

## NONASYMPTOTIC BOUNDS FOR AUTOREGRESSIVE TIME SERIES MODELING

BY ALEXANDER GOLDENSHLUGER AND ASSAF ZEEVI<sup>1</sup>

*University of Haifa and Stanford University*

The subject of this paper is autoregressive (AR) modeling of a stationary, Gaussian discrete time process, based on a finite sequence of observations. The process is assumed to admit an  $AR(\infty)$  representation with exponentially decaying coefficients. We adopt the nonparametric minimax framework and study how well the process can be approximated by a finite-order AR model. A lower bound on the accuracy of AR approximations is derived, and a nonasymptotic upper bound on the accuracy of the regularized least squares estimator is established. It is shown that with a “proper” choice of the model order, this estimator is minimax optimal in order. These considerations lead also to a nonasymptotic upper bound on the mean squared error of the associated one-step predictor. A numerical study compares the common model selection procedures to the minimax optimal order choice.

**1. Introduction.** The standard methods for estimating parameters of time series are based on the assumption that the observations come from an autoregressive (AR), moving average (MA), or mixed (ARMA) model of known orders. This assumption can rarely be justified in practice, and the less stringent assumption is that the time series data are observations from a linear stationary process. A common approach to modeling linear stationary processes is based on an AR approximation. In this framework a finite order AR model is fitted to the observations. The order of the AR model should provide an “optimal” finite AR approximation to the process, and it is usually chosen by selection procedures based on the data. This nonparametric AR approach to modeling linear stationary processes has been investigated by Shibata (1980), Bhansali (1981, 1986), An, Chen and Hannan (1982) and Hannan and Kavalieris (1986).

Shibata (1980) considered the problem of predicting a Gaussian infinite-order AR process by fitting a finite AR model. The notion of optimality for the model selection procedure proposed by Shibata (1980) is based on an asymptotic lower bound on the mean squared prediction error. Specifically, the procedure is asymptotically efficient if it attains the lower bound asymptotically. Shibata (1980) also established that the final prediction error (FPE) [Akaike (1970)] and the AIC [Akaike (1974)] criteria are asymptotically efficient in the above sense, provided that the linear process does not degenerate to a finite order autoregression. A similar result has been obtained

---

Received March 1998; revised November 2000.

<sup>1</sup>Supported in part by the NSF.

AMS 2000 *subject classifications*. Primary 62G05, 62M10, 62M20.

*Key words and phrases*. Autoregressive approximation, minimax risk, rates of convergence.

by Bhansali (1986) for the AR transfer function criterion (CAT) proposed by Parzen (1974).

Another motivation for fitting an AR model is the estimation of the spectral density function. Berk (1974) used AR approximation to estimate the spectral density of a linear process. It was shown there that the order of the approximating AR model should increase with the number of observations to ensure the consistency of the associated spectral density estimator. Shibata (1981) suggested another definition of selection procedures optimality which is based on an asymptotic lower bound for the relative integrated squared error in estimating the spectral density function. It was shown there that the FPE and AIC criteria are asymptotically efficient in this sense, provided that the linear process does not degenerate to a finite-order autoregression. Similar results for the CAT criterion have been obtained by Bhansali (1986). Some recent results on AR approximation can be found in Gerencsér (1992) and Bülmann (1995).

In spite of the fact that the FPE, AIC, and CAT criteria are asymptotically efficient as described above, the finite sample behavior of these selection procedures is not so clear. The definitions of optimality adopted in Shibata (1980, 1981) and Bhansali (1986) are essentially asymptotic. The assumption that the underlying linear process does not degenerate to a finite autoregression is also based on asymptotic considerations. If this assumption is violated, the AIC and FPE overestimate the true model order, and a different penalty term is called for. In particular, by penalizing each parameter by a factor of  $\ln n$ , with  $n$  being the sample size, one obtains the minimum description length (MDL) principle of Rissanen (1983), and the BIC criterion of Schwarz (1978). These criteria lead to consistent estimation of the model dimension in the case of an underlying finite-order autoregression. However, if the underlying process does not degenerate to a finite autoregression, they are not asymptotically efficient in the aforementioned sense [cf. the discussion in Shibata (1980), page 161]. Moreover, even if the underlying “true” linear process does not degenerate to a finite-order autoregression, the coefficients in its  $\text{AR}(\infty)$  representation can be small. In these situations, effectively the model is “close” to being finite-dimensional, and the behavior of the asymptotic efficient procedures can be poor even for “large” sample sizes. Several interesting questions arise in this context. Given a fixed number of observations from a linear process, how well can the underlying process be modeled using a finite-order autoregression? How can the finite sample behavior of selection procedures be assessed? It is evident that another notion of optimality is needed in order to address these questions.

In this paper we propose to use the nonparametric minimax approach to measuring the accuracy of an AR approximation. This framework is very common in nonparametric estimation problems such as nonparametric regression, density estimation and spectral density estimation. According to this methodology, we assume that the linear process belongs to a certain class, and the quality of an approximating AR model (and the associated one-step predictor) is measured by its worst-case modeling (respectively, prediction) risk over

the class. Establishing nonasymptotic upper and lower bounds on the risk, one can assess the accuracy of an estimator. Throughout the paper we consider the class of linear processes admitting an  $AR(\infty)$  representation with exponentially decaying coefficients. The practical importance of the class follows from the fact that it includes (but is not limited to) all causal invertible  $ARMA(p, q)$  processes. We derive a nonasymptotic lower bound on the accuracy of an AR approximation and show that the least squares estimator with a “proper” choice of the order is minimax optimal in order. These considerations lead also to a nonasymptotic upper bound on the mean squared error of the associated one-step predictor. Further, we present some numerical examples comparing common selection procedures (FPE, AIC and MDL) to the minimax optimal one. We note that our derivation is based on an exponential inequality on deviations of the sample covariances from their expectations; these results are of independent interest. The same technique has been used in Goldenshluger (1998) for derivation of nonasymptotic bounds in estimating impulse response sequences of linear dynamic systems.

The rest of the paper is organized as follows. In Section 2 we state formally the problem of nonparametric AR approximation in the minimax framework. Section 3 describes the construction of the estimator, and presents main results. In Section 4 we present our numerical examples. Some remarks are collected in Section 5. The proofs are given in Appendices A, B and C.

**2. Minimax framework and overview of results.** Let  $(X_t)_{t \in \mathbf{Z}}$  be a real-valued, purely nondeterministic, Gaussian stationary process with zero mean,  $E|X_t|^2 = 1$ , spectral density function  $f(\lambda)$ ,  $\lambda \in [-\pi, \pi]$  and covariance function  $\gamma(k)$ ,  $k \in \mathbf{Z}$ . According to the Wold decomposition theorem,  $(X_t)_{t \in \mathbf{Z}}$  can be represented as an  $MA(\infty)$  process,

$$(1) \quad X_t = \sum_{j=0}^{\infty} \psi_j \varepsilon_{t-j}, \quad \psi_0 = 1, \quad \sum_{j=0}^{\infty} \psi_j^2 < \infty,$$

where  $\{\varepsilon_t\}_{t \in \mathbf{Z}}$  is a sequence of independent Gaussian innovations with  $E\varepsilon_t = 0$  and  $E\varepsilon_t^2 = \sigma_\varepsilon^2$ . Assume that the MA transfer function  $\Psi(z) = \sum_{j=0}^{\infty} \psi_j z^j$  has no zeros in the unit disc  $|z| \leq 1$ ,  $z \in \mathbf{C}$  and  $\sum_{j=0}^{\infty} |\psi_j| < \infty$ ; then the linear process  $(X_t)_{t \in \mathbf{Z}}$  can also be represented as an invertible  $AR(\infty)$  process,

$$(2) \quad X_t = \sum_{j=1}^{\infty} \phi_j X_{t-j} + \varepsilon_t, \quad t \in \mathbf{Z},$$

where the coefficients  $\phi_j$ ,  $j = 1, \dots, \infty$  are given by  $1/\Psi(z) = 1 - \sum_{j=1}^{\infty} \phi_j z^j$ . Given observations  $X_1, \dots, X_n$  from the process  $(X_t)_{t \in \mathbf{Z}}$ , we are interested in modeling  $(X_t)_{t \in \mathbf{Z}}$  and predicting the future value  $X_{n+1}$ . The representation (2) motivates the use of AR approximation to approach the problems of modeling and prediction.

Assume that the process  $(X_t)_{t \in \mathbf{Z}}$  belongs to a certain family, and the quality of an approximating AR model (and the associated one step predictor) is measured by the worst-case modeling (respectively, prediction) error over

the family. The problem of modeling the process  $(X_t)_{t \in \mathbf{Z}}$  by a finite order AR model is identical to estimating the corresponding coefficient sequence  $\phi = (\phi_1, \phi_2, \dots)$  in the AR( $\infty$ ) representation of  $(X_t)_{t \in \mathbf{Z}}$ . More formally, let  $\mathcal{H}$  be a family of stationary Gaussian processes  $(X_t)_{t \in \mathbf{Z}}$  with zero mean and unit variance, admitting an AR( $\infty$ ) representation (2). Let  $\hat{\phi} = \hat{\phi}(X_1, \dots, X_n)$  be an estimator of the sequence  $\phi = (\phi_1, \phi_2, \dots)$ ; then the quality of the estimator  $\hat{\phi}$  is measured by its maximal risk over  $\mathcal{H}$ ,

$$\mathcal{R}_m[\hat{\phi}, \mathcal{H}] := \sup_{(X_t) \in \mathcal{H}} [E\|\hat{\phi} - \phi\|^2]^{1/2},$$

where  $\|\cdot\|$  is the standard  $l_2$  norm in the space of sequences. The *minimax* estimator  $\hat{\phi}_* = \hat{\phi}_*(X_1, \dots, X_n)$  is the one minimizing the maximal risk,

$$\mathcal{R}_m^*[n, \mathcal{H}] := \inf_{\hat{\phi}} \mathcal{R}_m[\hat{\phi}, \mathcal{H}] = \inf_{\hat{\phi}} \sup_{(X_t) \in \mathcal{H}} [E\|\hat{\phi} - \phi\|^2]^{1/2},$$

where the infimum is taken here over all possible estimators. Typically, the minimax estimators cannot be constructed; therefore, as usual in nonparametric estimation, we will be interested in *optimal in order* estimators for which

$$(3) \quad \mathcal{R}_m[\hat{\phi}, \mathcal{H}] \leq C(n) \mathcal{R}_m^*[n, \mathcal{H}], \quad \sup_n C(n) < \infty.$$

Similarly, in the problem of the prediction of  $X_{n+1}$  using observations  $X_1, \dots, X_n$  we will measure the accuracy of a prediction method  $\hat{X}_{n+1}(X_1, \dots, X_n)$  by its maximal prediction error over  $\mathcal{H}$ ,

$$\mathcal{R}_p[\hat{X}_{n+1}, \mathcal{H}] := \sup_{(X_t) \in \mathcal{H}} [E(\hat{X}_{n+1} - X_{n+1})^2 - \sigma_\varepsilon^2].$$

The minimax prediction error is defined as the infimum of the maximal prediction error, over all possible prediction methods,

$$\mathcal{R}_p^*[n, \mathcal{H}] := \inf_{\hat{X}_{n+1}} \mathcal{R}_p[\hat{X}_{n+1}, \mathcal{H}] = \inf_{\hat{X}_{n+1}} \sup_{(X_t) \in \mathcal{H}} [E(\hat{X}_{n+1} - X_{n+1})^2 - \sigma_\varepsilon^2].$$

In what follows, we will be interested in *optimal in order* predictors for which (3) holds with  $\mathcal{R}_m$  replaced by  $\mathcal{R}_p$ .

