

# Nonparametric regression with dependent errors

YUHONG YANG

*Department of Statistics, Iowa State University, Ames, IA 50011, USA.*

*E-mail: [yyang@iastate.edu](mailto:yyang@iastate.edu)*

We study minimax rates of convergence for nonparametric regression under a random design with dependent errors. It is shown that when the errors are independent of the explanatory variables, long-range dependence among the errors does not necessarily hurt regression estimation, which at first glance contradicts earlier results by Hall and Hart, Wang, and Johnstone and Silverman under a fixed design. In fact we show that, in general, the minimax rate of convergence under the square  $L_2$  loss is simply at the worse of two quantities: one determined by the massiveness of the class alone and the other by the severity of the dependence among the errors alone. The clear separation of the effects of the function class and dependence among the errors in determining the minimax rate of convergence is somewhat surprising. Examples of function classes under different covariance structures including both short- and long-range dependences are given.

*Keywords:* long-range dependent errors; minimax rate of convergence; nonparametric regression

## 1. Introduction

### 1.1. The problem of interest

Assume that we observe random variables  $(X_i, Y_i)_{i=1}^n$ , where  $Y_i$  takes values in  $\mathbb{R}$  and  $X_i$  takes values in  $\mathcal{X}$ , a subset of  $\mathbb{R}^d$  for some  $d \geq 1$ . The relationship between response variables  $Y_i$  and the explanatory or experimental variables  $X_i$  is modelled as

$$Y_i = u(X_i) + \varepsilon_i, \quad i \geq 1, \quad (1)$$

where  $u$  is an unknown regression function. The random errors  $\{\varepsilon_i, i \geq 1\}$  are assumed to have a joint normal distribution conditioned on  $\{X_i, i \geq 1\}$  with mean zero and a known covariance matrix. A goal is to estimate the regression function  $u$ , which is assumed a priori to be in a nonparametric function class  $\mathcal{U}$  (e.g., monotone or Lipschitz). In this paper, we study how well one can estimate  $u$  under a minimax consideration over the function class  $\mathcal{U}$ . The focus is on determination of minimax rates of convergence for the estimation problem when the errors are dependent. We will characterize how dependence of the errors as well as the function class affects the minimax rates of convergence under appropriate conditions.

## 1.2. Some background

In recent years, there has been an increasing interest in statistical estimation based on long-range dependent data – see Beran (1994) for a survey of work in this area. Long-range dependence has been observed in many applied scientific disciplines. Künsch *et al.* (1993) wrote: ‘Perhaps most unbelievable to many is the observation that high-quality measurement series from astronomy, physics, chemistry, generally regarded as prototypes of “i.i.d.” observations, are not independent but long-range correlated’. Based on the empirical evidence of long-range dependence in measurements and other applications, it becomes important to study how long-range dependence affects statistical estimation.

For parametric regression with fixed designs, asymptotic results for maximum likelihood and least-squares estimators under long-range dependence have been established by Yajima (see, for example, Yajima 1991). Künsch *et al.* (1993) show that for certain analysis of variance models with random designs, contrasts can be estimated at the same rate as that under independent errors. Asymptotic results for the estimation of long-range dependence parameters under parametric models are given, for example, in Beran (1986), Fox and Taqqu (1985), Dahlhaus (1989), Giraitis and Surgailis (1990) and Robinson (1995).

For nonparametric regression, the effect of long-range dependence on minimax rates of convergence is studied in a pioneering paper by Hall and Hart (1990a) for a differentiable function class, and later by Wang (1996) and Johnstone and Silverman (1997) for Besov classes, all under a fixed equally spaced design. These results show that a certain long-range dependence of errors damages the minimax rate of convergence for regression estimation. The latter two papers propose adaptive wavelet estimators. In addition, Wang shows that for some inhomogeneous Besov classes, linear estimators cannot achieve the minimax rate of convergence, and Johnstone and Silverman show that when an unknown dependence parameter is properly estimated, a wavelet threshold estimator is adaptive with respect to both the dependence parameter and the smoothness parameters. Robinson (1997) derives local asymptotic normality for kernel estimators under long-range dependence.

In this work, we study the effects of a general dependence among the errors on regression estimation for a general nonparametric function class, under a random design. The focus is on the theoretic determination of the minimax rate of convergence. We do not address issues of estimation of dependence and adaptive estimation.

Finally, we point out that, independently of our work, Efromovich (1999) obtains minimax rates of convergence for regression estimation for Hölder classes under a long-range dependence and a random design. He proposes a series expansion estimator and shows that it is adaptive with respect to a smoothness parameter. Our results on minimax rates of convergence apply to general classes of regression functions satisfying a mild richness assumption.

## 1.3. A summary of our findings

We summarize our results informally below. Our conclusions are in terms of minimax rates of convergence under the square  $L_2$  loss with a random design. The errors are assumed to be independent of the explanatory variables.

1. If the variances of the errors are uniformly upper-bounded, then the regression function up to a constant can be estimated as well as under independent and identically distributed (i.i.d.) errors.
2. Under some mild conditions, the minimax risk for estimating the regression function in a class converges at a rate of the maximum of two quantities: the minimax rate of the same function class but under i.i.d. errors, and the minimax rate for estimating the mean value of the regression function.

From the foregoing, the effect of dependence of serially correlated errors on regression is sort of ‘parametric’, in the sense that it does not affect the rate of convergence more than adding the risk for estimating a single parameter (the mean of the regression function). Similar phenomena have been observed earlier for some parametric models (see, for example, Künsch *et al.* 1993) and density estimation (Hall and Hart 1990b) both under long-range dependence.

The paper is organized as follows. Some preliminary considerations are addressed in Section 2. The main results on regression estimation are presented in Section 3. A key proposition on minimax risk bounds is presented in Section 4. The proofs of the main results as well as useful lemmas are given in Section 5.

## 2. Risks of interests and metric entropy

### 2.1. Risk for regression estimation

We assume that  $\{X_i, i \geq 1\}$  are i.i.d. with density  $h$  with respect to a measure  $\mu$ . For the nonparametric class  $\mathcal{U}$  supposed to contain  $u$ , we assume that  $\mathcal{U}$  is uniformly bounded throughout the paper.

For regression estimation, we obtain results when the errors are independent of  $X^n$ , that is, the conditional covariance matrix  $\Omega_n$  of  $\{\varepsilon_i, 1 \leq i \leq n\}$  given  $X^n = (X_1, \dots, X_n)$  does not depend on  $X^n$ . We assume that  $\Omega_n$  is known. Let  $\|u - v\|_{L_2(h)} = (\int (u - v)^2 h d\mu)^{1/2}$  be the  $L_2$  distance between two functions  $u$  and  $v$  with respect to the design density of  $X_1$ . Since  $\mathcal{U}$  is uniformly bounded, the distance is well defined within the function class.

The minimax risk we examine for estimating the regression function  $u$  is

$$R(\mathcal{U}; \Omega; n) = \min_{\hat{u}} \max_{u \in \mathcal{U}} E \|u - \hat{u}\|_{L_2(h)}^2,$$

where  $\hat{u}$  is over all estimators based on  $(X_i, Y_i)_{i=1}^n$  and the expectation is taken under the true regression function  $u$ . The minimax risk measures how well one can possibly estimate  $u$  uniformly over the function class.

The condition, that  $\text{tr}(\Omega_n^{-1})$  is of order  $n$  ( $\text{tr}(\cdot)$  denotes the trace of a square matrix), will be used for identifying minimax rates. It is satisfied by short- and long-range dependent cases as given in Section 3.3. It also holds for stationary invertible autoregressive errors as studied in Hall and Hart (1990a). Let  $\sigma_i^2 = \text{var}(\varepsilon_i)$ . A simple sufficient condition for  $\text{tr}(\Omega_n^{-1}) \asymp n$  is that  $\sup \sigma_i^2 < \infty$  and there is a white noise component in the errors, that is,  $\varepsilon_i = \varepsilon_i^{(1)} + \varepsilon_i^{(2)}$ , where  $\{\varepsilon_i^{(1)}, i \geq 1\}$  are i.i.d. and independent of  $\{\varepsilon_i^{(2)}\}$  (see Lemma 7 in

Section 5). When the trace condition is not satisfied, rates better than that under i.i.d. errors are possible. For instance, assume that the errors are independent with decreasing variances  $\sigma_i^2$  of order  $i^{-1}$ . Then it is intuitively clear that the rate of convergence can be faster compared with that under i.i.d. errors.

## 2.2. Estimation of the mean of the regression function

Related to the above problem of regression estimation is the problem of estimating the mean value of the regression function with respect to the design density. As will be seen, this 'parametric' problem characterizes the influence of serial dependence of errors on regression estimation.

Let  $\Delta = \{\eta(u) = \int uh \, d\mu: u \in \mathcal{U}\}$  be the set of all possible mean values of  $u(X)$  for the class  $\mathcal{U}$ . Let

$$r_n = \min_{\hat{\eta}} \max_{u \in \mathcal{U}} E(\hat{\eta} - \eta(u))^2 \quad (2)$$

be the minimax risk for estimating  $\eta(u)$ , where the minimization is over  $\hat{\eta}$  based on  $(X_i, Y_i)_{i=1}^n$ .

## 2.3. Estimation of the regression function up to a constant

Long-range dependence makes the estimation of the mean of the regression function harder and therefore may affect the rate for estimating the whole regression function. In some applications, it is the trend or change of the function that is of interest. Then it is appropriate to estimate the regression function up to a constant.

Let  $u_0(x) = u(x) - \eta(u)$  be a centred version of the regression function (centred according to the design density). The minimax risk for the estimation of  $u_0$  is

$$R_0(\mathcal{U}; \Omega; n) = \min_{\hat{u}_0} \max_{u \in \mathcal{U}} E\|u_0 - \hat{u}_0\|_{L_2(h)}^2,$$

where  $\hat{u}_0$  is over all estimators based on  $(X_i, Y_i)_{i=1}^n$ .

## 2.4. Metric entropy as a measure of massiveness of a function class

It is clear that as the function class  $\mathcal{U}$  grows larger, so the minimax risk increases (or at least does not decrease). For nonparametric regression with independent errors, it is known that massiveness of a target function class affects the minimax rate of convergence in terms of the metric entropy order of the function class (see, for example, Ibragimov and Hasminskii 1977; Bretagnolle and Huber 1979; Birgé 1983; 1986; Le Cam 1986, Chapter 16; Yatracos 1988; Yang and Barron 1999). Metric entropy as a measure of massiveness of a function class was intensively studied in Kolmogorov and Tihomirov (1959), and since then results have been obtained on the orders of metric entropy for the classical function classes and some others under various norms (see, for example, Lorentz *et al.*, 1996).

A finite subset  $N_\epsilon$  is called an  $\epsilon$ -packing set in  $\mathcal{U}$  under a distance  $d$  if  $d(u, v) > \epsilon$  for any  $u, v \in N_\epsilon$  with  $u \neq v$ . Let  $M_2(\epsilon) = M_2(\epsilon; \mathcal{U})$  be the maximal logarithm of the cardinality of any  $\epsilon$ -packing set under the  $L_2(h)$  distance. Clearly  $M_2(\epsilon)$  is non-increasing in  $\epsilon$ . The asymptotic behaviour of  $M_2(\epsilon)$  as  $\epsilon \rightarrow 0$  reflects how massive the class  $\mathcal{U}$  is under the given distance. We call  $M_2(\epsilon)$  the packing  $\epsilon$ -entropy or simply the metric entropy of  $\mathcal{U}$ .

