\frac{\kappa}{2\kappa-1}}, \exp \left\{ -\check{C}_1 (\log N)^{\frac{\kappa}{4\kappa-2}} N^{\frac{\kappa-1}{4\kappa-2}} \right\} \right) & \text{for } q = 0 \\ \left( N^{-\frac{\kappa}{2\kappa-1+q}}, \check{C}_1 N^{-\frac{(\kappa-1)q \leq 1+q}{q(2\kappa-1+q)}} \right) & \text{for } q > 0 \end{cases}$$

and  $b_N \triangleq \check{C}_2 (v_N)^{\frac{1}{\kappa}}$ .

When Assumption (MA3) does not hold, we define

$$(v_N, a_N) \triangleq \begin{cases} \left( \left[ \frac{\log(eN^{1/\kappa})}{N} \right]^{\frac{\kappa}{2\kappa-1}}, \exp \left\{ -\check{C}_1 (\log[eN^{1/\kappa}])^{\frac{\kappa}{4\kappa-2}} N^{\frac{\kappa-1}{4\kappa-2}} \right\} \right) & \text{for } q = 0 \\ \left( N^{-\frac{\kappa}{2\kappa-1+q}}, \check{C}_1 N^{-\frac{\kappa-1+q}{q(2\kappa-1+q)}} \right) & \text{for } 0 < q < 1 \\ \left( (\log N) N^{-\frac{1}{2}}, \check{C}_1 (\log N)^{-\frac{1}{2}} N^{-\frac{1}{2}} \right) & \text{for } q = 1 \\ \left( N^{-\frac{1}{1+q}}, \check{C}_1 N^{-\frac{1}{1+q}} \right) & \text{for } q > 1 \end{cases}$$

and  $b_N \triangleq \check{C}_2 v_N$ .

<sup>7</sup>The quantity  $C_0$  behaves as a constant only when  $\frac{h_+(u)}{Nu} \leq C$ , so when  $u \geq CN^{-\frac{1}{1+q}}$ .

For any classifier minimizing the empirical risk among a  $u_N$ -covering net  $\mathcal{N}_{u_N}$  such that

$$(3.10) \quad a_N \leq u_N \leq b_N$$

and

$$(3.11) \quad \log |\mathcal{N}_{u_N}| \leq \check{C}_3 h_q(u_N)$$

for some positive constants  $\check{C}_i, i = 1, \dots, 3$ , we have

$$\mathbb{P}^{\otimes N}[R(\hat{f}) - R(\tilde{f})] \leq C v_N$$

for some constant  $C > 0$  (depending on  $C''$ ,  $\check{C}_i, i = 1, \dots, 3$  [and also on  $c''$  under Assumption (MA3)]).

*Proof.* See Section 6.7. □

*Remark 3.8.* When  $q = 0$  and  $\kappa = +\infty$  (i.e. no margin assumption), the  $\log N$  factor in  $\log(eN^{1/\kappa})$  disappears. The suppression of the logarithmic factor, obtained by chaining, is similar to what occurs for VC classes (see Corollary 4.6). The difference is just that the complexity assumption concerns  $\mathbb{P}$ -nets here instead of empirical nets.

From Theorems 3.3 and the following theorem, these convergence rates are optimal (up to the logarithmic factor when we have  $q \in \{0; 1\}$ ).

**Theorem 3.11.** *There exist an input space  $(\mathcal{X}, \mathcal{B}_{\mathcal{X}})$ , a model  $\mathcal{F}$  and a set  $\mathcal{P}$  of probability distributions satisfying Assumptions (CA2) and (MA2) such that for any measurable estimator  $\hat{f} : \mathcal{Z}^N \rightarrow \mathcal{F}(\mathcal{X}, \mathcal{Y})$ ,*

$$\sup_{\mathbb{P} \in \mathcal{P}} \{\mathbb{P}^{\otimes N} R(\hat{f}) - R(\tilde{f})\} \geq CN^{-\frac{1}{1+q}}.$$

*Proof.* Apply Lemma 5.1 for a set  $\mathcal{P}$  equal to a  $(N^{\frac{q}{1+q}}, N^{-1}, \frac{1}{2})$ -constant hypercube and take  $\mathcal{F} \triangleq \{f_{\mathbb{P}}^* : \mathbb{P} \in \mathcal{P}\}$ . □

In Theorem 3.10, we consider classifiers which minimize the empirical risk on an almost minimal net  $\mathcal{N}$ . The following result just asserts that the same convergence rate holds for randomized estimators which “roughly” minimizes the empirical risk.

**Theorem 3.12.** *For any randomized classifier  $\hat{\rho} : \mathcal{Z}^N \rightarrow \mathcal{M}_+^1(\mathcal{N}_{u_N})$  such that there exists a function  $\check{f} \in \mathcal{F}$  satisfying*

- $\mathbb{P}_{\check{f}, \check{f}} \leq C u_N$ ,
- for any  $\epsilon > 0$ , with  $\mathbb{P}^{\otimes N}$ -probability at least  $1 - \epsilon$ ,  $\hat{\rho}r \leq r(\check{f}) + C \log(\epsilon^{-1}) v_N$ ,

we have

$$\mathbb{P}^{\otimes N} \hat{\rho}R - R(\check{f}) \leq \check{C} v_N.$$

*Proof.* It suffices to modify slightly the proof of Theorem 3.10. Let  $\tilde{f}_{\mathcal{N}}$  be the nearest neighbour of  $\check{f}$  in  $\mathcal{N}_{u_N}$ . We have  $\mathbb{P}_{\tilde{f}_{\mathcal{N}}, \check{f}} \leq \check{C} u_N$ . From inequality (2.3) with  $\mathcal{S}_1 = \{\check{f}\}$  and  $\mathcal{S}_2 = \{\tilde{f}_{\mathcal{N}}\}$ , with  $\mathbb{P}^{\otimes N}$ -probability at least  $1 - \epsilon$ , we have

$$r(\check{f}) \leq r(\tilde{f}_{\mathcal{N}}) + \check{C} \sqrt{\frac{u_N \log(\epsilon^{-1})}{N}} + \check{C} \frac{\log(\epsilon^{-1})}{N} + \sup_{\mathbb{P}_{\cdot, \check{f}} \leq \check{C} u_N} \Delta R,$$

hence  $r(\check{f}) \leq r(\check{f}_N) + \check{C} \log(e\epsilon^{-1})v_N$ . We obtain that  $\hat{r} \leq r(\check{f}_N) + \check{C} \log(e\epsilon^{-1})v_N$ . From this inequality and by using the last bound in Corollary 6.6, we obtain a new version of inequality (6.12) from which the convergence rate follows.  $\square$

As a consequence, the Gibbs estimators  $\pi_{-\lambda r}$  in which the prior distribution is the uniform distribution on a net  $\mathcal{N}$  perform as well as an (ERM,  $\mathcal{N}$ )-algorithm (i.e. a classifier which minimizes the empirical risk on the net  $\mathcal{N}$ ) as soon as the inverse temperature parameter  $\lambda$  is sufficiently large. This is not surprising to the extent that the Gibbs estimator  $\pi_{-\lambda r}$  when  $\lambda \rightarrow +\infty$  classifies, roughly speaking, as an (ERM,  $\mathcal{N}$ )-algorithm.

The following theorem completes Theorem 3.10.

**Theorem 3.13.** *Let  $\mathcal{N}_{u_N}$  be a  $u_N$ -covering net such that inequations (3.10) and (3.11) hold, let  $\lambda_N \geq \check{C}_4 \frac{h_q(u_N)}{v_N}$ , and let  $\pi$  be a probability distribution on the net  $\mathcal{N}_{u_N}$  satisfying  $(u_N, \check{C}_5)$ -(CA3) for some positive constants  $\check{C}_i, i = 1, \dots, 5$ . Then we have*

$$\mathbb{P}^{\otimes N} [\pi_{-\lambda_N r} R - R(\check{f})] \leq \check{C} v_N$$

for some constant  $\check{C} > 0$  (depending on  $C''$ ,  $\check{C}_i, i = 1, \dots, 5$  [and also on  $c''$  under Assumption (MA3)]).

*Proof.* Introduce the function  $\check{f}_N$  in the net  $\mathcal{N}_{u_N}$  such that  $\mathbb{P}_{\check{f}, \check{f}_N} \leq u_N$ . By Assumption  $(u_N, \check{C}_5)$ -(CA3) and inequality (3.11), we can choose the function  $\check{f}_N$  such that we also have  $\pi(\{\check{f}_N\}) \geq e^{-(\check{C}_5 + \check{C}_3)h_q(u_N)}$ . So we have

$$\pi_{-\lambda_N r} r - r(\check{f}_N) \leq \frac{\log[\pi(\check{f}_N)^{-1}]}{\lambda_N} \leq \check{C} \frac{h_q(u_N)}{\lambda_N} \leq \check{C} v_N.$$

The result then follows from Theorem 3.12.  $\square$

**3.5. Bracketing entropy.** To minimize the empirical risk over all the model  $\mathcal{F}$  can lead to inconsistency even for models with small covering entropy. For instance, define the set  $\mathcal{X} = [0; 1]$ , the functions  $f_0 \equiv 0$  and  $f_1 \equiv 1$ , and the probability distribution  $\mathbb{P}$  such that  $\mathbb{P}(dX) = \mathcal{U}([0; 1])(dX)$  (uniform law over  $\mathcal{X}$ ) and  $Y = \mathbb{1}_{X \geq \frac{3}{4}}$ . Consider the model formed by  $f_1$  and all the functions equal to  $f_0$  except on a finite number of points. For any  $u < 1$ , we have  $H(u, \mathcal{F}, \mathbb{P}, \cdot) = \log 2$ . However, in general, the ERM-algorithm will classify poorly<sup>8</sup>. (On the contrary, the classifier based on the ERM-principle over a  $(u, \mathcal{F}, \mathbb{P}, \cdot)$ -net for small  $u$  is efficient). This phenomenon occurs since the covering entropy does not suitably measures the complexity of models. In this section, we will see that the bracketing entropy does not suffer from this drawback.

Under polynomial bracketing entropy conditions, the empirical data contain what happens in expectation to the extent that two functions close for the distance  $\mathbb{P}_{\cdot, \cdot}$  are also close for the distance  $\bar{\mathbb{P}}_{\cdot, \cdot}$ .

Recall that if  $\mathcal{G}$  is a  $u$ -bracketing net of the set  $\mathcal{F}$ , then for any function  $f \in \mathcal{F}$ , there exist  $f_L, f_U \in \mathcal{G}$  satisfying  $f_L \leq f \leq f_U$  and  $\mathbb{P}_{f_L, f_U} \leq u$ . Let us define the mappings  $n_L, n_U : \mathcal{F} \rightarrow \mathcal{G}$  such that  $n_L(f) = f_L$  and  $n_U(f) = f_U$  (from the axiom of choice, they exist).

The following theorem, to be compared with Theorem 3.10, shows the influence of considering bracketing entropy assumptions instead of covering ones.

<sup>8</sup>That is why, in Theorem 3.10, we need to consider almost *minimal* nets (inequality (3.7)).



**Theorem 3.14.** *Let us define*

$$w_N \triangleq \begin{cases} \left[ \frac{\log(eN^{1/\kappa})}{N} \right]^{\frac{\kappa}{2\kappa-1}} & \text{under Assumptions (MA2)+(CA2) for } q = 0 \\ N^{-\frac{\kappa}{2\kappa-1+q}} & \text{under Assumptions (MA2)+(CA2) for } 0 < q < 1 \\ (\log N)N^{-\frac{1}{2}} & \text{under Assumptions (MA2)+(CA2) for } q = 1 \\ N^{-\frac{1}{1+q}} & \text{under Assumptions (MA2)+(CA2) for } q > 1 \end{cases} .$$

For any classifier  $\hat{f}_{\text{ERM},\mathcal{N}}$  minimizing the empirical risk in a  $u_N \triangleq \check{C}_1 w_N$ -covering net  $\mathcal{N}$  for some positive constant  $\check{C}_1$ , we have

$$\mathbb{P}^{\otimes N} [R(\hat{f}_{\text{ERM},\mathcal{N}}) - R(\tilde{f})] \leq \check{C} w_N$$

for some constant  $\check{C} > 0$  (depending on  $C', C''$  and  $\check{C}_1$ ).

*Proof.* Let  $\mathcal{N}'$  be a  $u_N$ -minimal bracketing net of the net  $\mathcal{N}$ . Let  $\tilde{f}_{\mathcal{N}'}$  be the nearest neighbour of the function  $\tilde{f}$  in the net  $\mathcal{N}'$ . By definition of the set  $\mathcal{N}'$ ,

- we have  $\log |\mathcal{N}'| \leq C' h_q(u_N)$ ,
- there exists a function  $f_{\mathcal{N}'}$  such that  $n_L(f_{\mathcal{N}'}) = \tilde{f}_{\mathcal{N}'}$  or  $n_U(f_{\mathcal{N}'}) = \tilde{f}_{\mathcal{N}'}$ ; consequently, we have  $r(f_{\mathcal{N}'}) \leq r(\tilde{f}_{\mathcal{N}'}) + u_N$ ,
- there exists a classifier  $\hat{f}_{\mathcal{N}'} : \mathcal{Z}^N \rightarrow \mathcal{N}'$  ( $\hat{f}_{\mathcal{N}'} \triangleq n_L(\hat{f}_{\text{ERM},\mathcal{N}'})$  for instance) such that we have  $r(\hat{f}_{\mathcal{N}'}) \leq r(\hat{f}_{\text{ERM},\mathcal{N}'}) + u_N$  and

$$(3.12) \quad R(\hat{f}_{\text{ERM},\mathcal{N}'}) \leq R(\hat{f}_{\mathcal{N}'}) + u_N.$$

So the estimator  $\hat{f}_{\mathcal{N}'} : \mathcal{Z}^N \rightarrow \mathcal{N}'$  satisfies

$$r(\hat{f}_{\mathcal{N}'}) \leq r(\hat{f}_{\text{ERM},\mathcal{N}'}) + u_N \leq r(f_{\mathcal{N}'}) + u_N \leq r(\tilde{f}_{\mathcal{N}'}) + 2u_N.$$

Then the result follows from Theorem 3.12 and inequality (3.12).  $\square$

*Remark 3.9.* Since we have  $(CA2) \Rightarrow (CA1)$ , Theorem 3.10 can be applied when the assumptions of Theorem 3.14 hold. We see that, under bracketing entropy assumptions, the ERM on nets containing a huge (possibly infinite) number of functions has also the optimal convergence rate. This was not the case under covering entropy assumptions.

*Remark 3.10.* The same convergence rate holds for classifiers minimizing the empirical risk up to an additive factor  $Cw_N$ .

The following theorem completes the previous one.

**Theorem 3.15.** *Let  $\lambda_N \geq \check{C}_1 \frac{h_q(w_N)}{w_N}$  and  $\pi$  be a probability distribution satisfying  $(\check{C}_2 w_N, \check{C}_3)$ -(CA3) for some positive constants  $\check{C}_i, i = 1, \dots, 3$ . Then we have*

$$\mathbb{P}^{\otimes N} [\pi_{-\lambda_N r} R - R(\tilde{f})] \leq \check{C} w_N$$

for some constant  $\check{C} > 0$  (depending on  $C'', \check{C}_i, i = 1, \dots, 3$ ).

*Proof.* See Section 6.8.  $\square$

*Remark 3.11.* From the previous theorem, the inverse temperature parameter  $\lambda_N$  should be taken as

$$\lambda_N \geq C \begin{cases} (\log N) N^{\frac{\kappa}{2\kappa-1}} & \text{under Assumptions (MA2)+(CA2) for } q = 0 \\ N^{\frac{\kappa(1+q)}{2\kappa-1+q}} & \text{under Assumptions (MA2)+(CA2) for } 0 < q < 1 \\ \frac{N}{(\log N)^2} & \text{under Assumptions (MA2)+(CA2) for } q = 1 \\ N & \text{under Assumptions (MA2)+(CA2) for } q > 1 \end{cases} .$$

The threshold value is all the smaller as the model is small<sup>9</sup> (i.e for small  $q$ ) and the margin assumption is weak<sup>10</sup> (i.e for large  $\kappa$ ).

Finally, the following theorem shows that, under polynomial bracketing entropy assumption, with high probability, the empirical covering nets are similar to the covering nets wrt the pseudo-distance  $\mathbb{P}(dX)$ .

**Theorem 3.16.** *Let  $\check{C}$  be positive constant and define*

$$(\alpha_q, \beta_q) = \begin{cases} \left( \frac{1}{N}, \frac{\log N}{N} \right) & \text{when } q = 0 \\ \left( \exp \left\{ -N^{\frac{q}{1+q}} \right\}, N^{-\frac{1}{1+q}} \right) & \text{when } q > 0 \end{cases}.$$

With  $\mathbb{P}^{\otimes N}$ -probability at least  $1 - (\alpha_q)^{\check{C}}$ , there exists  $\check{C}_1, \check{C}_2, \check{C}_3, \check{C}_4 > 0$  such that for any  $u \geq \check{C}_1 \beta_q$ ,

- a  $(u, \mathcal{F}, \mathbb{P}_{\cdot, \cdot})$ -covering net is a  $(\check{C}_3 u, \mathcal{F}, \bar{\mathbb{P}}_{\cdot, \cdot})$ -covering net,
- a  $(u, \mathcal{F}, \bar{\mathbb{P}}_{\cdot, \cdot})$ -covering net is a  $(\check{C}_2 u, \mathcal{F}, \mathbb{P}_{\cdot, \cdot})$ -covering net,
- $H(u, \mathcal{F}, \bar{\mathbb{P}}_{\cdot, \cdot}) \leq \check{C}_4 h_q(u)$ .

*Proof.* See Section 6.9. □

Therefore under polynomial bracketing entropy assumption, we can classify optimally by using the minimizer of the empirical risk on an *empirical* net of radius less than  $Cw_N$ . Note that another way of proving this result consists in saying that this classifier minimizes the empirical risk on the set  $\mathcal{F}$  up to an additive  $Cw_N$  factor.

#### 4. CLASSIFICATION UNDER EMPIRICAL COMPLEXITY ASSUMPTIONS

In this section, we will see that if we replace the complexity assumption concerning  $\mathbb{P}$ -entropies with a similar assumption on the empirical entropies, the same kind of convergence rates appear. VC-classes are a special case in which for any  $u > 0$  and any training set, we have  $H(u, \mathcal{F}, \bar{\mathbb{P}}_{\cdot, \cdot}) \leq CVh_0(u)$  where  $V$  is the VC-dimension of  $\mathcal{F}$ .

**4.1. Concentration of the empirical entropies.** In general, the link between the  $\mathbb{P}$ -entropies and  $\bar{\mathbb{P}}$ -entropies is not known. However, thanks to recent work by Boucheron, Bousquet, Lugosi and Massart, we are able to prove that the empirical entropies are concentrated.

**Theorem 4.1.** *For any  $\epsilon > 0$  and  $u \geq 0$*

- with  $\mathbb{P}^{\otimes N}$ -probability at least  $1 - \epsilon$ , we have

$$(4.1) \quad H(u, \mathcal{F}, \bar{\mathbb{P}}) \leq \mathbb{P}^{\otimes N} H(u, \mathcal{F}, \bar{\mathbb{P}}) + \frac{(\log 2) \log(\epsilon^{-1})}{3} \left( 1 + \sqrt{1 + \frac{18 \mathbb{P}^{\otimes N} H(u, \mathcal{F}, \bar{\mathbb{P}})}{(\log 2) \log(\epsilon^{-1})}} \right)$$

---

<sup>9</sup>This might be explained by looking at the size of the sets  $\{f \in \mathcal{F} : r(f) - \min_{\mathcal{F}} r = \frac{k}{N}\}$ . Indeed, when the model becomes larger and larger, the weight on these sets increases much more for small  $k$  than for very small  $k$ , hence we need to have larger  $\lambda$  to get rid of functions having a not-so-small empirical risk.

<sup>10</sup>This is not surprising since the stronger the margin assumption is, the smaller the optimal convergence rate is, and consequently the more selective we need to be.

equivalently

$$\mathbb{P}^{\otimes N} H(u, \mathcal{F}, \bar{\mathbb{P}}) \geq H(u, \mathcal{F}, \bar{\mathbb{P}}) + \frac{2(\log 2)\log(\epsilon^{-1})}{3} \left( 1 - \sqrt{1 + \frac{9H(u, \mathcal{F}, \bar{\mathbb{P}})}{2(\log 2)\log(\epsilon^{-1})}} \right)$$

- with  $\mathbb{P}^{\otimes N}$ -probability at least  $1 - \epsilon$ ,

$$H(u, \mathcal{F}, \bar{\mathbb{P}}) \geq \mathbb{P}^{\otimes N} H(u, \mathcal{F}, \bar{\mathbb{P}}) - \sqrt{2(\log 2)\log(\epsilon^{-1})\mathbb{P}^{\otimes N} H(u, \mathcal{F}, \bar{\mathbb{P}})}$$

equivalently

$$(4.2) \quad \mathbb{P}^{\otimes N} H(u, \mathcal{F}, \bar{\mathbb{P}}) \leq H(u, \mathcal{F}, \bar{\mathbb{P}}) + (\log 2)\log(\epsilon^{-1}) \left( 1 + \sqrt{1 + \frac{2H(u, \mathcal{F}, \bar{\mathbb{P}})}{(\log 2)\log(\epsilon^{-1})}} \right)$$

- with  $\mathbb{P}^{\otimes 2N}$ -probability at least  $1 - 2\epsilon$ ,

$$(4.3) \quad H(u, \mathcal{F}, \bar{\mathbb{P}}') \leq H(u, \mathcal{F}, \bar{\mathbb{P}}) + 2(\log 2)\log(\epsilon^{-1}) \left( \frac{6}{5} + \sqrt{1 + \frac{2H(u, \mathcal{F}, \bar{\mathbb{P}})}{(\log 2)\log(\epsilon^{-1})}} \right)$$

*Proof.* See Section 6.10.  $\square$

The previous result shows that the empirical entropies behave with high probability as the non empirical quantity  $\mathbb{P}^{\otimes N} H(u, \mathcal{F}, \bar{\mathbb{P}})$ . Specifically, by using a union bound on the different possible radius, we obtain that for any  $\check{C}' > 0$  there exists  $\check{C} > 0$  such that with probability at least  $1 - \frac{1}{N^{\check{C}'}}$ , for any  $u > 0$ , we have<sup>11</sup>

$$\begin{cases} \mathbb{P}^{\otimes N} H(u, \mathcal{F}, \bar{\mathbb{P}}_{\cdot, \cdot}) & \leq \check{C}'[H(u, \mathcal{F}, \bar{\mathbb{P}}_{\cdot, \cdot}) + \log N] \\ H(u, \mathcal{F}, \bar{\mathbb{P}}_{\cdot, \cdot}) & \leq \check{C}'[\mathbb{P}^{\otimes N} H(u, \mathcal{F}, \bar{\mathbb{P}}_{\cdot, \cdot}) + \log N] \\ H(u, \mathcal{F}, \bar{\mathbb{P}}_{\cdot, \cdot}) & \leq \check{C}'[H(u, \mathcal{F}, \bar{\mathbb{P}}_{\cdot, \cdot}) + \log N] \\ H(u, \mathcal{F}, \bar{\mathbb{P}}_{\cdot, \cdot}) & \leq H(u/2, \mathcal{F}, \bar{\mathbb{P}}_{\cdot, \cdot}) \end{cases}.$$

## 4.2. Chaining empirical quantities...

4.2.1. *...in the transductive learning.* In this section, we assume that we possess two samples of size  $N$ . The first sample is labeled:  $\{(X_1, Y_1), \dots, (X_N, Y_N)\}$ . The second one  $\{X_{N+1}, \dots, X_{2N}\}$  has to be labeled: the outputs  $\{Y_{N+1}, \dots, Y_{2N}\}$  are unknown. We will use the following notations:

$$\begin{cases} \bar{\mathbb{P}} & \triangleq \frac{1}{N} \sum_{i=1}^N \delta_{(X_i, Y_i)} \\ \bar{\mathbb{P}}' & \triangleq \frac{1}{N} \sum_{i=N+1}^{2N} \delta_{(X_i, Y_i)} \\ \bar{\bar{\mathbb{P}}} & \triangleq \frac{1}{2N} \sum_{i=1}^{2N} \delta_{(X_i, Y_i)} \\ r(f) & \triangleq \frac{1}{N} \sum_{i=1}^N \mathbb{1}_{Y_i \neq f(X_i)} = \bar{\mathbb{P}}[Y \neq f(X)] \\ r'(f) & \triangleq \frac{1}{N} \sum_{i=N+1}^{2N} \mathbb{1}_{Y_i \neq f(X_i)} = \bar{\mathbb{P}}'[Y \neq f(X)] \end{cases}$$

Let us start with a basic result which is not “chained”.