Throughout the paper we restrict attention to the following family  $\mathcal{H}_\rho(l, L)$  of stationary Gaussian processes  $(X_t)_{t \in \mathbf{Z}}$  satisfying  $EX_t = 0$ ,  $E|X_t|^2 = 1$ . For given finite real numbers  $\rho > 1$ ,  $0 < l < 1$  and  $L > 1$ , define  $\mathcal{H}_\rho(l, L)$  as

$$\mathcal{H}_\rho(l, L) := \{(X_t): 0 < l \leq |\Psi(z)| \leq L, \text{ for } |z| \leq \rho\},$$

where  $\Psi(\cdot)$  is the MA( $\infty$ ) transfer function. In words, the MA( $\infty$ ) transfer function of the process  $(X_t)_{t \in \mathbf{Z}} \in \mathcal{H}_\rho(l, L)$  is analytic in an open set containing the disc  $|z| \leq \rho$ , and bounded from above and below by constants  $L$  and  $l$ , respectively. The class  $\mathcal{H}_\rho(l, L)$  contains Gaussian stationary processes with spectral density function  $f(\lambda)$  bounded away from zero and infinity, which can be continued analytically over the interior of the strip  $\{(x + iy) \in \mathbf{C}: |y| < \ln \rho\}$  in the complex plane. The parameters  $l$  and  $L$  in the definition of  $\mathcal{H}_\rho(l, L)$  guarantee

uniform lower and upper bounds on the spectral density function. This, in turn, implies uniform bounds on the eigenvalues of the covariance matrices of all orders [cf. Grenander and Szegö (1984)]. The practical importance of the class  $\mathcal{H}_\rho(l, L)$  stems from the fact that it contains causal invertible ARMA  $(p, q)$  processes with proper restrictions on the magnitude of the coefficients. For example, all MA(1) processes with the coefficient  $|\psi_1| \leq \rho^{-1} \min\{1-l, L-1\}$  belong to  $\mathcal{H}_\rho(l, L)$ . The processes from  $\mathcal{H}_\rho(l, L)$  admit AR( $\infty$ ) representation with uniformly bounded exponentially decaying coefficients.

REMARK 1. We note here a simple imbedding relationship between classes  $\mathcal{H}_\rho(l, L)$  with different parameters:  $\mathcal{H}_\rho(l, L) \subseteq \mathcal{H}_r(l, L)$ ,  $\forall \rho \geq r > 1$ . As we shall see, the only important parameter for constructing a rate optimal AR estimator (predictor) is  $\rho$ .

REMARK 2. The classes of analytic functions are quite standard in non-parametric estimation problems. Our class is similar to those in Golubev and Levit (1996) and Golubev, Levit and Tsybakov (1997). In the context of spectral density estimation, a closely related class of processes was considered by Efromovich (1998).

The main contributions of this paper are the following. We study how well processes  $(X_t)_{t \in \mathbf{Z}} \in \mathcal{H}_\rho(l, L)$  can be approximated by a finite order AR model, obtaining a lower bound on the minimax risk  $\mathcal{R}_m^*[n, \mathcal{H}_\rho(l, L)]$ . We prove that if the sample size  $n$  is large enough then

$$\mathcal{R}_m^*[n, \mathcal{H}_\rho(l, L)] \geq K(l, L) \left( \frac{\rho - 1}{\rho} \right) \frac{1}{\sqrt{\ln \rho}} \sqrt{\frac{\ln n}{n}},$$

where the constant  $K(l, L)$  depends on  $l$  and  $L$  only. A nonasymptotic upper bound on the maximal risk of the regularized least squares estimator is derived. We show that the least squares estimator associated with the order  $d_* = \lfloor (2 \ln \rho)^{-1} \ln n \rfloor$  of the approximating AR model is optimal in order in the sense of inequality (3). These results have immediate implications for the prediction problem. In particular, we derive a nonasymptotic upper bound on  $\mathcal{R}_p[\hat{X}_{n+1}, \mathcal{H}_\rho(l, L)]$  for the corresponding one-step predictor and argue that the predictor associated with the order  $d_*$  is essentially minimax optimal in order. The nonasymptotic bounds we obtain are based on exponential inequalities on deviations of sample covariances from their expectations; these results are of independent interest. Further, through numerical examples we compare small samples behavior of some common order selection procedures to the minimax optimal choice  $d_*$ . In particular, simulating an MA(1) process, we found that for moderate values of  $\rho$ , that is, when the zeros are not “too close” to the unit disc, the AIC and FPE lead to an order selection that is comparable to the minimax optimal one. However, if  $\rho$  is close to unity then the AIC and FPE tend to select a smaller model order than the minimax optimal one. The MDL turns out to be slightly more conservative than the other methods, with the differences becoming marginal for larger values of  $\rho$ .

**3. Main results.** Consider the following estimate of the AR sequence  $\phi = (\phi_1, \phi_2, \dots)$ . Fix a natural number  $d$  and define

$$\theta^d = (\phi_1, \dots, \phi_d)', \quad Z_t = (X_{t-1}, \dots, X_{t-d})'.$$

We estimate  $\theta^d$  by the regularized least squares method,

$$(4) \quad \hat{\theta}^d = \left( \frac{1}{n} \sum_{t=1}^n Z_t Z_t' + n^{-1} I_d \right)^{-1} \left( \frac{1}{n} \sum_{t=1}^n X_t Z_t \right), \quad \hat{\theta}^d = (\hat{\phi}_1, \dots, \hat{\phi}_d)',$$

where  $I_d$  is the identity  $d \times d$  matrix. The corresponding estimate  $\hat{\phi}$  of the sequence  $\phi = (\phi_1, \phi_2, \dots)$  is given by

$$(5) \quad \hat{\phi} = (\hat{\phi}_1, \hat{\phi}_2, \dots, \hat{\phi}_d, 0, 0, \dots)$$

and the one-step predictor  $\hat{X}_{n+1}^d$  based on  $\hat{\phi}$  is defined by

$$(6) \quad \hat{X}_{n+1}^d = \sum_{j=1}^d \hat{\phi}_j X_{n+1-j}.$$

The reason why we consider a regularized version of the least squares estimate is that we are interested in a nonasymptotic upper bound on the expected value of the squared modeling (prediction) error. For this purpose we have to control the norm of the random matrix  $(n^{-1} \sum_{t=1}^n Z_t Z_t')^{-1}$ . Without the regularization term the matrix  $n^{-1} \sum_{t=1}^n Z_t Z_t'$  can be singular with nonzero probability for every fixed  $n$ . We note also that the vectors  $Z_t$ ,  $t = 1, \dots, n$  defined above can involve  $X_t$  with  $t \leq 0$ . In this case we suppose that the corresponding components of the vectors in (4) are replaced by zero. It should be stressed, however, that in our analysis we do not assume that  $X_t = 0$  for  $t \leq 0$ .

**3.1. Accuracy of AR approximation.** We are now ready to study the quality of an AR approximation of the stationary Gaussian process  $(X_t)_{t \in \mathbf{Z}} \in \mathcal{H}_\rho(l, L)$ .

**THEOREM 1.** *Let*

$$(7) \quad M = 1 + \frac{L\rho}{l(\rho-1)}, \quad r = 1 + \frac{1}{\ln \rho}.$$

*Suppose that  $n$  and  $d$  satisfy the following conditions:*

$$(8) \quad \frac{n}{(\ln n)^5} \geq c_1 d (rM)^5, \quad \sqrt{\frac{n}{\ln n}} \geq c_2 (L/l)^2 d^2 \sqrt{rM},$$

*where  $c_1$  and  $c_2$  are absolute constants which can be specified explicitly. Then for the estimate (4) and (5) one has*

$$(9) \quad \mathcal{R}_m[\hat{\phi}, \mathcal{H}_\rho(l, L)] \leq K_1(l, L) \left( \frac{1}{n(\rho-1)} + \frac{\sqrt{d}}{\rho^d(\rho-1)} + \sqrt{\frac{d}{n}} \right),$$

*where  $K_1(l, L)$  depends on  $l$  and  $L$  only.*

REMARK 3. Accuracy of the AR approximation is limited by two factors. First, we approximate the process by a finite-order AR model. The resulting *approximating error* [second term in the right-hand side of (9)] becomes smaller as the order  $d$  of the approximating AR model increases. Second, we estimate parameters of the approximating model. The resulting *estimating error* [third term in the right-hand side of (9)] grows as the order of the approximating model increases. The order  $d$  is viewed as a “smoothing parameter” that controls a trade-off between the approximation and estimation errors. The first term in the right-hand side of (9) is due to the use of a regularized version of the least squares estimator.

The following statement is an immediate consequence of Theorem 1.

COROLLARY 1. *Let  $n$  be large enough so that*

$$(10) \quad \frac{n}{(\ln n)^6} \geq c_1 r^6 M^5, \quad \frac{\sqrt{n}}{(\ln n)^{5/2}} \geq c_2 \left( \frac{L}{l \ln \rho} \right)^2 \sqrt{rM}.$$

*Then for the least squares estimator (4) and (5) associated with the choice  $d_* := \lfloor (2 \ln \rho)^{-1} \ln n \rfloor$  one has*

$$(11) \quad \mathcal{R}_m[\hat{\phi}_*, \mathcal{H}_\rho(l, L)] \leq K_2(l, L) \left( \frac{\rho}{\rho - 1} \right) \frac{1}{\sqrt{\ln \rho}} \sqrt{\frac{\ln n}{n}},$$

*where  $K_2(l, L)$  depends on  $l$  and  $L$  only.*

The next step in the analysis is to determine the limits of achievable accuracy for AR approximation. The following theorem gives a lower bound on approximation of  $(X_t)_{t \in \mathbb{Z}} \in \mathcal{H}_\rho(l, L)$  by a finite-order AR model.

THEOREM 2. *Let  $n$  be large enough so that for some constant  $K_3$  depending on  $l$  and  $L$  only,*

$$(12) \quad \ln n \geq K_3(l, L) \ln \rho.$$

*Then*

$$(13) \quad \mathcal{R}_m^*[n, \mathcal{H}_\rho(l, L)] \geq K_4(l, L) \left( \frac{\rho - 1}{\rho} \right) \frac{1}{\sqrt{\ln \rho}} \sqrt{\frac{\ln n}{n}}.$$

Theorem 2 and Corollary 1 imply that the least squares estimator (4) and (5) associated with the order  $d_* = \lfloor (2 \ln \rho)^{-1} \ln n \rfloor$  is optimal in order in the sense of inequality (3). It is interesting to note that for the class of spectral densities corresponding to processes closely related to the class  $\mathcal{H}_\rho(l, L)$ , this choice of the order leads to the asymptotically minimax spectral estimate [Efromovich (1998)].

