

# Model selection for weakly dependent time series forecasting

PIERRE ALQUIER<sup>1</sup> and OLIVIER WINTENBERGER<sup>2</sup>

<sup>1</sup>*Laboratoire de Probabilités et Modèles Aléatoires, Université Paris 7, site Chevaleret, 175, rue du Chevaleret, 75205 Paris Cedex 13, France, and CREST, Laboratoire de Statistique, 3, avenue Pierre Larousse, 92240 Malakoff, France. E-mail: [alquier@math.jussieu.fr](mailto:alquier@math.jussieu.fr)*

<sup>2</sup>*CEREMADE, Université Paris Dauphine, Place du Maréchal De Lattre De Tassigny, 75775 Paris Cedex 16, France. E-mail: [wintenberger@ceremade.dauphine.fr](mailto:wintenberger@ceremade.dauphine.fr)*

Observing a stationary time series, we propose a two-steps procedure for the prediction of its next value. The first step follows machine learning theory paradigm and consists in determining a set of possible predictors as randomized estimators in (possibly numerous) different predictive models. The second step follows the model selection paradigm and consists in choosing one predictor with good properties among all the predictors of the first step. We study our procedure for two different types of observations: causal Bernoulli shifts and bounded weakly dependent processes. In both cases, we give oracle inequalities: the risk of the chosen predictor is close to the best prediction risk in all predictive models that we consider. We apply our procedure for predictive models as linear predictors, neural networks predictors and nonparametric autoregressive predictors.

*Keywords:* adaptative inference; aggregation of estimators; autoregression estimation; model selection; randomized estimators; statistical learning; time series prediction; weak dependence

## 1. Introduction

When observing a time series, one crucial issue is to predict the (nonobserved) first future value using the observed past values. Since the seventies, different model selection procedures have been studied for inferring how many observed past values are needed for predicting the next value. Procedures as AIC [1], BIC (Schwarz [27]) and APE (Ing [18]) are used by practitioners to select a reasonable linear predictor. When the observations satisfy a linear model, those procedures are proved to be asymptotically efficient (see Ing [18] for more details).

In the same time, the progress of statistical learning theory in the i.i.d. setting brought new perspectives in model selection (see Vapnik [30] and Massart [20] among others). Machine-learning procedures allow to choose a predictor among a family, with the guarantee that this predictor performs almost as well as the best possible predictor of the family (called the oracle). Such results are called oracle inequalities; they provide guarantees on the quality of the prediction without any parametric assumption on the observations.

Few works have been done in the context of dependent observations. The machine learning theory was used successfully in the time series prediction context by Modha and Masry [23] and Meir [22]. However, their procedure relies on the knowledge of the  $\alpha$ -mixing coefficients. To our knowledge, there is no efficient estimation of this coefficients and their procedure seems difficult

to use in practice. Baraud *et al.* [5] use the model selection point of view to perform regression and auto-regression on dependent observations. They prove powerful oracle inequalities when the observations satisfy an additive auto-regressive model. When the observations are Harris recurrent Markov chains, Lacour [19] gives also oracle inequalities for a procedure completely free of the dependence properties. An alternative point of view is provided by the theory of individual sequences prediction (see Lugosi and Cesa-Bianchi [10] or Stoltz [29]). In these works, no assumption on the observations – not even a stochastic assumption – is done and oracle inequalities are given.

In this paper, our objectives are the following:

- (1) to build various predictors of different forms and using different numbers of past observations,
- (2) to select one of these predictors *without any assumption on the distribution of the observations*,
- (3) to prove oracle inequalities under weak assumptions on the observed time series.

In the end of this Introduction, let us fix the mathematical framework (see also Meir [22] for more details).

Let us observe  $(X_1, \dots, X_n)$  from a stationary time series  $X = (X_t)_{t \in \mathbb{Z}}$  distributed as  $\pi_0$  on  $\mathcal{X}^{\mathbb{Z}}$  where  $\mathcal{X}$  is an Hilbert space equipped with its usual norm  $\|\cdot\|$ . Fix a (possibly large) family of predictors  $\{f_\theta, \theta \in \Theta\}$ : for any  $\theta$  and any  $t$ ,  $f_\theta$  applied to the past values  $(X_{t-1}, X_{t-2}, \dots, X_1)$  is a possible prediction of  $X_t$ . We discretize the family of predictors by the number  $p$  of past values they use. Thus, we assume that

$$\Theta = \bigcup_{p=1}^{\lfloor n/2 \rfloor} \Theta_p,$$

where the  $\Theta_p$  are disjoint in order that for any  $\theta \in \Theta$ , there is only one  $p$  such that  $\theta \in \Theta_p$ . Now, for any  $\theta \in \Theta_p$ ,  $f_\theta$  is a function  $\mathcal{X}^p \rightarrow \mathcal{X}$  and at any time  $t$ ,  $f_\theta(X_{t-1}, \dots, X_{t-p})$  is a prediction of  $X_t$  according to  $\theta$  and denoted  $\hat{X}_t^\theta$ . As the predictor  $f_\theta$  may take different forms (linear functions, neural networks, ...), we write

$$\Theta_p = \bigcup_{\ell=1}^{m_p} \Theta_{p,\ell}$$

for a given  $m_p \in \{1, \dots, n\}$ . Finally, the risk of the prediction,  $R(\theta)$ , is defined by

$$R(\theta) = \pi_0[\|f_\theta(X_{t-1}, \dots, X_{t-p}) - X_t\|] = \pi_0[\|\hat{X}_t^\theta - X_t\|],$$

where here and all along the paper  $\pi[h] = \int h d\pi$  for any measure  $\pi$  and any integrable function  $h$ . Note that  $R(\theta)$  does not depend on  $t$  as  $X$  is stationary.

The mathematical counterparts of the points (1), (2) and (3) of our objectives are the following. The point (1) corresponds to build, on the basis of the observations, an estimator  $\hat{\theta}_{p,\ell}$  in each model  $\Theta_{p,\ell}$ , for  $1 \leq p \leq \lfloor n/2 \rfloor$  and  $1 \leq \ell \leq m_p$ . The point (2) consists in defining a procedure

to choose a  $\hat{\theta}$  among all the possible  $\hat{\theta}_{p,\ell}$ . Finally, point (3) is achieved by proving that  $R(\hat{\theta})$  is close to  $\inf_{\theta \in \Theta} R(\theta)$ . To attain these objectives, we use the PAC-Bayesian paradigm (introduced by Shawe-Taylor and Williamson [28] and McAllester [21]). Using this approach, Catoni [7–9], Audibert [4], Alquier [2], Tsybakov and Dalalyan [11] solve points (1), (2) and (3) simultaneously for various regression and classification problems in the i.i.d. setting. In this paper, we build a procedure that gives a predictor  $\hat{\theta}$  satisfying, under general conditions on  $X$  and with probability at least  $1 - \varepsilon$ ,

$$R(\hat{\theta}) \leq \inf_{d_{p,\ell} \leq n} \left\{ \min_{\theta \in \Theta_{p,\ell}} R(\theta) + \text{cst} \cdot \sqrt{\frac{d_{p,\ell}}{n}} \log^{5/2}(n) \right\} + \text{cst} \cdot \frac{\log(1/\varepsilon)}{\sqrt{n}},$$

where  $\text{cst} > 0$  is an explicit constant and  $d_{p,\ell}$  an estimate of the complexity of  $\Theta_{p,\ell}$ .

To obtain such oracle inequalities, we use sharp estimates (close to the ones in the i.i.d. case) on the Laplace transforms of the partial sums in dependent settings. For bounded observations, we use the  $\theta_\infty$ -coefficients (see [12]), introduced in Rio [25] as the  $\gamma$ -mixing coefficients. These coefficients generalize the uniform mixing ones. For unbounded observations, we use the causal Bernoulli shifts representation. It includes all classical linear ARMA models and also the more general chains with infinite memory introduced by Doukhan and Wintenberger [15]. These bounded and unbounded dependent frameworks are not comparable with the  $\beta$  or  $\alpha$ -mixing ones as they include some dynamical systems that are not mixing, see Andrews [3] and Dedecker and Prieur [13] for details. Finally, it is important to note that our prediction procedure is the same for the two dependence frameworks and does not depend on any unknown dependence coefficient. It is an advantage of our approach because it is impossible to estimate efficiently the dependence coefficients we use.

The paper is organized as follows: First, the prediction procedure is detailed in Section 2; Second, the assumptions on the observed time series and the corresponding oracle inequalities are given in Section 3. In Section 4, are given some examples of time series for which these oracle inequalities hold. Our procedure applied on some possible prediction models are given in Section 5. Linear predictors (with simulations), neural networks predictors and non-parametric predictors are considered. Finally, the complete proofs are collected in Section 6.

## 2. The prediction procedure

We observe  $(X_1, \dots, X_n)$  from a stationary time series  $X = (X_t)_{t \in \mathbb{Z}}$  distributed as  $\pi_0$  on  $\mathcal{X}^{\mathbb{Z}}$  where  $\mathcal{X}$  is an Hilbert space equipped with its usual norm  $\|\cdot\|$ . We fix a family of predictors  $\{f_\theta, \theta \in \Theta\}$  with

$$\Theta = \bigcup_{p=1}^{\lfloor n/2 \rfloor} \Theta_p = \bigcup_{p=1}^{\lfloor n/2 \rfloor} \left( \bigcup_{\ell=1}^{m_p} \Theta_{p,\ell} \right)$$

such that  $m_p \geq n$  and  $p(\theta)$  is the only  $p$  such that  $\theta \in \Theta_p$ . For any  $\theta \in \Theta$ , we denote  $\hat{X}_t^\theta = f_\theta(X_{t-1}, \dots, X_{t-p})$  and  $R(\theta) = \pi_0[\|\hat{X}_t^\theta - X_t\|]$ .

### 2.1. The Lipschitz predictors

Let  $M$  denotes the set of all possible pairs  $(p, \ell)$ :

$$M = \bigcup_{p=1}^{\lfloor n/2 \rfloor} \{p\} \times \{1, \dots, m_p\}.$$

Let  $\mathcal{T}$  be a  $\sigma$ -algebra on  $\Theta$  and  $\mathcal{T}_{p,\ell}$  be its restriction to  $\Theta_{p,\ell}$  for any  $(p, \ell) \in M$ . For any  $(p, \ell) \in M$ , we assume that  $\Theta_{p,\ell}$  is a compact subset of  $\mathbb{R}^q$  for some  $q < \infty$  ( $q$  depends on  $(p, \ell)$ ) and that there exists  $(a_j(\theta))_{j \in \{1, \dots, p\}}$  satisfying, for any  $(x_1, \dots, x_p), (y_1, \dots, y_p) \in \mathcal{X}^p$ , the relation

$$\|f_\theta(x_1, \dots, x_p) - f_\theta(y_1, \dots, y_p)\| \leq \sum_{j=1}^p a_j(\theta) \|x_j - y_j\|. \tag{2.1}$$

In order to bound the volatility of the predictors uniformly on  $M$ , we assume that

$$L := \sup_{(p,\ell) \in M} \sup_{\theta \in \Theta_{p,\ell}} \sum_{j=1}^p a_j(\theta) \quad \text{satisfies} \quad L \leq \log(n) - 1. \tag{2.2}$$

### 2.2. The complexity of $\Theta_{p,\ell}$

To control the complexity of each  $\Theta_{p,\ell}$  we assume that, for all  $(p, \ell) \in M$ , there exist a probability measure  $\pi_{p,\ell}$  on the measurable space  $(\Theta_{p,\ell}, \mathcal{T}_{p,\ell})$  and a constant  $1 \leq d_{p,\ell} < \infty$  satisfying

$$\sup_{\gamma > e} \left\{ \frac{-\log \int_{\Theta_{p,\ell}} [\exp(-\gamma(R(\theta) - R(\bar{\theta}_{p,\ell})))] d\pi_{p,\ell}(\theta)}{\log(\gamma)} \right\} \leq d_{p,\ell}. \tag{2.3}$$

Here  $\bar{\theta}_{p,\ell} = \arg \min_{\Theta_{p,\ell}} R$  for any  $(p, \ell) \in M$ . The parameter  $d_{p,\ell}$  is linked with classical complexities as the Vapnik dimension and entropy measures. In this paper, we only investigate the case where  $\pi_{p,\ell}$  is the Lebesgue measure on  $\Theta_{p,\ell}$ . We have the following result.

**Proposition 2.1.** *Let  $q \in \mathbb{N}^*$ ,  $x > 0$  and  $\mathcal{B}_x^q$  be the closed  $\ell^1$ -ball in  $\mathbb{R}^q$  of radius  $x > 0$  and centered at 0. If  $\Theta_{p,\ell} = \mathcal{B}_{c_{p,\ell}}^q$  for  $c_{p,\ell} > 0$  and  $\theta \rightarrow R(\theta)$  is a  $C$ -Lipschitz function then we have:*

$$d_{p,\ell} \leq q \times \left( 1 + \log \left( c_{p,\ell} \left( \frac{Ce}{q} \vee \frac{1}{c_{p,\ell} - \|\bar{\theta}_{p,\ell}\|} \right) \right) \right). \tag{2.4}$$

The proof of this result is given at the end of Section 6.4. Predictive models where complexity  $d_{p,\ell}$  is estimated are given in Section 5.