Throughout this paper, we assume  $M_2(\epsilon) < \infty$  for every  $\epsilon > 0$  (which necessarily requires  $\mathcal{U}$  to be bounded in  $L_2(h)$  norm) and  $M_2(\epsilon) \rightarrow \infty$  as  $\epsilon \rightarrow 0$  (which excludes trivial cases when  $\mathcal{U}$  is finite). These conditions are satisfied if  $\mathcal{U}$  is not finite, separable, and compact in  $L_2(h)$  norm.

For most function classes, the metric entropies are known only up to orders. For that reason, we assume that  $M(\epsilon)$  is an available non-increasing function known to be of order  $M_2(\epsilon)$ . We call a class  $\mathcal{U}$  rich if, for some constant  $0 < \tau < 1$ ,

$$\liminf_{\epsilon \rightarrow 0} \frac{M(\tau\epsilon)}{M(\epsilon)} > 1. \tag{3}$$

This condition is a characteristic of familiar nonparametric classes (except classes of analytic functions), for which the metric entropy is usually of order  $\epsilon^{-\alpha} \log(1/\epsilon)^\beta$  for some  $\alpha > 0$  and  $\beta \in \mathbb{R}$ .

### 3. Main results

In this paper, the expression  $a_n \leq b_n$  means that  $\limsup(a_n/b_n) < \infty$ . If  $a_n \leq b_n$  and  $b_n \leq a_n$  (i.e.,  $a_n$  and  $b_n$  are of the same order), we write  $a_n \asymp b_n$ .

#### 3.1. Regression estimation

The explanatory variables  $X_1, X_2, \dots$  are assumed to be i.i.d. with known density  $h$  with respect to a measure  $\mu$ . They are further assumed to be independent of the errors  $\epsilon_i$  in model (1). The following additional assumptions will be used for our results.

**Assumption 1.** *The class  $\mathcal{U}$  is uniformly bounded, that is, there exists a known constant  $L$  such that  $\sup_{u \in \mathcal{U}} \|u\|_\infty \leq L < \infty$ .*

**Assumption 2.** *The class  $\mathcal{U}$  is rich, as defined in (3).*

**Assumption 3.** *The class  $\mathcal{U}$  contains the constant functions  $u \equiv c$  with  $c \in \Delta$ .*

**Assumption 4.** *The mean value set  $\Delta$  contains an interval  $[a, b]$  with  $a < b$ .*

**Assumption 5.**  $\sup_{i \geq 1} \sigma_i^2 < \infty$ .

**Assumption 6.**  $\text{tr}(\Omega_n^{-1}) \asymp n$ .

Assumption 4 excludes cases where the estimation of  $\eta(u)$  is trivial.

Choose  $\epsilon_n$  such that

$$M(\epsilon_n) \asymp n\epsilon_n^2. \quad (4)$$

Under the richness assumption in (3), any two sequences of solution to the equation are of the same order, and  $\epsilon_n^2$  gives the minimax rate of convergence for estimating the regression function under i.i.d. errors (see, for example, Birgé 1983; Le Cam 1986; and Yang and Barron 1999). An interpretation of the equation is that if we discretize the function class  $\mathcal{U}$  using an  $\epsilon$ -net, then  $\epsilon_n$  balances the estimation error of order  $M(\epsilon)/n$  (due to identifying a good element in the  $\epsilon$ -net based on data) and the approximation error (bias squared, due to discretization)  $\epsilon^2$ . Throughout this paper, unless stated otherwise,  $\epsilon_n$  is defined as above.

**Theorem 1.** *If Assumptions 1–6 are satisfied, we have the following conclusions.*

(i) *The minimax risk for estimating  $u_0$  is of order  $\epsilon_n^2$ , that is,*

$$R_0(\mathcal{U}; \Omega; n) \asymp \epsilon_n^2. \quad (5)$$

(ii) *The minimax risk for regression function estimation is at the rate of the maximum (or equivalently, the sum) of two quantities: the minimax rate of the same class but under i.i.d. errors, and the rate for estimating the mean  $\eta = Eu(X)$  of the regression function under the correlated errors. That is,*

$$R(\mathcal{U}; \Omega; n) \asymp r_n + \epsilon_n^2. \quad (6)$$

**Remarks.** (i) Without assuming  $\text{tr}(\Omega_n^{-1}) \asymp n$  (Assumption 6), the above quantities  $\epsilon_n^2$  and  $r_n + \epsilon_n^2$  give valid upper rates respectively (see the proof of Theorem 1 in Section 5), but they are not necessarily optimal in general (see Section 2.1).

(ii) A parametric analogue of (5) is given in Künsch *et al.* (1993), where it is shown that the rate of convergence for estimating a contrast (similar in spirit to  $u_0$ ) remain unchanged for some ANOVA models.

From the foregoing, in particular, for stationary Gaussian errors independent of  $X^n$ , the regression function up to a constant can be estimated as well as under i.i.d. errors. For the estimation of the whole regression function, however, the minimax rate for estimating  $\eta(u)$  may hurt. Roughly speaking, the difficulty in estimating  $u$  is determined by the maximum of that caused by largeness of the function class  $\mathcal{U}$  and that caused by the dependence among the errors in estimating a constant. The separation of the roles of the function class and dependence is somewhat surprising. This separation may not hold when the random errors and the explanatory variables are not independent.

From Theorem 1, once we know the metric entropy order of a nonparametric class and the minimax rate for estimating  $\eta(u)$ , the minimax rate for regression is determined. The metric entropies for classical function classes are usually of order  $M(\epsilon) \asymp \epsilon^{-d/\alpha}(\log(1/\epsilon))^\beta$ , where  $d$  is the dimension of  $\mathcal{X}$ ,  $\alpha$  is a smoothness parameter of the class measured in some way (e.g., in terms of derivatives, or a modulus of continuity) and  $\beta \in \mathbb{R}$ . Then,

solving  $M(\epsilon_n) = n\epsilon_n^2$ , we have  $\epsilon_n^2$  of order  $n^{-2\alpha/(2\alpha+d)}(\log n)^{2\alpha\beta/(2\alpha+d)}$ . If  $r_n \asymp n^{-\gamma}$  for some  $0 < \gamma < 1$  (as for the long-range dependence case in Section 3.2), then

$$R(\mathcal{U}; \Omega; n) \asymp \begin{cases} n^{-2\alpha/(2\alpha+d)}(\log n)^{2\alpha\beta/(2\alpha+d)} & \text{if } \gamma > 2\alpha/(2\alpha+d), \text{ or } \gamma = 2\alpha/(2\alpha+d) \text{ and } \beta \geq 0 \\ n^{-\gamma} & \text{if } \gamma < 2\alpha/(2\alpha+d), \text{ or } \gamma = 2\alpha/(2\alpha+d) \text{ and } \beta < 0. \end{cases}$$

If for some reason  $Eu(X) = 0$  for all  $u \in \mathcal{U}$ , that is,  $\Delta = \{0\}$ , then there is no need to estimate  $\eta(u)$ . As a consequence of Theorem 1, the rate of convergence for estimating the regression function is of order  $\epsilon_n^2$  regardless of the dependence among the errors.

We now consider the rate of convergence of  $r_n$ . Under Assumption 3, the problem of estimating  $\eta \in \Delta$  based on  $Y_i = \eta + \epsilon_i$ ,  $1 \leq i \leq n$  (without  $X_i$ ,  $1 \leq i \leq n$ ) is an easier subproblem with smaller minimax risk than that of estimating  $\eta(u) = Eu(X)$  based on  $(X_i, Y_i)_{i=1}^n$  with  $Y_i = u(X_i) + \epsilon_i$ ,  $1 \leq i \leq n$  (see Lemma 6 in Section 5). That is,  $r_n \geq \tilde{r}_n$ , where  $\tilde{r}_n$  is the minimax mean square error of the easier problem. Since  $\{X_i\}_{i=1}^n$  is not involved,  $\tilde{r}_n$  is handled more easily. Some results on  $\tilde{r}_n$  were given in Hall and Hart (1990b). The following lemma gives useful bounds on  $r_n$  and  $\tilde{r}_n$ . Let  $\mathbf{1}^T = (1, 1, \dots, 1)$  be of dimension  $n$ .

**Lemma 1.** *Under Assumption 4, the minimax risk  $\tilde{r}_n$  satisfies*

$$(\mathbf{1}^T \Omega_n^{-1} \mathbf{1})^{-1} \leq \tilde{r}_n \leq (\mathbf{1}^T \Omega_n \mathbf{1})/n^2.$$

*If  $(\mathbf{1}^T \Omega_n^{-1} \mathbf{1})(\mathbf{1}^T \Omega_n \mathbf{1}) \asymp n^2$  and  $\mathbf{1}^T \Omega_n \mathbf{1} \asymp n$ , then under Assumptions 3 and 4,*

$$r_n \asymp \tilde{r}_n \asymp (\mathbf{1}^T \Omega_n \mathbf{1})/n^2.$$

**Remarks.** (i) The quantity  $(\mathbf{1}^T \Omega_n^{-1} \mathbf{1})^{-1}$  is the variance of the best linear unbiased estimator (BLUE) of  $\eta$  based on  $Y_1, \dots, Y_n$  with  $Y_i = \eta + \epsilon_i$ , where  $\{\epsilon_i, 1 \leq i \leq n\}$  have the covariance matrix  $\Omega_n$ . Adenstedt (1974) showed that for a wide range of stationary error sequences having a spectral density, the minimum variance  $(\mathbf{1}^T \Omega_n^{-1} \mathbf{1})^{-1}$  depends asymptotically only on the behaviour of the spectral density near the origin.

(ii) Note that  $(\mathbf{1}^T \Omega_n \mathbf{1})/n^2$  is the variance of the average of the errors, which determines the rate of convergence of  $\sum_{i=1}^n Y_i/n$  as a simple estimator of  $\eta$ . For the case of long-range dependence, it behaves as well as the BLUE in terms of rate of convergence (see Adenstedt 1974; Samarov and Taqqu 1988). The condition  $(\mathbf{1}^T \Omega_n^{-1} \mathbf{1})(\mathbf{1}^T \Omega_n \mathbf{1}) \asymp n^2$  is to say that BLUE and the simple estimator converge at the same speed (as in the case for the short- and long-range dependent cases in Section 3.3). The condition  $\mathbf{1}^T \Omega_n \mathbf{1} \asymp n$  excludes unusual situations (e.g., independent errors with  $\sigma_i^2 = i^{-1}$ ) where a better rate than  $\epsilon_n^2$  is possible for regression.

(iii) If  $|\sum_{i=1, j=1}^n \text{cov}(\epsilon_i, \epsilon_j)| \leq n$ , then the dependence is weak and  $\mathbf{1}^T \Omega_n \mathbf{1}/n^2 \asymp 1/n$ . As a result,  $r_n \asymp 1/n$  and, from Theorem 1, we have the same rate of convergence for regression estimation as in the case of i.i.d. errors. For another extreme with  $\mathbf{1}^T \Omega_n \mathbf{1} \asymp n^2$  (see Section 3.3, case 5), the minimax risk for estimating the regression function does not converge to zero at all under Assumptions 3 and 4, though the rate remains  $\epsilon_n^2$  for estimating  $u_0$ .