3.2. *Prediction via AR approximation.* In this section we establish a nonasymptotic upper bound on accuracy of the one-step predictor which is based on AR approximation. To simplify analysis we assume that the estimate  $\hat{\phi}$  of the sequence  $\phi$  is based on  $\lfloor n/2 \rfloor$  first observations  $(X_1, \dots, X_{\lfloor n/2 \rfloor})$  only. The assumption of this type is quite usual in investigating accuracy of prediction methods based on the estimated parameters. For instance, Shibata (1980) assumed the more stringent assumption that we have two independent realizations of the linear process: the first time series is used for estimating parameters, and then the estimated parameters are used to predict the second time series.

The associated one-step predictor is defined in (6). We note also that Theorem 1 remains unaltered for the estimate in question with  $n$  replaced by  $n/2$ .

**THEOREM 3.** *Let  $n > 4d$  and (8) hold with some absolute constants  $c_1$  and  $c_2$ . Then one has*

$$(14) \quad \mathcal{R}_p[\hat{X}_{n+1}^d, \mathcal{H}_\rho(l, L)] \leq K_5(l, L) \left( \frac{1}{n^2(\rho-1)^2} + \frac{d}{\rho^{2d}(\rho-1)^2} + \frac{d}{n} \right) \left( 1 + \frac{d\rho^{-2d}}{(\rho-1)^2} \right),$$

where  $K_5(l, L)$  depends on  $l$  and  $L$  only.

Now, choosing the model order  $d$  we obtain the prediction bounds.

**COROLLARY 2.** *Let (10) hold with some absolute constants  $c_1$  and  $c_2$ , and  $n(\ln n)^{-1} \geq 2(\ln \rho)^{-1}$ . Then for the one-step predictor  $\hat{X}_{n+1}^d$  associated with the choice*

$$d_* = \lfloor (2 \ln \rho)^{-1} \ln n \rfloor,$$

one has

$$(15) \quad \mathcal{R}_p[\hat{X}_{n+1}^{d_*}, \mathcal{H}_\rho(l, L)] \leq K_6(l, L) \left( \frac{\rho}{\rho-1} \right)^2 \left( \frac{1}{\ln \rho} \right) \frac{\ln n}{n},$$

where  $K_6(l, L)$  depends on  $l$  and  $L$  only.

Referring back to (9), we see that the upper bound on prediction accuracy given in (14) behaves as the square of accuracy of modeling. This is not surprising given the construction of the one-step predictor (6); clearly the resulting accuracy is determined by the quality of modeling via the AR approximation. In fact, one can argue that the predictor  $X_{n+1}^{d_*}$  is optimal in order. For the sake of simplicity assume as in Shibata (1980) that we have two independent copies  $Y_1, \dots, Y_n$  and  $X_1, \dots, X_n$  of the same linear process from  $\mathcal{H}_\rho(l, L)$ . Our goal is to predict  $X_{n+1}$ . Let  $\hat{X}_{n+1}$  be an arbitrary prediction method for  $X_{n+1}$  based on the observations  $Y_1, \dots, Y_n$ . Then  $\hat{X}_{n+1}$  can be decomposed into a sum of



two random variables  $\widehat{X}'_{n+1}$  and  $\widehat{X}''_{n+1}$  such that  $\widehat{X}'_{n+1}$  is the projection of  $\widehat{X}_{n+1}$  on  $\overline{\text{sp}}\{X_n, X_{n-1}, \dots\}$ , and  $\widehat{X}''_{n+1}$  is orthogonal to  $\overline{\text{sp}}\{X_n, X_{n-1}, \dots\}$ . Therefore,

$$\begin{aligned} \mathcal{R}_p^*[n, \mathcal{H}_\rho(l, L)] &\geq \sup_{(X_t) \in \mathcal{H}_\rho(l, L)} E \left| \widehat{X}'_{n+1} - \sum_{j=1}^{\infty} \phi_j X_{n+1-j} \right|^2 \\ &= \sup_{(X_t) \in \mathcal{H}_\rho(l, L)} E \left| \sum_{j=1}^{\infty} (\hat{\phi}_j - \phi_j) X_{n+1-j} \right|^2, \end{aligned}$$

where  $\hat{\phi}_j, j = 1, 2, \dots$  are measurable functions of  $Y_1, \dots, Y_n$ . Hence

$$\mathcal{R}_p^*[n, \mathcal{H}_\rho(l, L)] \geq K_7(l, L) \sup_{\phi \in \mathcal{H}_\rho(l, L)} E \|\hat{\phi} - \phi\|^2,$$

where the supremum is taken over all sequences  $\phi = (\phi_1, \phi_2, \dots)$  that define, through the AR( $\infty$ ) representation, the linear processes from  $\mathcal{H}_\rho(l, L)$ . Then the lower bound on the minimax prediction risk follows from Theorem 2.

**4. Numerical examples.** The choice of model order,  $d(n) = O(\ln n)$ , arises in Shibata (1980), and more recently in Hannan, Kavalieris (1986) and Gerencsér (1992). In particular, Shibata (1980) showed that the data-driven order selector based on the final prediction error (FPE) behaves asymptotically as  $O(\ln n)$  for the class of processes, similar to  $\mathcal{H}_\rho(m, L)$ . To investigate the practical impact of the above results, we compare the common model selection strategies (AIC, FPE and MDL) to the minimax optimal rule through a simple numerical example.

Consider the following MA(1) process

$$X_t = \varepsilon_t + \psi_1 \varepsilon_{t-1}$$

with  $\{\varepsilon_t\}$  a sequence of i.i.d. standard Gaussian random variables. We focus our attention on three particular cases, namely  $\psi_1 = 0.1, 0.5, 0.9$ , and the corresponding “margin of stability”  $\rho = 10, 2, 1.1111$ . This range of values will illustrate the sensitivity of the order selection methods to the moduli of the zeros of the transfer function  $\Psi(\cdot)$ . Suppose we are given  $n$  consecutive observations  $X_1, X_2, \dots, X_n$  from the process  $(X_t)$ . The selection procedures are defined [following the definitions in Shibata (1980)] as

$$\text{AIC}(d) := (n + 2d)\hat{\sigma}_d^2,$$

$$\text{FPE}(d) := n((n + d)/(n - d))\hat{\sigma}_d^2,$$

$$\text{MDL}(d) := (n + d \ln n)\hat{\sigma}_d^2$$

with

$$\hat{\sigma}_d^2 := \frac{1}{n - d} \sum_{t=d+1}^n \left( X_t - \sum_{j=1}^d \hat{\phi}_j X_{t-j} \right)^2.$$

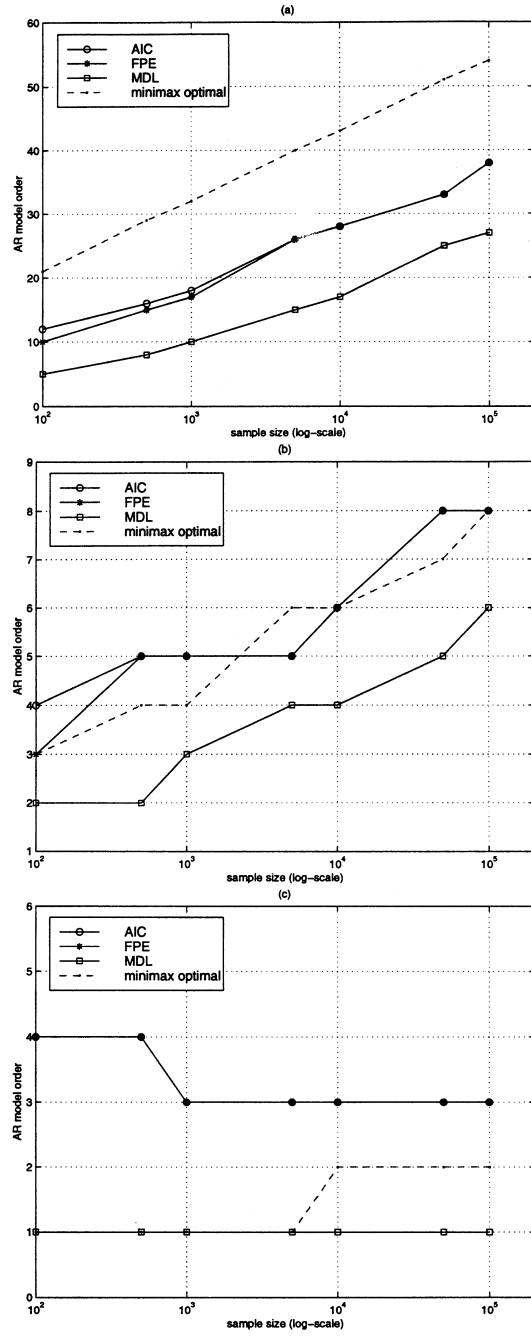


FIG. 1. Model order selected by different procedures plotted against the sample size (log scale); (a)  $\rho = 1.1111$ , (b)  $\rho = 2$  and (c)  $\rho = 10$ .

Recall also the minimax optimal order choice from Corollary 1:  $d_* = \lfloor (2 \ln \rho)^{-1} \ln n \rfloor$ .

The experiment was conducted by simulating 100 sample paths from the process, for each trial run a model order was selected using the three procedures, for sample sizes  $n = 100, 500, 1000, 5000, 10000, 50000$  and  $100000$ . Finally, we averaged out the selected orders over the 100 runs. The graphs in Figure 1 depict the behavior of the different order selection procedures.

A close look at Figure 1 reveals that the AIC and FPE, which are known to be asymptotically equivalent, behave in an almost identical way also for small values of sample size. The MDL leads to a choice that is more conservative than AIC and FPE, with this behavior being more pronounced for the case of small  $\rho$ . For the case of large  $\rho$ , the all three criteria are roughly the same as the minimax optimal choice. The case of moderate  $\rho$  depicts a behavior of AIC and FPE which is quite on a par with the minimax optimal choice. However, if  $\rho$  is close to unity, then the AIC and FPE tend to select a smaller model order than the minimax optimal one. It is interesting to note that all procedures lead to an order selection that exhibits logarithmic-like growth in the sample size, even for small sample sizes. This behavior is consistent with the asymptotic logarithmic growth of the order selected by AIC and FPE [cf. Shibata (1980), Example 4.1] and for MDL [cf. Gerencér (1992), Theorem 4].

To summarize the results, we observe that an infinite order AR model that is closer to a parametric (finite-dimensional) model gives rise to an order selection that is “close” to minimax optimal by all three methods. The case of more slowly decaying coefficients (larger  $\psi$  and  $\rho$  closer to unity, respectively) reveals that AIC and FPE “underestimate” with MDL being even more conservative. We note in passing that similar numerical results were obtained for more complicated ARMA structures.