**Lemma 4.2.** *Let  $\mathcal{S}_1$  and  $\mathcal{S}_2$  be two finite sets of functions from  $\mathcal{X}$  into  $\mathcal{Y}$  possibly depending on the data  $Z_1^{2N}$  in an exchangeable way. For any  $\epsilon > 0$ , with  $\mathbb{P}^{\otimes 2N}$ -probability at least  $1 - \epsilon$ , for any functions  $f_1 \in \mathcal{S}_1$  and  $f_2 \in \mathcal{S}_2$ , we have*

$$r'(f_2) - r'(f_1) + r(f_1) - r(f_2) \leq \sqrt{\frac{8\bar{\bar{\mathbb{P}}}_{f_1, f_2} \log(|\mathcal{S}_1||\mathcal{S}_2|\epsilon^{-1})}{N}}$$

<sup>11</sup>For the third inequality, we use the inequality  $H(u, \mathcal{F}, \bar{\mathbb{P}}_{\cdot, \cdot}) \leq H(u, \mathcal{F}, \bar{\mathbb{P}}_{\cdot, \cdot}) + H(u, \mathcal{F}, \bar{\mathbb{P}}'_{\cdot, \cdot})$  and inequality (4.3). The fourth inequality always holds.

*Proof.* The result comes from inequality (8.7) in [1] in which we take  $\nu$  equal to the uniform distribution on  $\mathcal{S}_1 \times \mathcal{S}_2$  and  $\mathcal{W}[(f_1, f_2), Z] = \mathbb{1}_{Y \neq f_2(X)} - \mathbb{1}_{Y \neq f_1(X)}$ .  $\square$

By chaining this inequality, we obtain:

**Theorem 4.3.** *Let  $U \in \mathbb{N}^*$  and  $u \triangleq 2^{-U}$ . Let  $\mathcal{N}$  be an  $u$ -minimal covering net. For any  $k \in \mathbb{N}^*$ , let  $H_k$  be an upper bound of  $H(2^{-k}, \mathcal{F}, \bar{\mathbb{P}}) \vee 1$ . For any  $\epsilon > 0$ , with  $\mathbb{P}^{\otimes 2N}$ -probability at least  $1 - \epsilon$ , for any  $f_1, f_2 \in \mathcal{N}$ ,*

$$\begin{aligned} r'(f_2) - r'(f_1) + r(f_1) - r(f_2) \\ \leq \sum_{k \in \mathbb{N}^*: u \leq 2^{-k} \leq \bar{\mathbb{P}}_{f_1, f_2} \vee u} 4 \sqrt{\frac{6 \times 2^{-k} \{2H_k + \log[3k(k+1)] + \log(\epsilon^{-1})\}}{N}}. \end{aligned}$$

*Proof.* See Section 6.11.  $\square$

*Remark 4.1.* The previous result can also be written in terms of integral. For instance, for  $H_k = H(2^{-k}, \mathcal{F}, \bar{\mathbb{P}}) \vee 1$ , the previous RHS is upper bounded by<sup>12</sup>

$$\frac{28}{\sqrt{N}} \int_{\frac{u}{2}}^{(\bar{\mathbb{P}}_{f_1, f_2} \vee u) \wedge \frac{1}{2}} \left( \sqrt{\frac{H(x, \mathcal{F}, \bar{\mathbb{P}}) \vee 1}{x}} + \sqrt{\frac{\log(4 \log x^{-1})}{x}} \right) dx + 34 \sqrt{\frac{(\bar{\mathbb{P}}_{f_1, f_2} \vee u) \log(3\epsilon^{-1})}{N}}.$$

4.2.2. ...in the inductive learning. The empirical bound for the inductive learning is derived from the one for the transductive learning and from the concentration properties of the pseudo-distances and the empirical entropies.

**Theorem 4.4.** *Let  $\epsilon > 0$  and  $H_\infty \triangleq 16 \log N(X_1^N) + 20 \log(5 \log N) + 12 \log(\epsilon^{-1})$ . With  $\mathbb{P}^{\otimes 2N}$ -probability  $1 - 3\epsilon$ , for any functions  $f_1$  and  $f_2$  in the set  $\mathcal{F}$ , we have*

$$\begin{aligned} r'(f_2) - r'(f_1) + r(f_1) - r(f_2) \\ \leq \frac{10}{\sqrt{N}} \sum_{\substack{k \in \mathbb{N}^*: \\ \frac{1}{2N} \leq 2^{-k} \leq \frac{5}{4} \bar{\mathbb{P}}_{f_1, f_2} + \frac{H_\infty}{N}}} \sqrt{8H(2^{-k}, \mathcal{F}, \bar{\mathbb{P}}_{f_1, f_2}) + 6 \log[k(k+1)\epsilon^{-1}] + 1}. \end{aligned}$$

*Proof.* See Section 6.12.  $\square$

The previous theorem gives, for instance, a guarantee of misclassification rate of the ERM-classifier on  $N$  new input data to classify. We recall that the leading term in the square root is generally the entropy one. Once more, we can upper bound the associated sum with the integral entropy

$$\frac{C}{\sqrt{N}} \int_{\frac{1}{4N}}^{\frac{5}{4} \bar{\mathbb{P}}_{f_1, f_2} + \frac{H_\infty}{N}} \sqrt{\frac{H(x, \mathcal{F}, \bar{\mathbb{P}}_{f_1, f_2})}{x}} dx$$

Note that this result is less general than the one for transductive learning since the integral starts from  $\frac{1}{4N}$ , which means that the largest complexity terms are taken into account. In Section 3, we have seen that for polynomial entropies with  $q \geq 1$ , the optimal convergence rate (which was of order  $N^{-\frac{1}{1+q}}$  up to the logarithmic factor) was proved since the largest complexities were not in the integral entropy.

On the contrary, for  $q < 1$ , we can recover the same convergence rates under the assumption  $H(u, \mathcal{F}, \bar{\mathbb{P}}) \leq C' h_q(u)$  for any  $u > 0$ , as under the polynomial bracketing entropy assumption. The following section deals with a special case of the case  $q = 0$ .

<sup>12</sup>Proof at the end of Section 6.11.

**4.3. Application to VC-classes.** At first sight, it is not obvious that Theorem 4.3 gives a tighter bound than Lemma 4.2 applied to  $(\mathcal{S}_1, \mathcal{S}_2) = (\mathcal{N}, \mathcal{N})$ . We will see in this section that for VC-classes, the two bounds gives the same convergence rate for the ERM-classifier, except when we have no margin assumption. In this last case, the chained result allows to get rid of a logarithmic factor.

Let us consider the binary classification setting:  $\mathcal{Y} = \{0; 1\}$ . Introduce the shattering number  $N(X_1^{2N}) \triangleq |\{[f(X_k)]_{k=1}^{2N} : f \in \mathcal{F}\}| = H(u, \mathcal{F}, \bar{\mathbb{P}})$  for any  $u < \frac{1}{2N}$ . Let  $V$  be the VC-dimension of the set  $\mathcal{F}$

$$V \triangleq \max \{|A| : A \in \mathcal{X}^{2N} \text{ such that } |\{A \cap f^{-1}(1) : f \in \mathcal{F}\}| = 2^{|A|}\}.$$

The empirical entropies satisfy<sup>13</sup>

$$H(u, \mathcal{F}, \bar{\mathbb{P}}) \leq \begin{cases} V \log \left( \frac{2Ne}{V} \right) \\ V \log \left( \frac{4e}{u} \right) \end{cases}.$$

Let  $\hat{f}_{\text{ERM}}$  be the minimizer of the empirical risk on the set  $\mathcal{F}$  and  $\tilde{f}'$  be the minimizer on  $\mathcal{F}$  of either  $r'$  or  $R$ . From Lemma 4.2, with  $\mathbb{P}^{\otimes 2N}$ -probability at least  $1 - \epsilon$ , we have

$$r'(\hat{f}_{\text{ERM}}) \leq \inf_{f \in \mathcal{F}} \left\{ r'(f) + 4 \sqrt{\frac{\bar{\mathbb{P}}_{\hat{f}_{\text{ERM}}, f} [V \log \left( \frac{2eN}{V} \right) + \frac{1}{2} \log(\epsilon^{-1})]}{N}} \right\},$$

and consequently, after some standard computations:

$$(4.4) \quad \mathbb{P}^{\otimes N} R(\hat{f}_{\text{ERM}}) - R(\tilde{f}) \leq 4 \sqrt{\frac{V \mathbb{P}^{\otimes 2N} \bar{\mathbb{P}}_{\hat{f}_{\text{ERM}}, \tilde{f}} \log \left( \frac{2eN}{V} \right)}{N}} + 2 \sqrt{\frac{2 \mathbb{P}^{\otimes 2N} \bar{\mathbb{P}}_{\hat{f}_{\text{ERM}}, \tilde{f}}}{N}}.$$

To compare, from Theorem 4.3, we obtain

**Corollary 4.5.** *For any  $\epsilon > 0$ , with  $\mathbb{P}^{\otimes 2N}$ -probability at least  $1 - \epsilon$ , we have*

$$r'(\hat{f}_{\text{ERM}}) \leq \inf_{f \in \mathcal{F}} \left\{ r'(f) + 47 \sqrt{\frac{(V+1) \bar{\mathbb{P}}_{\hat{f}_{\text{ERM}}, f} \log \left( \frac{8e}{\bar{\mathbb{P}}_{\hat{f}_{\text{ERM}}, f}} \right)}{N}} + 34 \sqrt{\frac{\bar{\mathbb{P}}_{\hat{f}_{\text{ERM}}, f} \log(\epsilon^{-1})}{N}} \right\}$$

and, consequently,

$$(4.5) \quad \mathbb{P}^{\otimes N} R(\hat{f}_{\text{ERM}}) - R(\tilde{f}) \leq 47 \sqrt{\frac{(V+1) \mathbb{P}^{\otimes 2N} \bar{\mathbb{P}}_{\hat{f}_{\text{ERM}}, \tilde{f}} \log \left( \frac{8e}{\mathbb{P}^{\otimes 2N} \bar{\mathbb{P}}_{\hat{f}_{\text{ERM}}, \tilde{f}}} \right)}{N}} + 34 \sqrt{\frac{\mathbb{P}^{\otimes 2N} \bar{\mathbb{P}}_{\hat{f}_{\text{ERM}}, \tilde{f}}}{N}}.$$

*Proof.* See Section 6.13. □

As a consequence, we obtain

**Corollary 4.6.** *Under assumption (MA2), for any set  $\mathcal{F}$  of VC-dimension  $V$ , the ERM-classifier satisfies*

$$\mathbb{P}^{\otimes N} R(\hat{f}_{\text{ERM}}) - R(\tilde{f}) \leq \check{C} \begin{cases} \left( \frac{V}{N} \log N \right)^{\frac{\kappa}{2\kappa-1}} & \text{when } 1 \leq \kappa < +\infty \\ \sqrt{\frac{V}{N}} & \text{when } \kappa = +\infty \end{cases}.$$

*Proof.* See section 6.14. □

<sup>13</sup>The first inequality is well-known consequence of Sauer's lemma; the second one comes from Haussler's formula ([7]), which asserts that for any  $u > 0$ ,  $H(u, \mathcal{F}, \bar{\mathbb{P}}) \leq V \log \left( \frac{2e}{u} \right) + \log[\epsilon(V+1)]$ .

*Remark 4.2.* This is an improvement of Massart and Nédélec results [12, Corollary 2.2] to the extent we do not have an extra additive term  $R(\tilde{f}) - R(f^*)$ , where  $f^*$  is the Bayes classifier. The second part of the corollary is a well-known result which is given with a simple proof in [5, p.31].

*Remark 4.3.* By comparing Inequalities (4.4) and (4.5), we see that the constants in chained inequalities are not satisfactory. The gap between the upper bound (6.23) and the lower bound (see Theorem 5.2) is the factor  $8 \times 83!$  We do not know how to chain inequalities with significantly tighter constants.

## 5. ASSOUD'S LEMMA

**Definition 5.1.** Let  $m \in \mathbb{N}^*$ ,  $w \in ]0; 1]$ ,  $b \in ]0; 1]$  and  $b' \in ]0; 1]$ . A  $(m, w, b, b')$ -hypercube of probability distributions is a family

$$\{\mathbb{P}_{\vec{\sigma}} \in \mathcal{M}_+^1(\mathcal{Z}) : \vec{\sigma} \triangleq (\sigma_1, \dots, \sigma_m) \in \{-1; +1\}^m\}$$

of  $2^m$  probability distributions having the same first marginal:

$$\mathbb{P}_{\vec{\sigma}}(dX) = \mathbb{P}_{(+1, \dots, +1)}(dX) \triangleq \mu \text{ for any } \vec{\sigma} \in \{-1; +1\}^m,$$

and such that there exists a partition  $\mathcal{X}_0, \dots, \mathcal{X}_m$  of  $\mathcal{X}$  satisfying

- for any  $j \in \{1, \dots, m\}$ , we have  $\mu(\mathcal{X}_j) = w$
- for any  $j \in \{0, \dots, m\}$ , for any  $X \in \mathcal{X}_j$ , we have

$$\mathbb{P}_{\vec{\sigma}}(Y = 1|X) = \frac{1 + \sigma_j \xi(X)}{2} = 1 - \mathbb{P}_{\vec{\sigma}}(Y = 0|X),$$

where  $\sigma_0 \triangleq 1$  and  $\xi : \mathcal{X} \rightarrow [0; 1]$  is such that for any  $j \in \{1, \dots, m\}$ ,

$$\begin{cases} b &= \sqrt{1 - (\mu[\sqrt{1 - \xi^2(X)} | X \in \mathcal{X}_j])^2} \\ b' &= \mu[\xi(X) | X \in \mathcal{X}_j] \end{cases}.$$

When  $\xi$  is constant on  $\mathcal{X}_j, j = 1, \dots, m$  (which implies  $\xi \equiv b' = b$  on  $\mathcal{X} - \mathcal{X}_0$ ), the hypercube will be said a  $(m, w, b)$ -constant hypercube. The hypercube will be said noiseless when  $\xi \equiv 1$  on  $\mathcal{X}_0$ .

The following lemma is Assouad's lemma adapted to the classification framework.

**Lemma 5.1.** *If a set  $\mathcal{P}$  of probability distributions contains a  $(m, w, b, b')$ -hypercube, then for any measurable estimator  $\hat{f} : \mathcal{Z}^N \rightarrow \mathcal{F}(\mathcal{X}, \mathcal{Y})$ , we have*

$$\sup_{\mathbb{P} \in \mathcal{P}} \{\mathbb{P}^{\otimes N} R_{\mathbb{P}}(\hat{f}) - R_{\mathbb{P}}(f_{\mathbb{P}}^*)\} \geq \frac{1 - b\sqrt{Nw}}{2} mwb'.$$

*Proof.* See Section 6.15. □

Lemma 5.1 gives a very simple strategy to obtain a lower bound for a given set  $\mathcal{P}$  of probability distributions: it consists in looking for the  $(m, w, b, b')$ -hypercube which is contained in the set  $\mathcal{P}$  and for which  $\frac{1 - b\sqrt{Nw}}{2} mwb'$  is maximized.

In general, the order of the bound is given by the quantity  $mwb'$  and  $w, b$  are taken such that  $\sqrt{Nwb} = \text{Cst} < 1$ . To obtain this order, we do not need the sophisticated computations detailed in the proof of the lemma. We can use two well-known lemmas instead (Birgé's lemma and Huber's lemma) as it is proved in Appendix E.

Lemma 5.1 implies lower bounds for VC-classes with decent constants. The following result is to be compared with Theorems 14.1 and 14.5 in [5].

**Theorem 5.2.** For any model  $\mathcal{F}$ , define  $\mathcal{P}_L$  as the set of probability distributions such that  $\inf_{f \in \mathcal{F}} R_{\mathbb{P}}(f) = L$  for a fixed  $L \in [0; 1/2]$ .

• When  $L = 0$ :

for any classification model  $\mathcal{F}$  of VC-dimension  $V \geq 2$ , for any measurable estimator  $\hat{f} : \mathcal{Z}^N \rightarrow \mathcal{F}(\mathcal{X}, \mathcal{Y})$ , we have

$$\sup_{\mathbb{P} \in \mathcal{P}_0} \left\{ \mathbb{P}^{\otimes N} R_{\mathbb{P}}(\hat{f}) - R_{\mathbb{P}}(f_{\mathbb{P}}^*) \right\} \geq \begin{cases} \frac{V-1}{2e^{(N+1)}} \text{ when } N \geq (V-2) \vee 1 \\ \frac{1}{2} \left(1 - \frac{1}{V-1}\right)^N \end{cases}.$$

• When  $0 < L \leq 1/2$ :

for any classification model  $\mathcal{F}$  of VC-dimension  $V \geq 2$ , for any measurable estimator  $\hat{f} : \mathcal{Z}^N \rightarrow \mathcal{F}(\mathcal{X}, \mathcal{Y})$ , we have

$$\sup_{\mathbb{P} \in \mathcal{P}_L} \left\{ \mathbb{P}^{\otimes N} R_{\mathbb{P}}(\hat{f}) - R_{\mathbb{P}}(f_{\mathbb{P}}^*) \right\} \geq \begin{cases} \sqrt{\frac{L(V-1)}{32N}} \vee \frac{V-1}{16N} \text{ when } \frac{(1-2L)^2 N}{V-1} \geq \frac{1}{4} \\ \left(\frac{1}{2} - L\right) \sqrt{\frac{L}{2}} \text{ otherwise} \end{cases}.$$

*Proof.* See Appendix F. □

*Remark 5.1.* It is a well known result that, when  $\inf_{f \in \mathcal{F}} R_{\mathbb{P}}(f)$  is of order  $1/N$  and when the complexity of the class is not too high, there exists an estimator such that  $\mathbb{P}^{\otimes N} R_{\mathbb{P}}(\hat{f}) - \inf_{f \in \mathcal{F}} R_{\mathbb{P}}(f) = O\left(\frac{1}{N}\right)$ . The previous theorem gives a corresponding lower bound.

## 6. PROOFS

**6.1. Proof of Lemma 3.1.** Let  $T_1(\pi) \triangleq -\log \pi \exp(-\lambda \Delta R)$  and

$$T_2(\pi) \triangleq 0 \vee \log \pi \exp\left(8.2 \frac{\lambda^2}{N} \mathbb{P}_{\cdot, \bar{f}} - \lambda \Delta R\right).$$

We start with the following lemma.

**Lemma 6.1.** For any  $\epsilon > 0$  and  $0 < \lambda \leq 0.19N$ , with  $\mathbb{P}^{\otimes N}$ -probability at least  $1 - \epsilon$ , we have  $\pi_{-\lambda r} R \leq \frac{\epsilon}{\lambda} [T_1(\pi) + T_2(\pi) + \log(4\epsilon^{-1})]$ .

*Proof.* Taking  $\chi = \frac{1}{2}$  and  $\gamma = \frac{1}{2}$  in Theorem 2.2, we get

$$\begin{aligned} \pi_{-\lambda r} \Delta R &\leq \pi_{-\frac{\lambda}{2} R} \Delta R + \frac{2}{\lambda} \left[ 16 \log \pi_{-\lambda R} \exp\left(\frac{8.2\lambda^2}{N} \mathbb{P}_{\cdot, \bar{f}}\right) + 5 \log(4\epsilon^{-1}) \right] \\ &\leq -\frac{2}{\lambda} \log \pi \exp\left(-\frac{\lambda}{2} \Delta R\right) + \frac{32}{\lambda} \log \pi_{-\lambda R} \exp\left(\frac{8.2\lambda^2}{N} \mathbb{P}_{\cdot, \bar{f}}\right) + \frac{10}{\lambda} \log(4\epsilon^{-1}) \\ &\leq -\frac{34}{\lambda} \log \pi \exp(-\lambda \Delta R) + \frac{32}{\lambda} \log \pi \exp\left(-\lambda \Delta R + \frac{8.2\lambda^2}{N} \mathbb{P}_{\cdot, \bar{f}}\right) \\ &\quad + \frac{10}{\lambda} \log(4\epsilon^{-1}). \end{aligned}$$

□

*Remark 6.1.* In order to explain the assumptions used in Lemma 3.1, let us give upper bounds for the quantities  $T_1$  and  $T_2$  using the strong complexity and margin Assumptions (CA1) and (MA3) for a well chosen distribution  $\pi$ . Under Assumption (CA1) (which is equivalent to Assumption (CA3)), there exists a distribution  $\pi^{(t)}$  such that for any  $f' \in \mathcal{F}$ ,  $\pi^{(t)}(\mathbb{P}_{\cdot, f'} \leq t) \geq e^{-C' t^{-q}}$ .