**Theorem 2.** Under Assumptions 1–6, if  $(\mathbf{1}^T \Omega_n^{-1} \mathbf{1})(\mathbf{1}^T \Omega_n \mathbf{1}) \asymp n^2$  and  $\mathbf{1}^T \Omega_n \mathbf{1} \succeq n$  then

$$R(\mathcal{L}; \Omega; n) \asymp (\mathbf{1}^T \Omega_n \mathbf{1})/n^2 + \epsilon_n^2. \tag{7}$$

### 3.2. Rates under long-range dependence

#### 3.2.1. Long-range dependence

Assume that the errors are stationary and that the spectral density, say  $f(\lambda)$ , of the serially correlated errors exists. Let  $r(i)$  denote the correlation between  $\epsilon_j$  and  $\epsilon_{j+i}$ . The error process is said to be long-range dependent if, for some  $c > 0$  and  $0 < \gamma < 1$ ,

$$f(\lambda) \sim c\lambda^{-(1-\gamma)} \quad \text{as } \lambda \rightarrow 0 \tag{8}$$

(see, for example, Cox 1984). Then  $r(j)$  is of order  $|j|^{-\gamma}$ .

**Corollary 1.** Assume that  $f(\lambda)$  satisfies (8), is continuous except at the origin and is bounded away from 0. Under Assumptions 1–4, we have  $\tilde{r}_n \asymp r_n \asymp n^{-\gamma}$  and the minimax rate of convergence for regression estimation is

$$R(\mathcal{L}; \Omega; n) \asymp n^{-\gamma} + \epsilon_n^2.$$

#### 3.2.2. An example with Besov classes

For  $1 \leq \sigma \leq \infty$ ,  $1 \leq q \leq \infty$ , and  $\alpha/d > 1/q - 1/2$ , let  $B_{\sigma,q}^\alpha(C)$  be the collections of all functions  $g \in L_q[0, 1]^d$  such that the Besov norm satisfies  $\|g\|_{B_{\sigma,q}^\alpha} \leq C$  (see, for example, DeVore and Lorentz 1993; Triebel 1975). Then the  $L_2$  metric entropy is of order  $\epsilon^{-d/\alpha}$  (see, for example, Triebel 1975; Lorentz *et al.*, 1996, Chapter 15). Assume the design density  $h(x)$  of  $X$  with respect to Lebesgue measure  $\mu$  is bounded above and away from zero. Then the metric entropy of the Besov class under  $L_2(h)$  distance is of order  $\epsilon^{-d/\alpha}$ . Application of Corollary 1 yields the minimax rate of convergence under the long-range dependence:

$$R(B_{\sigma,q}^\alpha(C); \Omega; n) \asymp n^{-\min(2\alpha/(2\alpha+d), \gamma)}. \tag{9}$$

#### 3.2.3. A comparison with an equally spaced fixed design

Results on minimax rates are obtained for long-range dependent errors with a one-dimensional equally spaced fixed design in Hall and Hart (1990a), Wang (1996) and Johnstone and Silverman (1997) for some concrete smoothness function classes. The model being considered is

$$Y_i = u(i/n) + \epsilon_i, \quad 1 \leq i \leq n,$$

where  $\text{corr}(\epsilon_i, \epsilon_j) \sim c|i - j|^{-\gamma}$  for some  $0 < \gamma < 1$ , and  $u$  is in Besov class  $B_{\sigma,q}^\alpha(C)$  (or a differentiable class in Hall and Hart 1990a). The minimax rate of convergence for estimating  $u$  under squared  $L_2$  loss is shown to be of order  $n^{-2\alpha\gamma/(2\alpha+\gamma)}$ .

Assume there are only measurement errors (independent of the sampling sites  $X_i$ ) in the



responses and the errors are long-range dependent in the order of measurements. For this case, if one uses an equally spaced fixed design, and if the order of measurements corresponds to the order of the sites, the rate of convergence is  $n^{-2\alpha\gamma/(2\alpha+\gamma)}$  from the foregoing. Alternatively, if one uses a random design, from (9), the rate of convergence is  $n^{-\min(2\alpha/(2\alpha+1),\gamma)}$ , which is faster than that with the fixed design. An explanation of the difference in rates is as follows. Under the fixed design, observations with  $x$  values close to each other are highly correlated. With the random design, however, the orders of the measurements of the observations at nearby  $x$  values are not necessarily adjacent but on average quite far away from each other, resulting in weaker correlations between observations that are close in terms of  $x$  values. Thus it is clear that the latter is preferred to the former design. A closer look suggests that the difference in rates is not due to the difference in random and fixed designs, but rather because the order of measurements is not randomized for the fixed design case. If one uses an equally spaced fixed design, one should randomize the order of measurements and we expect the same rate of convergence as under the random design. This example also illustrates importance of the randomization principle in statistical experimental design, as well demonstrated in Künsch *et al.* (1993) under some parametric settings with long-range dependence.

### 3.3. Examples of dependence

For simplicity, we focus on stationary errors.

1. *Exponentially decaying correlation.* Let  $r(j) = \sigma^2\theta^j, j \geq 0$ , for some constants  $\sigma^2 > 0$  and  $\theta$  with  $|\theta| < 1$ . Then it can be shown that  $\text{tr}(\Omega_n^{-1}) \asymp n$  and  $(\mathbf{1}^T \Omega_n \mathbf{1})/n^2$  is of order  $n^{-1}$ .
2. *Short-range dependence.* More generally than in case 1, we assume that the errors are weakly correlated or short-range dependent in the sense that  $\sum_{k=0}^m |r(k)|$  converges as  $m \rightarrow \infty$ . Then  $(\mathbf{1}^T \Omega_n \mathbf{1})/n^2$  is of order  $n^{-1}$ . A special case is finite memory dependence, where the errors are correlated only when they are not far away from each other, that is,  $r(j) = 0$  when  $j \geq j^*$  for some  $j^* > 1$ . Another example is  $r(k) \asymp |k|^{-\gamma}$  with  $\gamma > 1$ .
3. *Long-range dependence.* Assume  $f(\lambda) = f^*(\lambda)|1 - e^{i\lambda}|^{-(1-\gamma)}$  for some  $0 < \gamma < 1$ , where  $f^*(\lambda)$  is a strictly positive continuous function. This includes the spectral density of a fractional Gaussian noise model (Mandelbrot and Van Ness 1968) and a fractional ARIMA model (Granger and Joyeux 1980; Hosking 1981). For the first case,  $r(j) = (c/2)(|j+1|^{2-\gamma} - 2|j|^{2-\gamma} + |j-1|^{2-\gamma})$  (then  $r(j) \sim c'j^{-\gamma}$  for some constant  $c' > 0$ ). The fractional ARIMA( $p, d, q$ ) process has a spectral density  $f(\lambda; d, \phi, \theta) = c|\theta(e^{i\lambda})|^2/|\phi(e^{i\lambda})(1 - e^{i\lambda})^d|^2$ , where  $\theta(z) = 1 - \sum_{j=1}^q \theta_j z^j$  and  $\phi(z) = 1 - \sum_{j=1}^p \phi_j z^j$  are polynomials of order  $q$  and  $p$ , respectively. From Corollary 1,  $r_n \asymp n^{-\gamma}$  (see also Hall and Hart 1990b).
4. *Alternating dependence.* For the above long-range dependence, the errors are eventually positively correlated, that is,  $r(j) > 0$  when  $j$  is large enough. Now suppose  $r(j) \sim c(-1)^j |j|^{-\gamma}$  for some  $\gamma > 0$  as  $j \rightarrow \infty$ . One can obtain such a dependence from long-range dependent errors  $\{\varepsilon_i\}$  by considering  $\{(-1)^i \varepsilon_i\}$ . Then

because the covariances essentially cancel out even when  $0 < \gamma < 1$ , the rate of convergence for estimating  $\eta(u)$  under this correlation is still of order  $1/n$ .

5. *An excessively highly correlated case.* Let  $\Omega_n$  have diagonal elements  $\sigma^2$  and off-diagonal elements  $\sigma^2\theta$ . For  $0 < \theta < 1$ ,  $\Omega_n$  is positive definite for all  $n$ . For this case,  $(\mathbf{1}'\Omega_n\mathbf{1})/n^2 \asymp 1$ , and since  $\mathbf{1}$  is an eigenvector of  $\Omega_n$ , the product  $(\mathbf{1}'\Omega_n^{-1}\mathbf{1})(\mathbf{1}'\Omega_n\mathbf{1})$  is easily seen to be of order  $n^2$  as required in order to apply Theorem 2.

For cases 2–4, it is assumed that the spectral density of the errors is bounded away from 0. Then  $\text{tr}(\Omega_n^{-1}) \asymp n$  (see Lemma 8 in Section 5). Note that the trace condition is automatically satisfied for the other cases.

Take the Besov classes  $B_{\sigma,q}^\alpha(C)$ , for example. Based on Theorem 2, the minimax rate of convergence for estimating  $u$  is  $n^{-2\alpha/(2\alpha+d)}$  for cases 1, 2 and 4, and is worsened to  $n^{-\min(2\alpha/(2\alpha+d),\gamma)}$  for case 3 (as seen in the previous subsection). For case 5, by Theorem 1, the minimax rate for estimating  $u_0$  is still  $n^{-2\alpha/(2\alpha+d)}$ . However, since  $(\mathbf{1}'\Omega_n\mathbf{1})/n^2 \asymp 1$ , the minimax risk for estimating  $u$  does not converge at all.

## 4. A key proposition and its derivation

### 4.1. Minimax upper and lower bounds for regression

Assume that the errors are independent of  $X^n$ . Let  $\rho_n = \text{tr}(\Omega_n^{-1})$ .

Choose  $\tilde{\epsilon}_n$  such that

$$M_2(\tilde{\epsilon}_n) = \left(\frac{1}{2}\right)\rho_n\tilde{\epsilon}_n^2. \quad (10)$$

Let

$$\psi_n = \left(\frac{11}{2}\right)\rho_n\tilde{\epsilon}_n^2 + \log\left(8Ln^{1/2}/\tilde{\epsilon}_n\right)$$

and let  $\underline{\epsilon}_n$  be chosen to satisfy

$$M_2(\underline{\epsilon}_n) = 2\psi_n. \quad (11)$$

Let  $\bar{\epsilon}_n$  satisfy

$$M_2(\bar{\epsilon}_n) = n\bar{\epsilon}_n^2/2, \quad (12)$$

and define

$$\begin{aligned} \bar{\psi}_n &= \left(\frac{11}{2}\right)n\bar{\epsilon}_n^2 + \log\left(8Ln^{1/2}/\bar{\epsilon}_n\right), \\ \psi_n^* &= \min(\psi_n, \bar{\psi}_n). \end{aligned}$$

Typically (e.g., when  $\rho_n$  is of a polynomial order in  $n$ ), the component  $\rho_n\tilde{\epsilon}_n^2$  (or  $n\bar{\epsilon}_n^2$ ) dominates the other term in  $\psi_n$  (or  $\bar{\psi}_n$ ). Then under the richness condition in (3),  $\tilde{\epsilon}_n$  and  $\underline{\epsilon}_n$  are of the same order. If  $\rho_n \asymp n$ , then  $\tilde{\epsilon}_n$ ,  $\underline{\epsilon}_n$ ,  $\bar{\epsilon}_n$ ,  $\psi_n/n$ , and  $\bar{\psi}_n/n$  are all of the same order. They are also of the same order as  $\epsilon_n$  determined by  $M(\epsilon_n) = n\epsilon_n^2$  in (4) with  $M(\epsilon)$  of order  $M_2(\epsilon)$  (see Yang and Barron 1999). Let  $\bar{\sigma}^2 = \sup_{i \geq 1} \sigma_i^2$ .