**5. Discussion.** (i) The method of AR approximations is quite common for spectral density estimation in time series analysis [see, e.g., Berk (1974), Shibata (1981), Parzen (1983) among many others]. The minimax optimal model order ( $d_* = \lfloor \ln n / (2 \ln \rho) \rfloor$ ) for the AR approximation is also the optimal choice for spectral density estimation, and gives rise to the same convergence rates over the class  $\mathcal{H}_\rho(l, L)$ . It is worth noting, however, that spectral density estimation and AR approximation are not equivalent in the sense of comparison of experiments. Specifically, assume that the process belongs to the class of all invertible MA( $q$ ) processes whose MA-transfer function has no zeros inside the disc  $|z| \leq \rho$ ,  $\rho > 1$ . This class is a subset of  $\mathcal{H}_\rho(l, L)$  with proper  $l$  and  $L$ . The spectral density of such a process can be estimated with the parametric rate  $O(\sqrt{q/n})$ , while the accuracy of the AR approximation is  $O(\sqrt{\ln n / (n \ln \rho)})$ . An important implication of this fact is that even if a stationary process is approximated by an AR model with high accuracy, the corresponding spectral density estimate may be poor.

(ii) The nonparametric minimax approach, as applied to AR approximation, provides a useful criterion for assessment of finite sample behavior of selection methods. Within this approach, optimal selection methods are spec-

ified, and achievable lower bounds on the estimation accuracy are calculated. Note, however, that implementation of the minimax optimal rule requires a priori information on the parameter  $\rho$  of the class  $\mathcal{H}_\rho(l, L)$ . Developing adaptive selection rules with good minimax properties remains a challenging open problem. We conjecture that in the adaptive setting the rates of convergence for AR approximation remain unchanged.

(iii) Throughout the paper we assume that the process  $(X_t)_{t \in \mathbf{Z}}$  is Gaussian. This assumption is used to simplify the derivation of the exponential inequalities on the covariance estimates (Lemma 2 below). In addition, it facilitates the evaluation of higher order moments. The main results of the paper can be obtained under moment growth restrictions accompanied with some requirements ensuring exponential mixing properties of the process  $(X_t)_{t \in \mathbf{Z}}$ .

(iv) The family  $\mathcal{H}_\rho(l, L)$  allows for the processes admitting AR( $\infty$ ) representation with exponentially decaying coefficients. It seems that the exponential decay of the coefficients is essential for the exponential inequalities we derive. The techniques advocated in Lemma 6 and Lemma 7 below preclude polynomially decaying sequences. Thus, this restriction is a direct consequence of the limitations of our machinery. An interesting problem is to study rates of AR approximation for other classes of stationary, for example, with polynomially decaying AR coefficients.

## APPENDIX

**A. Preliminary results.** We collect here several preliminary results which will be used repeatedly in the subsequent proofs.

We start with establishing a relation between the properties of the sequences  $\gamma(k)$ ,  $\psi_j$  and  $\phi_j$  to the class  $\mathcal{H}_\rho(l, L)$ . Let us define  $\Gamma_d \equiv \{\gamma(i-j)\}_{i,j=1,\dots,d}$  for every natural number  $d$ .

LEMMA 1. *Let  $(X_t)_{t \in \mathbf{Z}} \in \mathcal{H}_\rho(l, L)$ ; then*

$$(16) \quad |\psi_j| \leq L\rho^{-j}, \quad |\phi_j| \leq l^{-1}\rho^{-j}, \quad j = 1, 2, \dots$$

*In addition, we have*

$$(17) \quad L^{-2} \leq \sigma_\epsilon^2 \leq l^{-2},$$

$$(18) \quad |\gamma(k)| \leq (L/l)^2 \frac{\rho^2}{\rho^2 - 1} \rho^{-|k|}, \quad k \in \mathbf{Z},$$

*and for any  $d$*

$$(19) \quad (l/L)^2 \leq \|\Gamma_d\| \leq (L/l)^2, \quad (l/L)^2 \leq \|\Gamma_d^{-1}\| \leq (L/l)^2,$$

*where  $\|\cdot\|$  stands for the standard Euclidean norm of a matrix.*

PROOF. By definition of the class  $\mathcal{H}_\rho(l, L)$ ,  $\Psi(z)$  is analytic in the open disc  $|z| < \rho$ , and  $|\Psi(z)| \leq L$ . Therefore the announced bound on  $|\psi_j|$  follows

immediately from the Cauchy estimates for the derivatives of  $\Psi(z)$  [see, e.g., Rudin (1974), page 229]. Further, note that

$$L^{-1} \leq |\Phi(z)| = 1/|\Psi(z)| \leq l^{-1} \quad \text{for } |z| < \rho.$$

Again applying the Cauchy estimates we obtain (16).

Note that  $f(\lambda) = (2\pi)^{-1} \sigma_\varepsilon^2 |\Psi(e^{-i\lambda})|^2$ , and therefore,

$$(20) \quad (2\pi)^{-1} \sigma_\varepsilon^2 l^2 \leq f(\lambda) \leq (2\pi)^{-1} \sigma_\varepsilon^2 L^2.$$

Taking into account that  $\gamma(0) = 1 = \int_{-\pi}^{\pi} f(\lambda) d\lambda$ , we obtain (17). The inequality (18) is an immediate consequence of the following evident inequalities:

$$|\gamma(k)| = \sigma_\varepsilon^2 \left| \sum_{j=0}^{\infty} \psi_j \psi_{j+|k|} \right| \leq L^2 \sigma_\varepsilon^2 \rho^{-|k|} \sum_{j=0}^{\infty} \rho^{-2j} = \frac{L^2 \sigma_\varepsilon^2 \rho^2}{(\rho^2 - 1)} \rho^{-|k|}$$

and (17) [here we have used the bound on  $\psi_j$  established in (16)]. The bounds on  $\|\Gamma_d\|$  and  $\|\Gamma_d^{-1}\|$  follow from the theorem on the eigenvalues of the Toeplitz forms [cf. Grenander and Szego (1984)]. In particular, we have

$$l^2 \sigma_\varepsilon^2 \leq \lambda_{\min}[\Gamma_d] \leq \lambda_{\max}[\Gamma_d] \leq L^2 \sigma_\varepsilon^2,$$

where  $\lambda_{\min}[\cdot]$  and  $\lambda_{\max}[\cdot]$  denote the minimal and maximal eigenvalues of a matrix, respectively. Applying (17) we obtain (19) which completes the proof.  $\square$

*A.1. An exponential inequality for sample covariances.* Here we establish an exponential inequality on the deviation of sample covariances from their expectations. This result is basic for our future developments; furthermore, it is interesting in its own right.

LEMMA 2. *Let  $(X_t)_{t \in \mathbf{Z}} \in \mathcal{H}_\rho(l, L)$ ; then there exist absolute constants  $C_1$  and  $C_2$  such that for every integer  $k$  one has*

$$(21) \quad P \left\{ \left| \frac{1}{n} \sum_{t=1}^n X_t X_{t+k} - \gamma(k) \right| > \delta \right\} \leq \begin{cases} \exp\left(-\frac{\delta^2 n}{4C_1 M k_* r}\right), & 0 \leq \delta \leq \left(\frac{k_* r}{n}\right)^{2/5} \left(\frac{C_1^3 M^3}{C_2}\right)^{1/5}, \\ \exp\left(-\frac{1}{4} \left[\frac{\delta n}{C_2 k_* r}\right]^{1/3}\right), & \delta \geq \left(\frac{k_* r}{n}\right)^{2/5} \left(\frac{C_1^3 M^3}{C_2}\right)^{1/5}, \end{cases}$$

where  $M$  and  $r$  are defined in (7) and  $k_* = |k|$  whenever  $k \neq 0$ , and  $k_* = 1$  whenever  $k = 0$ . The constants  $C_1$  and  $C_2$  are specified explicitly in the proof of the lemma.

REMARK 4. To establish the result of the lemma we use the general exponential inequalities for weakly dependent random sequences found in Saulis and Statulevičius (1991). Several other exponential-type inequalities for weakly dependent random sequences appear already in the literature [cf., e.g.,

Doukhan (1994), Bosq (1996)], of which Bosq [Theorem 1.4 (1996)] deals with conditions that are probably most akin to our set up. However, the machinery in Saulis and Statulevičius (1991) seems more suitable for our purposes, and leads to tighter bounds, in particular since we use the moderate deviations regime in (21).

For the proof, see Appendix C.1.

## B. Proofs of main results.

B.1. *Proof of Theorem 1.* In the proof below  $K_i$ ,  $i = 1, 2, \dots$  stand for absolute positive constants (unless otherwise specified), possibly different in different instances.

We first outline the main ideas in the proof. By straightforward algebra we have

$$(22) \quad \hat{\theta}^d - \theta^d = \mathbf{Q}^{-1} \left( -n^{-1} \theta^d + \frac{1}{n} \sum_{t=1}^n \mathbf{Z}_t \sum_{j=d+1}^{\infty} \phi_j X_{t-j} + \frac{1}{n} \sum_{t=1}^n \mathbf{Z}_t \varepsilon_t \right),$$

where  $\mathbf{Q} := n^{-1} \sum_{t=1}^n \mathbf{Z}_t \mathbf{Z}_t' + n^{-1} I_d$ . Thus, to prove a bound on the  $\ell_2$  distance between  $\theta^d$  and  $\hat{\theta}^d$ , we must bound the norm of the matrix  $\mathbf{Q}^{-1}$ , and of the vector multiplying it from the right in (22). The latter bound involves straightforward algebraic manipulations, therefore the real problem is to control the norm of  $\mathbf{Q}^{-1}$ . The key idea here is the following. Partition the sample space into two sets. One set corresponds to the samples of  $(X_t)_{t \in \mathbf{Z}}$ , for which the elements of  $\mathbf{Q}$  are uniformly “close” to their expectations. For the complement of this set,  $\|\mathbf{Q}^{-1}\|$  does not grow faster than  $n$ . Exponential inequalities on the uniform convergence of sample means to their expectations, in the spirit of Lemma 2, ensure that the “bad” set essentially does not contribute to the overall bound. We shall now make these statements rigorous.