### 2.3. The empirical risk

As the risk  $R(\theta)$  cannot be computed, we use its empirical counterpart  $r_n(\theta)$ :

$$r_n(\theta) = \frac{1}{n - p(\theta)} \sum_{t=p(\theta)+1}^n \|X_t - \widehat{X}_t^\theta\|.$$

### 2.4. The randomized estimators

For any  $(p, \ell) \in M$ , our randomized estimators  $\tilde{\theta}_{p,\ell}^\lambda$  is drawn randomly through a Gibbs measure

$$\tilde{\theta}_{p,\ell}^\lambda \sim \pi_{p,\ell}\{-\lambda r_n\}.$$

We recall that for any measure  $\pi$  and any measurable function  $h$  such that  $\pi[\exp(h)] < +\infty$ , the Gibbs measure denoted  $\pi\{h\}$  is defined by the relation:

$$\frac{d\pi\{h\}}{d\pi}(\theta) = \frac{\exp(h(\theta))}{\pi[\exp(h)]}. \tag{2.5}$$

Here the parameter  $\lambda$  is called the temperature (this terminology comes from the statistical thermodynamics). For  $n \geq 8e(1 + L)$ ,  $\lambda$  takes values in a finite grid  $\mathcal{G}_{p,\ell}$  defined as

$$\mathcal{G}_{p,\ell} = \left\{ g_1 \frac{\sqrt{d_{p,\ell} n \log(d_{p,\ell} n)}}{(1 + L) \log^{3/2}(n)}, \dots, g_{n_0} \frac{\sqrt{d_{p,\ell} n \log(d_{p,\ell} n)}}{(1 + L) \log^{3/2}(n)} \right\} \cap \left[ 2e, \frac{n}{4(1 + L)} \right],$$

where  $\check{c} \leq g_1 < \dots < g_{n_0} \leq \hat{c}$  with  $2 \leq n_0 \leq n$  and  $0 < \check{c} < 2/(1 + L) < 2e(1 + L) < \hat{c} < \infty$ . Remark that when  $\lambda$  grows,  $\pi_{p,\ell}\{-\lambda r_n\}$  tends to concentrate around the minimizer of the empirical risk.

### 2.5. The model selection

One way to select a predictor is to choose the minimizer of the penalized empirical risk  $\arg \min_{p,\ell} [r_n(\tilde{\theta}_{p,\ell}^\lambda) + \text{pen}(p, \ell, \lambda)]$ , for some well chosen penalization  $\text{pen}(p, \ell, \lambda)$ , see Massart [20]. Here we consider  $\hat{\theta} = \tilde{\theta}_{\hat{p},\hat{\ell}}^\lambda$  where

$$(\hat{p}, \hat{\ell}, \hat{\lambda}) = \arg \min_{\substack{(p,\ell) \in M \\ \lambda \in \mathcal{G}_{p,\ell}}} \hat{R}(p, \ell, \lambda).$$

The model criterion  $\hat{R}(p, \ell, \lambda)$  is given by the PAC-Bayesian approach:

$$\hat{R}(p, \ell, \lambda) = -\frac{1}{\lambda} \log \int_{\Theta_{p,\ell}} \exp(-\lambda r_n(\theta)) d\pi_{p,\ell}(\theta) + \frac{1}{\lambda} \log \left( n \left[ \frac{n}{2} \right] m_p \right) + \frac{\lambda(1 + L)^2 \log^3(n)}{n(1 - p/n)^2}.$$

### 3. Main results

In order to prove that  $R(\hat{\theta})$  is close to  $\inf_{\theta \in \Theta} R(\theta)$  with high probability, we restrict our study to two different contexts. Note that  $\hat{\theta}$  is defined independently of these contexts and that a practitioner may compute our predictor on any observed time series.

#### 3.1. Bounded weakly dependent processes (WDP)

In this case,  $X$  is bounded, that is,  $\|X\|_\infty := \sup_t \|X_t\| < \infty$ . We use the  $\theta_{\infty,n}(1)$ -coefficients in Dedecker *et al.* [12], a version of the  $\gamma$ -mixing of Rio [26]) adapted to stationary time series. If  $Z$  is a bounded variable in  $\mathcal{X}^q$  ( $q \geq 1$ ) defined on  $(\Omega, \mathcal{A}, \mathbb{P})$ , for any  $\sigma$ -algebra  $\mathfrak{G}$  of  $\mathcal{A}$  we have:

$$\theta_\infty(\mathfrak{G}, Z) = \sup_{f \in \Lambda_1} \|\mathbb{E}(f(Z)|\mathfrak{G}) - \mathbb{E}(f(Z))\|_\infty,$$

where  $\Lambda_1$  is the set of real 1-Lipschitz functions on  $\mathcal{X}^q$  equipped with the norm  $\|z\| = \sum_{i=1}^q \|z_i\|$ . Let us define the  $\sigma$ -algebra  $\mathfrak{G}_p = \sigma(X_t, t \leq p)$  for any  $p \in \mathbb{Z}$  and the coefficients

$$\theta_{\infty,k}(1) = \sup\{\theta_\infty(\mathfrak{G}_p, (X_{j_1}, \dots, X_{j_\ell})), p + 1 \leq j_1 < \dots < j_\ell, 1 \leq \ell \leq k\}.$$

Moreover, assume that there is a constant  $C > 0$  such that for any  $n$ ,  $\theta_{\infty,n}(1) < C$  (the short memory condition). Causal Bernoulli shifts with bounded innovations, uniform  $\varphi$ -mixing sequences and dynamical systems are classical  $\theta_\infty$  weakly-dependent examples, see Section 4 for more details. In this context, we prove the following oracle inequality.

**Theorem 3.1.** *Under (WDP) and condition (2.3), there are explicit constants*

$$(\text{cst}_1, \text{cst}_2) = \text{cst}(\check{c}, \hat{c}, L, C, \|X_0\|_\infty)$$

such that for all  $n \geq 8e(1 + L)$  with probability at least  $1 - \varepsilon$

$$R(\hat{\theta}) \leq \inf_{d_{p,\ell} \leq n} \left\{ \min_{\theta \in \Theta_{p,\ell}} R(\theta) + \text{cst}_1 \cdot \sqrt{\frac{d_{p,\ell}}{n}} \log^{5/2}(n) \right\} + \text{cst}_2 \cdot \frac{\log(1/\varepsilon)}{\sqrt{n}} + 4(1 + L) \left( \frac{(\|X_0\|_\infty + C)^2}{2} - \log^3(n) \right)_+.$$

The proof of this result is given in Section 6.2 page 899.

#### 3.2. Causal Bernoulli shifts (CBS)

Let  $\mathcal{X}'$  be some Banach space equipped with a norm also denoted  $\|\cdot\|$ . Let  $H : \mathcal{X}'^{\mathbb{N}} \mapsto \mathcal{X}$  be a satisfying, for some sequence  $(a_j(H))_{j \in \mathbb{N}}$ , and for any  $v = (v_j)_{j \in \mathbb{N}}$ ,  $v' = (v'_j)_{j \in \mathbb{N}} \in \mathcal{X}'^{\mathbb{N}}$ , the

relations:

$$\|H(v) - H(v')\| \leq \sum_{j=0}^{\infty} a_j(H) \|v_j - v'_j\|, \tag{3.1}$$

with

$$\sum_{j=0}^{\infty} ja_j(H) < +\infty. \tag{3.2}$$

We denote  $\sum_{j=0}^{\infty} a_j(H) := a(H)$ ,  $\sum_{j=0}^{\infty} ja_j(H) = \tilde{a}(H)$ . The causal Bernoulli shifts are defined by the relation

$$X_t = H(\xi_t, \xi_{t-1}, \xi_{t-2}, \dots) \quad \forall t \in \mathbb{Z},$$

where  $\xi_t$  for  $t \in \mathbb{Z}$  are i.i.d. variables called the innovations and distributed as  $\mu$ . We assume that we can choose, by quantile transformation, innovations that admit a finite Laplace transform  $\mu[\exp(c^* \|\xi_0\|)] := \Psi(c^*) < +\infty$  (the Cramer condition) for  $c^* \geq a(H)$ . Classical examples of such processes are causal linear ARMA models and chains with infinite memory with low-tail innovations, see Section 4 for more details. In this context, we prove the following oracle inequality

**Theorem 3.2.** *Under (CBS) and condition (2.3), there are explicit constants*

$$(cst'_1, cst'_2) = cst'(\check{c}, \hat{c}, L, a(H), \tilde{a}(H), \Psi(1))$$

such that for all  $n \geq 8e(1 + L)$  with probability at least  $1 - \varepsilon$

$$\begin{aligned} R(\hat{\theta}) &\leq \inf_{d_{p,\ell} \leq n} \left\{ \min_{\theta \in \Theta_{p,\ell}} R(\theta) + cst'_1 \cdot \sqrt{\frac{d_{p,\ell}}{n}} \log^{5/2}(n) \right\} + cst'_2 \cdot \frac{\log(1/\varepsilon)}{\sqrt{n}} \\ &\quad + \sqrt{\frac{d_{\hat{p},\hat{\ell}}}{n}} \log(d_{\hat{p},\hat{\ell}} n) 4(1 + L) \\ &\quad \times \hat{c} \left( 4a(H)\Psi(a(H)) + 2 \log^2(n) \left( 1 + \frac{\tilde{a}(H)}{a(H)} \right)^2 - \log^3(n) \right)_+. \end{aligned}$$

The proof of this result is given in Section 6.3 page 902.

### 3.3. Comments on the results

The constants are roughly (but explicitly) estimated in the proofs, see Sections 6.2 and 6.3. For example, we obtain

$$cst_1 \leq (1 + L) \left( \frac{6}{\check{c}} + 8\hat{c}(1 + \|X_0\|_{\infty} + C)^2 \right) \quad \text{and} \quad cst_2 \leq \frac{7(1 + L)}{\check{c}}.$$

For  $n$  sufficiently large, the last terms in the oracle inequalities vanish. Then it exists a constant  $C > 0$  such that under (WDP) or (CBS) for all  $n \geq 8e(1 + L)$  with probability at least  $1 - \varepsilon$ :

$$R(\hat{\theta}) \leq \inf_{d_{p,\ell} \leq n} \left\{ \min_{\theta \in \Theta_{p,\ell}} R(\theta) + C \sqrt{\frac{d_{p,\ell}}{n}} \log^{5/2}(n) \right\} + C \frac{\log(1/\varepsilon)}{\sqrt{n}}.$$

Similar oracles inequalities have already been proved by Modha and Masry [23] and Baraud *et al.* [5]. These inequalities are given in expectation while ours are true with high probability. Remark that integrating our oracle inequalities with respect to  $\varepsilon$  leads to a result in expectation: there exists a constant  $C > 0$  independent of  $n$  such that in both (WDP) and (CBS) cases

$$\pi_0[R(\hat{\theta})] \leq \inf_{d_{p,\ell} \leq n} \left\{ \min_{\theta \in \Theta_{p,\ell}} R(\theta) + C \sqrt{\frac{d_{p,\ell}}{n}} \log^{5/2}(n) \right\}.$$

The converse is not true: results in expectation do not lead to results that hold with high probability.

It is difficult to compare our oracle inequalities with the ones in [23] and [5]. Unlike our paper, those articles deal with the quadratic risk and ( $\beta$ - or  $\alpha$ -) mixing time series. However, remark that the additional terms in our oracle inequalities are proportional to  $\sqrt{d_{p,\ell}/n}$ , the rate in the i.i.d. case, times a term  $\log^{5/2}(n)$  term. Baraud *et al.* [5] obtain an oracle inequality for the quadratic risk with the same rate than in the i.i.d. case, while the one in Modha and Masry [23] suffers a loss  $(n/d_{p,\ell})^c$  for some  $c > 0$ .

## 4. Examples of time series satisfying (WDP) or (CBS)

We present several examples of time series satisfying (WDP) or (CBS).

### 4.1. Causal Bernoulli shifts

Causal Bernoulli shifts are stationary time series that admit the representation

$$X_t = H(\xi_t, \xi_{t-1}, \xi_{t-2}, \dots) \quad \forall t \in \mathbb{Z}, \tag{4.1}$$

where the  $\xi_t$  are i.i.d. variables called innovations. Almost all known stationary and ergodic processes have this form. However, we work here under the restrictive assumption (4.3). Remark that under this Lipschitz condition the existence of the stationary time series  $(X_t)$  follows from (4.1) and it satisfies the Cramer condition as soon as the innovations do. Some examples of causal Bernoulli shifts are presented below.

#### 4.1.1. Linear models

Let  $(X_t)$  be a real time series admitting the  $MA(\infty)$  representation

$$X_t = \sum_{j=0}^{\infty} a_j \xi_{t-j} \quad \text{with} \quad \sum_{j=0}^{\infty} j|a_j| < +\infty.$$

Then it satisfies (CBS) if the i.i.d. innovations  $\xi_t$  satisfy the Cramer condition. As an example, there is any causal AR( $\infty$ ) model  $X_t = \phi_0 + \sum_{j=1}^{\infty} \phi_j X_{t-j} + \xi_t$  with  $\phi(z) = 1 - \sum_{j=1}^{\infty} \phi_j z^j$  that have no root for  $|z| \leq 1$  (such that causal ARMA( $p, q$ ) models). Indeed, as  $\phi$  is a real analytic function on the unit disc,  $1/\phi$  is a well a real analytic function  $1/\phi(z) = \sum_{j=1}^{\infty} \psi_j z^j$  with the coefficients  $\psi_j$  that decrease exponentially fast (i.e., (3.2) is automatically satisfied).