For any  $0 < t < 1$ , we have

$$\begin{aligned} T_1 \left[ \pi(c'' t^{1/\kappa}) \right] &\leq -\log \left[ \pi(c'' t^{1/\kappa})(\Delta R \leq t) e^{-\lambda t} \right] && \text{(by Markov's inequality)} \\ &\leq -\log \left[ \pi(c'' t^{1/\kappa})(\mathbb{P}_{\cdot, \bar{f}} \leq c'' t^{\frac{1}{\kappa}}) \right] + \lambda t && \text{(according to (MA3))} \\ &\leq C' c''^{-q} t^{-\frac{q}{\kappa}} + \lambda t && \text{(by definition of } \pi^{(t)} \text{)}. \end{aligned}$$

Assumption (MA2) (recall that (MA3)  $\Rightarrow$  (MA2)) implies that for any  $\lambda > 0$

$$\begin{aligned} 8.2\frac{\lambda^2}{N}\mathbb{P}_{\cdot, \tilde{f}} - \lambda\Delta R &\leq 8.2C''\frac{\lambda^2}{N}(\Delta R)^{\frac{1}{\kappa}} - \lambda\Delta R \\ &\leq \lambda \sup_{x \geq 0} \left\{ 8.2C''\frac{\lambda}{N}x^{\frac{1}{\kappa}} - x \right\} \\ &= (\kappa - 1)\lambda \left( \frac{8.2C''\lambda}{\kappa N} \right)^{\frac{\kappa}{\kappa-1}}, \end{aligned}$$

hence  $T_2(\pi) \leq \check{C}\lambda\left(\frac{\lambda}{N}\right)^{\frac{\kappa}{\kappa-1}}$  for any distribution  $\pi$  and a constant  $\check{C} > 0$  depending on  $C''$  and  $\kappa$ . (Note that for the limit case  $\kappa = 1$ , we have  $T_2 = 0$  for any  $\lambda \leq \frac{N}{8.2C''}$ .) Therefore, with  $(\mathbb{P}^{\otimes N})_{*}$ -probability at least  $1 - 4\epsilon$ , we have

$$\left\{ \pi(c''t^{1/\kappa}) \right\}_{-\lambda r} R - R(\tilde{f}) \leq \check{C} \left[ t + \frac{t^{-\frac{q}{\kappa}} + \log(\epsilon^{-1})}{\lambda} + \left( \frac{\lambda}{N} \right)^{\frac{\kappa}{\kappa-1}} \right],$$

where the constant  $\check{C} > 0$  depends on  $C'$ ,  $c''$  and  $\kappa$ . The sum  $t + \frac{t^{-\frac{q}{\kappa}}}{\lambda} + \left(\frac{\lambda}{N}\right)^{\frac{\kappa}{\kappa-1}}$  has the minimal order  $N^{-\frac{\kappa}{2\kappa-1+q}}$  when  $\lambda$  has the order of  $N^{\frac{\kappa+q}{2\kappa-1+q}}$  and  $t$  has the order of  $N^{-\frac{\kappa}{2\kappa-1+q}}$ . This computation explains the choice of Assumptions (3.2) and (3.3).

From inequality (3.2), we have  $T_1(\pi) \leq \check{C}N^{\frac{q}{2\kappa-1+q}} + \lambda_N N^{-\frac{\kappa}{2\kappa-1+q}}$ . From Assumption (MA2), we have seen in the previous remark that  $T_2(\pi) \leq \check{C}\lambda_N\left(\frac{\lambda_N}{N}\right)^{\frac{\kappa}{\kappa-1}}$ . From inequality (3.3), we obtain the desired convergence rate.

Now let us prove the sharper result: inequality (3.4). Let  $a(\lambda) \triangleq \frac{\lambda}{N}g\left(\frac{\lambda}{N}\right)$ . From Theorem 6.2 in [1] and the same computations as for the quantity  $\mathcal{L}''$  in Section 9.12 of [1] to upper bound  $-\log \pi \exp\{-\lambda[r - r(\tilde{f})]\}$ , we obtain :

**Lemma 6.2.** *For any  $\epsilon > 0$ ,  $\lambda > 0$  and  $\xi > 0$ , with  $\mathbb{P}^{\otimes N}$ -probability at least  $1 - 2\epsilon$ , with  $\pi_{-\lambda r}$ -probability at least  $1 - \epsilon$ , we have*

$$\Delta R \leq a(\lambda)\mathbb{P}_{\cdot, \tilde{f}} + \frac{-\log \pi \exp\{-\lambda\Delta R - \lambda a\left(\frac{\lambda}{\xi}\right)\mathbb{P}_{\cdot, \tilde{f}}\} + (2+\xi)\log(\epsilon^{-1})}{\lambda}.$$

Taking  $\xi = 1$  and  $\lambda = \lambda_N$ , using the margin assumption  $\mathbb{P}_{\cdot, \tilde{f}} \leq C''(\Delta R)^{\frac{1}{\kappa}}$  and noting that  $a(\lambda_N) \leq g(\check{C}_4)\frac{\lambda_N}{N}$ , we get

$$\Delta R \leq \check{C}\frac{\lambda_N}{N}(\Delta R)^{\frac{1}{\kappa}} + \frac{3\log(\epsilon^{-1})}{\lambda_N} + \sup_{x \geq 0} \left\{ \check{C}\frac{\lambda_N}{N}x^{\frac{1}{\kappa}} - x \right\} - \frac{\log \pi \exp(-2\lambda_N \Delta R)}{\lambda_N},$$

where the constant  $\check{C} > 0$  only depends on  $C''$  and  $\check{C}_4$ . Now from the same computations as in Remark 6.1, when the Inequalities (3.2) and (3.3) hold, we get

$$\Delta R \leq \check{C} \left[ \frac{\lambda_N}{N}(\Delta R)^{\frac{1}{\kappa}} + N^{-\frac{\kappa}{2\kappa-1+q}} \right] + \frac{3\log(\epsilon^{-1})}{\lambda_N}.$$

We obtain successively

$$\Delta R \leq \check{C} \left[ N^{\frac{1-\kappa}{2\kappa-1+q}}(\Delta R)^{\frac{1}{\kappa}} + \log(\epsilon\epsilon^{-1})N^{-\frac{\kappa}{2\kappa-1+q}} \right]$$

and

$$\Delta R \leq \check{C} \log(\epsilon\epsilon^{-1})N^{-\frac{\kappa}{2\kappa-1+q}}.$$



**6.2. Proof of Theorem 3.3.** A standard idea to prove lower bounds is to consider an adequate hypercube of probability distributions and to use Assouad's lemma (see Section 5 for the definition of the hypercube of distributions).

Consider a  $(N^{\frac{q}{2\kappa-1+q}}, N^{-\frac{1+q}{2\kappa-1+q}}, aN^{-\frac{\kappa-1}{2\kappa-1+q}})$ -constant noiseless hypercube of probability distributions  $\{\mathbb{P}_{\vec{\sigma}} : \vec{\sigma} \in \{-1; +1\}^m\}$ , where  $a > 0$  is a constant which will be chosen later.

Had we replaced Assumption (MA3) with Assumption (MA2) in Theorem 3.3, the result would have been a direct consequence of Lemma 5.1 applied to this hypercube with  $a = \frac{1}{2}$ .

In this proof, we will not apply Assouad's lemma but Fano's lemma since Assumption (MA3) is not satisfied by the whole hypercube. First let us state the following classical result on the hypercube which is a refined version of Varshamov-Gilbert bound (1962).

**Lemma 6.3** (Huber, [8, p.256]). *Let  $\delta(\Sigma, \Sigma')$  denote the Hamming distance between  $\Sigma$  and  $\Sigma'$  in  $\{-1, 1\}^m$ :  $\delta(\Sigma, \Sigma') \triangleq \sum_{i=1}^m \mathbb{1}_{\Sigma_i \neq \Sigma'_i}$ . There exists a subset  $\mathcal{S}$  of the hypercube  $\{-1, 1\}^m$  such that*

- for any  $\Sigma \neq \Sigma'$  in  $\mathcal{S}$ , we have  $\delta(\Sigma, \Sigma') \geq \frac{m}{4}$
- $\log |\mathcal{S}| \geq \frac{m}{8}$ .

*Proof.* It suffices to upper bound the number of points in the ball centered at a point  $\sigma$  of the hypercube and of radius  $\frac{m}{4}$ . Consider the uniform distribution  $\nu(d\Sigma)$  on the hypercube  $\{-1, 1\}^m$ . Specifically, we have

$$\nu\left(\delta(\Sigma, \sigma) \leq \frac{m}{4}\right) \leq \nu e^{\frac{m}{4} - \delta(\Sigma, \sigma)} = e^{\frac{m}{4}} \left( \nu e^{-\mathbb{1}_{\Sigma_i \neq \sigma_i}} \right)^m = \left( \frac{e^{\frac{1}{4}}(1 + e^{-1})}{2} \right)^m \leq e^{-\frac{m}{8}},$$

which leads to the desired result.  $\square$

Let  $\mathcal{S} \subset \{-1; +1\}^m$  such that  $|\mathcal{S}| = \lfloor e^{\frac{m}{8}} \rfloor$  and for any  $\Sigma \neq \Sigma'$  in  $\mathcal{S}$ ,  $\delta(\Sigma, \Sigma') \geq \frac{m}{4}$ . From inequality (5.1) in [2], Birgé's version of Fano's lemma can be stated as

**Lemma 6.4.** *Given a non-trivial (i.e. cardinal  $\geq 2$ ) finite family  $\mathcal{D}$  of probability measures on some measurable set  $(E, \xi)$  and a random variable  $\bar{X}$  with an unknown distribution in the family, we have*

$$\inf_{\hat{T}} \sup_{\mathbb{P} \in \mathcal{D}} \mathbb{P}[\hat{T}(\bar{X}) \neq \mathbb{P}] \geq 0.36 \wedge \left( 1 - \frac{K_{\mathcal{D}}}{|\mathcal{D}| \log |\mathcal{D}|} \right),$$

where  $K_{\mathcal{D}} \triangleq \inf_{\mathbb{P} \in \mathcal{D}} \sum_{\mathbb{Q} \neq \mathbb{P}} K(\mathbb{Q}, \mathbb{P})$  and the infimum is taken over all measurable (possibly randomized) estimators based on  $\bar{X}$  with values in the finite set  $\mathcal{D}$ .

Define  $\mathcal{D}' \triangleq \{\mathbb{P}_{\vec{\sigma}} : \vec{\sigma} \in \mathcal{S}\}$ . Let us apply Birgé's lemma to the set of probability distributions

$$\mathcal{D} \triangleq \{\mathbb{P}^{\otimes N} : \mathbb{P} \in \mathcal{D}'\}.$$

With any estimator  $\hat{f} : \mathcal{Z}^N \rightarrow \mathcal{F}(\mathcal{X}, \mathcal{Y})$ , we can associate an estimator  $\hat{T} : \mathcal{Z}^N \rightarrow \mathcal{D}$  defined as  $\hat{T}(Z_1^N) = \mathbb{P}^{\otimes N}$ , where  $\mathbb{P} \in \mathcal{D}'$  minimizes  $\mu[\xi(X) \mathbb{1}_{f_{\mathbb{P}}^*(X) \neq \hat{f}(Z_1^N)(X)}]$ , where  $f_{\mathbb{P}}^*$  denotes the Bayes classifier associated with the distribution  $\mathbb{P}$ .

By Birgé's lemma, we have  $\sup_{\mathbb{Q} \in \mathcal{D}} \mathbb{Q}[\hat{T}(Z_1^N) \neq \mathbb{Q}] \geq 0.36 \wedge \left( 1 - \frac{K_{\mathcal{D}}}{|\mathcal{D}| \log |\mathcal{D}|} \right)$ . Now, when  $\hat{T}(Z_1^N) \neq \mathbb{P}^{\otimes N}$ , we have  $\mu(f_{\mathbb{P}}^* \neq \hat{f}) \geq \frac{m}{8} w$ , hence  $R_{\mathbb{P}}(\hat{f}) - R_{\mathbb{P}}(f_{\mathbb{P}}^*) \geq \frac{m}{8} w \beta$ .

Therefore, we get

$$(6.1) \quad \sup_{\mathbb{P}^{\otimes N} \in \mathcal{D}} R_{\mathbb{P}}(\hat{f}) - R_{\mathbb{P}}(f_{\mathbb{P}}^*) \geq \frac{m}{8} w \beta \left[ 0.36 \wedge \left( 1 - \frac{K_{\mathcal{D}}}{|\mathcal{D}| \log |\mathcal{D}|} \right) \right].$$

For any  $\mathbb{P} \neq \mathbb{Q} \in \mathcal{D}$ , we have  $K(\mathbb{P}, \mathbb{Q}) \leq N m w \beta \log \left( \frac{1+\beta}{1-\beta} \right)$ . Since we have  $|\mathcal{D}| = \lfloor e^{\frac{m}{8}} \rfloor$ , we obtain  $\frac{K_{\mathcal{D}}}{|\mathcal{D}| \log |\mathcal{D}|} \leq \frac{1}{\log \lfloor e^{\frac{m}{8}} \rfloor} N m w \beta \log \left( \frac{1+\beta}{1-\beta} \right) \leq 20 N w \beta^2$  for  $m$  large enough and  $\beta$  small enough. In our case, we have  $m = N^{-\frac{q}{2\kappa-1+q}}$ ,  $w = N^{-\frac{1+q}{2\kappa-1+q}}$  and  $\beta = a N^{-\frac{\kappa-1}{2\kappa-1+q}}$ . So for  $N$  large enough, we have  $\frac{K_{\mathcal{D}}}{|\mathcal{D}| \log |\mathcal{D}|} \leq 20 a^2$ . Let us choose  $a$  such that  $20 a^2 = 0.64$ . We obtain  $\sup_{\mathbb{P} \in \mathcal{D}'} R_{\mathbb{P}}(\hat{f}) - R_{\mathbb{P}}(f_{\mathbb{P}}^*) \geq 0.008 N^{-\frac{\kappa}{2\kappa-1+q}}$  for  $N$  sufficiently large.

Finally it remains to check that the set of distributions  $\mathcal{D}'$  is included in  $\mathcal{P}$ . For any  $\mathbb{P} \in \mathcal{D}'$ , the complexity Assumption (CA2) is satisfied since

- for  $u < mw$ ,  $H(u, \mathcal{F}, \mathbb{P}, \cdot) \leq \log |\mathcal{F}| \leq C u^{-q}$  for some constant  $C > 0$ .
- for  $u \geq mw$ ,  $H(u, \mathcal{F}, \mathbb{P}, \cdot) = 0 \leq C u^{-q}$ .

For any  $\mathbb{P} \in \mathcal{D}'$ , the margin Assumption (MA3) is satisfied since for any functions  $f \in \mathcal{F} - \{\tilde{f}\}$ ,  $\mathbb{P}_{f, \tilde{f}}$  has the order of  $mw = N^{-\frac{1}{2\kappa-1+q}}$  and  $\Delta R(f)$  has the order of  $mw\beta = a N^{-\frac{\kappa}{2\kappa-1+q}}$ . The margin Assumption (MA1) also holds since we have

$$\mathbb{P} \left( 0 < |\eta^*(x) - \frac{1}{2}| \leq t \right) = \begin{cases} 0 & \text{when } t < \beta \\ mw & \text{when } \beta \leq t < \frac{1}{2} \end{cases}$$

*Remark 6.2.* The proof also holds when  $q = 0$ . In this case, we take  $m = 1$ ,  $w = N^{-\frac{1}{2\kappa-1}}$  and  $\beta = \sqrt{\frac{0.64}{20}} N^{-\frac{\kappa-1}{2\kappa-1}}$ .

### 6.3. Proof of Theorem 3.5.

6.3.1. *First case:*  $\log \pi^{-1}(\Delta R \leq x) = -C' \log x + C''' + o(x^s)$ . Since we have<sup>14</sup>

$$(6.2) \quad \pi_{-\lambda R} \Delta R = \frac{C' + \underset{\lambda \rightarrow +\infty}{o}(\lambda^{-s})}{\lambda},$$

from Theorem 2.2, for any  $0 \leq \chi < 1$ , we get

$$\begin{cases} \pi_{-\lambda r} R \leq \frac{C' + \underset{N \rightarrow +\infty}{o}(\lambda^{-s}) + \underset{N \rightarrow +\infty}{O}(\chi)}{\lambda} + \frac{K(\pi_{-\lambda r}, \pi_{-\lambda R})}{\chi \lambda} \\ \pi_{-\lambda r} R \geq \frac{C' + \underset{N \rightarrow +\infty}{o}(\lambda^{-s}) + \underset{N \rightarrow +\infty}{O}(\chi)}{\lambda} - \frac{K(\pi_{-\lambda r}, \pi_{-\lambda R})}{\chi \lambda} \end{cases}.$$

Taking  $\chi = \sqrt{K(\pi_{-\lambda r}, \pi_{-\lambda R})}$ , we obtain

$$(6.3) \quad \lambda \pi_{-\lambda r} R = C' + \underset{N \rightarrow +\infty}{o}(\lambda^{-s}) + \underset{N \rightarrow +\infty}{O} \left( \sqrt{K(\pi_{-\lambda r}, \pi_{-\lambda R})} \right).$$

*First subcase:*  $\lambda = o(N^{\frac{\kappa}{2\kappa-1}})$ . Assume that  $\lambda = \underset{N \rightarrow +\infty}{o} \left( N^{\frac{\kappa}{2\kappa-1}} \right)$ . Then there exists

$\gamma \in ]0; \frac{1}{2}]$  such that  $\gamma = \underset{N \rightarrow +\infty}{o}(1)$  and  $\lambda \left( \frac{\lambda}{\gamma N} \right)^{\frac{\kappa}{\kappa-1}} = \underset{N \rightarrow +\infty}{o}(1)$ . We have<sup>15</sup>

$$(6.4) \quad \log \pi_{-\lambda R} \exp \left\{ C' \frac{\lambda^2}{\gamma N} (\Delta R)^{\frac{1}{\kappa}} \right\} = \underset{N \rightarrow +\infty}{o}(\lambda^{-s}) + \underset{N \rightarrow +\infty}{O} \left( \left[ \lambda \left( \frac{\lambda}{\gamma N} \right)^{\frac{\kappa}{\kappa-1}} \right]^{\frac{(\kappa-1)C'}{\kappa C' + \kappa - 1}} \right).$$

<sup>14</sup>See Appendix A.

<sup>15</sup>Proof in Appendix B.

Let  $L \triangleq \log(e\epsilon^{-1})$ . From Theorem 2.2, for any  $0 < \gamma < \frac{1}{2}$  and  $0 < \lambda \leq 0.39\gamma N$ , with  $\mathbb{P}^{\otimes N}$ -probability at least  $1 - \epsilon$ , we have

$$\begin{aligned} K(\pi_{-\lambda r}, \pi_{-\lambda R}) &\leq \log \pi_{-\lambda R} \exp \left\{ C \frac{\lambda^2}{\gamma N} (\Delta R)^{\frac{1}{\kappa}} \right\} + C\gamma L \\ &= \underset{N \rightarrow +\infty}{\mathfrak{o}}(\lambda^{-s}) + \underset{N \rightarrow +\infty}{\mathfrak{O}} \left( \left[ \lambda \left( \frac{\lambda}{N} \right)^{\frac{\kappa-1}{\kappa}} \right]^{\frac{(\kappa-1)C'}{\kappa C' + \kappa - 1}} \gamma^{-\frac{\kappa C'}{\kappa C' + \kappa - 1}} + \gamma L \right). \end{aligned}$$

Taking  $\gamma = \lambda^{\frac{(2\kappa-1)C'}{2\kappa C' + \kappa - 1}} N^{-\frac{\kappa C'}{2\kappa C' + \kappa - 1}} L^{-\frac{\kappa C' + \kappa - 1}{2\kappa C' + \kappa - 1}}$ , we obtain

$$(6.5) \quad K(\pi_{-\lambda r}, \pi_{-\lambda R}) = \underset{N \rightarrow +\infty}{\mathfrak{O}} \left( \lambda^{\frac{(2\kappa-1)C'}{2\kappa C' + \kappa - 1}} N^{-\frac{\kappa C'}{2\kappa C' + \kappa - 1}} L^{\frac{\kappa C'}{2\kappa C' + \kappa - 1}} \right) + \underset{N \rightarrow +\infty}{\mathfrak{o}}(\lambda^{-s}),$$

which, combined with equality (6.3), gives the desired result.