*Step 1.* First, proceed to bound  $\|\mathbf{Q}^{-1}\|$ , where  $\|\cdot\|$  denotes the standard Euclidean matrix norm. Note that the  $i, j$ -entry  $Q_{ij}$  of the matrix  $\mathbf{Q}$  with  $i \neq j$  may be expressed as follows:

$$\begin{aligned} Q_{ij} &= \frac{1}{n} \sum_{t=1}^n X_{t-i} X_{t-j} \\ &= \frac{1}{n} \sum_{\tau=1}^n X_{\tau} X_{\tau+j-i} - \frac{1}{n} \sum_{\tau=n-j+1}^n X_{\tau} X_{\tau+j-i} + \frac{1}{n} \sum_{\tau=1-j}^0 X_{\tau} X_{\tau+j-i} \\ &:= \widehat{Q}_{ij} - W_{ij} + V_{ij}. \end{aligned}$$

This, in turn, may be written as

$$(23) \quad \mathbf{Q} = \widehat{\mathbf{Q}} - \mathbf{W} + \mathbf{V} + n^{-1} I_d = \Gamma_d [I_d + \Gamma_d^{-1} (\widehat{\mathbf{Q}} + n^{-1} I_d)],$$

where  $\widehat{Q} = (\widehat{Q}_{ij})$ ;  $W = (W_{ij})$ ;  $V = (V_{ij})$ ;  $i, j = 1, \dots, d$ , and  $\widetilde{Q} := V - W + \widehat{Q} - \Gamma_d$ . Observe that  $\Gamma_d$  is nonsingular for every  $d$  (this follows from Lemma 1). Thus, the task of bounding  $\|Q^{-1}\|$  is reduced to establishing a bound on the norm of  $[I_d + \Gamma_d^{-1}(\widetilde{Q} + n^{-1}I_d)]^{-1}\Gamma_d^{-1}$ , where the only stochastic term is  $\widetilde{Q}$ . The main idea is the following. Write

$$\widetilde{Q} = (V - E[V]) - (W - E[W]) + (\widehat{Q} - \Gamma_d),$$

utilizing the fact that  $E[V] = E[W]$ . Note also that  $E[\widehat{Q}] = \Gamma_d$ . In addition, due to Lemma 2 we can evaluate how close  $\widehat{Q}$  to  $\Gamma_d$  is. Now, the key to bounding  $\|\widetilde{Q}\|$ , is to establish nonasymptotic exponential bounds on the probability that each one of the terms  $V$ ,  $W$  and  $\widehat{Q}$  deviate from their expectations.

LEMMA 3. *Let  $(X_t)_{t \in \mathbf{Z}} \in \mathcal{H}_\rho(l, L)$ . For any fixed  $i, j \in \{1, \dots, d\}$  we have*

$$(24) \quad \begin{aligned} &P\{|V_{ij} - EV_{ij}| > \delta\} \\ &\leq \begin{cases} \exp\left(-\frac{\delta^2 n}{4C_1 d}\right), & 0 \leq \delta \leq [dn^{-1}C_1^2 C_2^{-1}]^{1/3}, \\ \exp\left(-\frac{1}{4}\sqrt{\frac{\delta n}{C_2 d}}\right), & \delta \geq [dn^{-1}C_1^2 C_2^{-1}]^{1/3}, \end{cases} \end{aligned}$$

where  $C_1$  and  $C_2$  are as in Lemma 2. The same relations hold for  $W_{ij}$ .

For the proof, see Appendix C.2

Step 2. Recall that by definition  $EV_{ij} = EW_{ij}$ , so that

$$\widetilde{Q} = (V - E(V)) + (E(W) - W) + \widehat{Q} - \Gamma_d.$$

Applying the results of Lemma 2 and Lemma 3 we bound the norm of the matrix  $Q^{-1}$ . Let us fix  $\kappa \in (0, 1)$  and define the event

$$(25) \quad A_\kappa = \left\{ \omega \in \Omega : \max_{i, j=1, \dots, d} |\widetilde{Q}_{ij}| \leq C_\kappa \right\},$$

where

$$(26) \quad C_\kappa = 6\sqrt{C_1 r M} \sqrt{\frac{d}{n} \ln\left(\frac{6d^2}{\kappa}\right)}.$$

Here, and in the sequel,  $\Omega$  is the sample set of the underlying probability space  $(\Omega, \mathcal{F}, P)$ .

LEMMA 4. *Let  $(X_t)_{t \in \mathbf{Z}} \in \mathcal{H}_\rho(l, L)$  and for a fixed  $\kappa \in (0, 1)$  let  $d$  and  $n$  be such that*

$$(27) \quad d^{-1}n \geq (36C_1 r M \ln(6d^2/\kappa))^5$$

and

$$(28) \quad n^{-1} + dC_\kappa \leq \frac{1}{2}(l/L)^2.$$

Then  $P(A_\kappa^c) \geq 1 - \kappa$ , and  $\|\Gamma_d \mathbf{Q}^{-1}\| \leq 2$  if the event  $A_\kappa$  holds, and  $\|\mathbf{Q}^{-1}\| \leq n$  otherwise.

For the proof, see Appendix C.3.

*Step 4.* Now, recall for completeness (22),

$$\begin{aligned} \hat{\theta}^d - \theta^d &= \mathbf{Q}^{-1} \left( -n^{-1} \theta^d + \frac{1}{n} \sum_{t=1}^n \mathbf{Z}_t \sum_{j=d+1}^{\infty} \phi_j X_{t-j} + \frac{1}{n} \sum_{t=1}^n \mathbf{Z}_t \varepsilon_t \right) \\ &= \mathbf{Q}^{-1} (\mathcal{S}_1 + \mathcal{S}_2 + \mathcal{S}_3) = \mathbf{Q}^{-1} \mathcal{S}. \end{aligned}$$

Having established a bound on  $\|\mathbf{Q}^{-1} \Gamma_d\|$  we proceed to bound  $E\|\mathcal{S}\|$ .

LEMMA 5. *Let  $(X_t)_{t \in \mathbf{Z}} \in \mathcal{H}_\rho(l, L)$ . Then,*

$$\begin{aligned} \|\mathcal{S}_1\| &\leq \frac{1}{nl(\rho-1)}, \\ (E\|\mathcal{S}_2\|^4)^{1/2} &\leq \frac{K_1 d \rho^{-2d}}{l^2(\rho-1)^2}, \\ (E\|\mathcal{S}_3\|^4)^{1/2} &\leq \frac{K_2 d}{l^2 n}, \end{aligned}$$

where  $K_1$  and  $K_2$  are absolute constants.

For the proof, see Appendix C.4.

*Step 5.* Now we complete the proof of Theorem 1. We will proceed to bound  $E\|\hat{\theta}^d - \theta^d\|^2$  by evaluating the expectation over two disjoint subsets corresponding to the events  $A_\kappa$  and  $A_\kappa^c$ . Let  $\kappa = 6d^2 n^{-6}$ ,  $A_\kappa$  be given by (25) with  $C_\kappa$  defined by (26) with  $\kappa$  in question. It can be immediately checked that under (8) conditions of the Lemma 4 hold. Thus we can write

$$\begin{aligned} E[\|\hat{\theta}^d - \theta^d\|^2 \mathbf{1}_{\{A_\kappa\}}] &\leq \|\Gamma_d^{-1}\|^2 E[\|\Gamma_d \mathbf{Q}^{-1}\|^2 \|\mathcal{S}_1 + \mathcal{S}_2 + \mathcal{S}_3\|^2 \mathbf{1}_{\{A_\kappa\}}] \\ &\stackrel{(a)}{\leq} 4 \|\Gamma_d^{-1}\|^2 E(\|\mathcal{S}_1 + \mathcal{S}_2 + \mathcal{S}_3\|^2 \mathbf{1}_{\{A_\kappa\}}) \\ &\leq 16 \|\Gamma_d^{-1}\|^2 (\|\mathcal{S}_1\|^2 + E\|\mathcal{S}_2\|^2 + E\|\mathcal{S}_3\|^2) \\ &\stackrel{(b)}{\leq} K_3 \|\Gamma_d^{-1}\|^2 l^{-2} \left( \frac{1}{n^2(\rho-1)^2} + \frac{d}{\rho^{2d}(\rho-1)^2} + \frac{d}{n} \right), \end{aligned}$$



where (a) follows from Lemma 4, and (b) follows from the bounds established in Lemma 5. Similarly, we have

$$\begin{aligned} E[\|\hat{\theta}^d - \theta^d\|^2 \mathbf{1}_{\{A_\kappa^c\}}] &\leq 4E[\|Q^{-1}\|^2(\|\mathcal{J}_1\|^2 + \|\mathcal{J}_2\|^2 + \|\mathcal{J}_3\|^2) \mathbf{1}_{\{A_\kappa^c\}}] \\ &\leq 4n^2 \left[ \|\mathcal{J}_1\|^2 \mathbf{P}(A_\kappa^c) + \left( \sqrt{E\|\mathcal{J}_2\|^4} + \sqrt{E\|\mathcal{J}_3\|^4} \right) \sqrt{\mathbf{P}(A_\kappa^c)} \right] \\ &\leq K_4 n^2 l^{-2} \left[ \frac{\kappa}{n^2(\rho - 1)^2} + \frac{d\sqrt{\kappa}}{\rho^{2d}(\rho - 1)^2} + \frac{\sqrt{\kappa d}}{n} \right]. \end{aligned}$$

Substituting expression for  $\kappa$  and combining the two bounds above we have

$$(29) \quad [E\|\hat{\theta}^d - \theta^d\|^2]^{1/2} \leq K_5 \|\Gamma_d^{-1}\| l^{-1} \left( \frac{1}{n(\rho - 1)} + \frac{\sqrt{d}}{\rho^d(\rho - 1)} + \sqrt{\frac{d}{n}} \right),$$

whence

$$\begin{aligned} [E\|\hat{\phi} - \phi\|^2]^{1/2} &\leq [E\|\hat{\theta}^d - \theta^d\|^2]^{1/2} + \left( \sum_{j=d+1}^\infty |\phi_j|^2 \right)^{1/2} \\ &\leq K_5 \|\Gamma_d^{-1}\| l^{-1} \left( \frac{1}{n(\rho - 1)} + \frac{\sqrt{d}}{\rho^d(\rho - 1)} + \sqrt{\frac{d}{n}} \right) + \frac{1}{l\rho^d(\rho - 1)}. \end{aligned}$$

Applying (19) completes the proof.  $\square$

**B.2. Proof of Theorem 2.** Proof of the theorem rests upon the standard technique for deriving lower bounds in nonparametric estimation problems. In the proof below  $K_i, i = 1, 2, \dots$  denote positive constants depending on  $l$  and  $L$  only.