**Proposition 1.** *Under Assumptions 1 and 5, the minimax squared  $L_2(h)$  risk for regression function estimation is bounded as follows:*

$$\max(\underline{\epsilon}_n^2/8, r_n) \leq R(\mathcal{U}; \Omega; n) \leq r_n + C_{L, \bar{\sigma}^2} \psi_n^*/n,$$

where  $C_{L, \bar{\sigma}^2}$  is a constant depending on  $L$  and  $\bar{\sigma}^2$ .

**Remark.** Without the richness assumption (3), even under  $\rho_n \asymp n$ , the upper and lower bounds in the above proposition may not be of the same order. For example, for classes of analytic functions, the metric entropies are of polynomial orders of  $\log(1/\epsilon)$  (Kolmogorov and Tihomirov 1959) and the upper and lower bounds differ in a logarithmic term unless  $r_n$  dominates. It seems that the use of local entropy (instead of global entropy) as pioneered by Le Cam (1975) and Birgé (1983) in the construction of the upper bound may overcome the gap.

## 4.2. Proof of Proposition 1

In Yang and Barron (1999), minimax rates of convergence for regression under independent Gaussian errors are derived using a connection between density estimation and data compression. The Cesaro average of the Bayes predictive density estimators of the joint distribution of  $(X, Y)$  based on the uniform prior on a suitably chosen  $\epsilon$ -net in the regression function class  $\mathcal{U}$  is used to produce an estimator of the regression function to obtain a minimax upper bound. For regression with dependent errors, however, due to correlations, the Bayes predictive density ‘estimators’ are targeted at the conditional distributions of  $(X_i, Y_i)$ ,  $i \geq 1$ , given the past observations. They are no longer appropriate for estimating the distributions of  $(X_i, Y_i)$ . It becomes much harder to derive a rate-optimal estimator under general conditions on  $\mathcal{U}$  and  $\Omega$ . The difficulty is overcome through rather delicate adjustments of the Bayes predictive estimators, as will be seen.

Let  $Z = (X, Y)$ ,  $z = (x, y)$ ,  $z^n = (z_1, \dots, z_n)$ . Let  $U^n = (u(X_1), \dots, u(X_n))$  and  $u^n = (u(x_1), \dots, u(x_n))$ .

### 4.2.1. Lower bound

We prove  $R(\mathcal{U}; \Omega; n) \geq \underline{\epsilon}_n^2/8$  and  $R(\mathcal{U}; \Omega; n) \geq r_n$  separately. The second inequality follows basically from the observation that estimating the whole regression function is at least as difficult as estimating the mean of the regression function. The proof of the first one utilizes Fano’s inequality together with a suitable upper bound on the involved mutual information.

Let  $N_{\underline{\epsilon}_n}$  be an  $\underline{\epsilon}_n$ -packing set with the maximum cardinality in  $\mathcal{U}$  and let  $G_{\tilde{\epsilon}_n}$  be an  $\tilde{\epsilon}_n$ -net for  $\mathcal{U}$ , both under  $L_2(h)$  distance. Since an  $\epsilon$ -packing set with the maximum cardinality is automatically an  $\epsilon$ -covering set, we can find a  $G_{\tilde{\epsilon}_n}$  such that  $\log|G_{\tilde{\epsilon}_n}| = M_2(\tilde{\epsilon}_n)$ . Following now a standard argument using Fano’s inequality (see, for example, Birgé 1983, Proposition 2.8; Yu 1997, p. 427; Yang and Barron 1999, pp. 1570-1571), we have

$$\min_{\hat{u}} \max_{u \in \mathcal{U}} E_u \|u - \hat{u}\|_{L_2(h)}^2 \geq \frac{\epsilon_n^2}{4} \left( 1 - \frac{(U; Z^n) + \log 2}{\log |N_{\epsilon_n}|} \right),$$

where the Shannon mutual information  $I(U; Z^n)$  is equal to the average (with respect to the uniform prior  $w$ ) of the Kullback–Leibler (KL) divergence between  $p_u(z^n)$  and  $p^w(z^n) = \sum_{u \in N_{\epsilon_n}} p_u(z^n) / |N_{\epsilon_n}|$ . Here

$$p_u(z^n) = \left( \prod_{i=1}^n h(x_i) \right) (2\pi)^{-n/2} |\Omega_n|^{-1/2} \exp(-\frac{1}{2}(y^n - u^n)' \Omega_n^{-1} (y^n - u^n)).$$

Since the Bayes mixture density  $p^w(z^n)$  minimizes the average KL divergence over all choices of joint density  $q(z^n)$  on the sample space  $\mathcal{Z}^n$ , the mutual information is upper bounded by the maximum KL divergence between  $p_u(z^n)$  and any  $q(z^n)$ . That is,

$$I(U; Z^n) \leq \max_{u \in N_{\epsilon_n}} D(P_{Z^n, u} \| Q_{Z^n}).$$

We will choose  $q(z^n) = (1/|G|) \sum_{u \in G} p_u(z^n)$  for a certain appropriate covering set  $G$ .

Key to the analysis is the following expression for the KL divergence between  $P_{Z^n, u}$  and  $P_{Z^n, v}$  (see Lemma 2 in Section 5):

$$D(P_{Z^n, u} \| P_{Z^n, v}) = \frac{1}{2} \rho_n \|u - v\|_{L_2(h)}^2 + \frac{1}{2} \left( \sum_{i \neq j} \omega_{i,j}^{-1} \right) (Eu(X) - Ev(X))^2, \tag{13}$$

where  $\omega_{i,j}^{-1}$  denotes the  $(i, j)$ th element of  $\Omega_n^{-1}$ . When the errors are i.i.d., the second term in the above expression is zero and one can simply take  $G$  to be  $G_{\epsilon_n}$  and obtain the right order upper bound on  $\max_{u \in N_{\epsilon_n}} D(P_{Z^n, u} \| Q_{Z^n})$ , as shown in Yang and Barron (1999). For dependent errors,  $\sum_{i \neq j} \omega_{i,j}^{-1}$  might be large compared to  $\rho_n$  and the choice of  $G_{\epsilon_n}$ , together with the familiar bound  $(Eu(X) - Ev(X))^2 \leq \|u - v\|_{L_2(h)}^2$ , is not sufficient for the result. We instead construct a covering set carefully to handle this term  $(\sum_{i \neq j} \omega_{i,j}^{-1})(Eu(X) - Ev(X))^2$ . The idea is to enlarge  $G_{\epsilon_n}$  slightly by adding constants so that, for each  $u \in \mathcal{U}$ , we can find  $v$  in the enlarged covering set such that both terms in (13) are well behaved. Details are as follows.

Let  $A_n = \{a_1, a_2, \dots, a_m\}$ ,  $a_j = -2L + j\delta\tilde{\epsilon}_n$  be equally spaced points in  $[-2L, 2L]$  with width  $\delta\tilde{\epsilon}_n$  and  $m = \lfloor 4L/(\delta\tilde{\epsilon}_n) \rfloor$  (recall that  $L$  is an upper bound on the supremum norms of functions in  $\mathcal{U}$ ). Let us consider an enlarged net  $\tilde{G}_{\tilde{\epsilon}_n} = \{v + a : v \in G_{\tilde{\epsilon}_n} \text{ and } a \in A_n\}$ . Note that  $\log(|\tilde{G}_{\tilde{\epsilon}_n}|) \leq M_2(\tilde{\epsilon}_n) + \log(4L/(\delta\tilde{\epsilon}_n))$ . For any  $u \in \mathcal{U}$ , there exist  $\tilde{u} \in G_{\tilde{\epsilon}_n}$  and  $a^* \in A_n$  such that  $\|u - \tilde{u}\|_{L_2(h)} \leq \tilde{\epsilon}_n$  and  $|\int(\tilde{u} - u)h d\mu - a^*| \leq \delta\tilde{\epsilon}_n$ . Then  $|a^*| \leq \delta\tilde{\epsilon}_n + |\int(u - \tilde{u})h d\mu| \leq (1 + \delta)\tilde{\epsilon}_n$ . Let  $u^* = \tilde{u} - a^*$ ; then  $|\int(u - u^*)h d\mu| \leq \delta\tilde{\epsilon}_n$ , and  $\|u - u^*\|_{L_2(h)} \leq \|u - \tilde{u}\|_{L_2(h)} + \|u^* - \tilde{u}\|_{L_2(h)} \leq (2 + \delta)\tilde{\epsilon}_n$ . Clearly we have  $u^* \in \tilde{G}_{\tilde{\epsilon}_n}$ . From (13), we have  $D(P_{Z^n, u} \| P_{Z^n, u^*}) \leq \frac{1}{2}(2 + \delta)^2 \rho_n \epsilon_n^2 + \frac{1}{2} \max(0, \varpi_n) \delta^2 \tilde{\epsilon}_n^2$ , where  $\varpi_n = \sum_{i \neq j, 1 \leq i, j \leq n} \omega_{i,j}^{-1}$ . Now choose  $w_1$  to be the uniform prior on  $\tilde{G}_{\tilde{\epsilon}_n}$  and let  $q(z^n) = p^{w_1}(z^n) = \sum_{u \in \tilde{G}_{\tilde{\epsilon}_n}} w_1(u) p_u(z^n)$  and  $Q_{Z^n}$  be the corresponding Bayes mixture density and distribution, respectively. Let  $\lambda_{(1),n} \leq \lambda_{(2),n} \leq \dots \leq \lambda_{(n),n}$  be the eigenvalues of  $\Omega_n$ . Then  $\varpi_n \leq \mathbf{1}^T \Omega_n^{-1} \mathbf{1} \leq n \lambda_{(1),n}^{-1}$ . Since  $\rho_n = \sum_{i=1}^n \lambda_{(i),n}^{-1} \geq \lambda_{(1),n}^{-1}$ , we have  $\varpi_n / \rho_n \leq n$ . From the foregoing we have that, for any  $u \in \mathcal{U}$ ,

$$\begin{aligned}
 D(P_{Z^n, u} \| Q_{Z^n}) &= E \log \frac{p_u(z^n)}{(1/|\tilde{G}_{\tilde{\epsilon}_n}|) \sum_{u'} \in \tilde{G}_{\tilde{\epsilon}_n} p_{u'}(z^n)} \\
 &\leq E \log \frac{p_u(z^n)}{(1/|\tilde{G}_{\tilde{\epsilon}_n}|) p_{u^*}(z^n)} \\
 &= \log |\tilde{G}_{\tilde{\epsilon}_n}| + D(P_{Z^n, u} \| P_{Z^n, u^*}) \\
 &\leq M_2(\tilde{\epsilon}_n) + \log(4L/(\delta\tilde{\epsilon}_n)) + \frac{1}{2}(2 + \delta)^2 \rho_n \tilde{\epsilon}_n^2 + \frac{1}{2} n \rho_n \delta^2 \tilde{\epsilon}_n^2. \tag{14}
 \end{aligned}$$

Taking  $\delta = n^{-\frac{1}{2}}$ , together with our choice of  $\tilde{\epsilon}_n$  in (10), we have

$$D(P_{Z^n, u} \| Q_{Z^n}) \leq \log(4Ln^{1/2}/\tilde{\epsilon}_n) + \frac{11}{2} \rho_n \tilde{\epsilon}_n^2. \tag{15}$$

Thus we have shown that  $I(U; Z^n) \leq \log(4Ln^{1/2}/\tilde{\epsilon}_n) + \frac{11}{2} \rho_n \tilde{\epsilon}_n^2$ . By our choice of  $\tilde{\epsilon}_n$  in (11),  $(I(U; Z^n) + \log 2)/\log |N_{\tilde{\epsilon}_n}| \leq \frac{1}{2}$ . Thus  $\min_{\hat{u}} \max_{u \in \mathcal{U}} E \|u - \hat{u}\|_{L_2(h)}^2 \geq \underline{\epsilon}_n^2/8$ .