*Second subcase:  $\lambda = cN^{\frac{\kappa}{2\kappa-1}}$  for a small enough  $c$ .*

Then the previous computations can be adapted and we obtain that for any  $\beta > 0$  there exist  $c > 0$  and  $N_0 > 0$  such that for any  $N > N_0$  with  $\mathbb{P}^{\otimes N}$ -probability at least  $1 - \epsilon$ ,

$$\frac{C' - \beta}{\lambda} \leq \pi_{-\lambda r} R \leq \frac{C' + \beta}{\lambda}.$$

6.3.2. *Second case:  $\log \pi^{-1}(\Delta R \leq x) = C' x^{-\frac{q}{\kappa}} + C''' + \mathfrak{o}(1)$ . Since we have<sup>16</sup>*

$$(6.6) \quad \pi_{-\lambda R} \Delta R \underset{\lambda \rightarrow +\infty}{\sim} \left( \frac{qC'}{\kappa\lambda} \right)^{\frac{\kappa}{\kappa+q}},$$

from Theorem 2.2, for any  $0 \leq \chi < 1$ , we get

$$\begin{cases} \pi_{-\lambda r} R \leq \left( \frac{qC'}{\kappa\lambda} \right)^{\frac{\kappa}{\kappa+q}} \left[ 1 + \underset{N \rightarrow +\infty}{\mathfrak{o}}(1) + \underset{N \rightarrow +\infty}{\mathfrak{O}}(\chi) \right] + \frac{K(\pi_{-\lambda r}, \pi_{-\lambda R})}{\chi\lambda} \\ \pi_{-\lambda r} R \geq \left( \frac{qC'}{\kappa\lambda} \right)^{\frac{\kappa}{\kappa+q}} \left[ 1 + \underset{N \rightarrow +\infty}{\mathfrak{o}}(1) + \underset{N \rightarrow +\infty}{\mathfrak{O}}(\chi) \right] - \frac{K(\pi_{-\lambda r}, \pi_{-\lambda R})}{\chi\lambda} \end{cases}.$$

Taking  $\chi = \lambda^{-\frac{q}{2(\kappa+q)}} \sqrt{K(\pi_{-\lambda r}, \pi_{-\lambda R})}$ , we obtain

$$\pi_{-\lambda r} R = \left( \frac{qC'}{\kappa\lambda} \right)^{\frac{\kappa}{\kappa+q}} \left[ 1 + \underset{N \rightarrow +\infty}{\mathfrak{o}}(1) + \underset{N \rightarrow +\infty}{\mathfrak{O}} \left( \lambda^{-\frac{q}{2(\kappa+q)}} \sqrt{K(\pi_{-\lambda r}, \pi_{-\lambda R})} \right) \right].$$

*First subcase:  $\lambda = \mathfrak{o}(N^{-\frac{\kappa+q}{2\kappa-1+q}})$ . From Theorem 2.2, for any  $0 < \lambda \leq 0.19N$ , with  $\mathbb{P}^{\otimes N}$ -probability at least  $1 - \epsilon$ , we have*

$$K(\pi_{-\lambda r}, \pi_{-\lambda R}) \leq \log \pi_{-\lambda R} \exp \left\{ C \frac{\lambda^2}{N} (\Delta R)^{\frac{1}{\kappa}} \right\} + C \log(e\epsilon^{-1}).$$

We can prove<sup>17</sup> that for any  $\alpha \leq \check{c} \lambda^{-\frac{\kappa-1}{\kappa+q}}$  for  $\check{c}$  small enough, we have

$$(6.7) \quad \log \pi_{-\lambda R} \exp \left\{ \lambda \alpha (\Delta R)^{\frac{1}{\kappa}} \right\} = \underset{\lambda \rightarrow +\infty}{\mathfrak{O}} \left( \lambda^{\frac{q}{\kappa+q}} \alpha \lambda^{\frac{\kappa-1}{\kappa+q}} \right),$$

Let us assume that  $\lambda = \underset{N \rightarrow +\infty}{\mathfrak{o}} \left( N^{-\frac{\kappa+q}{2\kappa-1+q}} \right)$ . Then we have  $\frac{\lambda}{N} \lambda^{\frac{\kappa-1}{\kappa+q}} = \underset{N \rightarrow +\infty}{\mathfrak{o}}(1)$ , hence

$$K(\pi_{-\lambda r}, \pi_{-\lambda R}) \leq \underset{N \rightarrow +\infty}{\mathfrak{o}} \left( \lambda^{\frac{q}{\kappa+q}} \right) + C \log(e\epsilon^{-1}).$$

<sup>16</sup>See Appendix C.

<sup>17</sup>See Appendix D.

So we obtain that for  $\lambda = \underset{N \rightarrow +\infty}{o} \left( N^{-\frac{\kappa+q}{2\kappa-1+q}} \right)$ ,

$$\pi_{-\lambda r} R = \left( \frac{qC'}{\kappa\lambda} \right)^{\frac{\kappa}{\kappa+q}} \left[ 1 + \underset{N \rightarrow +\infty}{o} (1) \right].$$

*Second subcase:*  $\lambda = cN^{-\frac{\kappa+q}{2\kappa-1+q}}$  for a small enough  $c$ .

Once more, the previous computations can be adapted in order to obtain that for any  $\beta > 0$  there exist  $c > 0$  and  $N_0 > 0$  such that for any  $N > N_0$  with  $\mathbb{P}^{\otimes N}$ -probability at least  $1 - \epsilon$ ,

$$\left( \frac{qC' - \beta}{\kappa\lambda} \right)^{\frac{\kappa}{\kappa+q}} \leq \pi_{-\lambda r} R \leq \left( \frac{qC' + \beta}{\kappa\lambda} \right)^{\frac{\kappa}{\kappa+q}}.$$

**6.4. Proof of Theorem 3.6.** For any  $0 \leq j \leq \log N$ , introduce  $\lambda_j \triangleq 0.19\sqrt{N}e^{\frac{j}{2}}$ . Define  $L \triangleq \log[\log(eN)\epsilon^{-1}]$ . In [1, Section 3.4.2], an algorithm is proposed to choose the temperature of the standard Gibbs classifier. The associated generalization error is bounded by

$$\mathbb{G} \triangleq \min_{1 \leq j \leq \log N} \left\{ \pi_{-\lambda_{j-1}R} R + \frac{\sup_{0 \leq i \leq j} \left\{ \log \pi_{-\lambda_i R} \otimes \pi_{-\lambda_i R} \exp \left( \frac{C\lambda_i^2}{N} \mathbb{P}_{\cdot, \cdot} \right) \right\}}{\lambda_j} + C \frac{L}{\lambda_j} \right\}.$$

Under Assumptions (MA3) and (CA1), for any  $1 \leq j \leq \log N$  and  $t > 0$ , by Jensen's inequality, we have

$$\begin{aligned} \mathbb{G} &\leq -\frac{\log \pi \exp(-\lambda_{j-1}R)}{\lambda_{j-1}} + \frac{\sup_{0 \leq i \leq j} \left\{ \log \pi_{-\lambda_i R} \exp \left( \frac{C\lambda_i^2}{N} \mathbb{P}_{\cdot, f} \right) \right\}}{\lambda_j} + C \frac{L}{\lambda_j} \\ &\leq -\frac{\log \pi \exp(-\lambda_{j-1}R)}{\lambda_{j-1}} + \frac{\sup_{0 \leq i \leq j} \left\{ \log \pi \exp \left( -\lambda_i \Delta R + \frac{C\lambda_i^2}{N} (\Delta R)^{\frac{1}{\kappa}} \right) \right\}}{\lambda_j} \\ &\quad + \sup_{0 \leq i \leq j} \left\{ -\frac{\log \pi \exp(-\lambda_i \Delta R)}{\lambda_j} \right\} + C \frac{L}{\lambda_j} \\ &\leq R(\tilde{f}) - 2\sqrt{e} \frac{\log \pi \exp(-\lambda_j R)}{\lambda_j} + \frac{\sup_{0 \leq i \leq j; x \geq 0} \left\{ -\lambda_i x + \frac{C\lambda_i^2}{N} x^{\frac{1}{\kappa}} \right\}}{\lambda_j} + C \frac{L}{\lambda_j} \\ &\leq R(\tilde{f}) - 2\sqrt{e} \frac{\log[\pi(\Delta R \leq t) \exp(-\lambda_j t)]}{\lambda_j} + C \left( \frac{\lambda_j}{N} \right)^{\frac{\kappa}{\kappa-1}} + C \frac{L}{\lambda_j} \\ &\leq R(\tilde{f}) + C \frac{h_q(t^{1/\kappa})}{\lambda_j} + Ct + C \frac{L}{\lambda_j}. \end{aligned}$$

Taking  $j$  such that  $\lambda_j$  is of order  $N^{\frac{\kappa+q}{2\kappa-1+q}}$  and  $t$  minimizing  $C \frac{h_q(t^{1/\kappa})}{\lambda_j} + Ct$ , we obtain the desired rates (the ones given in Theorems 3.2 and 3.4). So the algorithm is adaptive wrt the margin parameter  $\kappa$ .

**6.5. Proof of Theorem 3.7.** We will prove the result for a minimal net. It is easy to generalize it to almost minimal nets. Let  $u > 0$ . Let  $\pi$  be the uniform distribution on a minimal  $(u, \mathcal{F}, \mathbb{P}_{\cdot, \cdot})$ -net denoted  $\mathcal{N}_u$ . Let  $\tilde{f}_u$  be the nearest neighbour of  $\tilde{f}$  in the net  $\mathcal{N}_u$ . Define  $a(\lambda) \triangleq \frac{\lambda}{N} g\left(\frac{\lambda}{N}\right)$ . From inequality (2.1) for  $(\rho_2, \pi_2, \rho_1, \pi_1) = (\delta_{\tilde{f}}, \pi, \delta_{\tilde{f}_u}, \delta_{\tilde{f}_u})$ , with  $(\mathbb{P}^{\otimes N})_*$ -probability at least  $1 - \epsilon$ , we have

$$R(\hat{f}) - R(\tilde{f}_u) + r(\tilde{f}_u) - r(\hat{f}) \leq a(\lambda) \mathbb{P}_{\tilde{f}, \tilde{f}_u} + \frac{H(u, \mathcal{F}, \mathbb{P}_{\cdot, \cdot}) + \log(\epsilon^{-1})}{\lambda},$$

When  $\hat{f} = \hat{f}_{\text{ERM}, u}$  minimizes the empirical risk over the net  $\mathcal{N}_u$ , we obtain

$$R(\hat{f}) - R(\tilde{f}_u) \leq a(\lambda) \mathbb{P}_{\tilde{f}, \tilde{f}_u} + \frac{H(u, \mathcal{F}, \mathbb{P}_{\cdot, \cdot}) + \log(\epsilon^{-1})}{\lambda},$$

hence

$$R(\hat{f}) - R(\tilde{f}) \leq R(\tilde{f}_u) - R(\tilde{f}) + a(\lambda)\mathbb{P}_{\tilde{f}, \tilde{f}} + a(\lambda)\mathbb{P}_{\tilde{f}, \tilde{f}_u} + \frac{H(u, \mathcal{F}, \mathbb{P}_{\cdot, \cdot}) + \log(\epsilon^{-1})}{\lambda}.$$

Let  $\check{C} > 0$  denote a constant (possibly depending on  $c'', C'', C', \check{C}_1$  and  $\check{C}_2$ ) whose value may differ from line to line. For any  $0 < \lambda \leq N$ , we have

$$(6.8) \quad \begin{aligned} \Delta R(\hat{f}) &\leq \check{C} \left( \mathbb{P}_{\tilde{f}, \tilde{f}_u}^\kappa + \frac{\lambda}{N} \Delta R_{\frac{1}{\kappa}}(\hat{f}) + \frac{\lambda}{N} \mathbb{P}_{\tilde{f}, \tilde{f}_u} + \frac{u^{-q} + \log(\epsilon^{-1})}{\lambda} \right) \\ &\leq \check{C} \left( \frac{\lambda}{N} \Delta R_{\frac{1}{\kappa}}(\hat{f}) + u^\kappa + \frac{\lambda}{N} u + \frac{u^{-q} + \log(\epsilon^{-1})}{\lambda} \right). \end{aligned}$$

Let us take  $u$  and  $\lambda$  such that  $\frac{\lambda}{N}u$ ,  $u^\kappa$  and  $\frac{u^{-q}}{\lambda}$  have the same orders. This is realized when Inequalities (3.6) hold and  $\lambda = N^{\frac{\kappa+q}{2\kappa-1+q}}$ . We obtain

$$\Delta R(\hat{f}) \leq \check{C} \left[ N^{-\frac{\kappa-1}{2\kappa-1+q}} \Delta R_{\frac{1}{\kappa}}(\hat{f}) + \log(\epsilon\epsilon^{-1}) N^{-\frac{\kappa}{2\kappa-1+q}} \right].$$

Simple computations lead to

$$\Delta R(\hat{f}) \leq \check{C} \log(\epsilon\epsilon^{-1}) N^{-\frac{\kappa}{2\kappa-1+q}}$$

and, then, to  $\mathbb{P}^{\otimes N} \Delta R(\hat{f}) \leq \check{C} N^{-\frac{\kappa}{2\kappa-1+q}}$ .

**6.6. Proof of Theorem 3.9.** The chaining idea comes from [6] and is well presented also, for instance, in [5, p.19-21]. Let  $\partial(\rho_1, \rho_2) \triangleq \rho_2 R - \rho_1 R + \rho_1 r - \rho_2 r$ . Let  $u_k = u2^k$ . Let  $c_k \triangleq h_+(u_k)$ . To shorten, denote  $\pi_{i,k} \triangleq \pi_{f_i, u_k}$ . Let  $K$  be the nonnegative integer such that  $\frac{\mathbb{P}_{f_1, f_2}}{2} \leq u_K < \mathbb{P}_{f_1, f_2}$ . The integer  $K$  exists as soon as  $\mathbb{P}_{f_1, f_2} > u$ . Let  $L'$  be the nonnegative integer such that  $\frac{1}{2} \leq u_{L'} < 1$ . Let  $\lambda_1, \dots, \lambda_{L'+1}$  be real positive parameters to be chosen. We apply inequality (2.1) for this  $L'+1$  parameters and for  $\pi_1 = \pi_2 = \pi$ .

With  $(\mathbb{P}^{\otimes N})_*$ -probability at least  $1 - (L'+1)\epsilon$ , we have

$$\begin{aligned} &\partial(\pi_{1,0}, \pi_{2,0}) \\ &= \partial(\pi_{1,K}, \pi_{2,K}) + \sum_{k=1}^K \left\{ \partial(\pi_{1,k-1}, \pi_{1,k}) + \partial(\pi_{2,k}, \pi_{2,k-1}) \right\} \\ &\leq \frac{\lambda_{K+1}}{N} g\left(\frac{\lambda_{K+1}}{N}\right) 4u_K + \frac{K(\pi_{1,K}, \pi) + K(\pi_{2,K}, \pi) + \log(\epsilon^{-1})}{\lambda_{K+1}} \\ &\quad + \sum_{k=1}^K \left\{ 2\frac{\lambda_k}{N} g\left(\frac{\lambda_k}{N}\right) (u_{k-1} + u_k) + \frac{2\log(\epsilon^{-1}) + \sum_{i=1}^2 \sum_{k'=k-1}^k K(\pi_{i,k'}, \pi)}{\lambda_k} \right\} \\ &\leq \frac{\lambda_{K+1}}{N} g\left(\frac{\lambda_{K+1}}{N}\right) 4u_K + \frac{2c_K + \log(\epsilon^{-1})}{\lambda_{K+1}} \\ &\quad + \sum_{k=1}^K \left\{ 6\frac{\lambda_k}{N} g\left(\frac{\lambda_k}{N}\right) u_{k-1} + \frac{2\log(\epsilon^{-1}) + 4c_{k-1}}{\lambda_k} \right\} \\ &\leq \sum_{k=1}^{K+1} \left\{ 6\frac{\lambda_k}{N} g\left(\frac{\lambda_k}{N}\right) u_{k-1} + \frac{2\log(\epsilon^{-1}) + 4c_{k-1}}{\lambda_k} \right\}. \end{aligned}$$

Let us choose the  $\lambda_k$ 's such that they do not depend on  $\epsilon$  and they roughly minimize the RHS of the last bound. Taking  $\lambda_k = \sqrt{\frac{4Nc_{k-1}}{3u_{k-1}}}$  for  $k \geq 1$ , we obtain

$$\begin{aligned} \partial(\pi_{1,0}, \pi_{2,0}) &\leq \sum_{k=1}^{K+1} \left[ 1 + 2g\left(\frac{\lambda_k}{N}\right) \right] \sqrt{\frac{12c_{k-1}u_{k-1}}{N}} + \sum_{k=1}^{K+1} \frac{2\log(\epsilon^{-1})}{\lambda_k} \\ &\leq \sum_{k=1}^{K+1} \left[ 1 + 2g\left(\frac{\lambda_k}{N}\right) \right] \sqrt{\frac{12c_{k-1}u_{k-1}}{N}} + 2\log(\epsilon^{-1}) \sqrt{\frac{3u}{4N}} \sum_{k=1}^{K+1} \sqrt{2}^{k-1}. \end{aligned}$$

For any  $k \in \{1, \dots, L'+1\}$ , we have  $\frac{\lambda_k}{N} \leq \sqrt{\frac{4c_0}{3Nu}} \triangleq C_1$ . Since  $C_0 = 2\sqrt{3}[1+2g(C_1)]$ , we obtain

$$\begin{aligned} \partial(\pi_{1,0}, \pi_{2,0}) &\leq \frac{2C_0}{\sqrt{N}} \sum_{k=1}^{K+1} (u_{k-1} - u_{k-2}) \sqrt{\frac{c_{k-1}}{u_{k-1}}} + \sqrt{\frac{6\mathbb{P}_{f_1, f_2}}{N} \frac{\log(\epsilon^{-1})}{\sqrt{2-1}}} \\ &\leq \frac{2C_0}{\sqrt{N}} \int_{u/2}^{\mathbb{P}_{f_1, f_2}} \sqrt{\frac{h_+(v)}{v}} dv + 6\sqrt{\frac{\mathbb{P}_{f_1, f_2}}{N}} \log(\epsilon^{-1}). \end{aligned}$$

*Remark 6.3.* We have used the ‘‘global bayesian entropy’’  $h(v) \triangleq \sup_{f \in \mathcal{F}} \log \pi^{-1}(B_{f,v})$

since it was convenient to have an (almost) optimal  $\lambda$ 's which do not depend on the functions  $f_1$  and  $f_2$ . Had we done a union bound on the parameters  $\lambda$ , we would have been able to make it depend on the functions  $f_1, f_2$ . Then the global bayesian entropy would have been replaced with the local ones  $h(v, f_1)$  and  $h(v, f_2)$  where  $h(v, f) \triangleq \log \pi^{-1}(B_{f,v})$ . In other words, the quantity  $\partial(\pi_{1,0}, \pi_{2,0})$  is mainly driven by the two integrals  $\int_{u/2}^{\mathbb{P}_{f_1, f_2}} \sqrt{\frac{h(v, f_1)}{vN}} dv$  and  $\int_{u/2}^{\mathbb{P}_{f_1, f_2}} \sqrt{\frac{h(v, f_2)}{vN}} dv$ .

### 6.7. Proof of Theorem 3.10.

6.7.1. *First step: upper bounds due to the chaining technique.* We start with the following chained result which is slightly different from Theorem 3.9 to the extent that we chained functions belonging to covering nets instead of chaining balls. Had we been interested in results for packing nets, Theorem 3.9 applied to an appropriate prior distribution<sup>18</sup> would have been sufficient. Let  $H(u) \triangleq H(u, \mathcal{F}, \mathbb{P}_{\cdot, \cdot})$ .

**Theorem 6.5.** *Let  $u > 0$ ,  $\mathcal{N}$  a minimal  $u$ -covering net and  $L \triangleq \frac{\log(2u^{-1})}{\log 2}$ . We have*

- for any  $\epsilon > 0$ , with  $\mathbb{P}^{\otimes N}$ -probability at least  $1 - \epsilon$ , for any  $f_1, f_2 \in \mathcal{N}_u$ ,

$$(6.9) \quad \begin{aligned} R(f_2) - R(f_1) + r(f_1) - r(f_2) &\leq 8\sqrt{\frac{3}{N}} \int_{u/2}^{\mathbb{P}_{f_1, f_2} \vee u} \sqrt{\frac{H(v)}{v}} dv + \frac{8}{3N} \int_{u/2}^{\mathbb{P}_{f_1, f_2} \vee u} \frac{H(v)}{v} dv \\ &\quad + 17\sqrt{\frac{\log(3L\epsilon^{-1})}{N}} \sqrt{\mathbb{P}_{f_1, f_2} \vee u} + \frac{2L \log(3L\epsilon^{-1})}{3N}, \end{aligned}$$

- for any  $\epsilon > 0$ , for any  $f_1 \in \mathcal{N}_u$ , with  $\mathbb{P}^{\otimes N}$ -probability at least  $1 - \epsilon$ , for any  $f_2 \in \mathcal{N}_u$ ,

$$\begin{aligned} R(f_2) - R(f_1) + r(f_1) - r(f_2) &\leq 4\sqrt{\frac{3}{N}} \int_{u/2}^{\mathbb{P}_{f_1, f_2} \vee u} \sqrt{\frac{H(v)}{v}} dv + \frac{4}{3N} \int_{u/2}^{\mathbb{P}_{f_1, f_2} \vee u} \frac{H(v)}{v} dv \\ &\quad + 8.5\sqrt{\frac{\log(2L\epsilon^{-1})}{N}} \sqrt{\mathbb{P}_{f_1, f_2} \vee u} + \frac{L \log(2L\epsilon^{-1})}{3N} \end{aligned}$$

*Proof.* The proof is similar to the one of Theorem 3.9. Instead of chaining balls, we will chain on covering nets. Let  $\partial(f_1, f_2) \triangleq R(f_2) - R(f_1) + r(f_1) - r(f_2)$ ,  $u_k = u2^k$  and  $c_k \triangleq H_+(u_k)$ . Introduce  $P \triangleq \mathbb{P}_{f_1, f_2} \vee u$  and let  $0 \leq K \leq L'$  be integers such that  $\frac{\mathbb{P}_{f_1, f_2}}{2} < u_K \leq P$  and  $\frac{1}{2} < u_{L'} \leq 1$ .

Consider the family  $(\mathcal{N}_k)_{k=\{0, \dots, L'\}}$  of minimal nets of radius  $u_k$ . For any  $(j, k) \in \{1, 2\} \times \{0, \dots, L'\}$ , introduce  $f_{j,k} \in \operatorname{argmin}_{\mathcal{N}_{u_{2^k}}} \mathbb{P}_{\cdot, f_j}$  a nearest neighbour of  $f_j$  in

<sup>18</sup>Let  $\mathcal{N}_p$  be a  $u$ -packing net. Using the notation of Section 6.6, an appropriate prior distribution is  $\pi = \frac{1}{L'+1} \sum_{k=0}^{L'} \pi_k$ , where  $\pi_k$  is the uniform distribution on a  $u_k$ -minimal packing net of the set  $\mathcal{F}$  built using points in  $\mathcal{N}_p$ . The log-cardinal of such a set is upper bounded by  $H(u_{k-1}, \mathcal{F}, \mathbb{P}_{\cdot, \cdot})$ , hence  $h(u_k) \leq H(u_{k-1}, \mathcal{F}, \mathbb{P}_{\cdot, \cdot}) + \log(L'+1)$ .

$\mathcal{N}_{u2^k}$ . Since  $f_1, f_2 \in \mathcal{N}_u$ , we have  $f_{1,0} = f_1$  and  $f_{2,0} = f_2$ . Let  $\pi_k$  be the uniform distribution on the net  $\mathcal{N}_k$ .