#### 4.1.2. Chains with infinite memory

Chains with infinite memory is a class of time series  $(X_t)$  introduced by Doukhan and Wintenberger [15] as the solution of the equation

$$X_t = F(X_{t-1}, X_{t-2}, \dots; \xi_t) \quad \text{almost surely} \tag{4.2}$$

for some function  $F: \mathcal{X}^{(\mathbb{N} \setminus \{0\})} \times \mathcal{X}' \rightarrow \mathcal{X}$ . Assume also that for some  $u > 0$ , for all  $x = (x_k)_{k \in \mathbb{N} \setminus \{0\}}$ ,  $x' = (x'_k)_{k \in \mathbb{N} \setminus \{0\}} \in \mathcal{X}^{(\mathbb{N} \setminus \{0\})}$  with  $x_k = x'_k = 0$  for all  $k > N$  for some  $N > 0$ , the following condition holds

$$\|F(x; y) - F(x'; y')\| \leq \sum_{j=1}^{\infty} a_j(F) \|x_j - x'_j\| + u \|y - y'\|, \tag{4.3}$$

with

$$\sum_{j=1}^{\infty} a_j(F) := a(F) < 1. \tag{4.4}$$

Many non linear econometrics time series are chains with infinite memory. The following proposition gives sufficient assumptions such that chains with infinite memory satisfy (CBS).

**Proposition 4.1.** *Under (4.3) and (4.4) there exists a unique solution  $(X_t)$  of equation (4.2) satisfying (CBS) if  $\xi_0$  satisfies the Cramer condition.*

The proof of Proposition 4.1 is given in Section 6.5.

## 4.2. Weakly dependent processes

### 4.2.1. Bounded causal Bernoulli shifts

Bounded causal Bernoulli shifts are examples of time series satisfying (WDP).

**Proposition 4.2.** *Under condition (4.3) and (3.2), any solution of the equation (4.1) is bounded by  $2a(H) \|\xi_0\|_{\infty}$  and is weakly dependent (WDP) with  $C = 2 \|\xi_0\|_{\infty} \tilde{a}(H)$ .*

The proof of this already known result is given in Section 6.5 for completeness. Below are presented two examples of time series satisfying (WDP) that are not bounded causal Bernoulli shifts.

4.2.2. Uniform  $\varphi$ -mixing processes

Let us recall the definition of the  $\varphi$ -mixing coefficients introduced in Ibragimov [17];

$$\varphi(r) = \sup_{(A,B) \in \mathfrak{G}_0 \times \mathfrak{F}_r} |\pi(B/A) - \pi(B)|,$$

where  $\mathfrak{F}_r = \sigma(Y_t, t \geq r)$ . The class of  $\varphi$ -mixing processes gives examples of time series that satisfied (WDP).

**Proposition 4.3.** *If  $(X_t)$  is a stationary bounded process, then it satisfies (WDP) with*

$$\theta_{\infty,n}(1) \leq 2 \|X_0\|_{\infty} \sum_{r=1}^n \varphi(r).$$

The proof of this already known result is given in Section 6.5 for completeness. Remark that  $(X_t)$  satisfies the short memory condition as soon as  $(\varphi(r))$  is summable. All uniform ergodic Markov chains are examples of  $\varphi$ -mixing processes with short memory, see Doukhan [14].

4.2.3. Dynamical systems on  $[0, 1]$

The AR(1) process  $X_t = 2^{-1}(X_{t-1} + \xi_t)$  with  $\xi_t$  Bernoulli distributed is not mixing, see [3] for more details. Through a reversion of the time, it can be viewed as a dynamical system  $X_t = T(X_{t+1})$  where  $T(x) = 2x$  if  $0 \leq x < 1/2$ ,  $T(x) = 2x - 1$  if  $1/2 \leq x \leq 1$ . Dedecker and Prieur [13] extended this counter-example to processes  $(X_t)$  such that  $X_t = T(X_{t+1})$  where  $T$  is an expanding map on  $[0, 1]$ , see Section 4.4 of [13] for a proper definition. Then  $(X_t)$  satisfies (WDP) with  $\mathcal{C} = K\sigma/(1 - \sigma)$  where  $K > 0, 0 \leq \sigma < 1$ , see Section 7.2 of [13].

## 5. Examples of predictors

We give some examples of Lipschitz predictors where we can estimate the complexity of the  $\Theta_{p,\ell}$  and then apply our main results. In this section,  $C > 0$  is a constant independent of  $\varepsilon$  and  $n$  that may be different from one inequality to another.

### 5.1. Linear predictors

Let  $\mathcal{X} = \mathbb{R}$  and we consider predictors of the form:

$$f_{\theta}(X_{t-1}, \dots, X_{t-p}) = \theta_0 + \sum_{i=1}^p \theta_i X_{t-i},$$

where  $\theta \in \Theta_p \subset \mathbb{R}^{p+1}$  with

$$\Theta_p = \Theta_{p,1} = \left\{ \theta \in \mathbb{R}^{p+1}, \|\theta\|_1 = \sum_{i=0}^p |\theta_i| \leq B \right\}$$

for some  $B > 0$  ( $m_p = 1$  for all  $p$  and we omit the index  $\ell$ ). Using Proposition 2.1 it follows that

$$d_p \leq (p + 1) \log \left( eB \left( \frac{e}{p + 1} \vee \frac{1}{B - \|\bar{\theta}_p\|} \right) \right),$$

where  $\bar{\theta}_p = \arg \min_{\Theta_p} R(\theta)$ . As a consequence of Theorems 3.1 and 3.2, we obtain the following corollary.

**Corollary 5.1.** *If  $\|\bar{\theta}_p\|_1 \leq B - e/(p + 1)$  for all  $p \geq 0$ , then, under (WDP) or (CBS), for all  $n \geq 8e(1 + L)$  with probability at least  $1 - \varepsilon$ :*

$$R(\hat{\theta}) \leq \inf_{p+1 \leq n/2} \left\{ \min_{\theta \in \Theta_p} R(\theta) + C \sqrt{\frac{p}{n}} \log^{5/2}(n) \right\} + C \frac{\log(1/\varepsilon)}{\sqrt{n}}.$$

Let us detail two examples:  $\text{AR}(p_0)$  and  $\text{AR}(\infty)$  models with innovations  $\xi_t$  i.i.d. satisfying the Cramer condition and  $\text{med}(\xi_0) = 0$ .

First, consider  $(X_t)$  a causal  $\text{AR}(p_0)$  process ( $0 \leq p_0 < \infty$ )

$$X_t = a_0 + \sum_{j=1}^{p_0} a_j X_{t-j} + \xi_j \quad \text{for all } t \in \mathbb{Z}.$$

If  $B \geq \sum_{j=0}^{p_0} |a_j| + e/(p + 1)$  for all  $0 \leq p \leq p_0$ , the error of the best linear predictor is  $\mu[|\varepsilon_j|]$ . Corollary 5.1 implies, for any  $0 < \varepsilon < 1$  and any  $n \geq 2(p_0 + 1)$ , the relation:

$$R(\hat{\theta}) - \mu[|\varepsilon_0|] \leq C \left( \sqrt{\frac{p_0}{n}} \log^{5/2}(n) + \frac{\log(1/\varepsilon)}{\sqrt{n}} \right) \quad \text{with probability at least } 1 - \varepsilon.$$

For  $\varepsilon > 0$  fixed independently of  $n$ , the rate of convergence of the excess risk is estimated by  $\sqrt{p_0/n} \log^{5/2}(n)$ . Note that  $\hat{\theta}$  achieves this rate even if  $p_0$  is unknown. One says that our procedure is adaptive in  $p_0$  and, using the terminology of [23], memory-universal.

Second, consider  $(X_t)$  a causal  $\text{AR}(\infty)$  process

$$X_t = a_0 + \sum_{i=1}^{\infty} a_i X_{t-i} + \xi_t \quad \text{for all } t \in \mathbb{Z}. \tag{5.1}$$

If  $B \geq \sum_{j=0}^p |a_j| + e/(p + 1)$  for all  $p \geq 0$ , we have  $\bar{\theta}_p = (a_0, \dots, a_p)$ . Then we roughly bound  $R(\bar{\theta}_p) = \pi_0[|\sum_{i>p} a_i X_{-i} + \xi_0|] \leq \mu[|\xi_0|] + \pi_0[|X_{-i}|] \sum_{i>p} |a_i|$  and with probability at least  $1 - \varepsilon$ :

$$R(\hat{\theta}) - \mu[|\xi_0|] \leq \inf_{p+1 \leq n/2} \left[ \pi_0[|X_0|] \sum_{i>p} |a_i| + C \sqrt{\frac{p}{n}} \log^{5/2}(n) \right] + C \frac{\log(1/\varepsilon)}{\sqrt{n}}.$$

In this nonparametric setting, to obtain a rate of convergence for the excess risk we have to specify the decay rate of the  $|a_i|$ . For example, if

$$\exists \gamma > 0, \exists \beta > 0, \forall p: \quad \sum_{i>p} |a_i| \leq \frac{\gamma}{p^\beta}$$

then the convergence rate is  $(\log^5(n)/n)^{\beta/(2\beta+1)}$  (consider the optimal  $p = n^{1/(2\beta+1)} \times \log^{5/(2\beta+1)}(n)$ ).

*Simulations*

We implement our linear prediction procedure using the R software [24]. We compare the results to the one obtained using the standard ARIMA procedure of R with the AIC criterion for model selection. Our theoretical penalization terms, driven by “the worst-case type” bounds, are necessarily pessimistic: our procedure systematically over-penalizes large models. Thus, for having an efficient procedure in practice, adjustments have been done. However, we aim with these simulations to show that

- (1) our linear prediction procedure is easily implementable;
- (2) its performances are reasonable when the implemented penalization term is smaller than the theoretical one.

We only consider observations from simulations of  $AR(p_0)$  models of the form

$$X_t = \sum_{i=1}^{p_0} a_i X_{t-i} + \xi_t,$$

where the  $\xi_t$  are i.i.d., either  $\mathcal{N}(0, \sigma^2)$ -distributed, either  $(\delta_0 + \mathcal{E}(\lambda))/2$  distributed, where  $\delta_0$  is the Dirac mass on 0 and  $\mathcal{E}(\lambda)$  the exponential distribution with parameter  $\lambda > 0$ . In both cases the Cramer condition is satisfied and  $\text{med}(\xi_0) = 0$ . Unlike the first case, mean and median are different in the second case. Thus, the minimizers of the  $\ell_1$  and the quadratic risks are the same in the first case and differ in the second one.

We use  $p_0 = 3, a_1 = 0.2, a_2 = 0.3, a_3 = 0.2, \sigma^2 \in \{1, 3\}, \lambda \in \{1, 1/\sqrt{12}\}$ , and  $n = 500$ ,

$$\Theta = \bigcup_{p=1}^8 \Theta_p = \bigcup_{p=1}^8 \{\theta \in \mathbb{R}^p: \|\theta\|_1 \leq 1\}$$

and

$$\lambda \in \mathcal{G} = \{2, 4, 8, \dots, 1024\}.$$

In view of our procedure, we compute the simplified penalized criterion

$$(\hat{\lambda}, \hat{p}) = \arg \min_{\substack{1 \leq p \leq 8 \\ \lambda \in \mathcal{G}}} -\frac{1}{\lambda} \log \int_{\Theta_{p,\ell}} \exp(-\lambda r_n(\theta)) \, d\pi_{p,\ell}(\theta) + \lambda \frac{K^2}{n}.$$

The theoretical value  $K = 2(\log n)^{3/2} \approx 9$  systematically over-penalizes the large models and always selects the simplest one ( $p = 1$ ). Thus, we fix in practice  $K = 0, 1$ . To compute the criteria, the integrand term is approximated using an acceptance-reject algorithm with gaussian proposal and 10,000 iterations. To compare one simulation of  $\hat{\theta} \sim \pi_{\hat{p}}\{-\hat{\lambda}_{r_n}\}$  with  $\hat{\theta}_{AIC}$  obtained by the classical R procedure, we simulate independently  $(X'_1, \dots, X'_{500})$  distributed as  $(X_1, \dots, X_{500})$  and we compare

$$\text{err}_1(\hat{\theta}) = \frac{1}{n-8} \sum_{i=9}^{500} \left| X'_i - \sum_{p=1}^{\hat{p}} (\hat{\theta})_p X'_{i-p} \right|$$

with  $\text{err}_1(\hat{\theta}_{AIC})$ . As the classical R procedure is based on least square estimators, we also compare the quadratic prevision error

$$\text{err}_2(\hat{\theta}) = \frac{1}{n-8} \sum_{i=9}^{500} \left( X'_i - \sum_{p=1}^{\hat{p}} (\hat{\theta})_p X'_{i-p} \right)^2$$

with  $\text{err}_2(\hat{\theta}_{AIC})$ . The results of 20 experiments are reported in Table 1.

The results are coherent with the theory: in the Gaussian cases, the optimal values of  $\theta$  for the  $\ell_1$  and the quadratic risks of prediction are the same. Both procedures estimate efficiently the same  $\theta$  and their prediction risks are the same. In the other cases, the optimal values of  $\theta$  for the  $\ell_1$  and the quadratic risks are not the same. We observe  $\text{err}_1(\hat{\theta}) < \text{err}_1(\hat{\theta}_{AIC})$  and  $\text{err}_2(\hat{\theta}) > \text{err}_2(\hat{\theta}_{AIC})$ . The choice between the two procedures only depends on the prediction risk considered.