The inequality  $R(\mathcal{U}; \Omega; n) \geq r_n$  follows from the simple fact that for any estimator  $\hat{u}$  based on  $Z^n$ , letting  $\hat{\eta} = \int \hat{u} h d\mu$ ,

$$E(\hat{\eta} - \eta)^2 = E \left( \int (\hat{u} - u) h d\mu \right)^2 \leq E \| \hat{u} - u \|_{L_2(h)}^2.$$

#### 4.2.2. Upper bound

We divide the proof of the upper bound into several steps. In step 1, as in the derivation of the lower bound, consider the covering set  $\tilde{G}_{\tilde{\epsilon}_n}$  with uniform prior. We show the resulting Bayes predictive densities (at different sample sizes) are good ‘estimators’ of the conditional densities of the observations  $Z_i$  given the past  $Z^{i-1}$ . The Bayes predictive densities are mixtures of Gaussian densities. In step 2, based on the Bayes predictive densities, we construct density estimators (of the same conditional densities) that have the form of a single Gaussian density (instead of a mixture), still with good risk bounds. Being a single Gaussian density is important in the later construction of the regression estimator. In step 3, the risk bounds on the estimators in Step 2 are shown to imply that the regression function can be estimated well up to a constant. In step 4, the estimation of the constant is shown to be determined by the correlations between the errors. Together with step 3, we have a good estimator of the regression function. In step 5, we consider the case when  $\rho_n$  is of higher order than  $n$ . A suitable modification improves the upper rate of convergence. This is why  $\psi_n^*$  is used instead of  $\psi_n$  in the upper bound in Proposition 1.

*Step 1.* As in the derivation of lower bounds, consider the covering set  $\tilde{G}_{\tilde{\epsilon}_n}$  with uniform prior  $w_1$ . Let the Bayes predictive density estimators be  $\hat{p}_i(z) = p(Z_{i+1}|Z^i)$  evaluated at  $Z_{i+1} = z$ , which equal  $p^{w_1}(Z^i, z)/p^{w_1}(Z^i)$  for  $i > 0$  and  $\hat{p}_i(z) = p^{w_1}(z) = (1/|\tilde{G}_{\tilde{\epsilon}_n}|) \sum_{u \in \tilde{G}_{\tilde{\epsilon}_n}} p_u(z)$  for  $i = 0$ . For  $n \geq 1$ , let

$$\Omega_n = \begin{pmatrix} \Omega_{n-1} & \beta_{n-1} \\ \beta'_{n-1} & \sigma_n^2 \end{pmatrix}$$

be the partition of  $\Omega_n$ . Under the Gaussian assumption, given  $X_{i+1} = x$  and  $(X_j, Y_j)_{j=1}^i, Y_{i+1}$  has a normal distribution with mean  $m_{i,u}(x|Z^i) = u(x) + \beta_i' \Omega_i^{-1}(Y^i - U^i)$  and variance  $\sigma_{i+1}^2 - \beta_i' \Omega_i^{-1} \beta_i$ . Let

$$p_{z_{i+1}|Z^i;u}(x_{i+1}, y_{i+1}) = h(x_{i+1}) \left( 2\pi \left( \sigma_{i+1}^2 - \beta_i' \Omega_i^{-1} \beta_i \right) \right)^{-1/2} \quad (16)$$

$$\times \exp \left( -1/2 \left( \sigma_{i+1}^2 - \beta_i' \Omega_i^{-1} \beta_i \right) (y_{i+1} - m_{i,u}(x_{i+1}|Z^i))^2 \right).$$

This is the conditional density of  $Z_{i+1}$  given  $Z^i$  under the regression function  $u$ . Then by the chain rule (see, for example, Barron 1987), for any  $u \in \mathcal{U}$ ,

$$\sum_{i=0}^{n-1} \mathbb{E} \log \frac{p_{z_{i+1}|Z^i;u}(Z_{i+1})}{\hat{p}_i(Z_{i+1})} = \mathbb{E} \log \frac{p_u(Z^n)}{p^{w_1}(Z^n)} = D(P_{Z^n, u} \| Q_{Z^n}) \leq \psi_n,$$

where the last inequality is as in (15). Thus

$$\max_{u \in \mathcal{U}} \sum_{i=0}^{n-1} \mathbb{E} D(p_{z_{i+1}|Z^i;u} \| \hat{p}_i) \leq \psi_n. \quad (17)$$

Since the squared Hellinger distance satisfies

$$d_{\text{H}}^2(p_1, p_2) = \int \left( p_1^{1/2} - p_2^{1/2} \right)^2 d\mu \leq D(p_1 \| p_2),$$

we have

$$\max_{u \in \mathcal{U}} \sum_{i=0}^{n-1} \mathbb{E} d_{\text{H}}^2(p_{z_{i+1}|Z^i;u}, \hat{p}_i) \leq \psi_n.$$

This means that we can estimate (or predict) well the conditional densities of  $Z_{i+1}$  given  $Z^i$  by  $\hat{p}_i$  in terms of the cumulative squared Hellinger risk.

*Step 2.* Note that  $\hat{p}_i(x_{i+1}, y_{i+1})$  takes the form of  $h(x_{i+1}) \hat{g}_i(y_{i+1}|x_{i+1})$ , where  $\hat{g}_i(y_{i+1}|x_{i+1})$  is an estimator of the conditional density of  $Y_{i+1}$  given  $X_{i+1}$  and  $Z^i$ . It is a mixture of Gaussians using a posterior based on the uniform prior on the  $\epsilon$ -net. We now construct an estimator taking the form of a single Gaussian density. The simplified form (instead of a mixture) is easier to work with in the next step. First, fix  $v^i \in \mathbb{R}^i$ . For given  $(X_j, Y_j)_{j=1}^i$  and  $v^i$ , for each  $x$ , let  $\tilde{m}_i(x) = \tilde{m}_i(x|v^i)$  be the minimizer of the Hellinger distance  $d_{\text{H}}(\hat{g}_i(\cdot|x), \phi_b)$  between  $\hat{g}_i(y|x)$  and the normal density  $\phi_b(y)$  with mean  $b$  and the variance  $\sigma_{i+1}^2 - \beta_i' \Omega_i^{-1} \beta_i$  over choices of  $b$  with  $|b - \beta_i' \Omega_i^{-1}(Y^i - v^i)| \leq L$ . Here  $\tilde{m}_i(x|v^i)$  and  $\tilde{u}_i(x) = \tilde{u}_i(x|v^i) = \tilde{m}_i(x) - \beta_i' \Omega_i^{-1}(Y^i - v^i)$  can be viewed as ‘estimators’ of the conditional mean  $m_{i,u}$  and of  $u$  respectively, based on  $(X_j, Y_j)_{j=1}^i$  except that  $v^i$  is used in place of  $U^i$  (unknown) in the second term of  $u(x) + \beta_i' \Omega_i^{-1}(Y^i - U^i)$ . Denote by  $p_{z_{i+1}|Z^i;v^i}$  the density function of  $(x_{i+1}, y_{i+1})$ :

$$h(x_{i+1}) \left( 2\pi \left( \sigma_{i+1}^2 - \beta_i' \Omega_i^{-1} \beta_i \right) \right)^{-1/2}$$

$$\times \exp \left( -1/2 \left( \sigma_{i+1}^2 - \beta_i' \Omega_i^{-1} \beta_i \right) \left( y_{i+1} - (s(x_{i+1}) + \beta_i' \Omega_i^{-1}(Y^i - v^i)) \right)^2 \right),$$

with given  $Z^i$ , function  $s(x)$ , and  $v^i$ . Let  $v_*^i$  be the minimizer of  $d_{\mathbb{H}}^2(\hat{p}_i, p_{z_{i+1}|Z^i; \bar{u}_i; v^i})$  over  $v^i \in \mathbb{R}^i$  and denote the corresponding  $\tilde{m}_i$  and  $\bar{u}_i$  by  $\tilde{m}_i^*$  and  $\bar{u}_i^*$ . Then, using the triangle inequality,

$$\begin{aligned} d_{\mathbb{H}}^2(p_{z_{i+1}|Z^i; u}, p_{z_{i+1}|Z^i; \bar{u}_i^*; v_*^i}) &\leq 2d_{\mathbb{H}}^2(p_{z_{i+1}|Z^i; u}, \hat{p}_i) + 2d_{\mathbb{H}}^2(p_{z_{i+1}|Z^i; \bar{u}_i^*; v_*^i}, \hat{p}_i) \\ &\leq 2d_{\mathbb{H}}^2(p_{z_{i+1}|Z^i; u}, \hat{p}_i) + 2d_{\mathbb{H}}^2(p_{z_{i+1}|Z^i; \bar{u}_i^0; U^i}, \hat{p}_i) \\ &\leq 2d_{\mathbb{H}}^2(p_{z_{i+1}|Z^i; u}, \hat{p}_i) + 2d_{\mathbb{H}}^2(p_{z_{i+1}|Z^i; u}, \hat{p}_i) \\ &= 4d_{\mathbb{H}}^2(p_{z_{i+1}|Z^i; u}, \hat{p}_i), \end{aligned}$$

where, in the second inequality,  $\bar{u}_i^0$  is  $\bar{u}_i(x|U^i)$  ( $v^i = U^i$ ), and, for the third inequality, we use the fact that

$$d_{\mathbb{H}}^2(p_{z_{i+1}|Z^i; \bar{u}_i^0; U^i}, \hat{p}_i) = \int h(x_{i+1}) d_{\mathbb{H}}^2(\hat{g}_i(\cdot|x_{i+1}), \phi_{\bar{u}_i^0 + \beta_i' \Omega_i^{-1}(Y^i - U^i)}) d\mu$$

is upper-bounded by

$$\int h(x_{i+1}) d_{\mathbb{H}}^2(\hat{g}_i(\cdot|x_{i+1}), \phi_{m_{i,u}}) d\mu = d_{\mathbb{H}}^2(\hat{p}_i, p_{z_{i+1}|Z^i; u}).$$

It follows that

$$\max_{u \in \mathcal{U}} \sum_{i=0}^{n-1} \mathbb{E} d_{\mathbb{H}}^2(p_{z_{i+1}|Z^i; u}, p_{z_{i+1}|Z^i; \bar{u}_i^*; v_*^i}) \leq 4 \max_{u \in \mathcal{U}} \sum_{i=0}^{n-1} \mathbb{E} d_{\mathbb{H}}^2(p_{z_{i+1}|Z^i; u}, \hat{p}_i) \leq 4\psi_n.$$

Thus the estimators  $p_{z_{i+1}|Z^i; \bar{u}_i^*; v_*^i}$  of a simpler form continue to have a good bound on the cumulative Hellinger risk.