Let us fix a natural number  $N$ , and consider the following family  $\mathcal{P}$  of the sequences  $\phi = (\phi_1, \phi_2, \dots)$ :  $\phi$  belongs to  $\mathcal{P}$  if and only if

$$\phi_j = \begin{cases} \pm\beta\rho^{-N}, & j = 1, \dots, N, \\ 0, & \text{otherwise,} \end{cases}$$

where  $\beta$  is a positive number to be chosen. We complement  $\mathcal{P}$  by the zero sequence  $\phi^{(0)} = (0, 0, \dots)$ . It is evident that there exists a choice of constant  $K_1$  (e.g., take  $K_1 \leq \min\{1 - L^{-1}, l^{-1} - 1\}$ ), such that with the choice  $\beta = K_1(1 - \rho^{-1})$  every  $\phi \in \mathcal{P}$  defines a process  $(X_t)_{t \in \mathbb{Z}}$  from  $\mathcal{H}_\rho(l, L)$ . In addition, cardinality of  $\mathcal{P}$  is equal to  $2^N + 1$ . According to the Varshamov–Gilbert lemma [see, e.g., Korostelev and Tsybakov (1993), page 79], one can choose a subfamily  $\mathcal{P}' \subset \mathcal{P}$  so that any two distinct sequences  $\phi', \phi''$  from  $\mathcal{P}'$  differ by at least  $N/16$  components, cardinality of  $\mathcal{P}'$  is equal to  $2^{\lfloor N/8 \rfloor} + 1$  and  $\phi^{(0)} \in \mathcal{P}'$ . Thus, for any  $\phi', \phi''$  one has

$$(30) \quad \|\phi' - \phi''\| \geq K_2 \sqrt{N}(1 - \rho^{-1})\rho^{-N} := s.$$

Let  $\hat{\phi}_n$  be an arbitrary estimate of  $\phi$  based on the data  $\{X_t\}_{t=1}^n$ ; then

$$(31) \quad \sup_{(X_t) \in \mathcal{H}_\rho(m, L)} E\|\hat{\phi}_n - \phi\| \geq \sup_{\phi \in \mathcal{P}'} E\|\hat{\phi}_n - \phi\| \geq \frac{s}{2} \sup_{\phi \in \mathcal{P}'} P\{\|\hat{\phi}_n - \phi\| \geq s/2\}.$$

Now consider the problem of testing between  $2^{\lfloor N/8 \rfloor} + 1$  hypotheses  $H_j$ :  $\phi = \phi^{(j)}$  using observations  $\{X_t\}_{t=1}^n$ ; here  $\phi^{(j)}$ ,  $j = 0, \dots, 2^{\lfloor N/8 \rfloor}$  stand for the sequences from  $\mathcal{P}'$ . Define the decision rule  $\tau: (X_1, \dots, X_n) \rightarrow \{0, \dots, 2^{\lfloor N/8 \rfloor}\}$  as follows. Given the observations, we compute  $\hat{\phi}_n$  and check to which of the sequences  $\phi^{(j)} \in \mathcal{P}'$  it is closer in  $\|\cdot\|$ -distance. Then we have

$$\sup_{\phi \in \mathcal{P}'} P\{\|\hat{\phi} - \phi\| \geq s/2\} = \sup_{j=0, \dots, 2^{\lfloor N/8 \rfloor}} P\{\tau \neq j | H_j\}$$

and we should evaluate from below the probability of error under the decision rule  $\tau$ . This can be done using the Fano inequality [see, e.g., Ibragimov and Has'minskii (1981), page 323]. Let  $g_j(y)$ ,  $j = 0, \dots, 2^{\lfloor N/8 \rfloor}$  denote joint density of observations  $X_1, \dots, X_n$  under the hypothesis  $H_j$ . Denote by  $\mathcal{K}(g_i, g_j)$  the Kullback–Leibler distance between the densities  $g_i$  and  $g_j$ . Then we have for  $i \neq j$ ,

$$\begin{aligned} \mathcal{K}(g_i, g_j) &\leq \sup_{i,j} E_i \ln \frac{g_i(X_1, \dots, X_n)}{g_j(X_1, \dots, X_n)} \\ &\stackrel{(a)}{=} \sup_{i,j} E_i \left[ -\frac{1}{2\sigma_\varepsilon^2} \sum_{t=1}^n \left( \left( X_t - \sum_{k=1}^N \phi_k^{(i)} X_{t-k} \right)^2 - \left( X_t - \sum_{k=1}^N \phi_k^{(j)} X_{t-k} \right)^2 \right) \right] \\ &\stackrel{(b)}{=} \sup_{i,j} \frac{n}{2\sigma_\varepsilon^2} E_i \left( \sum_{k=1}^N (\phi_k^{(i)} - \phi_k^{(j)}) X_{t-k} \right)^2 \\ &\leq \frac{n}{2\sigma_\varepsilon^2} \sup_{i,j} \sum_{k,l=1}^N (\phi_k^{(i)} - \phi_k^{(j)}) (\phi_l^{(i)} - \phi_l^{(j)}) E_i[X_{t-k} X_{t-l}] \\ &= \frac{n}{2\sigma_\varepsilon^2} \sup_{i,j} (\phi^{(i)} - \phi^{(j)})' \Gamma_N^{(i)} (\phi^{(i)} - \phi^{(j)}), \end{aligned}$$

where  $E_i$  denotes expectation with respect to the distribution related to the hypothesis  $H_i$ , and  $\Gamma_N^{(i)} = \{E_i[X_{t-k} X_{t-l}]\}_{k,l=1}^N$  is the  $N \times N$  covariance matrix under  $H_i$ . Here (a) follows from the fact that  $(X_t)$  is a Gaussian process, and (b) is obtained by taking expectation with respect to the density  $g_i$ . Using the bounds established in Lemma 1 on  $\sigma_\varepsilon^2$  and on the maximal eigenvalue of the covariance matrix  $\Gamma_N$  (which are uniform over the class  $\mathcal{H}_\rho(l, L)$  and  $N$ ) we have

$$\mathcal{K}(g_i, g_j) \leq \frac{nL^2}{2} (L/l)^2 \sup_{i,j} \|\phi^{(i)} - \phi^{(j)}\|^2 \leq K_3 n N \rho^{-2N}.$$

Now set

$$(32) \quad N = \left\lceil \frac{1}{2 \ln \rho} \ln(K_4 n) \right\rceil;$$

then due to the Fano inequality we can choose a constant  $K_4$  so that under (12) probability of the error under  $\tau$  will be at least, say, 1/4. Combining (30), (31) and (32) we come to the required statement.  $\square$

B.3. Proof of Theorem 3.

Step 1. We have the following decomposition of the prediction error:

$$\begin{aligned} X_{n+1} - \widehat{X}_{n+1}^d &= \sum_{j=1}^d (\phi_j - \hat{\phi}_j) X_{n+1-j} + \sum_{j=d+1}^{\infty} \phi_j X_{n+1-j} + \varepsilon_{n+1} \\ &:= \mathcal{E}_1 + \mathcal{E}_2 + \varepsilon_{n+1}. \end{aligned}$$

Therefore,

$$E(X_{n+1} - \widehat{X}_{n+1}^d)^2 = E(\mathcal{E}_1 + \mathcal{E}_2)^2 + \sigma_\varepsilon^2 \leq 2(E|\mathcal{E}_1|^2 + E|\mathcal{E}_2|^2) + \sigma_\varepsilon^2,$$

where we have used the fact that  $\varepsilon_{n+1}$  is independent of  $X_t$ , for  $t \leq n$ .

We first establish a bound on  $E|\mathcal{E}_1|^2$ . One clearly has

$$\begin{aligned} \mathcal{E}_1 &= \sum_{j=1}^d (\hat{\phi}_j - \phi_j) X_{n+1-j} \\ &= \sum_{j=1}^d (\hat{\phi}_j - \phi_j) \varepsilon_{n+1-j} + \sum_{k=1}^{\infty} \psi_k \sum_{j=1}^d (\hat{\phi}_j - \phi_j) \varepsilon_{n+1-j-k} \\ &:= \eta_{n+1} + \sum_{k=1}^{\infty} \psi_k \eta_{n+1-k}, \end{aligned}$$

where the second equality follows from the MA( $\infty$ ) representation of the process  $(X_t)_{t \in \mathbf{Z}}$ , and  $\eta_t := \sum_{j=1}^d (\hat{\phi}_j - \phi_j) \varepsilon_{t-j}$ . Thus, we have

$$\begin{aligned} E|\mathcal{E}_1|^2 &\leq 4 \left[ E|\eta_{n+1}|^2 + E \left( \sum_{k=1}^d \psi_k \eta_{n+1-k} \right)^2 + E \left( \sum_{k=d+1}^{\infty} \psi_k \eta_{n+1-k} \right)^2 \right] \\ &= 4 \left[ E|\eta_{n+1}|^2 + \sum_{k,l=1}^d \psi_k \psi_l E[\eta_{n+1-k} \eta_{n+1-l}] \right. \\ &\quad \left. + \sum_{k,l=d+1}^{\infty} \psi_k \psi_l E[\eta_{n+1-k} \eta_{n+1-l}] \right] \\ &= 4(E|\eta_{n+1}|^2 + \mathcal{E}_{11} + \mathcal{E}_{12}). \end{aligned}$$

Let  $\mathcal{F}_{-\infty}^n$  denote the  $\sigma$ -algebra on the common probability space  $(\Omega, \mathcal{F}, P)$  that is generated by the sequence  $(\varepsilon_n, \varepsilon_{n-1}, \dots)$ . We have

$$\begin{aligned} E|\eta_{n+1}|^2 &= E \sum_{i,j=1}^d (\hat{\phi}_i - \phi_i)(\hat{\phi}_j - \phi_j) \varepsilon_{n+1-i} \varepsilon_{n+1-j} \\ (33) \quad &= E \left( \sum_{i,j=1}^d (\hat{\phi}_i - \phi_i)(\hat{\phi}_j - \phi_j) E[\varepsilon_{n+1-i} \varepsilon_{n+1-j} | \mathcal{F}_{-\infty}^{[n/2]}] \right) \\ &= \sigma_\varepsilon^2 E \|\hat{\theta}^d - \theta^d\|^2, \end{aligned}$$

where the second equality follows from the fact that  $\hat{\phi}$  is  $\mathcal{F}_{-\infty}^{\lfloor n/2 \rfloor}$ -measurable and  $\varepsilon_{n+1-i}$ ,  $i = 1, \dots, d$  are independent of  $\mathcal{F}_{-\infty}^{\lfloor n/2 \rfloor}$  because  $n/2 > d$ . Further, applying the same reasoning for  $k, l = 1, \dots, d$  we obtain

$$\begin{aligned} E[\eta_{n+1-k}\eta_{n+1-l}] &= E \sum_{i,j=1}^d (\hat{\phi}_i - \phi_i)(\hat{\phi}_j - \phi_j)\varepsilon_{n+1-k-i}\varepsilon_{n+1-l-j} \\ &= E \left( \sum_{i,j=1}^d (\hat{\phi}_i - \phi_i)(\hat{\phi}_j - \phi_j) E[\varepsilon_{n+1-k-i}\varepsilon_{n+1-l-j} | \mathcal{F}_{-\infty}^{\lfloor n/2 \rfloor}] \right) \\ &= \sigma_\varepsilon^2 E \sum_{i,j=1, i=l+j-k}^d (\hat{\phi}_i - \phi_i)(\hat{\phi}_j - \phi_j) \\ &\leq K_1 \sigma_\varepsilon^2 E \|\hat{\theta}^d - \theta^d\|^2. \end{aligned}$$

Here we again have used the fact that  $\hat{\phi}$  is  $\mathcal{F}_{-\infty}^{\lfloor n/2 \rfloor}$ -measurable, and  $\varepsilon_{n+1-l-j}$ ,  $i, l = 1, \dots, d$  are independent of  $\mathcal{F}_{-\infty}^{\lfloor n/2 \rfloor}$  because  $n/4 > d$ . Therefore,

$$(34) \quad \mathcal{E}_{11} \leq K_1 \sigma_\varepsilon^2 E \|\hat{\theta}^d - \theta^d\|^2 \left( \sum_{k=1}^d |\psi_k| \right)^2 \leq K_1 \sigma_\varepsilon^2 \frac{L^2}{(\rho-1)^2} E \|\hat{\theta}^d - \theta^d\|^2$$

(here we have taken into account that  $|\psi_k| \leq L\rho^{-k}$ ).