**Table 1.** For each experiment, we report the median, mean and standard deviation of the  $\text{err}_i(\cdot)$  quantities on the 20 experiments realized. The best results, for both  $\text{err}_1(\cdot)$  and  $\text{err}_2(\cdot)$ , are bolded for each serie

$\xi_t$		$\text{err}_1(\hat{\theta})$	$\text{err}_1(\hat{\theta}_{AIC})$	$\text{err}_2(\hat{\theta})$	$\text{err}_2(\hat{\theta}_{AIC})$
$\mathcal{N}(0, 1)$	median	<b>0.790</b>	0.792	<b>0.975</b>	<b>0.975</b>
	mean	<b>0.797</b>	0.798	<b>0.985</b>	0.988
	s.d.	0.023	0.024	0.054	0.054
$\mathcal{N}(0, 3)$	median	2.433	<b>2.432</b>	0.918	<b>0.916</b>
	mean	<b>2.409</b>	2.412	<b>0.911</b>	0.912
	s.d.	0.078	0.065	0.496	0.412
$\frac{\delta_0 + \mathcal{E}(1)}{2}$	median	<b>0.567</b>	0.592	0.819	<b>0.813</b>
	mean	<b>0.580</b>	0.589	0.836	<b>0.813</b>
	s.d.	0.047	0.043	0.153	0.150
$\frac{\delta_0 + \mathcal{E}(1/\sqrt{12})}{2}$	median	<b>1.973</b>	2.000	9.525	<b>9.494</b>
	mean	<b>1.955</b>	1.997	9.733	<b>9.390</b>
	s.d.	0.158	0.162	1.656	1.522

### 5.2. Neural networks predictors

Similarly than in [23], we present a procedure that approximates the best possible predictor using the best possible number of past values  $p$  for the one-step prediction. Given  $p$ , the best possible predictor for the  $\mathbb{L}^1$ -risk is  $\text{med}(X_0|X_{-1}, \dots, X_{-p})$ . We denote  $R_p^*$  the corresponding risk. For  $\mathcal{X} = \mathbb{R}$ , we use the abstract neural networks predictors defined in Barron [6] by the relation

$$f_\theta = c_0 + \sum_{i=1}^{\ell} c_i \phi(a_i \cdot x + b_i) \quad \text{for all } x \in \mathbb{R}^p$$

for  $a_i \in \mathbb{R}^p$  and  $c_i, b_i \in \mathbb{R}$  for all  $1 \leq i \leq \ell$ , the sigmoidal function  $\phi(x) = (1 + \exp(-x))^{-1}$  for all  $x \in \mathbb{R}$  and  $\theta = (c_0, a_{1,1}, \dots, a_{1,p}, b_1, c_1, \dots, a_{\ell,1}, \dots, a_{\ell,p}, b_\ell, c_\ell)$  in  $\mathcal{B}_{c_p, \ell}^q$  for some  $c_{p, \ell} > 0$ ,  $q = \ell(p + 2) + 1$  and  $\ell \leq n$ . For any  $p \geq 1$ , we denote

$$r_p(x) = \text{med}(X_0|(X_{-1}, \dots, X_{-p}) = x) \quad \text{for all } x \in \mathbb{R}^p$$

and we assume that there exists a complex-valued function  $\tilde{r}_p$  on  $\mathbb{R}^p$  satisfying

$$\forall x \in \mathbb{R}^p \quad r_p(x) - r_p(0) = \int_{\mathbb{R}^p} (e^{iwx} - 1)\tilde{r}_p(w) dw \quad \text{and} \quad \int_{\mathbb{R}^p} \|w\|_1 |\tilde{r}_p(w)| dw \leq C' p^c$$

for some  $C', c > 0$ . Then

**Corollary 5.2.** *Under (WDP) if for any  $(p, \ell) \in M$*

$$\frac{q}{e} + 2\sqrt{\ell} \|X\|_\infty (C' p^c + \ell \log \ell) \leq c_{p, \ell} \tag{5.2}$$

*then, for all  $n \geq \max_M c_{p, \ell}$ , with probability at least  $1 - \varepsilon$ ,*

$$R(\hat{\theta}) \leq \inf_{10(1+\log n)^2 p^{1+2c} \leq n} \left\{ R_p^* + C \frac{p^{1/4+c/2} \log^3 n}{n^{1/4}} \right\} + C \frac{\log(1/\varepsilon)}{\sqrt{n}}.$$

If  $(X_t)$  satisfies the Markov condition of order  $p_0$ , then  $c = 0$  and for  $n$  sufficiently large

$$R(\hat{\theta}) - R_{p_0}^* \leq C \left( \frac{\log^3 n}{n^{1/4}} + \frac{\log(1/\varepsilon)}{\sqrt{n}} \right).$$

Compared to the i.i.d. case, the loss is  $\log^3 n$  and we do not need to know the order  $p_0$  (our procedure is memory-universal). Our loss is smaller than the one of the other memory-universal procedure given in [23].

### 5.3. Nonparametric auto-regressive predictors

As in Baraud, Comte and Viennet [5], we assume that  $(X_t)$  is a solution of the equation:

$$X_t = f_1(X_{t-1}) + \dots + f_{p_0}(X_{t-p_0}) + \xi_t \quad \text{for all } t \in \mathbb{Z},$$

where  $\xi_t \sim \mathcal{N}(0, \sigma^2)$ , the  $f_i$  are functions  $[-1; 1] \mapsto \mathbb{R}$  in Hölder class  $H(s_i, L_i)$ :  $f_i$  is derivable  $\lfloor s_i \rfloor$  times and

$$\exists \mathcal{L}_i > 0, \forall (x, x') \in [-1, 1]^2, \quad |f_i^{(\lfloor s_i \rfloor)}(x) - f_i^{(\lfloor s_i \rfloor)}(x')| \leq \mathcal{L}_i |x - x'|^{s_i - \lfloor s_i \rfloor}. \quad (5.3)$$

Consider the Fourier basis  $(\phi_j(\cdot))_{j \geq 1}$  on  $[-1, 1]$  composed by  $\phi_{2k}(x) = \sqrt{2} \cos(2\pi kx)$  and  $\phi_{2k+1}(x) = \sqrt{2} \sin(2\pi kx)$ . Assumption 5.3 implies the existence of  $\gamma_i > 0$  such that for any  $m \geq 0$  it holds

$$\min_{(\alpha_1, \dots, \alpha_m) \in \mathbb{R}^m} \left\{ \int_{-1}^1 \left[ f_i(t) - \sum_{j=1}^m \alpha_{i,j} \phi_j(t) \right]^2 ds \right\}^{1/2} \leq \gamma_i m^{-s_i}.$$

Natural predictors are given by

$$\widehat{X}_{n+1} = \sum_{i=1}^p \sum_{j=1}^{\ell} \theta_{i,j} \varphi_j(X_{n-i}) =: f_{\theta}(X_n, \dots, X_{n-p})$$

for any  $p \in \{1, \dots, \lfloor n/2 \rfloor\}$  and any  $\ell \in \{1, \dots, m_p = n\}$ . We restrict the procedure on  $\theta_{p,\ell}$  in the compact set

$$\Theta_{p,\ell} = \left\{ \theta \in \mathbb{R}^{p\ell}, \sum_{i=1}^p \sum_{j=1}^{\ell} \theta_{i,j}^2 (2\lfloor j/2 \rfloor)^2 \leq L^2 \right\}$$

such that any  $f_{\theta}$  is an  $L$ -Lipschitz function. We define also the coefficients  $\bar{\theta}_{p,\ell} \in \mathbb{R}^{p\ell}$  by the relation

$$\arg \min_{\theta \in \Theta_{p,\ell}} \pi_0 \left[ \left\| X_n - \sum_{i=1}^p \sum_{j=1}^{\ell} \theta_{i,j} \varphi_j(X_{n-i}) \right\|^2 \right].$$

As a consequence of Theorem 3.1, it holds

**Corollary 5.3.** *Under (CBS), if for any  $\ell \geq 1$  and any  $p \geq 1$*

$$\frac{\ell p}{e} + \left( \sum_{i=1}^{p_0} \sum_{j=1}^{\ell} (\bar{\theta}_{p_0,\ell})_{i,j}^2 (2\lfloor j/2 \rfloor)^2 \right)^{1/2} \leq L,$$

*then for all  $n \geq 8e(1 + L)$  with probability at least  $1 - \varepsilon$*

$$R(\hat{\theta}) - \mu[\|\xi_0\|] \leq C \left( \left( \frac{\log(n)}{n} \right)^{s/(2s+1)} + \frac{\log(1/\varepsilon)}{\sqrt{n}} \right),$$

*where  $s$  denotes  $\min\{s_1, \dots, s_{p_0}\}$ .*

The (i.i.d.) minimax rate of convergence with respect to  $s_1, \dots, s_{p_0}$  for the  $\ell^1$ -risk is achieved up to a logarithmic loss. In [5], the (i.i.d.) minimax rate of convergence for the quadratic risk is achieved for the empirical quadratic risk in expectation.

## 6. Proofs

To present the proofs in a unified version whether we work under (CBS) or (WDP), we truncate the observations if we are under (CBS):

$$\bar{X}_t = H(\bar{\xi}_t, \bar{\xi}_{t-1}, \bar{\xi}_{t-2}, \dots) \quad \text{for all } t \in \mathbb{Z},$$

where  $\bar{\xi}_t = (\xi_t \wedge C) \vee (-C)$ , under (WDP) we just take  $\bar{X}_t = X_t$ . We denote in the sequel  $\bar{X} = (\bar{X}_t)$  and  $\bar{r}, \bar{R}$  the risks associated with  $\bar{X}$  under (CBS) and with  $X$  under (WDP). To shorten the proofs, we denote  $K_n = (1 + L) \log^{3/2} n$  and  $w_{p,\ell} = 1/(m_p \lfloor n/2 \rfloor)$  in the sequel. The proof of our main theorem lies on estimates on Laplace transforms.

### 6.1. Preliminary lemmas: Estimates on Laplace transforms

The proofs of these lemmas are given in Section 6.4. The first lemma is an estimate of the Laplace transforms of the risk of  $\bar{X}$ ; it is a direct corollary of the result in Rio [25].

**Lemma 6.1 (Laplace transform of the risk).** *For any  $\lambda > 0$  and  $\theta \in \Theta$  we have:*

$$\pi_0[\exp(\lambda(\bar{R}(\theta) - \bar{r}_n(\theta)))] \leq \exp\left(\frac{\lambda^2 k_n^2}{n(1 - p/n)^2}\right),$$

where  $k_n = \sqrt{2}C(1 + L)(a(H) + \tilde{a}(H))$  under (CBS) and  $k_n = (1 + L)(\|X_0\|_\infty + \theta_{\infty,n}(1))/\sqrt{2}$  under (WDP).

Given a measurable space  $(E, \mathcal{E})$  we let  $\mathcal{M}_+^1(E)$  denote the set of all probability measures on  $(E, \mathcal{E})$ . The Kullback divergence is a pseudo-distance on  $\mathcal{M}_+^1(E)$  defined, for any  $(\pi, \pi') \in [\mathcal{M}_+^1(E)]^2$  by the equation

$$\mathcal{K}(\pi, \pi') = \begin{cases} \pi[\log(d\pi/d\pi')], & \text{if } \pi \ll \pi', \\ +\infty, & \text{otherwise.} \end{cases}$$

The proof of the following lemma is omitted as it can be found in [7] or [8].

**Lemma 6.2 (Legendre transform of the Kullback divergence function).** *For any  $\pi \in \mathcal{M}_+^1(E)$ , for any measurable function  $h : E \rightarrow \mathbb{R}$  such that  $\pi[\exp(h)] < +\infty$  we have:*

$$\pi[\exp(h)] = \exp\left(\sup_{\rho \in \mathcal{M}_+^1(E)} (\rho[h] - \mathcal{K}(\rho, \pi))\right), \tag{6.1}$$

with convention  $\infty - \infty = -\infty$ . Moreover, as soon as  $h$  is upper-bounded on the support of  $\pi$ , the supremum with respect to  $\rho$  in the right-hand side is reached for the Gibbs measure  $\pi\{h\}$  defined in (2.5).

Using Lemmas 6.1 and 6.2, we get an upper-bound for the Laplace transform of the mean risk of Gibbs estimators in all sub-models.

**Lemma 6.3.** *Under the assumptions of Theorem 3.1, we have for any  $\lambda > 0$  and  $(p, \ell) \in M$ :*

$$\pi_0 \left[ \exp \left( \sup_{\rho \in \mathcal{M}_+^1(\Theta_{p,\ell})} \{ \lambda \rho [\bar{R} - \bar{r}_n] - \mathcal{K}(\rho, \pi_{p,\ell}) \} - \frac{\lambda^2 k_n^2}{n(1-p/n)^2} \right) \right] \leq 1, \tag{6.2}$$

where  $k_n$  has the same expression than in Lemma 6.1.

Following the technique used by Catoni [7], we derive from Lemma 6.3 another upper-bound on the Laplace transform of the mean risk of any aggregation estimators of all Gibbs estimators.

**Lemma 6.4.** *For any measurable functions  $\hat{\rho}_{p,\ell} : \mathcal{X}^n \rightarrow \mathcal{M}_+^1(\Theta_{p,\ell})$  for  $(p, \ell) \in M$ , under the assumptions of Theorem 3.1, we have:*

$$\pi_0 \left[ \sum_{(p,\ell) \in M} \sum_{\lambda \in \mathcal{G}_{p,\ell}} \hat{\rho}_{p,\ell} \left[ \exp \left( \lambda (\bar{R} - \bar{r}_n) - \log \frac{d\hat{\rho}_{p,\ell}}{d\pi_{p,\ell}} - \frac{\lambda^2 k_n^2}{n(1-p/n)^2} + \log(w_{p,\ell}/n) \right) \right] \right] \leq 1$$

and

$$\pi_0 \left[ \sum_{(p,\ell) \in M} \sum_{\lambda \in \mathcal{G}_{p,\ell}} \exp \left( \lambda \hat{\rho}_{p,\ell} [\bar{r}_n - \bar{R}] - \mathcal{K}(\hat{\rho}_{p,\ell}, \pi_{p,\ell}) - \frac{\lambda^2 k_n^2}{n(1-p/n)^2} + \log(w_{p,\ell}/n) \right) \right] \leq 1,$$

where we remind that  $k_n$  is defined in Lemma 6.1.