Let  $l \triangleq \log[(3L' + 1)\epsilon^{-1}]$ . By applying  $3L' + 1$  times inequality (2.3), we obtain that with  $\mathbb{P}^{\otimes N}$ -probability at least  $1 - \epsilon$ , for any functions  $f_1, f_2$  in  $\mathcal{N}_u$ , we have

$$\begin{aligned}
(6.10) \quad \partial(f_1, f_2) &= \partial(f_{1,K}, f_{2,K}) + \sum_{k=1}^K \{ \partial(f_{1,k-1}, f_{1,k}) + \partial(f_{2,k}, f_{2,k-1}) \} \\
&\leq \sqrt{\frac{2[2H(u_K)+l]\mathbb{P}_{f_1, f_2, K}}{N} + \frac{2H(u_K)+l}{3N}} \\
&\quad + \sum_{k=1}^K \left\{ \sqrt{\frac{2[H(u_{k-1})+H(u_k)+l]\mathbb{P}_{f_{1,k-1}, f_{1,k}}}{N} + \frac{H(u_{k-1})+H(u_k)+l}{3N}} \right. \\
&\quad \left. + \sqrt{\frac{2[H(u_{k-1})+H(u_k)+l]\mathbb{P}_{f_{2,k-1}, f_{2,k}}}{N} + \frac{H(u_{k-1})+H(u_k)+l}{3N}} \right\} \\
&\leq 2 \sum_{k=1}^{K+1} \left\{ \sqrt{\frac{6[2H(u_{k-1})+l]u_{k-1}}{N} + \frac{2H(u_{k-1})+l}{3N}} \right\} \\
&\leq 2 \sum_{k=1}^{K+1} \left\{ \sqrt{\frac{12H(u_{k-1})u_{k-1}}{N}} + \sqrt{\frac{6lu_{k-1}}{N} + \frac{2H(u_{k-1})}{3N}} \right\} + \frac{2(K+1)l}{3N} \\
&\leq 2\sqrt{\frac{6lu}{N}} \frac{\sqrt{2^{K+1}}}{\sqrt{2-1}} + \frac{2(K+1)l}{3N} + 2 \sum_{k=1}^{K+1} \left\{ \sqrt{\frac{12H(u_{k-1})u_{k-1}}{N}} + \frac{2H(u_{k-1})}{3N} \right\}
\end{aligned}$$

Now the last sum can be upper bounded using integrals since the function  $v \mapsto H(v)$  is non increasing on  $\mathbb{R}_+^*$ . We obtain

$$\partial(f_1, f_2) \leq \frac{4}{\sqrt{2-1}} \sqrt{\frac{3lP}{N}} + \frac{2l}{3N} \frac{\log(\frac{2P}{u})}{\log 2} + 8\sqrt{\frac{3}{N}} \int_{u/2}^P \sqrt{\frac{H(v)}{v}} dv + \frac{8}{3N} \int_{u/2}^P \frac{H(v)}{v} dv.$$

For the second part of Theorem 6.5, it suffices to modify slightly the previous argument. This time, the functions  $f_{2,k}$  are defined as previously. The functions  $f_{1,k}$  are defined as  $f_{1,k} \triangleq f_1$ . Therefore we have  $\partial(f_{1,k-1}, f_{1,k}) = 0$ , hence the modification of the constants.  $\square$

Consider that Assumption (CA1) holds. Let  $c_q > 0$  such that for any  $0 < u \leq 1$ ,  $\sum_{k=0}^{L'} 3e^{-c_q h_q(u_k)} \leq 1$ . In the previous proof, we used an uniform union bound over the  $3L' + 1$  inequalities coming from (2.3). If we are just interested in the order of the bounds, we can weight the inequalities associated with  $\partial(f_{1,k-1}, f_{1,k})$  and  $\partial(f_{2,k}, f_{2,k-1})$  with  $e^{-c_q h_q(u_{k-1})}$  and those corresponding to  $\partial(f_{2,k}, f_{2,k-1})$  with at least weight  $e^{-c_q h_q(u_k)}$ .

Then, in Inequalities (6.10), we may replace  $2H(u_K) + l$  and  $H(u_{k-1}) + H(u_k) + l$  with respectively  $2H(u_K) + c_q h_q(u_K) + \log(\epsilon^{-1})$  and  $H(u_{k-1}) + H(u_k) + c_q h_q(u_{k-1}) + \log(\epsilon^{-1})$ , so that we obtain

$$\begin{aligned}
(6.11) \quad \partial(f_1, f_2) &\leq \check{C} \sum_{k=1}^{K+1} \left\{ \sqrt{\frac{[h_q(u_{k-1}) + \log(\epsilon^{-1})]u_{k-1}}{N} + \frac{h_q(u_{k-1}) + \log(\epsilon^{-1})}{N}} \right\} \\
&\leq \check{C} \sqrt{\frac{P \log(\epsilon^{-1})}{N}} + \check{C} \frac{\log(eu^{-1}) \log(\epsilon^{-1})}{N} + \check{C} \int_{u/2}^P \left( \sqrt{\frac{h_q(v)}{Nv}} + \frac{h_q(v)}{Nv} \right) dv.
\end{aligned}$$

**Corollary 6.6.** *Let  $\mathcal{N}$  denote a minimal  $u$ -net, where  $u$  is a positive real. Define  $U \triangleq \sup_{f: \mathbb{P}_{f, \tilde{f}} \leq u} \{R(f) - R(\tilde{f})\}$ . Introduce a function  $\tilde{f}_{\mathcal{N}} \in \mathcal{N}$  such that  $\mathbb{P}_{\tilde{f}_{\mathcal{N}}, \tilde{f}} \leq u$ . Let  $\gamma_u : [u; 1] \rightarrow \mathbb{R}$  and  $\Gamma_u : [u; 1] \rightarrow \mathbb{R}$  be non decreasing concave functions respectively upper bounding the functions  $\int_{\frac{u}{2}}^{\cdot} \sqrt{\frac{H(v)}{v}} dv$  and  $\int_{\frac{u}{2}}^{\cdot} \frac{H(v)}{v} dv$ .*

For any  $\epsilon > 0$ , with  $\mathbb{P}^{\otimes N}$ -probability at least  $1 - \epsilon$ , for any function  $f \in \mathcal{N}$ , we have

$$\begin{aligned} R(f) - R(\tilde{f}) &\leq r(f) - r(\tilde{f}_{\mathcal{N}}) + \frac{\check{C}}{\sqrt{N}} \left\{ \gamma_u(\mathbb{P}_{f, \tilde{f}} + u) + \sqrt{(\mathbb{P}_{f, \tilde{f}} + u) \log(\epsilon^{-1})} \right\} \\ &\quad + \frac{\check{C}}{N} \left\{ \Gamma_u(\mathbb{P}_{f, \tilde{f}} + u) + \log(eu^{-1}) \log(\epsilon^{-1}) \right\} + \mathcal{U} \end{aligned}$$

Consequently, for any  $\rho \in \mathcal{M}_+^1(\mathcal{N})$ , we have

$$\begin{aligned} \rho R - R(\tilde{f}) &\leq \rho r - r(\tilde{f}_{\mathcal{N}}) + \frac{\check{C}}{\sqrt{N}} \left\{ \gamma_u(\rho \mathbb{P}_{\cdot, \tilde{f}} + u) + \sqrt{(\rho \mathbb{P}_{\cdot, \tilde{f}} + u) \log(\epsilon^{-1})} \right\} \\ &\quad + \frac{\check{C}}{N} \left\{ \Gamma_u(\rho \mathbb{P}_{\cdot, \tilde{f}} + u) + \log(eu^{-1}) \log(\epsilon^{-1}) \right\} + \mathcal{U} \end{aligned}$$

*Proof.* The first inequality comes mainly from inequalities (6.11) and the decomposition:  $R(f) - R(\tilde{f}) = R(f) - R(\tilde{f}_{\mathcal{N}}) + R(\tilde{f}_{\mathcal{N}}) - R(\tilde{f})$ . The second inequality is then deduced from Jensen's inequality.  $\square$

6.7.2. *Second step: determining the radius of the nets.* Corollary 6.6 implies that for any  $\epsilon > 0$ , for any classifier  $\hat{f}$  minimizing the empirical risk over the net  $\mathcal{N}$ , with  $(\mathbb{P}^{\otimes N})_*$ -probability at least  $1 - \epsilon$ , we have

$$(6.12) \quad \begin{aligned} \Delta R(\hat{f}) &\leq \mathcal{U} + \frac{\check{C}}{\sqrt{N}} \left\{ \gamma_u(\mathbb{P}_{\hat{f}, \tilde{f}} + u) + \sqrt{(\mathbb{P}_{\hat{f}, \tilde{f}} + u) \log(\epsilon^{-1})} \right\} \\ &\quad + \frac{\check{C}}{N} \left\{ \Gamma_u(\mathbb{P}_{\hat{f}, \tilde{f}} + u) + \log(eu^{-1}) \log(\epsilon^{-1}) \right\}. \end{aligned}$$

Now we have

$$\mathcal{U} \leq \begin{cases} \check{C}u^\kappa & \text{under Assumption (MA}\beta\text{)} \\ 2u & \text{in any case} \end{cases},$$

and we can take

$$\gamma_u(x) \triangleq \begin{cases} \check{C} \sqrt{\log(e^2 x^{-1})} x & \text{under Assumption (CA1) for } q = 0 \\ \check{C} x^{\frac{1-q}{2}} & \text{under Assumption (CA1) for } 0 < q < 1 \\ \check{C} \log\left(\frac{2x}{u}\right) & \text{under Assumption (CA1) for } q = 1 \\ \check{C} u^{\frac{1-q}{2}} & \text{under Assumption (CA1) for } q > 1 \end{cases}$$

and

$$\Gamma_u(x) \triangleq \begin{cases} \check{C} [\log(eu^{-1})]^2 & \text{under Assumption (CA1) for } q = 0 \\ \check{C} u^{-q} & \text{under Assumption (CA1) for } q > 0 \end{cases}.$$

Then we have eight cases corresponding to the different complexity and margin assumptions. When we have  $q > 0$ , inequality (6.12) implies

$$\Delta R(\hat{f}) \leq \mathcal{U} + \check{C} \frac{\log(e\epsilon^{-1})}{\sqrt{N}} \gamma_u(\mathbb{P}_{\hat{f}, \tilde{f}} + u) + \check{C} \frac{\log(e\epsilon^{-1})}{N} u^{-q}.$$

Under Assumptions (MA2) and (CA1) for  $q = 0$

Let  $\Delta \triangleq \Delta R(\hat{f})$  to shorten. Inequality (6.12) becomes

$$\Delta \leq \check{C} \left[ \log(e\epsilon^{-1}) N^{-\frac{1}{2}} \left( \sqrt{\log(e^2 \Delta^{-\frac{1}{\kappa}})} \Delta^{\frac{1}{2\kappa}} + \sqrt{\log(e^2 u^{-1})} u \right) + u + \frac{[\log(eu^{-1})]^2}{N} \right].$$

We obtain  $\Delta \leq \check{C} \log(e\epsilon^{-1}) (\log[eN^{1/\kappa}])^{\frac{\kappa}{2\kappa-1}} N^{-\frac{\kappa}{2\kappa-1}}$  when<sup>19</sup>

$$\sqrt{\frac{\log(e^2 u^{-1})}{N}} + u + \frac{[\log(eu^{-1})]^2}{N} \leq \check{C} (\log[eN^{1/\kappa}])^{\frac{\kappa}{2\kappa-1}} N^{-\frac{\kappa}{2\kappa-1}},$$

<sup>19</sup>We use  $\log[eN^{1/\kappa}]$  since the logarithmic factor disappears for  $\kappa = +\infty$ . For  $\kappa < +\infty$ , the factor  $\log[eN^{1/\kappa}]$  can be simplified into  $\log N$  for  $N \geq 2$ .



hence when there exists  $\check{C}_1, \check{C}_2 > 0$  such that

$$\exp \left\{ -\check{C}_1 (\log[eN^{1/\kappa}])^{\frac{\kappa}{4\kappa-2}} N^{\frac{\kappa-1}{4\kappa-2}} \right\} \leq u \leq \check{C}_2 (\log[eN^{1/\kappa}])^{\frac{\kappa}{2\kappa-1}} N^{-\frac{\kappa}{2\kappa-1}}.$$

Under Assumptions (MA3) and (CA1) for  $q = 0$  and  $\kappa < +\infty$

Inequality (6.12) gives

$$\Delta \leq \check{C} \left[ \log(e\epsilon^{-1}) N^{-\frac{1}{2}} \left( \sqrt{\log(e^2 \Delta^{-\frac{1}{\kappa}}) \Delta^{\frac{1}{2\kappa}}} + \sqrt{\log(e^2 u^{-1}) u} \right) + u^\kappa + \frac{[\log(eu^{-1})]^2}{N} \right].$$

We obtain  $\Delta \leq \check{C} \log(e\epsilon^{-1}) (\log N)^{\frac{\kappa}{2\kappa-1}} N^{-\frac{\kappa}{2\kappa-1}}$  when

$$\sqrt{\frac{\log(e^2 u^{-1}) u}{N}} + u^\kappa + \frac{[\log(eu^{-1})]^2}{N} \leq \check{C} (\log N)^{\frac{\kappa}{2\kappa-1}} N^{-\frac{\kappa}{2\kappa-1}},$$

so when there exists  $\check{C}_1, \check{C}_2 > 0$  such that

$$\exp \left\{ -\check{C}_1 (\log N)^{\frac{\kappa}{4\kappa-2}} N^{\frac{\kappa-1}{4\kappa-2}} \right\} \leq u \leq \check{C}_2 (\log N)^{\frac{1}{2\kappa-1}} N^{-\frac{1}{2\kappa-1}}.$$

Under Assumptions (MA2) and (CA1) for  $0 < q < 1$

Inequality (6.12) becomes

$$\Delta \leq \check{C} \left[ u + \log(e\epsilon^{-1}) N^{-\frac{1}{2}} \left( \Delta^{\frac{1-q}{2\kappa}} + u^{\frac{1-q}{2}} \right) + \check{C} \frac{\log(e\epsilon^{-1})}{N} u^{-q} \right].$$

This leads to  $\Delta \leq \check{C} \log(e\epsilon^{-1}) N^{-\frac{\kappa}{2\kappa-1+q}}$  when the inequality

$$N^{-\frac{1}{2}} u^{\frac{1-q}{2}} + u + \check{C} \frac{\log(e\epsilon^{-1})}{N} u^{-q} \leq \check{C} N^{-\frac{\kappa}{2\kappa-1+q}}$$

holds, hence when there exist  $\check{C}_1, \check{C}_2$  such that  $\check{C}_1 N^{-\frac{\kappa-1+q}{q(2\kappa-1+q)}} \leq u \leq \check{C}_2 N^{-\frac{\kappa}{2\kappa-1+q}}$ .

Similarly, we deal with the five other cases. To finish the proof, we just have to notice that, when for any  $\epsilon > 0$  and some real function  $\phi$ , with  $\mathbb{P}^{\otimes N}$ -probability at least  $1 - \epsilon$ , we have  $\Delta \leq \log(e\epsilon^{-1}) \phi(N)$ , then we have  $\mathbb{P}^{\otimes N} \Delta \leq 2\phi(N)$ .

*Remark 6.4.* Once more, for sake of simplicity, we have done the proof for minimal nets without explicit values of the constants. It is easy to adapt the proof to almost minimal nets and to get an explicit constant  $\check{C}$  in terms of the other constants of the problem.

**6.8. Proof of Theorem 3.15.** Let  $u_N = \check{C}_2 w_N$ . Let  $\mathcal{N}'$  be a  $u_N$ -minimal bracketing net of the model  $\mathcal{F}$ . Let  $A \triangleq \{f \in \mathcal{F} : \mathbb{P}_{f, \bar{f}} \leq u_N\}$ . There exists a posterior distribution  $\hat{\rho}_{\mathcal{N}'} : \mathcal{Z}^N \rightarrow \mathcal{M}_+^1(\mathcal{N}')$  (for instance,  $\hat{\rho}_{\mathcal{N}'} \triangleq \pi_{-\lambda_{N'} r} \circ n_L^{-1}$ ) such that we have  $\hat{\rho}_{\mathcal{N}'} r \leq \pi_{-\lambda_{N'} r} r + u_N$  and

$$(6.13) \quad \pi_{-\lambda_{N'} r} R \leq \hat{\rho}_{\mathcal{N}'} R + u_N.$$

We have

$$(6.14) \quad \hat{\rho}_{\mathcal{N}'} r \leq \pi_{-\lambda_{N'} r} r + u_N \leq \pi|_A r + \frac{K(\pi|_A, \pi)}{\lambda_N} + u_N$$

and

$$(6.15) \quad K(\pi|_A, \pi) = \log[\pi(A)^{-1}] \leq \check{C} h_q(u_N) \leq \check{C} \lambda_N u_N.$$

From inequality (2.1) for  $(\rho_2, \pi_2, \rho_1, \pi_1) = (\delta_{\tilde{f}}, \delta_{\tilde{f}}, \pi|_A, \pi|_A)$  and  $\lambda = N$ , with  $\mathbb{P}^{\otimes N}$ -probability at least  $1 - \epsilon$ , we have  $\pi|_{Ar} - r(\tilde{f}) \leq \pi|_A R - R(\tilde{f}) + u_N + \frac{\log(\epsilon^{-1})}{N}$ , hence

$$(6.16) \quad \pi|_{Ar} - r(\tilde{f}) \leq 2u_N + \frac{\log(\epsilon^{-1})}{N} \leq C \log(\epsilon \epsilon^{-1}) w_N.$$

Combining Inequalities (6.14), (6.15) and (6.16), with  $\mathbb{P}^{\otimes N}$ -probability at least  $1 - \epsilon$ , we obtain  $\hat{\rho}_{\mathcal{N}'r} \leq r(\tilde{f}) + \check{C} \log(\epsilon \epsilon^{-1}) w_N$ . The result follows from Theorem 3.12 and inequality (6.13).

**6.9. Proof of Theorem 3.16.** Let  $\check{C}_1 > 0$ ,  $u \geq \check{C}_1 \beta_q$  and  $\mathcal{N}$  be a  $(u, \mathcal{F}, \mathbb{P})$ -minimal bracketing net. Let  $\pi$  be the uniform distribution on this net. From inequality (8.2) in [1] for  $\mathcal{W}((f_1, f_2), X) = \mathbb{1}_{f_1(X) \neq f_2(X)}$  and  $\nu = \pi \otimes \pi$ , we obtain that with  $\mathbb{P}^{\otimes N}$ -probability at least  $1 - \epsilon$ , for any function  $f'_1, f'_2$  in the net  $\mathcal{N}$  and any  $\lambda > 0$ , we have

$$\bar{\mathbb{P}}_{f'_1, f'_2} \leq \left[1 + \frac{\lambda}{N} g\left(\frac{\lambda}{N}\right)\right] \mathbb{P}_{f'_1, f'_2} + \frac{2 \log |\mathcal{N}| + \log(\epsilon^{-1})}{\lambda}$$

Recall that  $\check{C}$  is a constant (possibly depending on the other constants of the problem) which value may differ from line to line. Taking  $\epsilon = (\alpha_q)^C$  and  $\lambda = N$ , we obtain that with probability at least  $1 - (\alpha_q)^C$ , for any functions  $f_1, f_2$  in the set  $\mathcal{F}$ , we have

$$\begin{aligned} \bar{\mathbb{P}}_{f_1, f_2} &\leq \bar{\mathbb{P}}_{f_1, n_L(f_1)} + \bar{\mathbb{P}}_{n_L(f_1), n_L(f_2)} + \bar{\mathbb{P}}_{n_L(f_2), f_2} \\ &\leq \bar{\mathbb{P}}_{n_L(f_1), n_U(f_1)} + \bar{\mathbb{P}}_{n_L(f_1), n_L(f_2)} + \bar{\mathbb{P}}_{n_L(f_2), n_U(f_2)} \\ &\leq (e-1) \mathbb{P}_{n_L(f_1), n_L(f_2)} + \frac{6C' h_q(\check{C}_1 \beta_q) + 3C \log(\alpha_q^{-1})}{N} + 2u \\ &\leq (e-1) \mathbb{P}_{n_L(f_1), n_L(f_2)} + \check{C} \beta_q + 2u \\ &\leq (e-1) \mathbb{P}_{f_1, f_2} + \check{C} u. \end{aligned}$$

By applying inequality (8.2) in [1] to  $\mathcal{W}((f_1, f_2), X) = -\mathbb{1}_{f_1(X) \neq f_2(X)}$ , we can similarly prove that with probability at least  $1 - (\alpha_q)^C$ , for any functions  $f_1, f_2$  in the set  $\mathcal{F}$ , we have  $(3-e) \mathbb{P}_{f_1, f_2} \leq \bar{\mathbb{P}}_{f_1, f_2} + \check{C} u$ . (The constants  $e-1$  and  $3-e$  have nothing fundamental and we can make them as close as 1 as we want provided that we change the other constants.) These two inequalities allows to prove respectively the first two items of the theorem for one radius  $u \geq \check{C}_1 \beta_q$ . To get a result uniform wrt the radius, it suffices to make a union bound for radius in a geometric grid of  $[\check{C}_1 \beta_q; 1]$ .

For the last item, when we have  $u \geq \check{C}_1 \beta_q$  for a sufficiently large  $\check{C}_1$ , there exists a small constant  $\check{C}'$  satisfying

$$H(u, \mathcal{F}, \bar{\mathbb{P}}) \leq H_p(u, \mathcal{F}, \bar{\mathbb{P}}) \leq H_p(\check{C}' u, \mathcal{F}, \mathbb{P}) \leq H(\check{C}' u/2, \mathcal{F}, \mathbb{P}) \leq \check{C} h_q(u).$$

**6.10. Proof of Theorem 4.1.** The proof is adapted from the proof of the concentration of  $N(X_1^N)$  [10, p.42]. First, we prove that for any  $k \in \mathbb{N}$ , the quantity  $\frac{H(\frac{k}{N}, \mathcal{F}, \bar{\mathbb{P}}(\cdot, \cdot))}{\log 2}$  is a self-bounded quantity in the sense given in [10, p.23]. Let  $\mathcal{N}_k$  be a  $(\frac{k}{N}, \mathcal{F}, \bar{\mathbb{P}})$ -minimal net. Define the probability distribution

$$\bar{\mathbb{P}}^{(i)} \triangleq \frac{\delta_{z_1} + \dots + \delta_{z_{i-1}} + \delta_{z_{i+1}} + \dots + \delta_{z_N}}{N-1}.$$

Let  $H^{(i)}$  be the logarithm of the cardinal of the smallest  $(\frac{k}{N-1}, \mathcal{F}, \bar{\mathbb{P}}^{(i)})$ -net using only functions in the net  $\mathcal{N}_k$ . We have

$$0 \leq H\left(\frac{k}{N}, \mathcal{F}, \bar{\mathbb{P}}(\cdot, \cdot)\right) - H^{(i)} \leq H\left(\frac{k}{N}, \mathcal{F}, \bar{\mathbb{P}}(\cdot, \cdot)\right) - H\left(\frac{k}{N-1}, \mathcal{F}, \bar{\mathbb{P}}^{(i)}\right) \leq \log 2.$$

Let  $\mathbb{V} = (f(X_1), \dots, f(X_N))$  be the random vector induced by the uniform distribution on the net  $\mathcal{N}_k$ . The Shannon entropy of this vector is  $\log |\mathcal{N}_k| = H(\frac{k}{N}, \mathcal{F}, \bar{\mathbb{P}}_{\cdot, \cdot})$ . Define  $\mathbb{V}^{(i)} \triangleq (f(X_1), \dots, f(X_{i-1}), f(X_{i+1}), \dots, f(X_N))$ . Since the uniform distribution maximizes the entropy, we have  $H^{(i)} \geq H(\mathbb{V}^{(i)})$ . Then, using Han's inequality (see for instance [10, p.31]), we obtain that  $\frac{H(\frac{k}{N}, \mathcal{F}, \bar{\mathbb{P}}_{\cdot, \cdot})}{\log 2}$  is a self-bounded quantity.