*Step 3.* Now note that

$$\begin{aligned} &\mathbb{E} d_{\mathbb{H}}^2(p_{z_{i+1}|Z^i; u}, p_{z_{i+1}|Z^i; \bar{u}_i^*; v_*^i}) \\ &= 2\mathbb{E} \int h(x) \left( 1 - \exp\left(-\left((u(x) - \bar{u}_i^*(x)) - \beta_i' \Omega_i^{-1}(U^i - v_*^i)\right)^2 / (8(\sigma_{i+1}^2 - \beta_i' \Omega_i^{-1} \beta_i))\right)\right) d\mu. \end{aligned}$$

From Lemma 3 in Section 5,

$$\begin{aligned} &\int h(x) \left( 1 - \exp\left(-\left((u(x) - \bar{u}_i^*(x)) - \beta_i' \Omega_i^{-1}(U^i - v_*^i)\right)^2 / (8(\sigma_{i+1}^2 - \beta_i' \Omega_i^{-1} \beta_i))\right)\right) d\mu \\ &\geq c_{L, \bar{\sigma}^2} \int h(x) (u(x) - \bar{u}_i^*(x) - \tau_i)^2 d\mu, \end{aligned}$$

where  $\tau_i = \int h(x)(u(x) - \bar{u}_i^*(x)) d\mu$  and  $c_{L, \bar{\sigma}^2}$  is a constant depending only on  $L$  and  $\bar{\sigma}^2$ . Thus, for any  $u \in \mathcal{U}$ ,

$$\begin{aligned}
 & \sum_{i=0}^{n-1} \mathbb{E} \int h(x)(u(x) - \bar{u}_i^*(x) - \tau_i)^2 d\mu \tag{18} \\
 & \leq c_{L,\bar{\sigma}^2}^{-1} \sum_{i=0}^{n-1} \mathbb{E} \int h(x) \left( 1 - \exp \left( - \left( (u(x) - \bar{u}_i^*(x)) - \beta_i \Omega_i^{-1} (U^i - v_*^i)' \right)^2 / (8(\sigma_{i+1}^2 - \beta_i \Omega_i^{-1} \beta_i')) \right) \right) d\mu \\
 & = c_{L,\bar{\sigma}^2}^{-1} \sum_{i=0}^{n-1} 2^{-1} \mathbb{E} d_{\mathbb{H}}^2 \left( P_{z_{i+1}|Z^i;u}, P_{z_{i+1}|Z^i;\bar{u}_i^*,v_*^i} \right) \\
 & \leq 2c_{L,\bar{\sigma}^2}^{-1} \psi_n.
 \end{aligned}$$

This means that we have obtained a sequence of estimators  $\bar{u}_i^*$  of  $u$  with the variances  $\mathbb{E}(\int h(x)(u(x) - \bar{u}_i^*(x) - \tau_i)^2 d\mu)$  of  $u - \bar{u}_i^*$  well controlled on average. However, a possibly large bias remains. To obtain a final estimator of  $u$ , we estimate the mean  $\eta(u) = \int hu d\mu$  based on current data  $Z^i$ .

Step 4. For any  $\hat{\eta}_i$  based on  $Z^i$ , let  $\hat{u}_i(x) = \bar{u}_i^*(x) - \int \bar{u}_i^*(x)h(x) d\mu + \hat{\eta}_i$ . Then the new estimator satisfies

$$\int h(x)(u(x) - \hat{u}_i(x))^2 d\mu = \int h(x)(u(x) - \bar{u}_i^*(x) - \tau_i)^2 d\mu + (\hat{\eta}_i - \eta(u))^2.$$

It follows that

$$\begin{aligned}
 \sum_{i=0}^{n-1} \mathbb{E} \int h(x)(u(x) - \hat{u}_i(x))^2 d\mu &= \sum_{i=0}^{n-1} \mathbb{E} \int h(x) \left( u(x) - \bar{u}_i^*(x) - \tau_i \right)^2 d\mu + \sum_{i=0}^{n-1} \mathbb{E} (\hat{\eta}_i - \eta(u))^2 \\
 &\leq 2c_{L,\bar{\sigma}^2}^{-1} \psi_n + \sum_{i=0}^{n-1} \mathbb{E} (\hat{\eta}_i - \eta(u))^2.
 \end{aligned}$$

Taking  $\hat{\eta}_i$  to be the minimax estimator of  $\eta$  based on  $Z^i$ , we have

$$\sum_{i=0}^{n-1} \mathbb{E} \int h(x)(u(x) - \hat{u}_i(x))^2 d\mu \leq 2c_{L,\bar{\sigma}^2}^{-1} \psi_n + \sum_{i=0}^{n-1} r_i.$$

Here  $r_0 = \min_{\eta'} \max_{u \in \mathcal{U}} (\eta' - \eta(u))^2$ . As a consequence, we have the cumulative risk bound

$$\frac{1}{n} \sum_{i=0}^{n-1} \mathbb{E} \|u - \hat{u}_i\|_{L_2(h)}^2 \leq 2c_{L,\bar{\sigma}^2}^{-1} \frac{\psi_n}{n} + \bar{r}_n,$$

where  $\bar{r}_n = (1/n) \sum_{i=0}^{n-1} r_i$ . For the usual risk  $R(\mathcal{U}; \Omega; n)$ , we do not need to require  $\hat{\eta}_i$  to depend only on  $Z^i$ . Then we set  $\hat{\eta}_i = \hat{\eta}_n$  for all  $1 \leq i < n$ , where  $\hat{\eta}_n$  is the minimax estimator based on  $Z^n$ . Then the above risk bound becomes  $2c_{L,\bar{\sigma}^2}^{-1} \psi_n/n + r_n$ . From Lemma 4 in Section 5, we have an estimator  $\hat{u}_n$  based on  $Z^n$  such that

$$\max_{u \in \mathcal{U}} \mathbb{E} \|u - \hat{u}_n\|_{L_2(h)}^2 \leq \max_{u \in \mathcal{U}} \mathbb{E} \sum_{i=0}^{n-1} \|u - \hat{u}_i\|_{L_2(h)}^2 \leq 2c_{L,\bar{\sigma}^2}^{-1} \frac{\psi_n}{n} + r_n. \tag{4.19}$$



*Step 5.* When  $\rho_n$  is of higher order than  $n$ , the upper bound above may be suboptimal. For instance, suppose we have independent errors with  $\sigma_i^2 = i^{1-\delta}$  for some  $1 < \delta < 2$ , which implies that  $\rho_n \asymp n^\delta$ . Assume that  $M_2(\epsilon) \asymp \epsilon^{-d/\alpha}$  for some  $\alpha > 0$ . Then the upper bound rate given in terms of  $\psi_n$  is  $n^{\delta-1-2\alpha\delta/(2\alpha+d)}$ , which is worse than the rate  $n^{-2\alpha/(2\alpha+d)}$  obtained with i.i.d. errors. Clearly, this inferior rate is not because the problem is more difficult. It can be improved in general as follows. Let us generate i.i.d. random variables  $\tilde{\epsilon}_1, \tilde{\epsilon}_2, \dots, \tilde{\epsilon}_n$  from a standard normal distribution. Let  $\tilde{Y}_i = Y_i + \tilde{\epsilon}_i$ ,  $1 \leq i \leq n$ . Then the random errors  $\epsilon_i + \tilde{\epsilon}_i$  in  $\tilde{Y}_i$  have covariance matrix  $\tilde{\Omega}_n = I_n + \Omega_n$  ( $I_n$  is the  $n \times n$  identity matrix). Then  $\tilde{\rho}_n = \text{tr}(\tilde{\Omega}_n^{-1}) \leq \text{tr}(I_n^{-1}) = n$  because  $I_n + \Omega_n \geq I_n$  implies  $(I_n + \Omega_n)^{-1} \leq I_n^{-1}$  (here the symbol ' $\geq$ ' for matrix comparison means the difference is non-negative definite). Note also that the variances of the new errors  $\epsilon_i + \tilde{\epsilon}_i$  are uniformly upper-bounded by  $\tilde{\sigma}^2 = \bar{\sigma}^2 + 1$ . Applying similar analysis to  $(X_i, \tilde{Y}_i)$  replacing  $\rho_n$  by  $n$  yields

$$\sum_{i=0}^{n-1} \mathbb{E} \int h(x)(u(x) - \tilde{u}_i^*(x) - \tau_i)^2 d\mu \leq 2c_{L, \tilde{\sigma}^2}^{-1} \tilde{\psi}_n, \tag{4.20}$$

where the  $\tilde{u}_i^*$  are obtained with the new data  $(X_j, \tilde{Y}_j)_{j=1}^i$ . Estimating  $\eta(u)$  the same way as before, we obtain a randomized estimator  $\hat{u}_n$  with risk bounded by  $2c_{L, \tilde{\sigma}^2}^{-1} \tilde{\psi}_n/n + r_n$ . The estimator depends on both  $Z^n$  and the generated random variables  $\tilde{\epsilon}_i$ ,  $1 \leq i \leq n$ . One could average out the randomness in  $\tilde{\epsilon}_i$  to obtain a non-randomized estimator with no bigger risk since the loss being considered is convex. Thus  $R(\mathcal{B}; \Omega; n) \leq 2c_{L, \tilde{\sigma}^2}^{-1} \tilde{\psi}_n/n + r_n$ . This completes the proof of Proposition 1.

## 5. Proofs

### 5.1. Main results

**Proof of Lemma 1.** For the upper rate on  $\tilde{r}_n$ , taking  $\hat{\eta} = \sum_{i=1}^n Y_i/n$ , we obtain  $\tilde{r}_n \leq (\mathbf{1}^T \Omega_n \mathbf{1})/n^2$ . For the lower bound, consider  $2^m$  equally spaced points in  $\Delta_n = [a_n, b_n] \subset \Delta$ . Denote the set of these points by  $D_n$  and let  $\Theta$  take values in  $D_n$  with equal probability. Let  $\delta_n = (b_n - a_n)2^{-m}$ . Then as in the proof of Proposition 1, we have

$$\tilde{r}_n \geq \frac{\delta_n^2}{4} \left( 1 - \frac{I(\Theta; Y^n) + \log 2}{m \log 2} \right).$$

Similarly to the analysis there, consider a rougher net in  $\Delta_n$ . Let  $D'_n$  be the set of  $2^{m'}$  equally spaced points in  $\Delta_n$  and let  $\delta'_n = (b_n - a_n)2^{-m'}$ . Then it can be shown similarly that  $I(\Theta; Y^n) \leq m' \log 2 + (1/2)(\delta'_n)^2 (\mathbf{1}^T \Omega_n^{-1} \mathbf{1})$ . Take  $b_n - a_n$  of order  $(\mathbf{1}^T \Omega_n^{-1} \mathbf{1})^{-1/2}$  and  $m' = 1$  to have  $I(\Theta; Y^n) \leq 1$  (note that  $(\mathbf{1}^T \Omega_n^{-1} \mathbf{1})^{-1}$  is the variance of the BLUE and thus  $(\mathbf{1}^T \Omega_n^{-1} \mathbf{1})^{-1} \leq 1$ ). Thus there exists a constant  $C$  such that  $I(\Theta; Y^n) \leq C$  for all  $n$ . Take  $m$  suitably large (independent of  $n$ ) such that  $(C + \log 2)/(m \log 2) \leq \frac{1}{2}$ . Then  $\tilde{r}_n \geq \delta_n^2/8$ . This establishes the lower bound rate  $(\mathbf{1}^T \Omega_n^{-1} \mathbf{1})^{-1}$ .