Finally, we use a lemma that quantify the error in the risk due to the truncation under (CBS).

**Lemma 6.5.** *Under (CBS), for any truncation level  $C > 0$  and any  $0 \leq \lambda \leq n/(4(1+L))$ , we have*

$$\pi_0 \left[ \exp \left( \lambda \sup_{\theta \in \Theta} |r_n(\theta) - \bar{r}_n(\theta)| - \lambda 2(1+L) \Psi(a(H)) \left( \frac{a(H)^2 C}{\exp(a(H)C) - 1} + \lambda \frac{4(1+L)}{n} \right) \right) \right] \leq 1.$$

### 6.2. Proof of Theorem 3.1

Remark that (WDP) is satisfied, so  $\bar{R} = R$  and  $\bar{r} = r$ . We apply the first inequality of Lemma 6.4 to  $\hat{\rho}_{p,\ell}^\lambda = \pi_{p,\ell}\{-\lambda r_n\}$ . Remembering that  $(\hat{\rho}, \hat{\ell}, \hat{\lambda}) = \arg \min \hat{R}(p, \ell, \lambda)$ , we obtain in particu-

lar:

$$\pi_{0\hat{\rho}_{\hat{p},\hat{\ell}}^\lambda} \left[ \exp \left( \hat{\lambda} (R - r_n) - \log \left( \frac{d\hat{\rho}_{\hat{p},\hat{\ell}}^\lambda}{d\pi_{\hat{p},\hat{\ell}}} \right) - \frac{\hat{\lambda}^2 k_n^2}{n(1 - \hat{p}/n)^2} + \log \left( \frac{w_{\hat{p},\hat{\ell}}}{n} \right) \right) \right] \leq 1. \tag{6.3}$$

Remark that  $\pi_{0\hat{\rho}_{\hat{p},\hat{\ell}}^\lambda}$  is a well defined probability measure as  $\hat{\rho}$  are defined conditionally on the observations. Remark also that  $\hat{\theta} \sim \hat{\rho}_{\hat{p},\hat{\ell}}^\lambda$  by definition, then using the classical Chernov bound we derive that with probability  $1 - \varepsilon$  it holds:

$$R(\hat{\theta}) \leq r_n(\hat{\theta}) + \frac{\hat{\lambda} k_n^2}{n(1 - \hat{p}/n)^2} + \frac{1}{\hat{\lambda}} \log \left( \frac{d\hat{\rho}_{\hat{p},\hat{\ell}}^\lambda}{d\pi_{\hat{p},\hat{\ell}}} \right) + \frac{1}{\hat{\lambda}} \log \left( \frac{n}{w_{\hat{p},\hat{\ell}}} \right) + \frac{1}{\hat{\lambda}} \log \frac{1}{\varepsilon}. \tag{6.4}$$

In order that the term  $\hat{R}$  appears, we notice that (6.4) is equivalent to

$$\begin{aligned} R(\hat{\theta}) &\leq -\frac{1}{\hat{\lambda}} \log \int_{\Theta_{\hat{p},\hat{\ell}}} \exp(-\hat{\lambda} r_n(\theta)) \pi_{\hat{p},\hat{\ell}}(d\theta) + \frac{\hat{\lambda} k_n^2}{n(1 - \hat{p}/n)^2} + \frac{1}{\hat{\lambda}} \log \left( \frac{n}{w_{\hat{p},\hat{\ell}}} \right) + \frac{1}{\hat{\lambda}} \log \frac{1}{\varepsilon} \\ &\leq \inf_{p,\ell,\lambda} \hat{R}(p, \ell, \lambda) + \frac{\hat{\lambda}(k_n^2 - K_n^2)}{n(1 - \hat{p}/n)^2} - \frac{1}{\hat{\lambda}} \log \varepsilon \end{aligned}$$

(remind that  $K_n = (1 + L) \log^{3/2} n$ ). Now, we upper bound the term  $\hat{R}(p, \ell, \lambda)$ , for any  $p, \ell$  and  $\lambda$ . Using the second inequality of Lemma 6.4, we obtain for any  $(p, \ell) \in M, \lambda \in \mathcal{G}$  and  $\rho \in \mathcal{M}_+^1(\Theta_{p,\ell})$ ,

$$\begin{aligned} \int_{\Theta_{p,\ell}} r_n(\theta) \rho(d\theta) &\leq \int_{\Theta_{p,\ell}} R(\theta) \rho(d\theta) + \frac{\lambda k_n^2}{n(1 - p/n)^2} + \frac{1}{\lambda} \mathcal{K}(\rho, \pi_{p,\ell}) \\ &\quad + \frac{1}{\lambda} \log \frac{n}{w_{p,\ell}} + \frac{1}{\lambda} \log \frac{1}{\varepsilon}. \end{aligned} \tag{6.5}$$

From (6.5) and using Lemma 6.2 two times, we derive that

$$\begin{aligned} &-\frac{1}{\lambda} \log \int_{\Theta_{p,\ell}} \exp(-\lambda r_n(\theta)) \pi_{p,\ell}(d\theta) \\ &= \inf_{\rho \in \mathcal{M}_+^1(\Theta_{p,\ell})} \left\{ \int_{\Theta_{p,\ell}} r_n(\theta) \rho(d\theta) + \frac{1}{\lambda} \mathcal{K}(\rho, \pi_{p,\ell}) \right\} \\ &\leq \inf_{\rho \in \mathcal{M}_+^1(\Theta_{p,\ell})} \left\{ \int_{\Theta_{p,\ell}} R(\theta) \rho(d\theta) + \frac{2}{\lambda} \mathcal{K}(\rho, \pi_{p,\ell}) \right\} + \frac{\lambda k_n^2}{n(1 - p/n)^2} + \frac{1}{\lambda} \log \frac{n}{\varepsilon w_{p,\ell}} \\ &= -\frac{2}{\lambda} \log \int_{\Theta_{p,\ell}} \exp \left( -\frac{\lambda}{2} R(\theta) \right) \pi_{p,\ell}(d\theta) + \frac{\lambda k_n^2}{n(1 - p/n)^2} + \frac{1}{\lambda} \log \frac{n}{\varepsilon w_{p,\ell}}. \end{aligned}$$

Finally, we obtain:

$$\hat{R}(p, \ell, \lambda) \leq -\frac{2}{\lambda} \log \int_{\Theta_{p,\ell}} \exp\left(-\frac{\lambda}{2} R(\theta)\right) \pi_{p,\ell}(d\theta) + \frac{\lambda(k_n^2 + K_n^2)}{n(1-p/n)^2} + \frac{1}{\lambda} \log \frac{n}{\varepsilon w_{p,\ell}}. \quad (6.6)$$

Under Assumption (2.3), as soon as  $\lambda > 2e$  it holds

$$-\log \pi_{p,\ell} \left[ \exp\left(-\frac{\lambda}{2} (R - R(\bar{\theta}_{p,\ell}))\right) \right] \leq d_{p,\ell} \log \frac{\lambda}{2}$$

and it easily follows that

$$-\log \pi_{p,\ell} \left[ \exp\left(-\frac{\lambda}{2} \bar{R}\right) \right] \leq d_{p,\ell} \log \frac{\lambda}{2} + \frac{\lambda}{2} R(\bar{\theta}_{p,\ell}).$$

We plug this result into the inequality (6.6) to obtain:

$$\hat{R}(p, \ell, \lambda) \leq R(\bar{\theta}_{p,\ell}) + \frac{1}{\lambda} \left( 2d_{p,\ell} \log \frac{\lambda}{2} + \log \frac{n}{\varepsilon w_{p,\ell}} \right) + \frac{\lambda(k_n^2 + K_n^2)}{n(1-p/n)^2}. \quad (6.7)$$

Collecting the inequalities (6.4) and (6.7), we obtain:

$$\begin{aligned} R(\hat{\theta}) \leq & \inf_{p,\ell,\lambda \in \mathcal{G}_{p,\ell}} \left\{ R(\bar{\theta}_{p,\ell}) + \frac{1}{\lambda} \left( 2d_{p,\ell} \log \frac{\lambda}{2} + \log \frac{n}{\varepsilon w_{p,\ell}} \right) + \frac{\lambda(k_n^2 + K_n^2)}{n(1-p/n)^2} \right\} \\ & + \frac{\hat{\lambda}(k_n^2 - K_n^2)}{n(1-\hat{p}/n)^2} - \frac{1}{\hat{\lambda}} \log \varepsilon. \end{aligned} \quad (6.8)$$

As for  $\lambda \in \mathcal{G}_{p,\ell}$ , we have, by definition of  $\mathcal{G}_{p,\ell}$  that

$$\lambda \in \left[ \check{c} \frac{\sqrt{d_{p,\ell} n} \log(d_{p,\ell} n)}{K_n}, \dots, \hat{c} \frac{\sqrt{d_{p,\ell} n} \log(d_{p,\ell} n)}{K_n} \right] \cap [2e, n]$$

then it holds

$$\begin{aligned} R(\hat{\theta}) \leq & \inf_{d_{p,\ell} \leq n} \left\{ R(\bar{\theta}_{p,\ell}) + \frac{K_n}{\check{c} \sqrt{d_{p,\ell} n} \log(d_{p,\ell} n)} \left( 2d_{p,\ell} \log \frac{n}{2} + \log \frac{n}{\varepsilon w_{p,\ell}} \right) \right. \\ & \left. + 4\hat{c}(k_n^2 + K_n^2) \sqrt{\frac{d_{p,\ell} \log(nd_{p,\ell})}{n} \frac{1}{K_n}} \right\} + 4(k_n^2 - K_n^2)_+ + \frac{(1+L) \log(1/\varepsilon)}{\check{c} \sqrt{n}}. \end{aligned} \quad (6.9)$$

For the sake of simplicity, we use rough estimates ( $1 \leq d_{p,\ell}$ ,  $1 \leq 1/\varepsilon$ ,  $m_p \leq n$ , ...) to obtain

$$\begin{aligned} R(\hat{\theta}) \leq & \inf_{d_{p,\ell} \leq n} \left\{ R(\bar{\theta}_{p,\ell}) + (1+L) \left( \frac{6}{\check{c}} + 8\hat{c}(1 + \|X_0\|_\infty + \theta_{\infty,n}(1))^2 \right) \sqrt{\frac{d_{p,\ell}}{n} \log^{5/2}(n)} \right\} \\ & + 4(k_n^2 - K_n^2)_+ + \frac{7(1+L) \log(n/\varepsilon)}{\check{c} \sqrt{n}}. \end{aligned}$$

This ends the proof as

$$k_n^2 - K_n^2 = (1 + L) \left( \frac{(\|X_0\|_\infty + \theta_{\infty,n}(1))^2}{2} - \log^3(n) \right).$$

### 6.3. Proof of Theorem 3.2

As we work under (CBS), we have to deal with the error of approximation of  $r$  and  $R$  by  $\bar{R}$ . To quantify it, we use Lemma 6.5. First, remark that as  $R = \pi_0[r]$  it holds

$$\exp\left(\lambda \sup_{\theta \in \Theta} |R(\theta) - \bar{R}(\theta)| - \lambda\phi(C, \lambda)\right) \leq 1,$$

where

$$\phi(C, \lambda) = 2(1 + L)\Psi(a(H)) \left( \frac{a(H)^2 C}{\exp(a(H)C) - 1} + \lambda \frac{4(1 + L)}{n} \right).$$

An immediate consequence is that

$$\pi_0 \left[ \exp\left(\lambda \sup_{\theta \in \Theta} |(r_n - R)(\theta) - (\bar{r}_n - \bar{R})(\theta)| - 2\lambda\phi(C, \lambda)\right) \right] \leq 1.$$

As  $R - r_n = \bar{r}_n - \bar{R} + (r_n - R) - (\bar{r}_n - \bar{R})$ , for any measurable function  $\rho_{p,\ell} : \mathcal{X}^n \rightarrow \mathcal{M}_+^1(\Theta_{p,\ell})$  the Cauchy–Schwarz inequality gives

$$\begin{aligned} & \pi_0 \rho \left[ \exp(\lambda/2(R - r_n)) \right] \\ & \leq \sqrt{\pi_0 \rho \left[ \exp(\lambda(\bar{R} - \bar{r}_n)) \right] \pi_0 \rho \left[ \exp\left(\lambda \sup_{\theta \in \Theta} |(r_n - R)(\theta) - (\bar{r}_n - \bar{R})(\theta)|\right) \right]}. \end{aligned}$$

Using this remark and the same reasoning than in the proof of Theorem 3.1 that gives (6.3) from Lemma 6.4, we get the inequality

$$\begin{aligned} & \pi_0 \hat{\rho}_{\hat{p}, \hat{\ell}}^{\hat{\lambda}} \left[ \exp\left(\frac{\hat{\lambda}}{2}(R - r_n) - 0, 5 \log\left(\frac{d\hat{\rho}_{\hat{p}, \hat{\ell}}^{\hat{\lambda}}}{d\pi_{\hat{p}, \hat{\ell}}}\right) - 0, 5 \frac{\hat{\lambda}^2 k_n^2}{n(1 - \hat{p}/n)^2} \right. \right. \\ & \left. \left. + 0, 5 \log\left(\frac{w_{\hat{p}, \hat{\ell}}}{n}\right) - \lambda\phi(C, \lambda)\right) \right] \leq 1. \end{aligned}$$