Therefore, we can apply Corollary 5 in [10, p.43]. To shorten, we write until the end of the proof  $H(k)$  for  $H(\frac{k}{N}, \mathcal{F}, \bar{\mathbb{P}}_{\cdot, \cdot})$ . For any  $\eta > 0$ , we have

$$(6.17) \quad \mathbb{P}^{\otimes N} \left[ H(k) \geq \mathbb{P}^{\otimes N} H(k) + (\log 2)\eta \right] \leq e^{-\frac{\eta^2}{\frac{2\mathbb{P}^{\otimes N} H(k)}{\log 2} + \frac{2\eta}{3}}}$$

and

$$(6.18) \quad \mathbb{P}^{\otimes N} \left[ H(k) \leq \mathbb{P}^{\otimes N} H(k) - (\log 2)\eta \right] \leq e^{-\frac{\eta^2}{2\mathbb{P}^{\otimes N} H(k)}}.$$

Introducing  $\epsilon = e^{-\frac{\eta^2}{\frac{2\mathbb{P}^{\otimes N} H(k)}{\log 2} + \frac{2\eta}{3}}}$ , equivalently

$$\begin{aligned} \eta^2 - \left( \frac{2\mathbb{P}^{\otimes N} H(k)}{\log 2} + \frac{2\eta}{3} \right) \log(\epsilon^{-1}) &= 0, \\ \eta &= \frac{\log(\epsilon^{-1})}{3} \left( 1 + \sqrt{1 + \frac{18\mathbb{P}^{\otimes N} H(k)}{(\log 2)\log(\epsilon^{-1})}} \right), \end{aligned}$$

we obtain that for any  $\epsilon > 0$ ,

$$\mathbb{P}^{\otimes N} \left\{ H(k) \geq \mathbb{P}^{\otimes N} H(k) + \frac{(\log 2)\log(\epsilon^{-1})}{3} \left( 1 + \sqrt{1 + \frac{18\mathbb{P}^{\otimes N} H(k)}{(\log 2)\log(\epsilon^{-1})}} \right) \right\} \leq \epsilon,$$

which is the first assertion of the lemma. The second inequality of the lemma is a direct consequence of inequality (6.18). The following two inequalities in the lemma comes similarly from inequality (6.17). Finally, inequality (4.3) comes from combining Inequalities (4.1) and (4.2).

**6.11. Proof of Theorem 4.3.** For any  $k \in \{0, \dots, U\}$ , let  $\mathcal{N}_k$  be a  $2^{-k}$ -minimal covering net of  $\mathcal{F}$  for the pseudo-distance  $\bar{\mathbb{P}}$ . For any  $(i, k) \in \{1, 2\} \times \{0, \dots, U\}$ , let  $f_{i,k}$  be a nearest neighbour of  $f_i$  in the set  $\mathcal{N}_k$ . Let  $0 \leq K \leq U$  be the integer satisfying  $\frac{\bar{\mathbb{P}}_{f_1, f_2} \vee u}{2} < 2^{-K} \leq \bar{\mathbb{P}}_{f_1, f_2} \vee u$ .

Since we have

$$\partial(f_1; f_2) = \partial(f_{1,K}; f_{2,K}) + \sum_{k=K+1}^U \left\{ \partial(f_{1,k}; f_{1,k-1}) + \partial(f_{2,k-1}; f_{2,k}) \right\},$$

we need to apply Lemma 4.2 to  $(\mathcal{S}_1, \mathcal{S}_2) \in \cup_{1 \leq k \leq U} \{(\mathcal{N}_{k-1}, \mathcal{N}_k) \cup (\mathcal{N}_k, \mathcal{N}_{k-1}) \cup (\mathcal{N}_k, \mathcal{N}_k)\}$  and to do a union bound on the associated  $3U$  inequalities. Let  $w_k, k \in \mathbb{N}^*$  be positive integers such that  $\sum_{k \geq 1} w_k = 1$ . With probability at least  $1 - \epsilon$ , for any  $k \in \mathbb{N}^*$ , for any  $(f'_1, f'_2) \in (\mathcal{N}_{k-1} \times \mathcal{N}_k) \cup (\mathcal{N}_k \times \mathcal{N}_{k-1}) \cup (\mathcal{N}_k \times \mathcal{N}_k)$ , we have

$$\partial(f'_1; f'_2) \leq \sqrt{\frac{8\bar{\mathbb{P}}_{f'_1, f'_2} \log(3|\mathcal{N}_k|^2 w_k^{-1} \epsilon^{-1})}{N}}$$

For any  $(i, k) \in \{1, 2\} \times \{1, \dots, U\}$ , we have  $\mathbb{P}_{f_{i,k-1}, f_{i,k}} \leq 3 \times 2^{-k}$ . Denote

$$B_k \triangleq \sqrt{\frac{24 \times 2^{-k} \log(3|\mathcal{N}_k|^2 w_k^{-1} \epsilon^{-1})}{N}}$$

We have  $\mathbb{P}_{f_1, K, f_2, K} \leq 4 \times 2^{-K}$  and  $f_{1,0} = f_{2,0}$ . Chaining the inequalities, we obtain that, with  $\mathbb{P}^{\otimes 2N}$ -probability at least  $1 - \epsilon$ ,

$$\partial(f_1; f_2) \leq 2B_K \mathbf{1}_{K>0} + 2 \sum_{k=K+1}^U B_k \leq 2 \sum_{k=K \vee 1}^U B_k,$$

hence

$$\partial(f_1; f_2) \leq 4 \sum_{k=K \vee 1}^U \sqrt{\frac{6 \times 2^{-k} [2H_k + \log(3w_k^{-1}) + \log(\epsilon^{-1})]}{N}}.$$

We want that the union bound term  $\log(3w_k^{-1})$  remains negligible wrt the complexity term  $2H_k$ . This leads to choose, for instance,  $w_k = \frac{1}{k(k+1)}$  since  $\sum_{k \geq 1} \frac{1}{k(k+1)} = 1$  and for small classes of functions (i.e. VC-classes), the entropy  $H_k$  has the order of  $k$ , hence  $\log(3w_k^{-1}) \ll H_k$ . We obtain

$$\partial(f_1; f_2) \leq \sum_{k \in \mathbb{N}^*: u \leq 2^{-k} \leq \bar{\mathbb{P}}_{f_1, f_2} \vee u} 4 \sqrt{\frac{6 \times 2^{-k} \{2H_k + \log[3k(k+1)] + \log(\epsilon^{-1})\}}{N}}.$$

*Remark 6.5.* The previous result can also be written in terms of integral. Introduce the set  $E \triangleq \{k \in \mathbb{N} : 2^{-k} \leq (\bar{\mathbb{P}}_{f_1, f_2} \vee u) \wedge \frac{1}{2}\}$  and take  $H_k = H(2^{-k}, \mathcal{F}, \bar{\mathbb{P}}) \vee 1$ . We get

$$\begin{aligned} \partial(f_1; f_2) &\leq 16 \sqrt{\frac{3}{N}} \sum_{k \in E} \sqrt{\frac{H_k}{2^{-k}}} (2^{-k} - 2^{-k-1}) + 4 \sqrt{\frac{6 \log(3\epsilon^{-1})}{N}} \sum_{k \in E} (\sqrt{2})^{-k} \\ &\quad + 4 \sqrt{\frac{6}{N}} \sum_{k \in E} (\sqrt{2})^{-k} \sqrt{\log[k(k+1)]} \\ &\leq 16 \sqrt{\frac{3}{N}} \int_{\frac{u}{2}}^{(\bar{\mathbb{P}}_{f_1, f_2} \vee u) \wedge \frac{1}{2}} \sqrt{\frac{H(x, \mathcal{F}, \bar{\mathbb{P}}) \vee 1}{x}} dx \\ &\quad + 4 \sqrt{\frac{6 \log(3\epsilon^{-1})}{N}} \frac{\sqrt{2}}{\sqrt{2}-1} \sqrt{\bar{\mathbb{P}}_{f_1, f_2} \vee u} + 4 \sqrt{\frac{6}{N}} \sum_{k \in E} (\sqrt{2})^{-k} \sqrt{2 \log(k+1)}. \end{aligned}$$

Let  $\varphi(x) \triangleq \frac{1}{\sqrt{x}} \sqrt{\log(e \frac{\log x^{-1}}{\log 2})}$  for any  $0 < x \leq \frac{1}{2}$ . The function  $\varphi$  is decreasing on  $[0; \frac{1}{2}]$ . The last term can be written as

$$8 \sqrt{\frac{3}{N}} \sum_{k \in E} \frac{2(2^{-k} - 2^{-k-1})}{\sqrt{2^{-k}}} \sqrt{\log(k+1)} \leq 16 \sqrt{\frac{3}{N}} \int_{\frac{u}{2}}^{(\bar{\mathbb{P}}_{f_1, f_2} \vee u) \wedge \frac{1}{2}} \varphi(x) dx.$$

**6.12. Proof of Theorem 4.4.** Let us take  $U \in \mathbb{N}$  such that  $2^{-U} < \frac{1}{2N}$ . From Theorem 4.3, for any  $\epsilon > 0$ , with  $\mathbb{P}^{\otimes 2N}$ -probability at least  $1 - \epsilon$ , for any functions  $f_1$  and  $f_2$  in the set  $\mathcal{F}$ ,

(6.19)

$$\begin{aligned} r'(f_2) - r'(f_1) + r(f_1) - r(f_2) \\ \leq \sum_{k \in \mathbb{N}: \frac{1}{2N} \leq 2^{-k} \leq \bar{\mathbb{P}}_{f_1, f_2} \wedge \frac{1}{2}} 4 \sqrt{\frac{6 \times 2^{-k} \{2H_k + \log[3k(k+1)] + \log(\epsilon^{-1})\}}{N}}. \end{aligned}$$

Let  $\mathcal{N}$  be a  $(\frac{1}{3N}, \mathcal{F}, \bar{\mathbb{P}}, \cdot)$ -minimal covering net. From Theorem 8.4 in [1] apply to  $\mathcal{W}((f_1, f_2), X) = \mathbf{1}_{f_1(X) \neq f_2(X)}$ , we obtain that with  $\mathbb{P}^{\otimes N}$ -probability at least  $1 - \epsilon$ , for any function  $f_1, f_2$  in the net  $\mathcal{F}$ , we have

$$\bar{\mathbb{P}}'_{f_1, f_2} \leq \bar{\mathbb{P}}_{f_1, f_2} + 2 \sqrt{\frac{\bar{\mathbb{P}}_{f_1, f_2} \log[N^2 (X_1^{2N}) \epsilon^{-1}]}{N}},$$

hence

$$\bar{\mathbb{P}}_{f_1, f_2} \leq \bar{\mathbb{P}}_{f_1, f_2} + \sqrt{\frac{\bar{\mathbb{P}}_{f_1, f_2} \log[N^2 (X_1^{2N}) \epsilon^{-1}]}{N}}.$$

Let  $\check{\mathcal{K}} \triangleq \frac{2 \log N(X_1^{2N}) + \log(\epsilon^{-1})}{N}$ . By solving the previous inequation, we obtain

$$(6.20) \quad \bar{\mathbb{P}}_{f_1, f_2} \leq \bar{\mathbb{P}}_{f_1, f_2} + \sqrt{\check{\mathcal{K}} \bar{\mathbb{P}}_{f_1, f_2} + \frac{\check{\mathcal{K}}^2}{4}} + \frac{\check{\mathcal{K}}}{2}.$$

From inequality (4.3), we have

$$\begin{aligned} H(u, \mathcal{F}, \bar{\mathbb{P}}') &\leq H(u, \mathcal{F}, \bar{\mathbb{P}}) + (\log 2) \log(\epsilon^{-1}) \left( \frac{12}{5} + 1 + \frac{2H(u, \mathcal{F}, \bar{\mathbb{P}})}{(\log 2) \log(\epsilon^{-1})} \right) \\ &= 3H(u, \mathcal{F}, \bar{\mathbb{P}}) + \frac{17 \log 2}{5} \log(\epsilon^{-1}), \end{aligned}$$

hence

$$H(u, \mathcal{F}, \bar{\mathbb{P}}) \leq 4H(u, \mathcal{F}, \bar{\mathbb{P}}) + \frac{17 \log 2}{5} \log(\epsilon^{-1}).$$

Taking an union bound with weight  $\frac{1}{k(k+1)}$ , we obtain that with  $\mathbb{P}^{\otimes 2N}$ -probability at least  $1 - \epsilon$ , for any  $k \geq 1$ , we have

$$(6.21) \quad H(2^{-k}, \mathcal{F}, \bar{\mathbb{P}}) \leq 4H(2^{-k}, \mathcal{F}, \bar{\mathbb{P}}_{\cdot, \cdot}) + 2.4 \log[k(k+1)\epsilon^{-1}].$$

Let  $H'_k \triangleq 4H(2^{-k}, \mathcal{F}, \bar{\mathbb{P}}_{\cdot, \cdot}) + 2.4 \log[k(k+1)\epsilon^{-1}]$ . Rigorously, we cannot apply Theorem 4.3 for  $H_k = H'_k$  since  $H'_k$  is not always an upper bound of  $H(2^{-k}, \mathcal{F}, \bar{\mathbb{P}}) \vee 1$ . However we can modify the proof of Theorem 4.3 to take into a “probably approximately correct” inequality. Therefore, combining Inequalities (6.19), (6.20) and (6.21), letting  $\bar{\mathcal{K}} \triangleq \frac{2H'_k + \log(\epsilon^{-1})}{N}$  and

$$\bar{E} \triangleq \left\{ k \in \mathbb{N}^* : \frac{1}{2N} \leq 2^{-k} \leq \bar{\mathbb{P}}_{f_1, f_2} + \sqrt{\bar{\mathcal{K}} \bar{\mathbb{P}}_{f_1, f_2} + \frac{\bar{\mathcal{K}}^2}{4}} + \frac{\bar{\mathcal{K}}}{2} \right\},$$

we obtain that with  $\mathbb{P}^{\otimes 2N}$ -probability at least  $1 - 3\epsilon$ , for any functions  $f_1, f_2$  in  $\mathcal{F}$ , we have

$$r'(f_2) - r'(f_1) + r(f_1) - r(f_2) \leq \sum_{k \in \bar{E}} 4 \sqrt{\frac{6 \times 2^{-k} \{2H'_k + \log[3k(k+1)] + \log(\epsilon^{-1})\}}{N}}.$$

To obtain the announced result, we simplify this formula by using

$$(6.22) \quad \begin{aligned} \bar{\mathbb{P}}_{f_1, f_2} + \sqrt{\bar{\mathcal{K}} \bar{\mathbb{P}}_{f_1, f_2} + \frac{\bar{\mathcal{K}}^2}{4}} + \frac{\bar{\mathcal{K}}}{2} &\leq \bar{\mathbb{P}}_{f_1, f_2} + \sqrt{\bar{\mathcal{K}} \bar{\mathbb{P}}_{f_1, f_2}} + \bar{\mathcal{K}} \leq \frac{5}{4} \bar{\mathbb{P}}_{f_1, f_2} + 2\bar{\mathcal{K}}, \\ \log[U(U+1)] &\leq \log \left[ \left( 2 + \frac{\log N}{\log 2} \right) \left( 3 + \frac{\log N}{\log 2} \right) \right] \leq \log 6 + 2 \log \left( \frac{e}{2 \log 2} \log N \right) \end{aligned}$$

and

$$2H'_k + \log[3k(k+1)\epsilon^{-1}] \leq 8H(2^{-k}, \mathcal{F}, \bar{\mathbb{P}}_{\cdot, \cdot}) + 6 \log[k(k+1)\epsilon^{-1}] + 1.$$

**6.13. Proof of Corollary 4.5.** Let  $f \in \mathcal{F}$ . If  $\bar{\mathbb{P}}_{\hat{f}_{\text{ERM}}, f} = 0$ , then we trivially have  $r'(\hat{f}_{\text{ERM}}) \leq r'(f)$ . Otherwise, we have  $\bar{\mathbb{P}}_{\hat{f}_{\text{ERM}}, f} \geq \frac{1}{2N}$ . Let  $K \in \mathbb{N}$  such that  $\frac{\bar{\mathbb{P}}_{\hat{f}_{\text{ERM}}, f}}{2} < 2^{-K} \leq \bar{\mathbb{P}}_{\hat{f}_{\text{ERM}}, f}$ . From Theorem 4.3, with  $\mathbb{P}^{\otimes 2N}$ -probability at least  $1 - \epsilon$ , we have

$$\begin{aligned} r'(\hat{f}_{\text{ERM}}) - r'(f) &\leq \sum_{k \geq K} 4 \sqrt{\frac{6 \times 2^{-k} \{2V \log(e2^{k+2}) + \log[3k(k+1)] + \log(\epsilon^{-1})\}}{N}} \\ &\leq \sum_{k \geq K} 4 \sqrt{\frac{6 \times 2^{-k} \{(2V+1) \log(e2^{k+2}) + \log(\epsilon^{-1})\}}{N}} \\ &\leq 4 \sqrt{\frac{6(2V+1)}{N}} \sum_{k \geq K} \sqrt{2^{-k} \log(e2^{k+2})} + 4 \sqrt{\frac{6 \log(\epsilon^{-1})}{N}} \sum_{k \geq K} (\sqrt{2})^{-k} \end{aligned}$$

Now, for any  $k \geq K \geq 0$  and  $V \geq 1$ , we have  $\frac{\log(e2^{k+2})}{\log(e2^{K+2})} \leq \frac{(k-K)(\log 2)}{1+2\log 2} + 1$ . Therefore we get

$$\begin{aligned} r'(\hat{f}_{\text{ERM}}) - r'(f) &\leq 4\sqrt{\frac{6(2V+1)2^{-K}\log(e2^{K+2})}{N}} \sum_{k \geq 0} \sqrt{2^{-k} \left( \frac{k \log 2}{1+2\log 2} + 1 \right)} \\ &\quad + 4\sqrt{\frac{6\log(\epsilon^{-1})}{N}} (\sqrt{2})^{-K} \frac{\sqrt{2}}{\sqrt{2-1}} \\ &\leq 4\sqrt{\frac{6(2V+1)2^{-K}\log(e2^{K+2})}{N}} \sum_{k \geq 0} \sqrt{2^{-k} \left( \frac{k \log 2}{1+2\log 2} + 1 \right)} \\ &\quad + 4\sqrt{\frac{6\log(\epsilon^{-1})}{N}} (\sqrt{2})^{-K} \frac{\sqrt{2}}{\sqrt{2-1}} \\ &\leq 47\sqrt{\frac{V+1}{N}} \sqrt{\bar{\mathbb{P}}_{\hat{f}_{\text{ERM}},f} \log\left(\frac{8e}{\bar{\mathbb{P}}_{\hat{f}_{\text{ERM}},f}}\right)} + 34\sqrt{\frac{\bar{\mathbb{P}}_{\hat{f}_{\text{ERM}},f} \log(\epsilon^{-1})}{N}} \end{aligned}$$

For the second assertion of the corollary, we use Jensen's inequality and the concavity of  $x \mapsto \sqrt{x \log(8ex^{-1})}$  in order to obtain that for any function  $f \in \mathcal{F}$ ,

$$\begin{aligned} \mathbb{P}^{\otimes N} R(\hat{f}_{\text{ERM}}) &\leq R(f) + 47\sqrt{\frac{(V+1)\mathbb{P}^{\otimes 2N} \bar{\mathbb{P}}_{\hat{f}_{\text{ERM}},f}}{N}} \log\left(\frac{8e}{\mathbb{P}^{\otimes 2N} \bar{\mathbb{P}}_{\hat{f}_{\text{ERM}},f}}\right) \\ &\quad + 34\sqrt{\frac{\mathbb{P}^{\otimes 2N} \bar{\mathbb{P}}_{\hat{f}_{\text{ERM}},f}}{N}}. \end{aligned}$$

**6.14. Proof of Corollary 4.6.** For  $\kappa = +\infty$  (i.e. no margin assumption), the result comes from inequality (4.5) since the function  $x \mapsto x \log(8e/x)$  is an increasing function on  $[0; 1]$ , hence upper bounded by its value for  $x = 1$ . Specifically, we obtain

$$(6.23) \quad \mathbb{P}^{\otimes N} R(\hat{f}_{\text{ERM}}) - R(\tilde{f}) \leq 83\sqrt{\frac{V+1}{N}} + \frac{34}{\sqrt{N}}.$$

Note that it is thanks to the chaining that we get rid of the  $\log N$  factor. For  $\kappa < +\infty$ , chained and unchained results lead to the same convergence rate:  $\left(\frac{V}{N} \log N\right)^{\frac{\kappa}{2\kappa-1}}$ .

To obtain this rate from the previous bounds, we just need to link the variance term  $\mathbb{P}^{\otimes 2N} \bar{\mathbb{P}}_{\hat{f}_{\text{ERM}},\tilde{f}}$  with  $\mathbb{P}^{\otimes N} \mathbb{P}_{\hat{f}_{\text{ERM}},\tilde{f}}$  in order to use the margin assumption.

Combining Inequalities (6.20) and (6.22), we obtain

$$\begin{aligned} \bar{\mathbb{P}}_{\hat{f}_{\text{ERM}},\tilde{f}} &\leq \frac{5}{4} \bar{\mathbb{P}}_{\hat{f}_{\text{ERM}},\tilde{f}} + \frac{4\log N(X_1^{2N}) + 2\log(\epsilon^{-1})}{N} \\ &\leq \frac{5}{4} \bar{\mathbb{P}}_{\hat{f}_{\text{ERM}},\tilde{f}} + \frac{4V \log\left(\frac{2eN}{V}\right) + 2\log(\epsilon^{-1})}{N}, \end{aligned}$$

hence

$$(6.24) \quad \mathbb{P}^{\otimes 2N} \bar{\mathbb{P}}_{\hat{f}_{\text{ERM}},\tilde{f}} \leq \frac{5}{4} \mathbb{P}^{\otimes N} \mathbb{P}_{\hat{f}_{\text{ERM}},\tilde{f}} + \frac{4V \log\left(\frac{2eN}{V}\right) + 2}{N}.$$

Now, by the margin assumption and Jensen's inequality, we have

$$(6.25) \quad \mathbb{P}^{\otimes N} \mathbb{P}_{\hat{f}_{\text{ERM}},\tilde{f}} \leq C'' \mathbb{P}^{\otimes N} (\Delta^{\frac{1}{\kappa}}) \leq C''' (\mathbb{P}^{\otimes N} \Delta)^{\frac{1}{\kappa}}.$$

The convergence rate then follows from (6.24), (6.25) and either (4.4) or (4.5).