For an upper bound on  $r_n$  in the second statement, consider  $\hat{\eta}_n = \bar{Y} = (1/n) \sum_{j=1}^n Y_j$ . Then

$$\begin{aligned}
 \mathbb{E}(\hat{\eta}_n - \eta(u))^2 &= \mathbb{E}\left(\frac{1}{n} \sum_{i=1}^n (u(X_i) - \eta(u)) + \frac{1}{n} \sum_{i=1}^n \varepsilon_i\right)^2 \\
 &= \mathbb{E}\left(\frac{1}{n} \sum_{i=1}^n (u(X_i) - \eta(u))\right)^2 + \mathbb{E}\left(\frac{1}{n} \sum_{i=1}^n \varepsilon_i\right)^2 \\
 &= \frac{1}{n} \int (u(x) - \eta(u))^2 h(x) \, d\mu + \frac{\mathbf{1}^\top \Omega_n \mathbf{1}}{n^2} \\
 &\leq \frac{4L^2}{n} + \frac{\mathbf{1}^\top \Omega_n \mathbf{1}}{n^2}.
 \end{aligned}$$

Under the given conditions, together with Lemma 6 later in this section, we have

$$\tilde{r}_n \leq r_n \leq (\mathbf{1}^\top \Omega_n \mathbf{1})/n^2.$$

If  $(\mathbf{1}^\top \Omega_n^{-1} \mathbf{1})(\mathbf{1}^\top \Omega_n \mathbf{1}) \asymp n^2$ , then clearly  $\tilde{r}_n \asymp r_n \asymp (\mathbf{1}^\top \Omega_n \mathbf{1})/n^2$ . This completes the proof of Lemma 1.  $\square$

**Proof of Theorem 1.** The upper bound part for the first conclusion follows from (20) in the proof of Proposition 1 using  $\hat{u}_0 = (1/n) \sum_{i=0}^{n-1} (\bar{u}_i^*(x) - \int \bar{u}_i^*(x)h(x) \, d\mu)$  as an estimator of  $u_0$ . From (20) and using Lemma 4, we have that

$$\begin{aligned}
 \mathbb{E} \int h(x)(u_0(x) - \hat{u}_0(x))^2 \, d\mu &\leq \frac{1}{n} \sum_{i=0}^{n-1} \mathbb{E} \int h(x) \left( u_0(x) - \left( \bar{u}_i^*(x) - \int \bar{u}_i^*(x)h(x) \, d\mu \right) \right)^2 \, d\mu \\
 &\leq 2c_{L, \bar{\sigma}^2}^{-1} \bar{\psi}_n \leq \epsilon_n^2.
 \end{aligned}$$

Note that Assumption 6 is not needed for the above upper rate of convergence for estimating  $u_0$ .

To prove  $\epsilon_n^2$  is also a lower rate for  $R_0(\mathcal{U}; \Omega; n)$ , consider the distance  $d_0$  defined as  $d_0(u, v) = \int (u_0 - v_0)^2 h \, d\mu$ , where  $u_0 = u - \int uh \, d\mu$  and  $v_0 = v - \int vh \, d\mu$ . Replacing  $L_2(h)$  distance by  $d_0$  in the derivation of the lower bound in the proof of Proposition 1, we have

$$R_0(\mathcal{U}; \Omega; n) \geq \eta_n^2/8,$$

where  $\eta_n$  is determined by  $M_0(\eta_n) = 2\psi_n$  with  $M_0(\epsilon)$  being the packing entropy of  $\mathcal{U}$  under  $d_0$ . It is straightforward to show that  $M_0(\epsilon)$  is of the same order as  $M_2(\epsilon)$  for a uniformly bounded rich class. As a consequence, under Assumption 6,  $\eta_n \asymp \epsilon_n$ .

The second conclusion in (6) follows directly from Proposition 1 using that  $\underline{\epsilon}_n^2$  and  $\psi_n^*/n$  are both of order  $\epsilon_n^2$  under the condition  $\text{tr}(\Omega_n^{-1}) \asymp n$ . Note that the upper bound in Proposition 1 always satisfies  $\psi_n^*/n \leq \epsilon_n^2$ , regardless of the trace condition. This completes the proof of Theorem 1.  $\square$

**Proof of Theorem 2.** The conclusion follows directly from Theorem 1 and Lemma 1.  $\square$

**Proof of Corollary 1.** Assumption 5 is obviously satisfied. From Lemma 8 later in this section, Assumption 6 is satisfied. It remains to verify  $(\mathbf{1}^\top \Omega_n^{-1} \mathbf{1})(\mathbf{1}^\top \Omega_n \mathbf{1}) \asymp n^2$ ,  $(\mathbf{1}^\top \Omega_n \mathbf{1}) \succeq n$  and  $(\mathbf{1}^\top \Omega_n \mathbf{1})/n^2 \asymp n^{-\gamma}$ . Since  $r(j) \sim |j|^{-\gamma}$ , it is straightforward to show that  $(\mathbf{1}^\top \Omega_n \mathbf{1}) \asymp n^{2-\gamma}$ . Under our assumptions on the spectral density, Adenstedt (1974, Theorem 5.2) shows that  $(\mathbf{1}^\top \Omega_n^{-1} \mathbf{1})^{-1}$  is of order  $n^{-\gamma}$  – note that  $(\mathbf{1}^\top \Omega_n^{-1} \mathbf{1})^{-1}$  is the variance of the BLUE. This completes the proof of Corollary 1.  $\square$

## 5.2. Technical lemmas

Let  $P_{Z^n, u}$  denote the distribution of  $Z^n = (X_i, Y_i)_{i=1}^n$  when the regression function is  $u$ . The density of  $P_{Z^n, u}$  is

$$p_u(z^n) = \left( \prod_{i=1}^n h(x_i) \right) (2\pi)^{-n/2} |\Omega_n|^{-1/2} \exp\left(-\frac{1}{2}(y^n - u^n)' \Omega_n^{-1} (y^n - u^n)\right).$$

Let  $\omega_{i,j}^{-1}$  denote the  $(i, j)$ th element of  $\Omega_n^{-1}$ . Recall that the Kullback–Leibler divergence  $D(P\|Q)$  between two distributions  $P$  and  $Q$  with densities  $p$  and  $q$  with respect to  $\mu$  is defined as  $D(P\|Q) = \int p \log(p/q) d\mu$ .

**Lemma 2.** *The KL divergence between  $P_{Z^n, u}$  and  $P_{Z^n, v}$  is*

$$D(P_{Z^n, u}\|P_{Z^n, v}) = \frac{1}{2} \text{tr}(\Omega_n^{-1}) \|u - v\|_{L_2(h)}^2 + \frac{1}{2} \left( \sum_{i \neq j} \omega_{i,j}^{-1} \right) (\mathbb{E}u - \mathbb{E}v)^2.$$

**Proof.** We have

$$2 \log \frac{p_u(z^n)}{p_v(z^n)} = 2(u^n - v^n)' \Omega_n^{-1} y^n - (u^n)' \Omega_n^{-1} u^n + (v^n)' \Omega_n^{-1} v^n.$$

Given  $X^n$ ,

$$\begin{aligned} 2\mathbb{E}_{Z^n|X^n; u} \log \frac{p_u(Z^n)}{p_v(Z^n)} &= 2(u^n - v^n)' \Omega_n^{-1} (u^n)' - u^n' \Omega_n^{-1} (u^n)' + v^n' \Omega_n^{-1} (v^n)' \quad (21) \\ &= (u^n - v^n)' \Omega_n^{-1} (u^n - v^n)'. \end{aligned}$$

Then

$$\begin{aligned} 2\mathbb{E}_{Z^n, u} \log \frac{p_u(Z^n)}{p_v(Z^n)} &= \mathbb{E} \left( \sum_{i,j} \omega_{i,j}^{-1} (u(X_i) - v(X_i))(u(X_j) - v(X_j)) \right) \\ &= \sum_{i=1}^n \omega_{i,i}^{-1} \|u - v\|_{L_2(h)}^2 + \sum_{i \neq j} \omega_{i,j}^{-1} \mathbb{E}(u(X_i) - v(X_i))(u(X_j) - v(X_j)). \end{aligned}$$

Under the i.i.d. assumption on  $X_1, \dots, X_n$ , we have

$$2E_{Z^n, u} \log \frac{p_u(Z^n)}{p_v(Z^n)} = \sum_{i=1}^n \omega_{i,i}^{-1} \|u - v\|_{L_2(h)}^2 + \left( \sum_{i \neq j} \omega_{i,j}^{-1} \right) (E(u(X) - v(X)))^2.$$

This completes the proof of Lemma 2. □

**Lemma 3.** Assume  $\sup_x |g(x)| \leq A$  for some constant  $A$  and  $\sigma^2 \leq \sigma_0^2$ . Let  $h(x)$  be a probability density function. Then

$$\min_{\theta \in \mathbb{R}} \int h(x) \left( 1 - e^{-(g(x)-\theta)^2/\sigma^2} \right) d\mu \geq c \int h(x) \left( g(x) - \int h(x)g(x) d\mu \right)^2 d\mu,$$

where the constant  $c$  depends only on  $A$  and  $\sigma_0^2$ .

**Proof.** It is easy to prove that, for  $|g| \leq A$ ,

$$1 - e^{-(g(x)-\theta)^2/\sigma^2} \geq \begin{cases} c(g - \theta)^2, & |\theta| \leq 2A, \\ cg^2, & |\theta| > 2A, \end{cases}$$

for some constant  $c$  depending only on  $A$  and  $\sigma_0^2$ . It follows that

$$\int h(x) \left( 1 - e^{-(g-\theta)^2/\sigma^2} \right) d\mu \geq \begin{cases} c \int h(x)(g(x) - \theta)^2 d\mu, & |\theta| \leq 2A, \\ c \int h(x)g(x)^2 d\mu, & |\theta| > 2A. \end{cases}$$

Since  $\int h(x)(g(x) - a)^2 d\mu$  is minimized when  $a = \int h(x)g(x) d\mu$ , the lemma follows. □

**Lemma 4.** Let  $\hat{u}_1, \dots, \hat{u}_k$  be  $k$  estimators of  $u$ . Then the estimator  $\hat{u}_k = (1/k) \sum_{i=1}^k \hat{u}_i$  satisfies

$$E \|u - \hat{u}_k\|_{L_2(h)}^2 \leq \frac{1}{k} \sum_{i=1}^k E \|u - \hat{u}_i\|_{L_2(h)}^2.$$

**Proof.** The result follows from the fact that  $\|u - v\|_{L_2(h)}^2$  is convex in  $v$ . □

**Lemma 5.** Let  $\Omega_n$  be the  $n \times n$  finite section of the covariance matrix of a stationary process. Assume  $\Omega_n$  is invertible for  $n \geq 1$ . Then  $\text{tr}(\Omega_n^{-1})$  is at least of order  $n$ . More generally, if  $\sup_{i \geq 1} \sigma_i^2 < \infty$ , then  $\text{tr}(\Omega_n^{-1}) \succeq n$ .