As in the proof of Theorem 3.1, we derive an equivalent of (6.4), that is, with probability  $1 - \varepsilon$  it holds:

$$R(\hat{\theta}) \leq r_n(\hat{\theta}) + \frac{\hat{\lambda} k_n^2}{n(1 - \hat{p}/n)^2} + \frac{1}{\hat{\lambda}} \log\left(\frac{d\hat{\rho}_{\hat{p}, \hat{\ell}}^{\hat{\lambda}}}{d\pi_{\hat{p}, \hat{\ell}}}\right) + \frac{1}{\hat{\lambda}} \log\left(\frac{n}{w_{\hat{p}, \hat{\ell}}}\right) + 2\phi(C, \hat{\lambda}) + \frac{2}{\hat{\lambda}} \log \frac{1}{\varepsilon}.$$

With similar arguments, we derive an equivalent of (6.5):

$$\int_{\Theta_{p,\ell}} r_n(\theta)\rho(d\theta) \leq \int_{\Theta_{p,\ell}} R(\theta)\rho(d\theta) + \frac{\lambda k_n^2}{n(1-p/n)^2} + \frac{1}{\lambda}\mathcal{K}(\rho, \pi_{p,\ell}) + \frac{1}{\lambda} \log \frac{n}{w_{p,\ell}} + 2\phi(C, \lambda) + \frac{2}{\lambda} \log \frac{1}{\varepsilon}$$

and also

$$R(\hat{\theta}) \leq \inf_{p,\ell,\lambda} \left\{ R(\bar{\theta}_{p,\ell}) + \frac{1}{\lambda} \left( 2d_{p,\ell} \log \frac{\lambda}{2} + \log \frac{n}{\varepsilon w_{p,\ell}} \right) + \frac{\lambda(k_n^2 + K_n^2)}{n(1-p/n)^2} + 2\phi(C, \lambda) \right\} + \frac{\hat{\lambda}(k_n^2 - K_n^2)}{n(1-p/n)^2} + 2\phi(C, \hat{\lambda}) - \frac{2}{\hat{\lambda}} \log \varepsilon. \tag{6.10}$$

We still have

$$-\frac{2}{\hat{\lambda}} \log \varepsilon \leq \frac{2(1+L)}{\check{c}\sqrt{n}} \log \frac{1}{\varepsilon}$$

so we now have to upper bound  $2\phi(C, \hat{\lambda})$ . As  $\hat{\lambda} \leq n/(4(1+L))$  by definition of the  $\mathcal{G}_{p,\ell}$ , fixing  $C = a(H)^{-1} \log n$  we obtain:

$$\phi(C, \hat{\lambda}) \leq \frac{4a(H)(1+L)\Psi(a(H))[2\hat{\lambda}(1+L) + a(H) \log(n)]}{n}.$$

As  $\hat{\lambda} \leq \hat{c}d_{\hat{p},\hat{\ell}} \log(d_{\hat{p},\hat{\ell}}n)/(1+L)$  by definition of  $\mathcal{G}_{p,\ell}$ , we obtain

$$\begin{aligned} \frac{\hat{\lambda}(k_n^2 - K_n^2)}{n(1-p/n)^2} + 2\phi(C, \hat{\lambda}) &\leq \frac{8a(H)^2(1+L)\Psi(a(H)) \log(n)}{n} \\ &\quad + \sqrt{\frac{d_{\hat{p},\hat{\ell}}}{n}} \log(d_{\hat{p},\hat{\ell}}n)4(1+L) \\ &\quad \times \hat{c}(4a(H)\Psi(a(H)) + 2 \log^2(n)(1 + \tilde{a}(H)/a(H))^2 - \log^3(n))_+. \end{aligned}$$

We now plug this result into (6.10) to end the proof.

### 6.4. Proofs of Lemmas 6.1, 6.3, 6.4, 6.5 and of Proposition 2.1

**Proof of Lemma 6.1.** The proof of this lemma is based on the following result of Rio [25] on  $\bar{\mathcal{X}}$ .

**Theorem 6.6.** Let  $Y = (Y_t)_{t \in \mathbb{Z}}$  be a bounded stationary time series bounded distributed as  $\pi_0$  on  $\mathcal{X}^{\mathbb{Z}}$ . Let  $h$  be a 1-Lipschitz function of  $\mathcal{X}^n \rightarrow \mathbb{R}$ , that is, such that:

$$\forall (x_1, y_1, \dots, x_n, y_n) \in \mathcal{X}^{2n}, \quad |h(x_1, \dots, x_n) - h(y_1, \dots, y_n)| \leq \sum_{i=1}^n \|x_i - y_i\|. \tag{6.11}$$

Then for any  $t \geq 0$  we have:

$$\pi_0[\exp(t(\pi_0[h(X_1, \dots, X_n)] - h(X_1, \dots, X_n)))] \leq \exp(t^2 n (\|X_0\|_\infty + \theta_{\infty, n}(1))^2 / 2).$$

**Proof.** This version of Theorem 1 of [25] comes rewriting the inequality (3) in [25] as, for any 1-Lipschitz function  $g$ :

$$\Gamma(g) = \|\mathbb{E}(g(X_{\ell+1}, \dots, X_n) | \mathcal{F}_\ell) - \mathbb{E}(g(X_{\ell+1}, \dots, X_n))\|_\infty \leq \theta_{\infty, n-\ell}(1).$$

The result is proved as  $\sup_{1 \leq r \leq n} \theta_{\infty, r}(1) \leq \theta_{\infty, n}(1)$ . □

We now apply the result of Theorem 6.6 on  $Y = \bar{X}$  to obtain the result of Lemma 6.1. Let us fix  $\lambda > 0$ ,  $(p, \ell) \in M$ ,  $\theta \in \Theta_{p, \ell}$  and  $t = (1 + L)\lambda / [n - p(\theta)]$  and the function  $h$  defined by:

$$h(x_1, \dots, x_n) = \frac{1}{1 + L} \sum_{i=p(\theta)+1}^n \|x_i - f_\theta(x_{i-1}, \dots, x_{i-p(\theta)})\|.$$

We easily check that  $h$  satisfies condition (6.11):

$$\begin{aligned} & |h(x_1, \dots, x_n) - h(y_1, \dots, y_n)| \\ & \leq \frac{1}{1 + L} \sum_{i=p(\theta)+1}^n \left| \|x_i - f_\theta(x_{i-1}, \dots, x_{i-p(\theta)})\| - \|y_i - f_\theta(y_{i-1}, \dots, y_{i-p(\theta)})\| \right| \\ & \leq \frac{1}{1 + L} \sum_{i=p(\theta)+1}^n \|x_i - y_i - f_\theta(x_{i-1}, \dots, x_{i-p(\theta)}) + f_\theta(y_{i-1}, \dots, y_{i-p(\theta)})\| \\ & \leq \frac{1}{1 + L} \sum_{i=p(\theta)+1}^n \|x_i - y_i\| \\ & \quad + \frac{1}{1 + L} \sum_{i=p(\theta)+1}^n \|f_\theta(x_{i-1}, \dots, x_{i-p(\theta)}) - f_\theta(y_{i-1}, \dots, y_{i-p(\theta)})\| \\ & \leq \frac{1}{1 + L} \sum_{i=p(\theta)+1}^n \|x_i - y_i\| + \frac{1}{1 + L} \sum_{i=p(\theta)+1}^n \sum_{j=1}^{p(\theta)} a_j(\theta) \|x_{i-j} - y_{i-j}\| \\ & \leq \frac{1}{1 + L} \sum_{i=p(\theta)+1}^n \|x_i - y_i\| + \frac{L}{1 + L} \sum_{i=1}^n \|x_i - y_i\| \\ & \leq \sum_{i=1}^n \|x_i - y_i\|. \end{aligned}$$

The direct application of Theorem 6.6 ends the proof under (WDP). Under (CBS),  $k_n$  follows from the estimates of  $\|X_0\|_\infty$  and  $\theta_{\infty,n}(1)$  obtained in Proposition 4.1.  $\square$

**Proof of Lemma 6.3.** Integrate the inequality in Lemma 6.1 with respect  $\pi_{p,\ell}$  on  $\Theta_{p,\ell}$  (then  $p(\theta) = p$ ) for any  $(p, \ell) \in M$  in order to obtain:

$$\pi_{p,\ell}[\pi_0[\exp(\lambda(\bar{R} - \bar{r}_n))]] \leq \exp\left(\frac{\lambda^2 k_n^2}{n(1 - p/n)^2}\right).$$

Fubini's theorem implies that

$$\pi_0\left[\pi_{p,\ell}\left[\exp\left(\lambda(\bar{R} - \bar{r}_n) - \frac{\lambda^2 k_n^2}{n(1 - p/n)^2}\right)\right]\right] \leq 1.$$

Applying Lemma 6.2 for  $\pi = \pi_{p,\ell}$  and  $h = \lambda(\bar{R} - \bar{r}_n) - \lambda^2 k_n^2 / (n(1 - p/n)^2)$  on  $\mathcal{M}_+^1(\Theta_{p,\ell})$  leads to the inequality:

$$\pi_0\left[\exp\left(\sup_{\rho \in \mathcal{M}_+^1(\Theta_{p,\ell})} \{\lambda\rho[\bar{R} - \bar{r}_n] - \mathcal{K}(\rho, \pi_{p,\ell})\} - \frac{\lambda^2 k_n^2}{n(1 - p/n)^2}\right)\right] \leq 1.$$

This ends the proof.  $\square$

**Proof of Lemma 6.4.** First, let us choose  $\lambda \in \Lambda$ . Let  $h_{p,\ell}^\lambda$  denotes, for any  $(p, \ell) \in M$ :

$$h_{p,\ell}^\lambda = \sup_{\rho_{p,\ell} \in \mathcal{M}_+^1(\Theta_{p,\ell})} \{\lambda\rho_{p,\ell}[\bar{R} - \bar{r}_n] - \mathcal{K}(\rho_{p,\ell}, \pi_{p,\ell})\} - \frac{\lambda^2 k_n^2}{n(1 - p/n)^2}.$$

From Lemma 6.3 applied on the different  $\mathcal{M}_+^1(\Theta_{p,\ell})$  we have, for any  $(p, \ell) \in M$ :

$$\pi_0\left[\sum_{(p,\ell) \in M} w_{p,\ell} \exp(h_{p,\ell}^\lambda)\right] \leq 1.$$

Now we apply Inequality (6.1) in Lemma 6.2 for  $\pi = \sum_{(p,\ell) \in M} w_{p,\ell} \delta_{(p,\ell)}$  and  $h = \sum_{(p,\ell) \in M} h_{p,\ell}^\lambda \mathbb{1}_{\Theta_{p,\ell}}$  and we obtain

$$\pi_0\left[\exp\left(\sup_{\sum_{(p,\ell) \in M} w'_{p,\ell} = 1} \left\{ \sum_{(p,\ell) \in M} w'_{p,\ell} h_{p,\ell} - \sum_{(p,\ell) \in M} w'_{p,\ell} \log(w'_{p,\ell}/w_{p,\ell}) \right\}\right)\right] \leq 1$$

and, by Jensen's inequality, and replacing  $h_{p,\ell}^\lambda$  by its definition,

$$\begin{aligned} \pi_0\left[\sup_{\sum_{(p,\ell) \in M} w'_{p,\ell} = 1} \left\{ \sum_{(p,\ell) \in M} w'_{p,\ell} \sup_{\rho_{p,\ell} \in \mathcal{M}_+^1(\Theta_{p,\ell})} \exp\left(\lambda\rho_{p,\ell}\left[\lambda(\bar{R} - \bar{r}_n) - \log \frac{d\rho_{p,\ell}}{d\pi_{p,\ell}}\right] \right. \right. \right. \\ \left. \left. \left. - \frac{\lambda^2 k_n^2}{n(1 - p/n)^2} + \log \frac{w_{p,\ell}}{w'_{p,\ell}}\right)\right\}\right] \leq 1. \end{aligned} \tag{6.12}$$

By Jensen again, we obtain a bound for the first term in the sum bounded in Lemma 6.4:

$$\pi_0 \left[ \sup_{\sum_{(p,\ell) \in M} w'_{p,\ell} = 1} \left\{ \sum_{(p,\ell) \in M} w'_{p,\ell} \sup_{\rho_{p,\ell} \in \mathcal{M}_+^1(\Theta_{p,\ell})} \rho_{p,\ell} \left[ \exp \left( \lambda(\bar{R} - \bar{r}_n) - \log \frac{d\rho_{p,\ell}}{d\pi_{p,\ell}} - \frac{\lambda^2 k_n^2}{n(1-p/n)^2} + \log \frac{w_{p,\ell}}{w'_{p,\ell}} \right) \right] \right\} \right] \leq 1.$$

Finally, we sum this inequality over all  $\lambda \in \mathcal{G}$  to bound the first expectation.