**6.15. Proof of Lemma 5.1.** Let  $\vec{\sigma}_{j,r} \triangleq (\sigma_1, \dots, \sigma_{j-1}, r, \sigma_{j+1}, \dots, \sigma_m)$  for any  $r \in \{-1, 0, +1\}$ . The distribution  $\mathbb{P}_{\vec{\sigma}_{j,0}}$  is such that  $\mathbb{P}_{\vec{\sigma}_{j,0}}(dX) = \mu(dX)$  and

$$\mathbb{P}_{\vec{\sigma}_{j,0}}(Y = 1|X) = \begin{cases} \frac{1}{2} & \text{for any } X \in \mathcal{X}_j \\ \mathbb{P}_{\vec{\sigma}}(Y = 1|X) & \text{otherwise} \end{cases}.$$

Introduce the quantities  $\pi_{r,j} \triangleq \frac{\mathbb{P}_{\sigma_j,r}^{\otimes N}}{\mathbb{P}_{\sigma_j,0}^{\otimes N}}(Z_1^N) = \prod_{i=1}^N [1 + r \mathbb{1}_{X_i \in \mathcal{X}_j} (2Y_i - 1) \xi(X_i)]$  for any  $r \in \{-1; +1\}$ . Let  $\nu$  denote the distribution of a Rademacher variable:

$$\nu(\sigma = +1) = \nu(\sigma = -1) = \frac{1}{2}.$$

The variational distance between two probability distributions  $\mathbb{P}_1$  and  $\mathbb{P}_2$  is defined as  $V(\mathbb{P}_1, \mathbb{P}_2) \triangleq \sup_{A \text{ measurable set}} \{\mathbb{P}_1(A) - \mathbb{P}_2(A)\}$ . When the distributions  $\mathbb{P}_1$  and  $\mathbb{P}_2$  are absolutely continuous wrt a probability distribution  $\mathbb{P}_0$ , we have

$$V(\mathbb{P}_1, \mathbb{P}_2) = \frac{1}{2} \int \left| \frac{\mathbb{P}_1}{\mathbb{P}_0} - \frac{\mathbb{P}_2}{\mathbb{P}_0} \right| d\mathbb{P}_0 = 1 - \int \left( \frac{\mathbb{P}_1}{\mathbb{P}_0} \wedge \frac{\mathbb{P}_2}{\mathbb{P}_0} \right) d\mathbb{P}_0.$$

We have successively

$$\begin{aligned}
& \sup_{\mathbb{P} \in \mathcal{P}} \left\{ \mathbb{P}^{\otimes N} \mathbb{P}[\hat{f}(X) \neq Y] - \mathbb{P}[f_{\mathbb{P}}^*(X) \neq Y] \right\} \\
& \geq \sup_{\bar{\sigma} \in \{-1; +1\}^m} \left\{ (\mathbb{P}_{\bar{\sigma}}^{\otimes N}) \mathbb{P}_{\bar{\sigma}}[\hat{f}(X) \neq Y] - \mathbb{P}_{\bar{\sigma}}[f_{\mathbb{P}_{\bar{\sigma}}}^* \neq Y] \right\} \\
& = \sup_{\bar{\sigma} \in \{-1; +1\}^m} \left\{ (\mathbb{P}_{\bar{\sigma}}^{\otimes N}) \mathbb{P}_{\bar{\sigma}} \left( \xi(X) \mathbb{1}_{\hat{f}(X) \neq f_{\mathbb{P}_{\bar{\sigma}}}^*(X)} \right) \right\} \\
& = \sup_{\bar{\sigma} \in \{-1; +1\}^m} \left\{ \mathbb{P}_{\bar{\sigma}}^{\otimes N} \left( \sum_{j=1}^m \mu \left[ \xi(X) \mathbb{1}_{\hat{f}(X) \neq \frac{1+\sigma_j}{2}; X \in \mathcal{X}_j} \right] \right) \right\} \\
& \geq \mathbb{E}_{\nu^{\otimes m}} \sum_{j=1}^m \mathbb{P}_{\bar{\sigma}}^{\otimes N} \left[ \mu \left[ \xi(X) \mathbb{1}_{\hat{f}(X) \neq \frac{1+\sigma_j}{2}; X \in \mathcal{X}_j} \right] \right] \\
& = \mathbb{E}_{\nu^{\otimes m}} \sum_{j=1}^m \mathbb{P}_{\bar{\sigma}_j,0}^{\otimes N} \left( \frac{\mathbb{P}_{\bar{\sigma}_j}^{\otimes N}}{\mathbb{P}_{\bar{\sigma}_j,0}^{\otimes N}} \mu \left[ \xi(X) \mathbb{1}_{\hat{f}(X) \neq \frac{1+\sigma_j}{2}; X \in \mathcal{X}_j} \right] \right) \\
(6.26) \quad & = \mathbb{E}_{\nu^{\otimes(m-1)}(d\sigma_1, \dots, d\sigma_{j-1}, d\sigma_{j+1}, \dots, d\sigma_m)} \sum_{j=1}^m \mathbb{P}_{\bar{\sigma}_j,0}^{\otimes N} \mathbb{E}_{\nu}(d\sigma_j) \\
& \quad \left( \frac{\mathbb{P}_{\bar{\sigma}_j}^{\otimes N}}{\mathbb{P}_{\bar{\sigma}_j,0}^{\otimes N}} \mu \left[ \xi(X) \mathbb{1}_{\hat{f}(X) \neq \frac{1+\sigma_j}{2}; X \in \mathcal{X}_j} \right] \right) \\
& \geq \mathbb{E}_{\nu^{\otimes(m-1)}(d\sigma_1, \dots, d\sigma_{j-1}, d\sigma_{j+1}, \dots, d\sigma_m)} \sum_{j=1}^m \mathbb{P}_{\bar{\sigma}_j,0}^{\otimes N} \\
& \quad \left[ \left( \pi_{-1,j} \wedge \pi_{+1,j} \right) \mathbb{E}_{\nu}(d\sigma_j) \mu \left[ \xi(X) \mathbb{1}_{\hat{f}(X) \neq \frac{1+\sigma_j}{2}; X \in \mathcal{X}_j} \right] \right] \\
& = \mathbb{E}_{\nu^{\otimes(m-1)}(d\sigma_1, \dots, d\sigma_{j-1}, d\sigma_{j+1}, \dots, d\sigma_m)} \sum_{j=1}^m \\
& \quad \frac{1}{2} \mu \left[ \xi(X) \mathbb{1}_{X \in \mathcal{X}_j} \right] \left[ 1 - V \left( \mathbb{P}_{\bar{\sigma}_j,-1}^{\otimes N}, \mathbb{P}_{\bar{\sigma}_j,+1}^{\otimes N} \right) \right] \\
& = \frac{mw}{2} \mu \left[ \xi(X) \mid X \in \mathcal{X}_j \right] \left[ 1 - V \left( \mathbb{P}_{-1,1,\dots,1}^{\otimes N}, \mathbb{P}_{1,1,\dots,1}^{\otimes N} \right) \right].
\end{aligned}$$

Now let us prove

$$(6.27) \quad V \left( \mathbb{P}_{-1,1,\dots,1}^{\otimes N}, \mathbb{P}_{1,1,\dots,1}^{\otimes N} \right) \leq b\sqrt{Nw}.$$

First, we have

$$(6.28) \quad V \left( \mathbb{P}_{-1,1,\dots,1}^{\otimes N}, \mathbb{P}_{1,1,\dots,1}^{\otimes N} \right) = \sum_{l=1}^N \binom{N}{l} w^l (1-w)^{N-l} \mathcal{V}_l,$$

where  $\mathcal{V}_l \triangleq V \left( \mathbb{P}_{-1}^{\otimes l}, \mathbb{P}_{+1}^{\otimes l} \right)$  and  $\mathbb{P}_{\sigma} \triangleq \mathbb{P}_{\sigma,1,\dots,1}(\cdot | X \in \mathcal{X}_1)$  for any  $\sigma \in \{-1, +1\}$ . By simple computations, we get  $\mathcal{V}_1 = \mu[\xi(X) | X \in \mathcal{X}_j]$ . From Jensen's inequality, by the concavity of  $x \mapsto \sqrt{1-x^2}$ , we have

$$\sqrt{1-b^2} = \mu \left[ \sqrt{1-\xi^2(X)} \mid X \in \mathcal{X}_j \right] \leq \sqrt{1 - \{\mu[\xi(X) | X \in \mathcal{X}_j]\}^2},$$

hence  $\mathcal{V}_1 \leq b$ .

For  $l \geq 2$ , we upper bound the variational distance by the Hellinger distance. By definition, the Hellinger distance  $H(\mathbb{P}, \mathbb{Q})$  satisfies  $1 - \frac{H^2(\mathbb{P}, \mathbb{Q})}{2} = \int \sqrt{d\mathbb{P}} \sqrt{d\mathbb{Q}}$ . Hence the tensorization equality is  $1 - \frac{H^2(\mathbb{P}^{\otimes l}, \mathbb{Q}^{\otimes l})}{2} = \left(1 - \frac{H^2(\mathbb{P}, \mathbb{Q})}{2}\right)^l$ . We have

$$\mathcal{V}_l \leq H\left(\mathbb{P}_{-1}^{\otimes l}, \mathbb{P}_{+1}^{\otimes l}\right) = \sqrt{2\left(1 - \left[1 - \frac{H^2(\mathbb{P}_{-1}, \mathbb{P}_{+1})}{2}\right]^l\right)}$$

and  $1 - \frac{H^2(\mathbb{P}_{-1}, \mathbb{P}_{+1})}{2} = \mu[\sqrt{1 - \xi^2(X)} | X \in \mathcal{X}_1] = \sqrt{1 - b^2}$ , by definition of  $b$ . Now, for any  $l \geq 2$  and  $x \geq 0$ ,  $2\left(1 - [1 - x^2]^{\frac{l}{2}}\right) \leq lx^2$ . Finally, for any  $l \geq 1$ , we have  $\mathcal{V}_l \leq b\sqrt{l}$ . Putting this result in equality (6.28), we get

$$V\left(\mathbb{P}_{-1,1,\dots,1}^{\otimes N}, \mathbb{P}_{1,1,\dots,1}^{\otimes N}\right) \leq b \sum_{l=0}^N \mathbb{P}\left(\sum_{i=1}^N \epsilon_i = l\right) \sqrt{l},$$

where the  $\epsilon_i$  are i.i.d. random variables such that  $\mathbb{P}(\epsilon_i = 1) = w = 1 - \mathbb{P}(\epsilon_i = 0)$ . So we have  $V\left(\mathbb{P}_{-1,1,\dots,1}^{\otimes N}, \mathbb{P}_{1,1,\dots,1}^{\otimes N}\right) \leq b\mathbb{P}\sqrt{\sum_{i=1}^N \epsilon_i} \leq b\sqrt{\mathbb{P}\sum_{i=1}^N \epsilon_i} = b\sqrt{Nw}$ .

*Remark 6.6.* The last inequality in (6.26) is an equality when for any  $j \in \{1, \dots, m\}$ ,  $\hat{f} = \operatorname{argmax}_{r \in \{-1; +1\}} \mathbb{P}^{\otimes N} \bar{\sigma}_{j,r}$  on  $\mathcal{X}_j$ , i.e. when  $\hat{f}$  is the maximum likelihood estimator on the set  $\mathcal{X} - \mathcal{X}_0$ .

#### APPENDIX A. PROOF OF INEQUALITY (6.2)

For any  $r > 0$ , define  $\Gamma(r) \triangleq \int_0^{+\infty} u^{r-1} \exp(-u) du$ . Integrating by parts, we obtain the well-known property  $\Gamma(r+1) = r\Gamma(r)$ .

• We have

$$\begin{aligned} \pi \exp(-\lambda \Delta R) &= \int_0^{+\infty} \pi \{ \exp(-\lambda \Delta R) \geq u \} du \\ (A.1) \quad &= \exp(-\lambda) + \int_{\exp(-\lambda)}^1 \pi \{ \exp(-\lambda \Delta R) \geq u \} du \\ &= \exp(-\lambda) + \int_0^1 \lambda \exp(-\lambda x) \pi(\Delta R \leq x) dx \end{aligned}$$

Let us introduce  $A' \triangleq \exp(-C''')$ . Since we have  $\pi(\Delta R \leq x) = x^{C'} [A' + \eta(x)]$  with  $\eta(x) = \underset{x \rightarrow 0}{\circ} (x^s)$  and  $\eta(x) \leq x^{-C'}$  and  $\int_0^{+\infty} \lambda \exp(-\lambda x) x^{C'} dx = \frac{\Gamma(C'+1)}{\lambda^{C'}}$ , we get

$$\begin{aligned} &\left| \pi \exp(-\lambda \Delta R) - A' \frac{\Gamma(C'+1)}{\lambda^{C'}} - \exp(-\lambda) \right| \\ &= \int_0^1 \lambda \exp(-\lambda x) x^{C'} \eta(x) dx + A' \int_1^{+\infty} \lambda \exp(-\lambda x) x^{C'} dx \end{aligned}$$

Since we have

$$\begin{aligned} &\int_0^1 \lambda \exp(-\lambda x) x^{C'} \eta(x) dx \\ &= \int_0^{\frac{1}{\sqrt{\lambda}}} \lambda \exp(-\lambda x) x^{C'+s} \underset{x \rightarrow 0}{\circ} (1) dx + \int_{\frac{1}{\sqrt{\lambda}}}^1 \lambda \exp(-\lambda x) x^{C'} \eta(x) dx \\ &\leq \underset{\lambda \rightarrow +\infty}{\circ} (\lambda^{-(C'+s)}) + \int_{\frac{1}{\sqrt{\lambda}}}^1 \lambda \exp(-\lambda x) dx \\ &= \underset{\lambda \rightarrow +\infty}{\circ} (\lambda^{-(C'+s)}) \end{aligned}$$

and  $\int_1^{+\infty} \lambda \exp(-\lambda x) x^{C'} dx = \underset{\lambda \rightarrow +\infty}{\circ} (\lambda^{-(C'+s)})$ , we obtain

$$(A.2) \quad \pi \exp(-\lambda \Delta R) = A' \frac{\Gamma(C'+1) + \underset{\lambda \rightarrow +\infty}{\circ} (\lambda^{-s})}{\lambda^{C'}}$$



- From equalities (A.1), we have

(A.3)

$$\pi[\Delta R \exp(-\lambda \Delta R)] = \exp(-\lambda) + \int_0^1 (\lambda x - 1) \exp(-\lambda x) \pi(\Delta R \leq x) dx.$$

We have just seen that  $\int_0^1 \exp(-\lambda x) \pi(\Delta R \leq x) dx = A' \frac{\Gamma(C'+1) + \frac{0}{\lambda \rightarrow +\infty}(\lambda^{-s})}{\lambda^{C'+1}}$ . Besides, from the same argument as above, we have

$$\int_0^1 \lambda x \exp(-\lambda x) \pi(\Delta R \leq x) dx = A' \frac{\Gamma(C'+2) + \frac{0}{\lambda \rightarrow +\infty}(\lambda^{-s})}{\lambda^{C'+1}}.$$

Since  $\Gamma(C'+2) = (C'+1)\Gamma(C'+1)$ , we obtain

$$\pi[\Delta R \exp(-\lambda \Delta R)] = A' \frac{C'\Gamma(C'+1) + \frac{0}{\lambda \rightarrow +\infty}(\lambda^{-s})}{\lambda^{C'+1}}.$$

- Combining the previous results, we obtain  $\pi_{-\lambda R} \Delta R = \frac{C' + \frac{0}{\lambda \rightarrow +\infty}(\lambda^{-s})}{\lambda}$ .

#### APPENDIX B. PROOF OF INEQUALITY (6.4)

Let  $\alpha > 0$  depend on  $\lambda$  such that  $\lambda \alpha^{\frac{\kappa}{\kappa-1}} \rightarrow 0$ . Then there exists  $0 < \zeta < 1$  depending on  $\lambda$  such that  $\zeta \rightarrow 1$  and  $\lambda \left(\frac{\alpha}{1-\zeta}\right)^{\frac{\kappa}{\kappa-1}} \rightarrow 0$ . Let  $h_\alpha(x) \triangleq x - \alpha x^{\frac{1}{\kappa}}$  and  $x_0 \triangleq \left(\frac{\alpha}{\kappa}\right)^{\frac{\kappa}{\kappa-1}}$ . The function  $h$  decreases on  $[0; x_0]$  and increases on  $[x_0; +\infty]$ . We have

$$\begin{aligned} \pi \exp[-\lambda h_\alpha(\Delta R)] &= \pi \left\{ \exp[-\lambda h_\alpha(\Delta R)] \mathbb{1}_{h_\alpha(\Delta R) \leq \zeta \Delta R} \right\} \\ &\quad + \pi \left\{ \exp[-\lambda h_\alpha(\Delta R)] \mathbb{1}_{h_\alpha(\Delta R) > \zeta \Delta R} \right\} \\ &\leq \exp[-\lambda h_\alpha(x_0)] \pi \left\{ \Delta R \leq \left(\frac{\alpha}{1-\zeta}\right)^{\frac{\kappa}{\kappa-1}} \right\} + \pi \exp[-\lambda \zeta \Delta R] \\ &\leq \exp\left[(\kappa-1)\lambda \left(\frac{\alpha}{\kappa}\right)^{\frac{\kappa}{\kappa-1}} \left(\frac{\alpha}{1-\zeta}\right)^{\frac{C'\kappa}{\kappa-1}} \left[1 + \frac{0}{\lambda \rightarrow +\infty}(1)\right]\right] \\ &\quad + \pi \exp[-\lambda \zeta \Delta R] \end{aligned}$$

From equality (A.2), we get

$$\begin{aligned} \pi_{-\lambda R} \exp(\lambda \alpha \Delta R) &= \exp\left[(\kappa-1)\lambda \left(\frac{\alpha}{\kappa}\right)^{\frac{\kappa}{\kappa-1}}\right] \lambda^{C'} \left(\frac{\alpha}{1-\zeta}\right)^{\frac{C'\kappa}{\kappa-1}} \left[\frac{1}{\Gamma(C'+1)} + \frac{0}{\lambda \rightarrow +\infty}(1)\right] + \frac{1 + \frac{0}{\lambda \rightarrow +\infty}(\lambda^{-s})}{\zeta^{C'}} \\ &= \underset{\lambda \rightarrow +\infty}{\mathbf{O}} \left(\lambda^{C'} \left[\frac{\alpha}{1-\zeta}\right]^{\frac{\kappa C'}{\kappa-1}}\right) + 1 + \frac{0}{\lambda \rightarrow +\infty}(\lambda^{-s}) + \underset{\lambda \rightarrow +\infty}{\mathbf{O}}(1-\zeta) \end{aligned}$$

Taking  $\zeta = 1 - \left(\lambda \alpha^{\frac{\kappa}{\kappa-1}}\right)^{\frac{(\kappa-1)C'}{\kappa C' + \kappa - 1}}$ , we obtain

$$\log \pi_{-\lambda R} \exp(\lambda \alpha \Delta R) = \underset{\lambda \rightarrow +\infty}{\mathbf{O}} \left( \left[ \lambda \alpha^{\frac{\kappa}{\kappa-1}} \right]^{\frac{(\kappa-1)C'}{\kappa C' + \kappa - 1}} + \frac{0}{\lambda \rightarrow +\infty}(\lambda^{-s}) \right).$$

#### APPENDIX C. PROOF OF INEQUALITY (6.6)

We start with the following lemma.

**Lemma C.1.** *Let  $h : \mathbb{R}^* \rightarrow \mathbb{R}$  be a  $C^3$  convex function such that there exists  $u_0 > 0$  satisfying  $h'(u_0) = 0$  and  $h''(u_0) > 0$ . Let  $\phi : \mathbb{R}^* \rightarrow \mathbb{R}$  be a continuous non negative*

function such that  $\phi(u_0 > 0)$  and  $u \mapsto \phi(u) \exp(-t_0 u)$  integrable for some  $t_0 > 0$ . Then for any  $A > u_0$ , we have

$$\int_0^A \phi(u) \exp[-th(u)] du \underset{t \rightarrow +\infty}{\sim} \phi(u_0) \exp[-th(u_0)] \sqrt{\frac{2\pi}{th''(u_0)}}.$$

*Proof.* • Since the function  $h''$  is non negative, continuous,  $h''(u_0) > 0$  and  $h'(u_0) = 0$ , there exists  $c > 0$  such that for any  $u \in [0; A]$ ,  $h(u) - h(u_0) \geq c(u - u_0)^2$ . Let  $\alpha_t \triangleq t^{-p}$  with  $\frac{1}{3} < p < \frac{1}{2}$ . We have

$$\begin{aligned} & \int_{[0; u_0 - \alpha_t] \cup [u_0 + \alpha_t; A]} \phi(u) \exp[-th(u)] du \\ & \leq \exp[-th(u_0)] \int_{[0; u_0 - \alpha_t] \cup [u_0 + \alpha_t; A]} \phi(u) \exp[-tc(u - u_0)^2] du \\ & = \exp[-th(u_0)] \underset{t \rightarrow +\infty}{\mathbf{O}}(\exp[-ct\alpha_t^2]). \end{aligned}$$

• From Taylor's theorem, for any  $u \in [u_0 - \alpha_t; u_0 + \alpha_t]$ , there exists  $u^* \in [u_0 - \alpha_t; u_0 + \alpha_t]$  such that

$$h(u) = h(u_0) + \frac{h''(u_0)}{2}(u - u_0)^2 + \frac{h'''(u^*)}{6}(u - u_0)^3$$

Let  $A'' \triangleq \sup_{[u_0/2; A]} |h'''(u)|$  and  $I_t \triangleq \int_{[u_0 - \alpha_t; u_0 + \alpha_t]} \phi(u) \exp[-t \frac{h''(u_0)}{2}(u - u_0)^2] du$ . We get

$$\int_{[u_0 - \alpha_t; u_0 + \alpha_t]} \phi(u) \exp[-th(u)] du \leq \exp[A''t\alpha_t^3] \exp[-th(u_0)] I_t$$

and

$$\int_{[u_0 - \alpha_t; u_0 + \alpha_t]} \phi(u) \exp[-th(u)] du \geq \exp[-A''t\alpha_t^3] \exp[-th(u_0)] I_t.$$

We have

$$\begin{aligned} & \left| I_t - \int_{-\infty}^{+\infty} \phi(u_0) \exp\left[-t \frac{h''(u_0)}{2}(u - u_0)^2\right] du \right| \\ & \leq \int_{[u_0 - \alpha_t; u_0 + \alpha_t]} |\phi(u) - \phi(u_0)| \exp\left[-t \frac{h''(u_0)}{2}(u - u_0)^2\right] du \\ & \quad + \int_{]-\infty; u_0 - \alpha_t] \cup [u_0 + \alpha_t; +\infty[} \phi(u_0) \exp\left[-t \frac{h''(u_0)}{2}(u - u_0)^2\right] du \\ & \leq \underset{t \rightarrow +\infty}{\mathbf{O}}\left(\int_{-\infty}^{+\infty} \exp\left[-t \frac{h''(u_0)}{2}(u - u_0)^2\right] du\right) \\ & \quad + \underset{t \rightarrow +\infty}{\mathbf{O}}\left(\exp\left[-\frac{h''(u_0)}{2}t\alpha_t^2\right]\right) \end{aligned}$$

Since we have  $\int_{-\infty}^{+\infty} \exp\left[-t \frac{h''(u_0)}{2}(u - u_0)^2\right] du = \sqrt{\frac{2\pi}{th''(u_0)}}$ , we obtain

$$I_t = [\phi(u_0) + \underset{t \rightarrow +\infty}{\mathbf{O}}(1)] \sqrt{\frac{2\pi}{th''(u_0)}}.$$

• Combining the previous results, we obtain

$$\int_0^A \phi(u) \exp[-th(u)] du = [\phi(u_0) + \underset{t \rightarrow +\infty}{\mathbf{O}}(1)] \exp[-th(u_0)] \sqrt{\frac{2\pi}{th''(u_0)}}.$$

□

By assumption, we may write  $\pi(\Delta R \leq x) = \exp(-C'x^{-\frac{q}{\kappa}} - C''') [1 + \eta(x)]$  with  $\eta(x) = \underset{x \rightarrow 0}{\mathbf{O}}(1)$ . Let  $A' \triangleq \exp(-C''')$ ,  $u_0 \triangleq \operatorname{argmin}_{x>0} (x + C'x^{-\frac{q}{\kappa}})$ ,  $H \triangleq u_0 + C'u_0^{-\frac{q}{\kappa}}$  and  $\theta \triangleq 2H\lambda^{-\frac{\kappa}{\kappa+q}}$ .