**Proof.** Let  $\Omega_n$  be as in step 1 of the proof of the upper bound of Proposition 1. Then simple linear algebra gives

$$\Omega_n^{-1} = \begin{pmatrix} \Omega_{n-1}^{-1} + (\sigma_n^2 - \beta'_{n-1} \Omega_{n-1}^{-1} \beta_{n-1})^{-1} \Omega_{n-1}^{-1} \beta_{n-1} \beta'_{n-1} \Omega_{n-1}^{-1} & -(\sigma_n^2 - \beta'_{n-1} \Omega_{n-1}^{-1} \beta_{n-1})^{-1} \Omega_{n-1}^{-1} \beta_{n-1} \\ -(\sigma_n^2 - \beta_{n-1} \Omega_{n-1}^{-1} \beta'_{n-1})^{-1} \beta'_{n-1} \Omega_{n-1}^{-1} & (\sigma_n^2 - \beta'_{n-1} \Omega_{n-1}^{-1} \beta_{n-1})^{-1} \end{pmatrix}.$$

It follows that

$$\text{tr}(\Omega_n^{-1}) \geq \text{tr}(\Omega_{n-1}^{-1}) + \left(\sigma_n^2 - \beta_{n-1}' \Omega_{n-1}^{-1} \beta_{n-1}\right)^{-1} \geq \text{tr}(\Omega_{n-1}^{-1}) + \sigma_n^{-2}.$$

The conclusion follows by induction. □

Let  $r_n$  be defined as in (2) and let  $\tilde{r}_n = \min_{\hat{\eta}} \max_{\eta \in \Delta} E(\hat{\eta} - \eta)^2$  be the minimax risk for estimating  $\eta$  based on  $(Y_i)_{i=1}^n$  under the model  $Y_i = \eta + \varepsilon_i$ ,  $1 \leq i \leq n$ .

**Lemma 6.** *Under Assumption 3, we have  $r_n \geq \tilde{r}_n$ .*

**Proof.** Under Assumption 3,  $r_n$  decreases when  $u \in \mathcal{U}$  is instead restricted to the set of constant functions  $\{\eta, \eta \in \Delta\}$ . For the restricted model, it is easy to see by the factorization theorem that  $(Y_1, \dots, Y_n)$  is a sufficient statistic for  $\eta$ . Then for any estimator  $\hat{\eta}$  based on  $(X_i, Y_i)_{i=1}^n$ , we may take  $\hat{\eta} = E(\hat{\eta} | Y_1, \dots, Y_n)$  to obtain an estimator based only on  $Y_1, \dots, Y_n$  with no bigger mean squared error. The conclusion follows. □

The following two lemmas give sufficient conditions for  $\text{tr}(\Omega_n^{-1}) \asymp n$  as used in Section 3.

**Lemma 7.** *Assume that  $\sup_i \sigma_i^2 < \infty$  and that  $\Omega_n$  can be expressed as the sum of two components  $\Omega_n = \Omega_n^{(1)} + \Omega_n^{(2)}$ , where  $\Omega_n^{(1)} = \text{diag}(\omega_{1,n}, \dots, \omega_{n,n})$  with  $\min_{1 \leq i \leq n} \omega_{i,n} \geq c > 0$  for some constant  $c > 0$  independent of  $n$ , and  $\Omega_n^{(2)}$  is non-negative definite. Then  $\text{tr}(\Omega_n^{-1}) \asymp n$ .*

**Proof.** By Lemma 5, under the condition  $\sup_i \sigma_i^2 < \infty$ , we have  $\text{tr}(\Omega_n^{-1}) \geq n$ . Under the other condition, we have  $\Omega_n \geq \Omega_n^{(1)}$  and hence  $\Omega_n^{-1} \leq (\Omega_n^{(1)})^{-1}$ . So  $\text{tr}(\Omega_n^{-1}) \leq \text{tr}(\Omega_n^{(1)})^{-1} \leq n$ . This completes the proof of Lemma 7. □

**Lemma 8.** *For stationary serially correlated errors with spectral density bounded away from zero,  $\text{tr}(\Omega_n^{-1}) \asymp n$ .*

**Proof.** From Lemma 5,  $\text{tr}(\Omega_n^{-1})$  is at least of order  $n$ . From Grenander and Szegö (1958, p. 64), the minimum eigenvalue of  $\Omega_n$  is uniformly bounded away from zero for  $n \geq 1$ . Since  $\text{tr}(\Omega_n^{-1})$  is the sum of the reciprocals of the eigenvalues of  $\Omega_n$ , we have  $\text{tr}(\Omega_n^{-1}) \leq n$ . This completes the proof of Lemma 8. □

## Acknowledgements

I am very grateful to Andrew Barron for his comments, suggestions and encouragement. I thank Wayne Fuller, Cun-Hui Zhang, Zhiliang Ying, and Nicolas Hengartner for helpful discussions. I also thank the reviewers and the editor for their valuable comments to improve the presentation.

## References

- Adenstedt, R.K. (1974) On large sample estimation for the mean of a stationary sequence. *Ann. Statist.*, **2**, 1095–1107.
- Barron, A.R. (1987) Are Bayes rules consistent in information? In T.M. Cover and B. Gopinath (eds), *Open Problems in Communication and Computation*, pp. 85–91. Berlin: Springer-Verlag.
- Beran, J. (1986) Estimation, testing and prediction for self-similar and related processes. Doctoral thesis, ETH, Zurich.
- Beran, J. (1994) *Statistics for Long-Memory Processes*. London: Chapman & Hall.
- Birgé, L. (1983) Approximation dans les espaces métriques et théorie de l'estimation. *Z. Wahrscheinlichkeitstheorie Verw. Geb.*, **65**, 181–237.
- Birgé, L. (1986) On estimating a density using Hellinger distance and some other strange facts. *Probab. Theory Related Fields*, **71**, 271–291.
- Bretagnolle, J. and Huber, C. (1979) Estimation des densités: risque minimax. *Z. Wahrscheinlichkeitstheorie Verw. Geb.*, **47**, 119–137.
- Cox, D.R. (1984) Long-range dependence: a review. In H.A. David and H.T. David (eds), *Statistics: An Appraisal. Proceedings of a Conference Marking the 50th Anniversary of the Statistical Laboratory, Iowa State University*, pp. 55–74. Ames: Iowa State University Press.
- Dahlhaus, R. (1989) Efficient parameter estimation for self-similar processes. *Ann. Statist.*, **17**, 1749–1766.
- DeVore, R.A. and Lorentz, G.G. (1993) *Constructive Approximation*. Berlin: Springer-Verlag.
- Efromovich, S. (1999) How to overcome the curse of long-memory errors. *IEEE Trans. Inform. Theory*, **45**, 1735–1741.
- Fox, R. and Taqqu, M.S. (1985) Non-central limit theorems for quadratic forms in random variables having long-range dependence. *Ann. Probab.*, **13**, 428–446.
- Giraitis, L. and Surgailis, D. (1990) A central limit theorem for quadratic forms in strongly dependent linear variables and application to asymptotical normality of Whittle's estimate. *Probab. Theory Related Fields*, **86**, 87–104.
- Granger, C.W.J. and Joyeux, R. (1980) An introduction to long-range time series models and fractional differencing. *J. Time Ser. Anal.*, **1**, 15–30.
- Grenander, U. and Szegő, G. (1958) *Toeplitz Forms and Their Applications*. Berkeley: University of California Press.
- Hall, P. and Hart, J.D. (1990a) Nonparametric regression with long-range dependence. *Stochastic Process. Appl.*, **36**, 339–351.
- Hall, P. and Hart, J.D. (1990b) Convergence rates in density estimation for data from infinite-order moving average processes. *Probab. Theory Related Fields*, **87**, 253–274.
- Hosking, J.R.M. (1981) Fractional differencing. *Biometrika*, **68**, 165–176.
- Ibragimov, I.A. and Hasminskii, R.Z. (1977) On the estimation of an infinite-dimensional parameter in Gaussian white noise. *Soviet Math. Dokl.*, **18**, 1307–1309.
- Johnstone, I.M. and Silverman, B.W. (1997) Wavelet threshold estimators for data with correlated noise. *J. Roy. Statist. Soc. Ser. B*, **59**, 319–351.
- Kolmogorov, A.N. and Tihomirov, V.M. (1959)  $\epsilon$ -entropy and  $\epsilon$ -capacity of sets in functional spaces. *Uspekhi Mat. Nauk*, **14**, 3–86. English translation (1961): S.N. Cernikov *et al.*, *Twelve Papers on Algebra and Real Functions*, Amer. Math. Soc. Transl. (2), **17**, pp. 277–364. Providence, RI: American Mathematical Society.
- Künsch, H., Beran, J. and Hampel, F. (1993) Contrasts under long-range correlations. *Ann. Statist.*, **21**, 943–964.

- Le Cam, L.M. (1975) On local and global properties in the theory of asymptotic normality of experiments. In M. Puri (ed.), *Stochastic Processes and Related Topics*, Vol. 1, pp. 13–54. New York: Academic Press.
- Le Cam, L.M. (1986) *Asymptotic Methods in Statistical Decision Theory*. Berlin: Springer-Verlag.
- Lorentz, G.G., Golitschek, M. v. and Makovoz, Y. (1996) *Constructive Approximation: Advanced Problems*. Berlin: Springer-Verlag.
- Mandelbrot, B.B. and Van Ness, J.W. (1968) Fractional Brownian motions, fractional noises and applications. *SIAM Rev.*, **10**, 422–437.
- Robinson, P.M. (1995) Gaussian semiparametric estimation of long-range dependence. *Ann. Statist.*, **23**, 1630–1661.
- Robinson, P.M. (1997) Large-sample inference for nonparametric regression with dependent errors. *Ann. Statist.*, **25**, 2054–2083.
- Samarov, A. and Taqqu, M.S. (1988) On the efficiency of the sample mean in long-memory noise. *J. Time Ser. Anal.*, **9**, 191–200.
- Triebel, H. (1975) Interpolation properties of  $\epsilon$ -entropy and diameters. Geometric characteristics of embedding for function spaces of Sobolev–Besov type. *Mat. Sb.*, **98**, 27–41. English translation (1977): *Math. USSR Sb.*, **27**, 23–37.
- Wang, Y. (1996) Function estimation via wavelet shrinkage for long-memory data. *Ann. Statist.*, **24**, 466–484.
- Yang, Y. and Barron, A.R. (1999) Information-theoretic determination of minimax rates of convergence. *Ann. Statist.*, **27**, 1564–1599.
- Yajima, Y. (1991) Asymptotic properties of LSE in a regression model with long-memory stationary errors. *Ann. Statist.*, **19**, 158–177.
- Yatracos, Y.G. (1988) A lower bound on the error in nonparametric regression type problems. *Ann. Statist.*, **16**, 1180–1187.
- Yu, B. (1997) Assouad, Fano, and Le Cam. In D. Pollard, E. Torgersen and G.L. Yang (eds), *Research Papers in Probability and Statistics: Festschrift in Honor of Lucien Le Cam*. New York: Springer-Verlag.

Received March 1999 and revised April 2001