The second expectation is bounded by choosing specific weights  $w'_{p,\ell}$  in the supremum in inequality (6.12) such that  $w'_{p,\ell} = 1$  for  $(p, \ell) = \arg \max_M \{h_{p,\ell}\}$ :

$$\pi_0 \left[ \sup_{\substack{(p,\ell) \in M \\ \rho_{p,\ell} \in \mathcal{M}_+^1(\Theta_{p,\ell})}} \left\{ \exp \left( \lambda \rho_{p,\ell} [\bar{R} - \bar{r}_n] - \mathcal{K}(\rho_{p,\ell}, \pi_{p,\ell}) - \frac{\lambda^2 k_n^2}{n(1-p/n)^2} + \log w_{p,\ell} \right) \right\} \right] \leq 1.$$

Again a summation over all  $\lambda \in \mathcal{G}$  leads to the result. This ends the proof. □

**Proof of Lemma 6.5.** From the proof of the Lemma 6.1, we already know that  $|\bar{r}_n(\theta) - r_n(\theta)| \leq (1 + L)/(n - p) \sum_{i=1}^n \|X_i - \bar{X}_i\|$ . This bound holds uniformly on  $\Theta$ . As  $p \leq n/2$  it remains to estimate  $\pi_0[\exp(\lambda 2(1 + L)/n \sum_{i=1}^n \|X_i - \bar{X}_i\|)]$ . From the assumption (4.3), the stationarity of  $X$  and as the  $\xi_i$ s are i.i.d. we have:

$$\begin{aligned} & \pi_0 \left[ \exp \left( \lambda 2(1 + L)/n \sum_{i=1}^n \|X_i - \bar{X}_i\| \right) \right] \\ & \leq \pi_0 \left[ \exp \left( \lambda 2(1 + L)/n \sum_{i=1}^n \sum_{j=0}^{\infty} a_j(H) \|\xi_{i-j} - \bar{\xi}_{i-j}\| \right) \right] \\ & \leq \pi_0 \left[ \exp \left( \lambda 2(1 + L)/n \sum_{j=0}^{\infty} \sum_{i=1 \vee (n-j)}^n a_{n-i+j}(H) \|\xi_{n-j} - \bar{\xi}_{n-j}\| \right) \right] \\ & \leq \prod_{j=0}^{\infty} \pi_0 \left[ \exp \left( \lambda 2(1 + L)/n \sum_{i=1 \vee (n-j)}^n a_{n-i+j}(H) \|\xi_0\| \mathbb{1}_{\|\xi_0\| > C} \right) \right]. \end{aligned}$$

Denoting  $c_j = \lambda 2(1 + L) \sum_{i=1 \vee (n-j)}^n a_{n-i+j}(H)/n$ , we develop for all  $j \geq 0$

$$\pi_0[\exp(c_j \|\xi_0\| \mathbb{1}_{\|\xi_0\| > C})] = 1 + c_j \pi_0[\|\xi_0\| \mathbb{1}_{\|\xi_0\| > C}] + \sum_{k \geq 2} \frac{c_j^k \pi_0[\|\xi_0\|^k \mathbb{1}_{\|\xi_0\| > C}]}{k!}.$$

As  $\Psi(a(H)) = \pi_0[\exp(a(H)\|\xi_0\|)] = \sum_{k \geq 0} a(H)^k \pi_0[\|\xi_0\|^k]/k!$  is a convergent series of sequence of positive numbers, one gets

$$\pi_0[\|\xi_0\|^k \mathbb{1}_{\|\xi_0\| > C}] \leq \pi_0[\|\xi_0\|^k] \leq \frac{k! \Psi(a(H))}{a(H)^k} \quad \forall k \geq 2.$$

As  $\lambda < n/(4(1 + L))$  then  $2c_j \leq a(H)$  for all  $j \geq 0$  and then we derive that for all  $j \geq 0$ :

$$\begin{aligned} \pi_0[\exp(c_j \|\xi_0\| \mathbb{1}_{\|\xi_0\| > C})] &\leq 1 + c_j \pi_0[\|\xi_0\| \mathbb{1}_{\|\xi_0\| > C}] + \Psi(a(H)) \sum_{k \geq 2} (c_j/a(H))^k \\ &\leq 1 + c_j \pi_0[\|\xi_0\| \mathbb{1}_{\|\xi_0\| > C}] + \frac{\Psi(a(H))c_j^2}{a(H)(a(H) - c_j)} \\ &\leq 1 + c_j \pi_0[\|\xi_0\| \mathbb{1}_{\|\xi_0\| > C}] + c_j^2 \frac{2\Psi(a(H))}{a(H)^2}. \end{aligned}$$

As  $\phi(x) = (\exp(x) - 1)/x$  is an increasing function for  $x > 0$ , then  $\mathbb{1}_{\|\xi_0\| > C} \leq \phi(a(H)\|\xi_0\|)/\phi(a(H)C)$  and the Markov formula gives for all  $j \geq 0$

$$\pi_0[\exp(c_j \|\xi_0\| \mathbb{1}_{\|\xi_0\| > C})] \leq 1 + c_j \frac{\Psi(a(H))a(H)C}{\exp(a(H)C) - 1} + c_j^2 \frac{2\Psi(a(H))}{a(H)^2}.$$

Collecting those bounds, we obtain

$$\pi_0\left[\exp\left(\lambda \sup_{\theta \in \Theta} |\bar{r}_n(\theta) - r_n(\theta)|\right)\right] \leq \prod_{j=0}^{\infty} \left(1 + c_j \frac{\Psi(a(H))a(H)C}{\exp(a(H)C) - 1} + c_j^2 \frac{2\Psi(a(H))}{a(H)^2}\right).$$

Using that  $\log(1 + x) \leq x$  for all  $x > 0$ , we finally obtain:

$$\log\left(\pi_0\left[\exp\left(\lambda \sup_{\theta \in \Theta} |\bar{r}_n(\theta) - r_n(\theta)|\right)\right]\right) \leq \sum_{j=0}^{\infty} c_j \frac{\Psi(a(H))a(H)C}{\exp(a(H)C) - 1} + \sum_{j=0}^{\infty} c_j^2 \frac{2\Psi(a(H))}{a(H)^2}.$$

The desired result follows from the estimates  $\sum_{j=0}^{\infty} c_j \leq \lambda a(H)2(1 + L)$  and  $\sum_{j=0}^{\infty} c_j^2 \leq \lambda^2 a(H)^2 4(1 + L)^2/n$ . □

Now give the proof of the useful Proposition 2.1.

**Proof of Proposition 2.1.** Let us introduce a parameter  $\zeta > 0$  then we have

$$\begin{aligned} -\frac{1}{\gamma} \log \pi_{p,\ell}[\exp(-\gamma(R - R(\bar{\theta}_{p,\ell})))] - \zeta &= -\frac{1}{\gamma} \log \pi_{p,\ell}[\exp(-\gamma(R - R(\bar{\theta}_{p,\ell}) - \zeta))] \\ &\leq -\frac{1}{\gamma} \log \pi_{p,\ell}(R(\theta) - R(\bar{\theta}_{p,\ell}) \leq \zeta). \end{aligned}$$

Then we directly derive from the definition of  $d_{p,\ell}$  that

$$d_{p,\ell} \leq \sup_{\gamma > e} \frac{\inf_{\zeta > 0} \{ \zeta \gamma - \log \pi_{p,\ell}(R(\theta) - R(\bar{\theta}_{p,\ell}) \leq \zeta) \}}{\log \gamma}.$$

So

$$\zeta \gamma - q \log \frac{\zeta}{C c_{p,\ell}} \leq q \wedge \gamma C(c_{p,\ell} - \|\bar{\theta}_{p,\ell}\|) + q \log \left( \frac{C c_{p,\ell} \gamma}{q} \vee \frac{c_{p,\ell}}{c_{p,\ell} - \|\bar{\theta}_{p,\ell}\|} \right).$$

Now if  $q \leq \gamma C(c_{p,\ell} - \|\bar{\theta}_{p,\ell}\|)$  then we get the estimate  $q(1 + \log(C c_{p,\ell} \gamma / q)) / \log \gamma$  which decreases with  $\gamma$ . We then get the desired bound when the supremum is established for  $\gamma = e \vee q / (C(c_{p,\ell} - \|\bar{\theta}_{p,\ell}\|))$ . If  $q \geq \gamma C(c_{p,\ell} - \|\bar{\theta}_{p,\ell}\|)$ , then we get the estimate  $(\gamma C(c_{p,\ell} - \|\bar{\theta}_{p,\ell}\|) + q \log(c_{p,\ell} / (c_{p,\ell} - \|\bar{\theta}_{p,\ell}\|))) / \log \gamma$  which increases with  $\gamma$ . We have to consider  $\gamma$  as large as possible, that is, when  $q = \gamma C(c_{p,\ell} - \|\bar{\theta}_{p,\ell}\|)$  and we are going back to the case treated above. □

### 6.5. Proofs of the results given in Section 4

After proving Proposition 4.1, we give Lemma 6.7 that introduces a coupling argument used to estimate the coefficients  $\theta_{\infty,n}(1)$  in Propositions 4.2 and 4.3.

**Proof of Proposition 4.1.** The Theorem 3.1 of Doukhan and Wintenberger [15] gives the existence of a unique stationary solution and the existence of an  $H$  such that  $X_t = H(\xi_t, \xi_{t-1}, \xi_{t-2}, \dots)$ . We prove that conditions (3.1) and (3.2) are automatically satisfied. Let  $(x_i)$  and  $(y_i)$  be two sequences such that there exists  $j \in \mathbb{N}$  with  $x_i = y_i$  for all  $i \neq j$ . Then  $H(x) = u_0^\infty$  where  $u_0^\infty = \lim_{k \rightarrow \infty} u_0^k$  for  $(u_{-i}^k)_{i \in \mathbb{N}}$  defined recursively by

$$u_{-i}^k = F(u_{-i-1}^k, u_{-i-2}^k, \dots, u_{1-k}^k, u_{-k}^k, 0, \dots; x_i) \quad \forall 0 \leq i \leq k.$$

Similarly, we denote  $H(y) = v_0^\infty$  such that  $\|H(x) - H(y)\| = \|u_0^\infty - v_0^\infty\|$ . For  $j = 0$ , using (4.3)  $\|u_0^k - v_0^k\| \leq u \|x_j - y_j\|$  for all  $k$ . For  $j \geq 1$ , as  $x_i = y_i$  for  $i > j$ , for  $k$  sufficiently large it holds (with the convention  $\sum_{\ell=1}^{-k} = 0$  for  $k \geq 0$ ):

$$\|u_0^k - v_0^k\| \leq \sum_{\ell_1=1}^j a_{\ell_1}(F) \sum_{\ell_2=1}^{j-\ell_1} a_{\ell_2}(F) \dots \sum_{\ell_j=1}^{j-\ell_1-\dots-\ell_{j-1}} a_{\ell_j}(F) \|u_{-j}^k - v_{-j}^k\|.$$

By definition  $\|u_{-j}^k - v_{-j}^k\| \leq u \|x_j - y_j\|$  and we obtain  $\|u_0^k - v_0^k\| \leq u a(F)^{j-1} \|x_j - y_j\|$  for sufficiently large  $k$ . As the estimate does not depends on  $k$ , we derive that (3.1) holds with  $a_j(H) = u a(F)^{j-1}$  and that (3.2) follows from the condition (4.4). □

Now we state a useful coupling lemma;  $(X_t^*)$  is said to be a coupling version of  $(X_t)$  if it is similarly distributed and such that  $(X_t^*)_{t>0}$  is independent of  $\mathfrak{S}_0 = \sigma(X_t, t \leq 0)$ . From a version

of the Kantorovitch–Rubinstein duality, see Dedecker and Prieur [13] for more details, we obtain an estimate of  $\theta_{\infty,n}(1)$ .

**Lemma 6.7.** *For any version  $(X_t^*)$ , we have*

$$\theta_{\infty,n}(1) \leq \sum_{i=1}^n \left\| \mathbb{E}(\|X_i - X_i^*\|/\mathfrak{G}_0) \right\|_{\infty}.$$

For the sake of completeness, we give the proof of this already known result.

**Proof of Lemma 6.7.** As we equipped  $\mathcal{X}^n$  with the norm  $\|(x_1, \dots, x_n)\| = \sum_{i=1}^n \|x_i\|$ , we immediately get the inequality

$$\theta_{\infty,n}(1) \leq \left\| \mathbb{E}(\|(X_1, \dots, X_n) - (X_1^*, \dots, X_n^*)\|/\mathfrak{G}_0) \right\|_{\infty} \leq \sum_{i=1}^n \left\| \mathbb{E}(\|X_i - X_i^*\|/\mathfrak{G}_0) \right\|_{\infty}. \quad \square$$

The proof of Propositions 4.2 and 4.3 are simple applications of this lemma.

**Proof of Proposition 4.2.** Let us consider the coupling version of the causal Bernoulli shift  $(X_t)$  given by

$$X_t^* = H(\xi_t, \xi_{t-1}, \dots, \xi_1, \xi_0^*, \xi_{-1}^*, \dots) \quad \forall t \in \mathbb{Z},$$

where  $(\xi_t^*)$  is similarly distributed than  $(\xi_t)$  and the two processes are independent. Then from Lemma 6.7 and condition (3.1), we obtain:

$$\theta_{\infty,n}(1) \leq \sum_{i=1}^n \left\| \sum_{j=i}^{\infty} a_j(H) \mathbb{E}(\|\xi_{i-j} - \xi_{i-j}^*\|/\mathfrak{G}_0) \right\|_{\infty} \leq \sum_{j=i}^{\infty} j a_j(H) \left\| \mathbb{E}(\|\xi_{i-j} - \xi_{i-j}^*\|/\mathfrak{G}_0) \right\|_{\infty}$$

and the desired result follows. □

**Proof of Proposition 4.3.** Here we will consider the maximal coupling scheme of Goldstein [16]: there exists a version  $(X_t^*)$  such that

$$\left\| \mathbb{P}(X_t \neq X_t^* \text{ for some } t \geq r/\mathfrak{G}_0) \right\|_{\infty} = \sup_{(A,B) \in \mathfrak{G}_0 \times \mathfrak{F}_r} |\mathbb{P}(A/B) - P(B)| = \varphi(r).$$

As  $\|Y - Z\| \leq 2\|X_0\|_{\infty} \mathbb{1}_{Y \neq Z}$  for any variables  $Y, Z$  bounded by  $\|X_0\|_{\infty}$ , we have:

$$\left\| \mathbb{E}(\|X_i - X_i^*\|/\mathfrak{G}_0) \right\|_{\infty} \leq 2\|X_0\|_{\infty} \left\| \mathbb{E}(\mathbb{1}_{X_i \neq X_i^*}/\mathfrak{G}_0) \right\|_{\infty} \leq 2\|X_0\|_{\infty} \left\| \mathbb{P}(X_i \neq X_i^*/\mathfrak{G}_0) \right\|_{\infty}.$$

As  $\mathbb{P}(X_i \neq X_i^*/\mathfrak{G}_0) \leq \mathbb{P}(X_t \neq X_t^* \text{ for some } t \geq r/\mathfrak{G}_0)$ , we conclude using Lemma 6.7. □

### 6.6. Proofs of the results given in Section 5

We proof the Corollaries 5.2 and 5.3 of Theorem 3.1 applied in the context of Neural Networks and projection in the Fourier basis predictors.