To bound from above  $\mathcal{E}_{12}$  note first that for  $k \geq d+1$  by the Cauchy–Schwarz inequality,

$$\begin{aligned} E|\eta_{n+1-k}|^2 &= E \left( \sum_{i=1}^d (\hat{\phi}_i - \phi_i)\varepsilon_{n+1-k-i} \right)^2 \leq \sum_{i=1}^d E \left[ \|\hat{\theta}^d - \theta^d\|^2 \varepsilon_{n+1-k-i}^2 \right] \\ &\leq K_2 d \sigma_\varepsilon^2 (E \|\hat{\theta}^d - \theta^d\|^4)^{1/2}. \end{aligned}$$

Thus, one has

$$(35) \quad \begin{aligned} \mathcal{E}_{12} &\leq K_2 d \sigma_\varepsilon^2 (E \|\hat{\theta}^d - \theta^d\|^4)^{1/2} \left( \sum_{k=d+1}^{\infty} |\psi_k| \right)^2 \\ &\leq K_2 d \sigma_\varepsilon^2 (E \|\hat{\theta}^d - \theta^d\|^4)^{1/2} \rho^{-2d} \frac{L^2}{(\rho-1)^2}. \end{aligned}$$

Combining (35), (34) and (33) we come to the bound on  $E|\mathcal{E}_1|^2$ ,

$$(36) \quad E|\mathcal{E}_1|^2 \leq K_3 \sigma_\varepsilon^2 \left[ E \|\hat{\theta}^d - \theta^d\|^2 + \frac{d\rho^{-2d}L^2}{(\rho-1)^2} (E \|\hat{\theta}^d - \theta^d\|^4)^{1/2} \right].$$

*Step 2.* Our next step is to bound from above  $E\|\hat{\theta}^d - \theta^d\|^4$ . Choose  $\kappa = 6d^2n^{-10}$ , and let  $A_\kappa$  be given by (25) with  $C_\kappa$  for  $\kappa$  in question. Write

$$E\|\hat{\theta}^d - \theta^d\|^4 = E\|\hat{\theta}^d - \theta^d\|^4 \mathbf{1}_{\{A_\kappa\}} + E\|\hat{\theta}^d - \theta^d\|^4 \mathbf{1}_{\{A_\kappa^c\}}.$$

It can be easily verified that under premise of the theorem the conditions of Lemma 4 hold. Therefore,

$$\begin{aligned} E\|\hat{\theta}^d - \theta^d\|^4 \mathbf{1}_{\{A_\kappa\}} &\leq 2^4 \|\Gamma_d^{-1}\|^4 E(\|\mathcal{J}_1 + \mathcal{J}_2 + \mathcal{J}_3\|^4 \mathbf{1}_{\{A_\kappa\}}) \\ &\leq K_4 \|\Gamma_d^{-1}\|^4 (E\|\mathcal{J}_1\|^4 + E\|\mathcal{J}_2\|^4 + E\|\mathcal{J}_3\|^4). \end{aligned}$$

Applying Lemma 5 we obtain

$$E\|\hat{\theta}^d - \theta^d\|^4 \mathbf{1}_{\{A_\kappa\}} \leq K_5 \|\Gamma_d^{-1}\|^4 l^{-4} \left( \frac{1}{n^4(\rho - 1)^4} + \frac{d}{\rho^{4d}(\rho - 1)^4} + \frac{d^2}{n^2} \right).$$

For the other term, involving the indicator of the event  $A_\kappa^c$ , the result follows the derivation in the proof of Theorem 1. In particular, we now require bounds in Lemma 5 to hold for  $(E\|\mathcal{J}_j\|^8)^{1/2}$  for  $j = 2, 3$ . It is straightforward to extend the results of the lemma; the details are omitted. Thus, we obtain

$$E\|\hat{\theta}^d - \theta^d\|^4 \mathbf{1}_{\{A_\kappa^c\}} \leq K_6 n^4 l^{-4} \left( \frac{\kappa}{n^4(\rho - 1)^4} + \frac{d^2 \sqrt{\kappa}}{\rho^{4d}(\rho - 1)^4} + \frac{d^2 \sqrt{\kappa}}{n^2} \right)$$

and finally, substituting  $\kappa = 4d^2 n^{-10}$  and combining the above bounds we have

$$(E\|\hat{\theta}^d - \theta^d\|^4)^{1/2} \leq K_7 (L^4/l^6) \left( \frac{1}{n^2(\rho - 1)^2} + \frac{d}{\rho^{2d}(\rho - 1)^2} + \frac{d}{n} \right).$$

Thus, it follows from (36) and (29) that

$$E|\mathcal{E}_1|^2 \leq K_8 \sigma_\varepsilon^2 (L^4/l^6) \left( \frac{1}{n^2(\rho - 1)^2} + \frac{d}{\rho^{2d}(\rho - 1)^2} + \frac{d}{n} \right) \left( 1 + \frac{d\rho^{-2d}L^2}{(\rho - 1)^2} \right).$$

*Step 3.* Now we complete the proof of the theorem. We have the following upper bound on  $E|\mathcal{E}_2|^2$ :

$$E|\mathcal{E}_2|^2 = E \left| \sum_{j=d+1}^\infty \phi_j X_{n+1-j} \right|^2 \leq \left( \sum_{j=d+1}^\infty \phi_j \right)^2 \leq l^{-2} \rho^{-2d} (\rho - 1)^{-2},$$

where we have used the fact that  $|\phi_k| \leq l^{-1} \rho^{-k}$ . Combining the above bounds on  $E|\mathcal{E}_1|^2$  and  $E|\mathcal{E}_2|^2$  we come to (14). This completes the proof of the theorem.  $\square$

### C. Proofs of auxiliary results.

#### C.1. Proof of Lemma 2.

*Step 1.* First observe that the process  $(X_t)_{t \in \mathbf{Z}}$  is *strongly mixing*. We recall the definition of the strong mixing condition [cf. Bradley (1986), page 169]. For  $-\infty \leq s \leq k \leq \infty$  let  $\mathcal{F}_s^k$  denote the  $\sigma$ -algebra generated by  $(X_s, X_{s+1}, \dots, X_k)$ . The process  $(X_t)_{t \in \mathbf{Z}}$  is said to be strongly mixing if

$$\alpha_{\mathbf{X}}(\tau) = \sup_{s \in \mathbf{Z}} \alpha(\mathcal{F}_{-\infty}^s, \mathcal{F}_{s+\tau}^\infty) \rightarrow 0 \quad \text{as } \tau \rightarrow \infty,$$

where

$$\alpha(\mathcal{F}_{-\infty}^s, \mathcal{F}_{s+\tau}^\infty) = \sup_{A \in \mathcal{F}_{-\infty}^s, B \in \mathcal{F}_{s+\tau}^\infty} |P(AB) - P(A)P(B)|.$$

Since  $(X_t)_{t \in \mathbf{Z}}$  is Gaussian and stationary, we have  $\alpha_X(\tau) = \alpha(\mathcal{F}_{-\infty}^0, \mathcal{F}_\tau^\infty)$ . The strong mixing coefficient  $\alpha_X(\tau)$  is bounded from above by the maximal correlation coefficient,

$$(37) \quad \alpha_X(\tau) \leq \sup_{\zeta_1, \zeta_2} E(\zeta_1 \zeta_2),$$

where the supremum in (37) is taken over all pairs of zero mean random variables  $(\zeta_1, \zeta_2)$  such that  $\zeta_1 \in \mathcal{F}_{-\infty}^0$ ,  $\zeta_2 \in \mathcal{F}_\tau^\infty$  and  $E|\zeta_1|^2 = E|\zeta_2|^2 = 1$ . Further, let  $\mathcal{E}_{\tau-1}(f)$  denote the error of the best approximation of the spectral density  $f(\lambda)$  by trigonometric polynomials of the degree  $\leq \tau-1$  on the interval  $[-\pi, \pi]$  in the uniform norm. We have

$$(38) \quad \begin{aligned} \mathcal{E}_{\tau-1}(f) &\leq \frac{1}{\pi} \max_{\lambda \in [-\pi, \pi]} \left| \sum_{k=\tau}^{\infty} \gamma(k) \cos(\lambda k) \right| \leq \frac{1}{\pi} \sum_{k=\tau}^{\infty} |\gamma(k)| \\ &\leq \frac{L^2 \sigma_\varepsilon^2 \rho^2}{\pi(\rho-1)^2} \rho^{-\tau}, \end{aligned}$$

where the last inequality follows from (18). It is well known [cf. Ibragimov and Rozanov (1978), page 146] that for a stationary process with continuous and strictly positive spectral density the maximal correlation coefficient does not exceed  $[\min_{\lambda \in [-\pi, \pi]} f(\lambda)]^{-1} \mathcal{E}_{\tau-1}(f)$ . Taking into account (20) we obtain

$$(39) \quad \alpha_X(\tau) \leq 2(L/l)^2 \left( \frac{\rho}{\rho-1} \right)^2 \rho^{-\tau}.$$

Now fix integer number  $k$  and define

$$U_{t,k} = \frac{1}{n} X_t X_{t+k} - \frac{\gamma(k)}{n}, \quad t \in \mathbf{Z}.$$

Without loss of generality we assume that  $k$  is a nonnegative integer number. Let  $\mathcal{U}_{t,k}^s$  be the  $\sigma$ -algebra generated by  $(U_{t,k}, U_{t+1,k}, \dots, U_{s,k})$ . Observe that  $\mathcal{U}_{-\infty,k}^t \subseteq \mathcal{F}_{-\infty}^{t+k}$  and  $\mathcal{U}_{t+\tau,k}^\infty \subseteq \mathcal{F}_{t+\tau}^\infty$ . This implies that the process  $(U_{t,k})_{t \in \mathbf{Z}}$  is also strongly mixing with the rate

$$\alpha_U(\tau) \leq \alpha_X(\tau - k) \quad \forall \tau > k.$$

For  $\tau \leq k$  we have the following trivial inequality  $\alpha_U(\tau) \leq 1$ .

*Step 2.* To complete the proof of the lemma we need the following two auxiliary statements, adapted from Saulis and Statulevičius [(1991), Theorem 4.17, Lemma 2.4].

LEMMA 6. Let  $(Y_t)_{t \in \mathbb{Z}}$  be a strongly mixing random process,  $S_n = \sum_{t=1}^n Y_t$ , and  $\text{cum}_p(S_n)$  be the  $p$ th order cumulant of the sum  $S_n$ . For  $\nu > 0$  define the function

$$\Lambda_n[\alpha_Y, \nu] = \max \left\{ 1; \sum_{\tau=0}^n [\alpha_Y(\tau)]^{1/\nu} \right\}.$$

If for some  $\mu \geq 0, H > 0$

$$E|Y_t|^p \leq (p!)^{\mu+1} H^p, \quad t = 1, \dots, n, \quad p = 2, 3, \dots,$$

then  $|\text{cum}_p(S_n)| \leq 2^{p(1+\mu)+1} 12^{p-1} (p!)^{2+\mu} H^p \{\Lambda_n[\alpha_X, 2(p-1)]\}^{p-1} n$ .

For definition of the cumulants see, for example, Brillinger [(1975), page 19].