From inequality (A.1), we have

$$(C.1) \quad \begin{aligned} \pi \exp(-\lambda \Delta R) &= \exp(-\lambda) + \int_0^1 \lambda \exp(-\lambda x) \pi(\Delta R \leq x) dx \\ &\leq \exp(-\lambda \theta) + \int_0^\theta \lambda \exp(-\lambda x) \pi(\Delta R \leq x) dx \end{aligned}$$

Besides, we have

$$\begin{aligned} \int_0^\theta \exp(-\lambda x) \pi(\Delta R \leq x) dx \\ &= A' \int_0^\theta \exp(-\lambda x - C' x^{-\frac{q}{\kappa}}) [1 + \eta(x)] dx \\ &= A' \int_0^{2H} \exp(-\lambda^{\frac{q}{\kappa+q}} [x + C' x^{-\frac{q}{\kappa}}]) [1 + \eta(\lambda^{-\frac{\kappa}{\kappa+q}} x)] dx \end{aligned}$$

For any  $\beta > 0$ , there exists  $\lambda_0$  such that for any  $\lambda > \lambda_0$  and any  $x \leq \theta$ , we have  $|\eta(\lambda^{-\frac{\kappa}{\kappa+q}} x)| \leq \beta$ . We obtain

$$\left| \frac{\int_0^\theta \exp(-\lambda x) \pi(\Delta R \leq x) dx}{A' \int_0^{2H} \exp(-\lambda^{\frac{q}{\kappa+q}} [x + C' x^{-\frac{q}{\kappa}}]) dx} - 1 \right| \leq \beta.$$

Using Lemma C.1, we get

$$(C.2) \quad \int_0^\theta \exp(-\lambda x) \pi(\Delta R \leq x) dx \underset{\lambda \rightarrow +\infty}{\sim} A' \int_0^{2H} \exp(-\lambda^{\frac{q}{\kappa+q}} [x + C' x^{-\frac{q}{\kappa}}]) dx \\ \underset{\lambda \rightarrow +\infty}{\sim} \exp(-\lambda^{\frac{q}{\kappa+q}} H)$$

So inequality (C.1) implies

$$(C.3) \quad \pi \exp(-\lambda \Delta R) \underset{\lambda \rightarrow +\infty}{\sim} \lambda \exp(-\lambda^{\frac{q}{\kappa+q}} H).$$

From equality (A.3), we have

$$\begin{aligned} \pi[\Delta R \exp(-\lambda \Delta R)] &= \exp(-\lambda) - \int_0^1 \exp(-\lambda x) \pi(\Delta R \leq x) dx \\ &\quad + \lambda \int_0^1 x \exp(-\lambda x) \pi(\Delta R \leq x) dx. \end{aligned}$$

Using similar computations to the one used to prove (C.2) and from the equality

$$\begin{aligned} \int_0^\theta x \exp(-\lambda x - C' x^{-\frac{q}{\kappa}}) [1 + \eta(x)] \\ &= \int_0^{2H} \lambda^{-\frac{\kappa}{\kappa+q}} x \exp(-\lambda^{\frac{q}{\kappa+q}} [x + C' x^{-\frac{q}{\kappa}}]) [1 + \eta(\lambda^{-\frac{\kappa}{\kappa+q}} x)] dx, \end{aligned}$$

we obtain

$$\pi[\Delta R \exp(-\lambda \Delta R)] \underset{\lambda \rightarrow +\infty}{\sim} \lambda^{\frac{q}{\kappa+q}} u_0 \exp(-\lambda^{\frac{q}{\kappa+q}} H).$$

Consequently, we have proved  $\pi_{-\lambda R} \Delta R \underset{\lambda \rightarrow +\infty}{\sim} u_0 \lambda^{-\frac{\kappa}{\kappa+q}}$ . By definition of  $u_0$ , we get  $\pi_{-\lambda R} \Delta R \underset{\lambda \rightarrow +\infty}{\sim} \left(\frac{qC'}{\kappa\lambda}\right)^{\frac{\kappa}{\kappa+q}}$ .

#### APPENDIX D. PROOF OF INEQUALITY (6.7)

Let  $0 < \alpha \leq \check{\alpha} \lambda^{-\frac{\kappa-1}{\kappa+q}}$  for some constant  $\check{\alpha} > 0$  to be determined, and  $\frac{1}{2} < \zeta < 1$ . We use once more the function  $h_\alpha(x) \triangleq x - \alpha x^{\frac{1}{\kappa}}$  which is minimum at  $x_0 \triangleq \left(\frac{\alpha}{\kappa}\right)^{\frac{\kappa}{\kappa-1}}$ . Let  $\nu \triangleq \eta\left\{\left(\frac{\alpha}{1-\zeta}\right)^{\frac{\kappa}{\kappa-1}}\right\}$ . We have

$$\begin{aligned} \pi \exp[-\lambda h_\alpha(\Delta R)] &= \pi \left\{ \exp[-\lambda h_\alpha(\Delta R)] \mathbb{1}_{h_\alpha(\Delta R) \leq \zeta \Delta R} \right\} \\ &\quad + \pi \left\{ \exp[-\lambda h_\alpha(\Delta R)] \mathbb{1}_{h_\alpha(\Delta R) > \zeta \Delta R} \right\} \\ &\leq \exp[-\lambda h_\alpha(x_0)] \pi \left\{ \Delta R \leq \left(\frac{\alpha}{1-\zeta}\right)^{\frac{\kappa}{\kappa-1}} \right\} + \pi \exp[-\lambda \zeta \Delta R] \\ &\leq \exp \left\{ \lambda(\kappa-1) \left(\frac{\alpha}{\kappa}\right)^{\frac{\kappa}{\kappa-1}} - C' \left(\frac{1-\zeta}{\alpha}\right)^{\frac{q}{\kappa-1}} - C''' \right\} (1 + \nu) \\ &\quad + \pi \exp[-\lambda \zeta \Delta R] \end{aligned}$$

Using (C.3), we obtain

$$\begin{aligned} \pi_{-\lambda R} \exp(\lambda \alpha \Delta R) &\leq \lambda^{-1} \exp\left(H \lambda^{\frac{q}{\kappa+q}} + \lambda(\kappa-1) \left(\frac{\alpha}{\kappa}\right)^{\frac{\kappa}{\kappa-1}} - C' \left(\frac{1-\zeta}{\alpha}\right)^{\frac{q}{\kappa-1}}\right) \left[\check{C} + \underset{\lambda \rightarrow +\infty}{\mathcal{O}}(1)\right] \\ &\quad + \zeta \exp\left(\lambda^{\frac{q}{\kappa+q}} H [1 - \zeta^{\frac{q}{\kappa+q}}]\right) \left[1 + \underset{\lambda \rightarrow +\infty}{\mathcal{O}}(1)\right] \end{aligned}$$

Let  $\zeta = 1 - \left(\frac{2H}{C'}\right)^{\frac{\kappa-1}{q}} \alpha \lambda^{\frac{\kappa-1}{\kappa+q}}$  so that  $C' \left(\frac{1-\zeta}{\alpha}\right)^{\frac{q}{\kappa-1}} = 2H \lambda^{\frac{q}{\kappa+q}}$  and let  $\check{c} > 0$  such that  $(\kappa-1) \left(\frac{\check{c}}{\kappa}\right)^{\frac{\kappa}{\kappa-1}} \leq H$ . Then we have

$$\begin{aligned} \pi_{-\lambda R} \exp(\lambda \alpha \Delta R) &\leq \lambda^{-1} \left[\check{C} + \underset{\lambda \rightarrow +\infty}{\mathcal{O}}(1)\right] + \exp\left\{\lambda^{\frac{q}{\kappa+q}} \frac{q}{\kappa+q} H \mathcal{O}(1 - \zeta)\right\} \left[1 + \underset{\lambda \rightarrow +\infty}{\mathcal{O}}(1)\right] \\ &= \underset{\lambda \rightarrow +\infty}{\mathcal{O}}\left(\exp\left\{\check{C} \lambda^{\frac{q}{\kappa+q}} \alpha \lambda^{\frac{\kappa-1}{\kappa+q}}\right\}\right), \end{aligned}$$

hence

$$\log \pi_{-\lambda R} \exp(\lambda \alpha \Delta R) = \underset{\lambda \rightarrow +\infty}{\mathcal{O}}\left(\lambda^{\frac{q}{\kappa+q}} \alpha \lambda^{\frac{\kappa-1}{\kappa+q}}\right).$$

#### APPENDIX E. ANOTHER WAY OF GETTING THE RIGHT ORDER

This section proves that by using well-known results, we can obtain a lower bound having the same spirit as Lemma 5.1 but without proper constants.

Applying Lemma 6.4 to the set of probability distributions

$$\mathcal{D} \triangleq \{\mathbb{P}^{\otimes N} : \mathbb{P} \in \mathcal{D}'\}$$

where  $\mathcal{D}' \triangleq \{\mathbb{P}_{\sigma^m} : \sigma_1^m \in \mathcal{S} \subset \{-1; +1\}^m\}$  and  $\mathcal{S}$  satisfies  $\delta(\Sigma, \Sigma') \geq \frac{m}{4}$  for any  $\Sigma \neq \Sigma' \in \mathcal{S}$  and  $|\mathcal{S}| = \lfloor e^{\frac{m}{8}} \rfloor$ . From Lemma 6.3, such a set  $\mathcal{S}$  exists. With any estimator  $\hat{f} : \mathcal{Z}^N \rightarrow \mathcal{F}(\mathcal{X}, \mathcal{Y})$ , we can associate an estimator  $\hat{T} : \mathcal{Z}^N \rightarrow \mathcal{D}$  defined as  $\hat{T}(Z_1^N) = \mathbb{P}^{\otimes N}$ , where  $\mathbb{P} \in \mathcal{D}'$  minimizes  $\mu[\xi(X) \mathbb{1}_{f_{\mathbb{P}}^*(X) \neq \hat{f}(Z_1^N)(X)}]$ .

By Birgé's lemma, we have  $\sup_{\mathbb{P} \in \mathcal{D}'} \mathbb{P}^{\otimes N}[\hat{T}(Z_1^N) \neq \mathbb{P}] \geq 0.36 \wedge \left(1 - \frac{K_{\mathcal{D}}}{|\mathcal{D}| \log |\mathcal{D}|}\right)$ .

Now, when  $\hat{T}(Z_1^N) \neq \mathbb{P}$ , we have  $R_{\mathbb{P}}(\hat{f}) - R_{\mathbb{P}}(f_{\mathbb{P}}^*) \geq \frac{m}{8} w b'$ . Therefore, we get

$$(E.1) \quad \sup_{\mathbb{P} \in \mathcal{D}} \{\mathbb{P}^{\otimes N} R_{\mathbb{P}}(\hat{f}) - R_{\mathbb{P}}(f_{\mathbb{P}}^*)\} \geq \frac{m}{8} w b' \left[0.36 \wedge \left(1 - \frac{K_{\mathcal{D}}}{|\mathcal{D}| \log |\mathcal{D}|}\right)\right].$$

For any  $\mathbb{P} \neq \mathbb{Q} \in \mathcal{D}'$ , we have

$$K(\mathbb{P}, \mathbb{Q}) \leq N \mu \left[ \xi(X) \log \left( \frac{1 + \xi(X)}{1 - \xi(X)} \right) \mathbb{1}_{X \notin \mathcal{X}_0} \right] \leq N m w b' \log \left( \frac{1 + B}{1 - B} \right),$$

where  $B \triangleq \sup_{x \in \mathcal{X} - \mathcal{X}_0} \xi(x)$ . If we assume that  $B \leq C b < 1$ , we get

$$K(\mathbb{P}, \mathbb{Q}) \leq C N m w b^2$$

for some constant  $C > 0$ . Since we have  $|\mathcal{D}| = \lfloor e^{\frac{m}{8}} \rfloor$ , we obtain

$$\frac{K_{\mathcal{D}}}{|\mathcal{D}| \log |\mathcal{D}|} \leq \frac{C}{\log \lfloor e^{\frac{m}{8}} \rfloor} N m w b^2 \leq C' N w b^2$$

for  $m$  large enough and some constant  $C' > 0$ . So we obtain the right order to the extent that when the quantity  $N w b^2$  is small enough, the order of the bound is given by the product  $m w b$ .

## APPENDIX F. PROOF OF THEOREM 5.2

• When  $L = 0$ : let  $x_0, x_1, \dots, x_{V-1}$  denote the  $V$  points shattered by the model. Let us take

$$\begin{cases} m = V - 1 \\ \mathcal{X}_0 = \mathcal{X} - \{x_1, \dots, x_{V-1}\} \\ \mathcal{X}_j = \{x_j\} \\ \mu(\mathcal{X}_j) = w \quad \text{for any } j \in \{1, \dots, m\} \\ \mu(\{x_0\}) = 1 - mw \\ b = 1 \quad (\xi \equiv 1) \end{cases} ,$$

where  $w$  is a free positive parameter which satisfies  $mw \leq 1$  (since  $\mu$  is a probability distribution). By noticing that

$$1 - V \left( \mathbb{P}_{-1,1,\dots,1}^{\otimes N}, \mathbb{P}_{1,1,\dots,1}^{\otimes N} \right) = \mu^{\otimes N}(\text{for any } i \in \{1, \dots, N\}, X_i \notin \mathcal{X}_j) = (1 - w)^N$$

and using inequality (6.26), we obtain

$$\sup_{\mathbb{P} \in \mathcal{P}} \left\{ \mathbb{P}^{\otimes N} R_{\mathbb{P}}(\hat{f}) - R_{\mathbb{P}}(f^*) \right\} \geq \frac{V-1}{2} \sup_{w \leq \frac{1}{V-1}} \left\{ w(1-w)^N \right\}$$

This supremum is attained for  $w = \frac{1}{N+1}$  when  $N \geq (V-2) \vee 1$  and for  $w = \frac{1}{V-1}$  otherwise.

• When  $0 \leq L \leq \frac{1}{2}$ : once more,  $x_0, x_1, \dots, x_{V-1}$  denote the  $V$  points shattered by the model. This time, we take

$$\begin{cases} m = V - 1 \\ \mathcal{X}_0 = \mathcal{X} - \{x_1, \dots, x_{V-1}\} \\ \mathcal{X}_j = \{x_j\} \\ \mu(\mathcal{X}_j) = w \quad \text{for any } j \in \{1, \dots, m\} \\ \mu(\{x_0\}) = 1 - mw \\ \xi(x) = \begin{cases} b_0 & \text{when } x \in \mathcal{X}_0 \\ b & \text{otherwise} \end{cases} \end{cases} ,$$

where  $w$  is a free positive parameter which satisfies  $mw \leq 1$  (since  $\mu$  is a probability distribution) and  $b$  and  $b_0$  belong to  $[0; 1]$ . Since we have

$$L = \frac{1}{2}mw(1-b) + \frac{1}{2}(1-mw)(1-b_0)$$

and  $b_0 \in [0; 1]$ , the parameters  $m, w$  and  $b$  should satisfy

$$mw(1-b) \leq 2L \leq 1 - mw.$$

Since this condition implies that  $mw \leq 1$ , we have the following lower bound

$$\sup_{\mathbb{P} \in \mathcal{P}} \left\{ \mathbb{P}^{\otimes N} R_{\mathbb{P}}(\hat{f}) - R_{\mathbb{P}}(f^*) \right\} \geq \sup_{\substack{w \geq 0 \\ 0 \leq b \leq 1 \\ mw(1-b) \leq 2L \leq 1 - mw}} \frac{1}{2}mw b [1 - b\sqrt{Nw}].$$

From this lower bound, one can recover the first assertion of Theorem 5.2 with a constant slightly worsened (due to the upper bound (6.27)). We will now slightly weaken this result in order to get a simple lower bound. Introduce  $x = b^2wN$ . The

previous supremum can be written as

$$\begin{aligned}
& \sup_{\substack{x>0 \\ 0<b\leq 1 \\ \frac{(V-1)x}{bN} \frac{1-b}{b} \leq 2L \\ \frac{(V-1)x}{bN} \leq 1-2L}} \frac{1}{2} \frac{(V-1)x}{bN} [1 - \sqrt{x}] \\
& \geq \sup_{\substack{x>0 \\ 0<b<1 \\ \frac{(V-1)x}{bN} \frac{1-b}{b} = 2L \\ 2L \frac{1-b}{b} \leq 1-2L}} \frac{1}{2} \frac{(V-1)x}{bN} [1 - \sqrt{x}] \\
& = \sup_{\substack{x>0 \\ b = \frac{1}{1 + \sqrt{1 + \frac{8LN}{(V-1)x}}} \\ b \leq 1-2L}} \frac{1}{2} \frac{(V-1)x}{bN} [1 - \sqrt{x}] \\
& = \sup_{0 < x \leq \frac{(1-2L)^2 N}{V-1}} \frac{(V-1)x}{4N} \left( 1 + \sqrt{1 + \frac{8LN}{(V-1)x}} \right) [1 - \sqrt{x}] \\
& > \sup_{0 < x \leq \frac{(1-2L)^2 N}{V-1}} \sqrt{\frac{L(V-1)x}{2N}} [1 - \sqrt{x}] \tag{A} \\
& = \begin{cases} \sqrt{\frac{L(V-1)}{32N}} & \text{when } \frac{(1-2L)^2 N}{V-1} \geq \frac{1}{4} \\ \sqrt{\frac{L(1-2L)^2}{8}} & \text{otherwise.} \end{cases}
\end{aligned}$$

Note that the step (A) prevents us to have a good lower bound when  $L = o\left(\frac{V-1}{N}\right)$ . In this last case, the lower bound (A) can be replaced with  $\frac{V-1}{2N}x(1 - \sqrt{x})$  which, by taking  $x = \frac{1}{4}$ , leads to the desired bound  $\frac{V-1}{16N}$ .

#### REFERENCES

1. J.-Y. Audibert, *Data-dependent generalization error bounds for (noisy) classification: the PAC-Bayesian approach*, Preprint, Laboratoire de Probabilité et Modèles Aléatoires, 2004.
2. L. Birgé, *A new look at an old result: Fano's lemma*, Preprint, Laboratoire de Probabilité et Modèles Aléatoires, 2001.
3. O. Catoni, *Statistical learning theory and stochastic optimization*, Lecture notes, Saint-Flour summer school on Probability Theory, 2001, Springer, to appear.
4. ———, *A PAC-Bayesian approach to adaptive classification*, Preprint, Laboratoire de Probabilité et Modèles Aléatoires, 2003.
5. L. Devroye and G. Lugosi, *Combinatorial methods in density estimation*, Springer-Verlag, 2000.
6. R.M. Dudley, *Central limit theorems for empirical measures*, Ann. Probab. **6** (1978), 899–929.
7. D. Haussler, *Sphere packing numbers for subsets of the boolean  $n$ -cube with bounded vapnik-chervonenkis dimension*, Journal of Combinatorial Theory **69** (1995), 217–232, Series A.
8. C. Huber, *Lower bounds for function estimation*, Research paper in probability and statistics: Festschrift for Lucien Le Cam (1996), 245–258.
9. V.I. Koltchinskii, *On the central limit theorem for empirical measures*, Theory Probab. Math. Stat. **24** (1981), 71–82.
10. G. Lugosi, *Concentration-of-measure inequalities*, 2003, Lecture notes, Machine Learning Summer School, Canberra.
11. E. Mammen and A.B. Tsybakov, *Smooth discrimination analysis*, Ann. Stat. **27** (1999), 1808–1829.
12. P. Massart and E. Nédélec, *Risk bounds for statistical learning*, Available from <http://www.math.u-psud.fr/~massart/margin.pdf>, 2003.
13. D. A. McAllester, *PAC-Bayesian model averaging*, Morgan Kaufmann Publishers, 1999.
14. D. Pollard, *A central limit theorem for empirical measures*, J. Aust. Math. Soc., Ser. A **33** (1982), 235–248.

15. A.B. Tsybakov, *Optimal aggregation of classifiers in statistical learning*, Ann. Stat. **32** (2004), no. 1.
16. A.B. Tsybakov and S. van de Geer, *Square root penalty: adaptation to the margin in classification and in edge estimation*, Preprint, Laboratoire de Probabilité et Modèles Aléatoires, 2004.
17. A. van der Vaart and J. Wellner, *Weak convergence and empirical processes with application to statistics*, John Wiley & Sons, New York, 1996.