**Proof of Corollary 5.2.** Let us check that all the predictors are  $L$ -Lipschitz functions of the observations. For any  $x, y \in \mathbb{R}^p$ , as the function  $\phi$  is 1-Lipschitz, we have

$$\begin{aligned} |f_\theta(x) - f_\theta(y)| &\leq \left| \sum_{k=1}^{\ell} c_k (\phi(a_k \cdot x + b_k) - \phi(a_k \cdot y + b_k)) \right| \\ &\leq \sum_{k=1}^{\ell} |c_k| |a_k \cdot (x - y)| \leq \sum_{k=1}^{\ell} |c_k| \|a_k\|_1 \|x - y\|_\infty \\ &\leq \|a_k\|_1 \sum_{k=1}^{\ell} |c_k| \sum_{i=1}^p |x_i - y_i|. \end{aligned}$$

For  $\theta \in \mathcal{B}_{c_{p,\ell}}^q$  then  $L = (c_{p,\ell} \vee 1)^3$  is convenient. Next, using Jensen to estimate  $\mathbb{L}_1$ -risk by  $\mathbb{L}_2$ -risk, we obtain from the Theorem 1 of Barron [6] the existence of  $C > 0$  such that

$$\pi_0[|\text{med}(X_0|X_{-1}, \dots, X_{-p}) - f_{\bar{\theta}_{p,\ell}}(X_{-1}, \dots, X_{-p})|] \leq C \frac{p^c \|X_0\|_\infty}{\sqrt{\ell}},$$

where  $\bar{\theta}_{p,\ell}$  belongs to the compact set

$$\mathcal{B}'_{p,\ell} = \left\{ \theta \in \mathbb{R}^{\ell(p+2)+1}; \sum_{i=1}^{\ell} |c_i| \leq C' c^p; \max_{1 \leq i \leq \ell} \|a_i\| \leq \sqrt{\ell} \log \ell; \max_{1 \leq i \leq \ell} |b_i| \leq \|X_0\|_\infty \sqrt{\ell} \log \ell \right\}.$$

Remark that under the assumptions of Corollary 5.2, we have  $c_{p,\ell} - \|\bar{\theta}_{p,\ell}\| \geq q/e$ . It implies by Proposition 2.1 that  $d_{p,\ell} \leq 3q(1 + \log(c_{p,\ell}))$  when  $c_{p,\ell} \geq 1$ . From Theorem 3.1 there exists  $C > 0$  satisfying

$$R(\hat{\theta}) \leq \inf_{d_{p,\ell} \leq n} \left\{ R_p^* + C \left( \frac{p^c}{\sqrt{\ell}} + \log^3(n) \sqrt{\frac{p\ell}{n}} \right) \right\} + C \frac{\log(1/\varepsilon)}{\sqrt{n}}.$$

The result follows from considering  $\ell = \sqrt{n} p^{c-1/2}$ . □

**Proof.** Proof of Proposition 5.3 Let us apply Theorem 3.2: there exists  $C > 0$  such that

$$\begin{aligned} R(\hat{\theta}) &\leq \inf_{p,\ell:d_{p,\ell} \leq n} \left\{ \min_{\theta \in \Theta_{p,\ell}} R(\theta) + C \sqrt{\frac{d_{p,\ell}}{n}} \log^{5/2}(n) \right\} + C \frac{\log(1/\varepsilon)}{\sqrt{n}} \\ &\leq \inf_{\ell:d_{p_0,\ell} \leq n} \left\{ \min_{\theta \in \Theta_{p_0,\ell}} R(\theta) + C \sqrt{\frac{d_{p_0,\ell}}{n}} \log^{5/2}(n) \right\} + C \frac{\log(1/\varepsilon)}{\sqrt{n}}. \end{aligned}$$

Remarking that

$$\begin{aligned}
 R(\bar{\theta}_{p_0, \ell}) &= \inf_{\theta \in \Theta} \pi_0[|X_{p+1} - f_{\bar{\theta}_{p_0, \ell}}(X_p, \dots, X_1)|] \\
 &\leq \pi_0 \left[ \left| X_{p+1} - \sum_{i=1}^{p_0} f_i(X_{p-i}) \right| \right] + \inf_{\theta \in \Theta} \pi_0 \left[ \left| \sum_{i=1}^{p_0} f_i(X_{p-i}) - \sum_{i=1}^{p_0} \sum_{j=1}^n \theta_{i,j} \varphi_j(X_{p-i}) \right| \right] \\
 &\leq \mu[|\xi_0|] + \inf_{\theta \in \Theta} \sum_{i=1}^{p_0} \pi_0 \left[ \left| f_i(X_1) - \sum_{j=1}^n \theta_{i,j} \varphi_j(X_1) \right| \right].
 \end{aligned}$$

Note also that under our hypothesis  $X_1$  has a density upper bounded by  $1/\sqrt{2\pi\sigma^2}$ . It then holds

$$\begin{aligned}
 R(\bar{\theta}_{p_0, \ell}) &\leq \mu[|\xi_0|] + \frac{1}{\sqrt{2\pi\sigma^2}} \inf_{\theta \in \Theta} \sum_{i=1}^{p_0} \int \left| f_i(x) - \sum_{j=1}^n \theta_{i,j} \varphi_j(x) \right| dx \\
 &\leq \mu[|\xi_0|] + \frac{1}{\sqrt{2\pi\sigma^2}} \inf_{\theta \in \Theta} \sum_{i=1}^{p_0} \left( \int \left[ f_i(x) - \sum_{j=1}^n \theta_{i,j} \varphi_j(x) \right]^2 dx \right)^{1/2} \\
 &\leq \mu[|\xi_0|] + \frac{1}{\sqrt{2\pi\sigma^2}} \sum_{i=1}^{p_0} \gamma_i \ell^{-s_i} \leq \mu[|\xi_0|] + \frac{\sum_{i=1}^{p_0} \gamma_i}{\sqrt{2\pi\sigma^2}} \ell^{-s}.
 \end{aligned}$$

Then we have

$$\pi_0[R(\hat{\theta})] \leq \mu[|\xi_0|] + \inf_{\ell} \left\{ \ell^{-s} \frac{\sum_{i=1}^{p_0} \gamma_i}{\sqrt{2\pi\sigma^2}} + C \sqrt{\frac{d_{p_0, \ell}}{n}} \log^{5/2}(n) \right\} + C \frac{\log(1/\varepsilon)}{\sqrt{n}}. \tag{6.13}$$

The estimate of  $d_{p_0, \ell}$  from Proposition 2.1 is plugged into (6.13) to obtain for some  $C > 0$

$$\pi_0[R(\hat{\theta})] \leq \mu[|\xi_0|] + \inf_{\ell} \left\{ \ell^{-s} \frac{\sum_{i=1}^{p_0} \gamma_i}{\sqrt{2\pi\sigma^2}} + C \sqrt{\frac{p_0 \ell}{n}} \log^{5/2}(n) \right\} + C \frac{\log(1/\varepsilon)}{\sqrt{n}}.$$

In particular, fixing  $\ell$  proportional to  $n^{1/(2s+1)}$  leads to the result. □

## Acknowledgements

We would like to thank the anonymous referees for the various corrections and improvements they suggested.

## References

- [1] Akaike, H. (1973). Information theory and an extension of the maximum likelihood principle. In *Second International Symposium on Information Theory (Tsahkadsor, 1971)* 267–281. Budapest: Akadémiai Kiadó. [MR0483125](#)
- [2] Alquier, P. (2008). PAC-Bayesian bounds for randomized empirical risk minimizers. *Math. Methods Statist.* **17** 279–304. [MR2483458](#)
- [3] Andrews, D.W.K. (1984). Nonstrong mixing autoregressive processes. *J. Appl. Probab.* **21** 930–934. [MR0766830](#)
- [4] Audibert, J.Y. (2004). Aggregated estimators and empirical complexity for least square regression. *Ann. Inst. Henri Poincaré Probab. Stat.* **40** 685–736. [MR2096215](#)
- [5] Baraud, Y., Comte, F. and Viennet, G. (2001). Adaptive estimation in autoregression or  $\beta$ -mixing regression via model selection. *Ann. Statist.* **29** 839–875. [MR1865343](#)
- [6] Barron, A.R. (1994). Approximation and estimation bounds for artificial neural networks. *Machine Learning* **14** 115–133.
- [7] Catoni, O. (2003). A PAC-Bayesian approach to adaptative classification. Preprint, Laboratoire de Probabilités et Modèles Aléatoires.
- [8] Catoni, O. (2004). *Statistical Learning Theory and Stochastic Optimization. Lecture Notes in Math.* **1851**. Berlin: Springer. Lecture notes from the 31st Summer School on Probability Theory held in Saint-Flour, July 8–25, 2001. [MR2163920](#)
- [9] Catoni, O. (2007). *Pac-Bayesian Supervised Classification: The Thermodynamics of Statistical Learning. Institute of Mathematical Statistics Lecture Notes – Monograph Series* **56**. Beachwood, OH: IMS. [MR2483528](#)
- [10] Cesa-Bianchi, N. and Lugosi, G. (2006). *Prediction, Learning, and Games*. Cambridge: Cambridge Univ. Press. [MR2409394](#)
- [11] Dalalyan, A. and Tsybakov, A. (2008). Aggregation by exponential weighting, sharp oracle inequalities and sparsity. *Machine Learning* **72** 39–61.
- [12] Dedecker, J., Doukhan, P., Lang, G., León R., J.R., Louhichi, S. and Prieur, C. (2007). *Weak Dependence: With Examples and Applications. Lecture Notes in Statistics* **190**. New York: Springer. [MR2338725](#)
- [13] Dedecker, J. and Prieur, C. (2005). New dependence coefficients. Examples and applications to statistics. *Probab. Theory Related Fields* **132** 203–236. [MR2199291](#)
- [14] Doukhan, P. (1994). *Mixing: Properties and Examples. Lecture Notes in Statistics* **85**. New York: Springer. [MR1312160](#)
- [15] Doukhan, P. and Wintenberger, O. (2008). Weakly dependent chains with infinite memory. *Stochastic Process. Appl.* **118** 1997–2013. [MR2462284](#)
- [16] Goldstein, S. (1978/79). Maximal coupling. *Z. Wahrsch. Verw. Gebiete* **46** 193–204. [MR0516740](#)
- [17] Ibragimov, I. (1962). Some limit theorems for stationary processes. *Theory Probab. Appl.* **7** 349–382.
- [18] Ing, C.K. (2007). Accumulated prediction errors, information criteria and optimal forecasting for autoregressive time series. *Ann. Statist.* **35** 1238–1277. [MR2341705](#)
- [19] Lacour, C. (2008). Nonparametric estimation of the stationary density and the transition density of a Markov chain. *Stochastic Process. Appl.* **118** 232–260. [MR2376901](#)
- [20] Massart, P. (2007). *Concentration Inequalities and Model Selection. Lecture Notes in Math.* **1896**. Berlin: Springer. Lectures from the 33rd Summer School on Probability Theory held in Saint-Flour, July 6–23, 2003. [MR2319879](#)
- [21] McAllester, D.A. (1998). Some PAC-Bayesian theorems. In *Proceedings of the Eleventh Annual Conference on Computational Learning Theory (Madison, WI, 1998)* 230–234 (electronic). New York: ACM. [MR1811587](#)

- [22] Meir, R. (2000). Nonparametric model selection through adaptive model selection. *Machine Learning* **39** 5–34.
- [23] Modha, D.S. and Masry, E. (1998). Memory-universal prediction of stationary random processes. *IEEE Trans. Inform. Theory* **44** 117–133. [MR1486652](#)
- [24] R Development Core Team (2008). *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing. ISBN 3-900051-07-0.
- [25] Rio, E. (2000). Inégalités de Hoeffding pour les fonctions lipschitziennes de suites dépendantes. *C. R. Acad. Sci. Paris Sér. I Math.* **330** 905–908. [MR1771956](#)
- [26] Rio, E. (2000). *Théorie Asymptotique des Processus Aléatoires Faiblement Dépendants. Mathématiques & Applications (Berlin) [Mathematics & Applications]* **31**. Berlin: Springer. [MR2117923](#)
- [27] Schwarz, G. (1978). Estimating the dimension of a model. *Ann. Statist.* **6** 461–464. [MR0468014](#)
- [28] Shawe-Taylor, J. and Williamson, R. (1997). A pac analysis of a Bayes estimator. In *Proceedings of the Tenth Annual Conference on Computational Learning Theory, COLT'97* 2–9. New York: ACM.
- [29] Stoltz, G. (2005). Information incomplète et regret interne en prédiction de suites individuelles. Ph.D. thesis, Univ. Paris Sud.
- [30] Vapnik, V.N. (1995). *The Nature of Statistical Learning Theory*. New York: Springer. [MR1367965](#)

*Received February 2009 and revised January 2011*