LEMMA 7. Let  $\xi$  be an arbitrary random variable with  $E\xi = 0$ . If there exist  $\mu_1 \geq 0, H_1 > 0$  and  $\Delta > 0$  such that

$$|\text{cum}_p(\xi)| \leq \left(\frac{p!}{2}\right)^{1+\mu_1} \frac{H_1}{\Delta^{p-2}}, \quad p = 2, 3, \dots,$$

then

$$P(|\xi| \geq x) \leq \begin{cases} \exp\{-x^2/(4H_1)\}, & 0 \leq x \leq (H_1^{1+\mu_1} \Delta)^{1/(2\mu_1+1)}, \\ \exp\{-(x\Delta)^{1/(1+\mu_1)}/4\}, & x \geq (H_1^{1+\mu_1} \Delta)^{1/(2\mu_1+1)}. \end{cases}$$

Step 3. Using Lemma 6 we will derive the upper bound on the cumulants of the sum  $\sum_{t=1}^n U_{t,k}$ , and then applying Lemma 7 we will obtain the required exponential inequality. First, we verify the conditions of Lemma 6 in order to apply it to the process  $(U_{t,k})_{t \in \mathbb{Z}}$ . Observe that for any natural  $p$  we have

$$\begin{aligned} E|U_{t,k}|^p &\leq \frac{2^{p-1}}{n^p} (|\gamma(k)|^p + [E|X_t|^{2p} E|X_{t+k}|^{2p}]^{1/2}) \\ &\leq \frac{2^{p-1}}{n^p} (1 + p!2^p) \leq \frac{2^{2p} p!}{n^p}, \end{aligned}$$

where the second inequality follows from the fact that  $X_t$  is a Gaussian random variable,  $|\gamma(k)| \leq E|X_t|^2 = 1$ , and  $E|X_t|^{2p} \leq p!2^p$ . Further,

$$\begin{aligned} \Lambda_n[\alpha_U, 2(p-1)] &\leq k + \sum_{\tau=k}^n [\alpha_X(\tau-k)]^{1/(2p-2)} \\ &\stackrel{(a)}{\leq} k + \left(\frac{2L\rho}{l(\rho-1)}\right)^{2/(2p-2)} \sum_{\tau=0}^{n-k} \rho^{-\tau/(2p-2)} \\ &\leq k + \left(\frac{2L\rho}{l(\rho-1)}\right)^{1/(p-1)} \frac{\rho^{1/(2p-2)}}{\rho^{1/(2p-2)} - 1} \\ &\stackrel{(b)}{\leq} k + \left(\frac{2L\rho}{l(\rho-1)}\right)^{1/(p-1)} \left(1 + \frac{2p-2}{\ln \rho}\right) \end{aligned}$$

where (a) follows from the bound in (39), and (b) follows from the elementary inequality  $\exp(x) - 1 \geq x$  for  $x \geq 0$ . Thus, one has

$$\begin{aligned} \{\Lambda_n[\alpha_U, 2(p-1)]\}^{p-1} &\leq 2^{p-2} \left[ k^{p-1} + \frac{2L\rho}{l(\rho-1)} (p-1)^{p-1} \left(1 + \frac{2}{\ln \rho}\right)^{p-1} \right] \\ &\leq 2^{p-2} k_*^{p-1} (p-1)! e^{p-1} \left(1 + \frac{2}{\ln \rho}\right)^{p-1} \left(1 + \frac{2L\rho}{l(\rho-1)}\right) \\ &\leq (4e)^{p-1} p! (k_* r)^{p-1} M, \end{aligned}$$

where the inequality  $(p-1)^{p-1} \leq (p-1)! e^{p-1}$  has been used, and  $k_*$ ,  $r$  and  $M$  are defined in (7). Setting  $\mu = 0$  and  $H = 4n^{-1}$  we see that Lemma 6 applies for  $(U_{t,k})_{t \in \mathbf{Z}}$ , and thus

$$(40) \quad \left| \text{cum}_p \left( \sum_{t=1}^n U_{t,k} \right) \right| \leq 2^{3p+1} 12^{2p-2} (p!)^3 (k_* r)^{p-1} M n^{-p+1}.$$

Now, to apply Lemma 7, put  $\mu_1 = 2$ ,  $H_1 = C_1 M k_* r n^{-1}$  and  $\Delta = n(C_2 k_* r)^{-1}$ , where  $C_1$ , and  $C_2$  are absolute constants ( $C_1 = 2^{10} 12^2$ ,  $C_2 = 2^3 12^2$ ). It is immediately seen that the conditions of Lemma 7 hold for the parameters in question. Applying Lemma 7 completes the proof.  $\square$

*C.2. Proof of Lemma 3.* The basis is the same argument as the one in Theorem 2. Recall the definition of  $V_{ij}$ :

$$V_{ij} = \frac{1}{n} \sum_{\tau=1-j}^0 X_\tau X_{\tau+j-i};$$

we have  $EV_{ij} = jn^{-1}\gamma(j-i)$ , whence  $E|V_{ij}| \leq dn^{-1}$ . Fix  $i, j \in \{1, \dots, d\}$  and define

$$U_t = \frac{1}{n} [X_t X_{t+j-i} - \gamma(j-i)],$$

then  $V_{ij} - EV_{ij} = \sum_{t=1-j}^0 U_t$ . For any natural number  $p$  one has

$$\begin{aligned} E|U_t|^p &\leq \frac{2^{p-1}}{n^p} (|\gamma(j-i)|^p + [E|X_t|^{2p} E|X_{t+j-i}|^{2p}]^{1/2}) \\ &\leq p! 2^{2p} n^{-p}, \end{aligned}$$

where the second inequality follows from the bound on  $E|X_t|^{2p}$  established in Step 4 of the proof of Theorem 2 and the fact that  $|\gamma(k)| \leq 1$ ,  $\forall k$ .

Taking into account the strong mixing property of the sequence  $(U_t)_{t \in \mathbf{Z}}$  and the fact that  $\Lambda_j[\alpha_U, 2(p-1)] \leq j$  for  $1 \leq j \leq d$ , we can apply Lemma 6 with  $\mu = 0$  and  $H = 4n^{-1}$ . Thus,

$$(41) \quad \left| \text{cum}_p \left( \sum_{t=1-j}^0 U_t \right) \right| \leq 2^{3p+1} 12^{p-1} (p!)^2 \frac{d^{p-1}}{n^{p-1}}.$$



It is immediately seen that the conditions of Lemma 7 hold with  $\mu_1 = 1$ ,  $\Delta = n(C_2d)^{-1}$ , and  $H_1 = C_1dn^{-1}$ , and  $C_1$  and  $C_2$  may be chosen as in Theorem 2. The same argument is valid for  $W_{ij}$ . This completes the proof.  $\square$

C.3. *Proof of Lemma 4.* (i) First we establish that  $P(A_\kappa) \geq 1 - \kappa$ . We have

$$\begin{aligned} P(A_\kappa^c) &\leq P\left\{\max_{i,j=1,\dots,d} (|V_{ij} - EV_{ij}| + |W_{ij} - EW_{ij}| + |(\widehat{Q} - \Gamma_d)_{ij}|) > C_\kappa\right\} \\ &\leq P\left\{\max_{i,j=1,\dots,d} |V_{ij} - EV_{ij}| > C_\kappa/3\right\} + P\left\{\max_{i,j=1,\dots,d} |W_{ij} - EW_{ij}| > C_\kappa/3\right\} \\ &\quad + P\left\{\max_{i,j=1,\dots,d} |(\widehat{Q} - \Gamma_d)_{ij}| > C_\kappa/3\right\} \\ &:= P_1 + P_2 + P_3. \end{aligned}$$

It can be easily verified that under the condition of (27),  $C_\kappa \leq (d/n)^{2/5}$ . Thus we may apply the results of Lemma 2 and Lemma 3 in the range of “moderate” deviations. Note that  $P_3$  can be bounded using the first inequality in (21) and the Toeplitz structure of the matrix  $\widehat{Q} - \Gamma_d$ ,

$$P_3 \leq 2d \exp\left\{-\frac{C_\kappa^2 n}{36C_1 drM}\right\}.$$

The probabilities  $P_1$  and  $P_2$  are bounded, in turn, using the first inequality in (24),

$$P_i \leq 2d^2 \exp\left\{-\frac{C_\kappa^2 n}{36C_1 d}\right\}, \quad i = 1, 2.$$

Thus using the fact that  $r \geq 1$  and  $M \geq 1$ , we have

$$P(A_\kappa^c) \leq 6d^2 \exp\left\{-\frac{C_\kappa^2 n}{36C_1 drM}\right\}.$$

Now, it is straightforward to verify that the choice of  $C_\kappa$  is made so as to satisfy  $P(A_\kappa^c) \leq \kappa$ .

(ii) Suppose that the event  $A_\kappa$  holds. Since  $\widetilde{Q}$  is a symmetric  $d \times d$  matrix we have

$$(42) \quad \|\widetilde{Q}\| = \lambda_{\max}(\widetilde{Q}) \leq \max_i \left\{ \sum_j |\widetilde{Q}_{ij}| \right\} \leq dC_\kappa.$$

Therefore (23) and the definition of  $\widetilde{Q}$  together imply that

$$(43) \quad \|\mathcal{Q}^{-1}\Gamma_d\| \leq \frac{1}{1 - \|\Gamma_d^{-1}(\widetilde{Q} + n^{-1}I_d)\|},$$

provided that  $\|\Gamma_d^{-1}(\widetilde{Q} + n^{-1}I_d)\| < 1$ . This condition will subsequently be verified. Taking into account (19) we have

$$(44) \quad \begin{aligned} \|\Gamma_d^{-1}(\widetilde{Q} + n^{-1}I_d)\| &\leq \|\Gamma_d^{-1}\|(n^{-1} + dC_\kappa) \\ &\leq (L/l)^2(n^{-1} + dC_\kappa) \leq 1/2, \end{aligned}$$

where the last inequality follows from the condition imposed in (28). Thus, (44) along with (43) imply the statement of the lemma for the case where the event  $A_\kappa$  holds.

(iii) Now consider the case of  $\omega \in A_\kappa^c$ . Independently of the event  $A_\kappa$ , the matrix  $Q$  is positive-definite,  $\lambda_{\min}[Q] \geq n^{-1}$  and whence  $\lambda_{\max}[Q^{-1}] \leq n$ . Since  $Q^{-1}$  is symmetric, we obtain immediately that  $\|Q^{-1}\| \leq n$ . This completes the proof of the lemma.  $\square$

C.4. *Proof of Lemma 5.* The upper bound on  $\|\mathcal{S}_1\|$  follows immediately from (16):

$$\|\mathcal{S}_1\| = n^{-1}\|\theta^d\| = n^{-1}\left(\sum_{j=1}^d |\phi_j|^2\right)^{1/2} \leq \frac{1}{nl(\rho-1)}.$$

Let us denote the  $k$ th component of  $\mathcal{S}_2$  as

$$\mathcal{S}_{2,k} := \frac{1}{n} \sum_{t=1}^n X_{t-k} \sum_{j=d+1}^{\infty} \phi_j X_{t-j}, \quad k = 1, 2, \dots, d.$$

We have

$$\begin{aligned} E|\mathcal{S}_{2,k}|^4 &= \sum_{j_1, \dots, j_4=d+1}^{\infty} \phi_{j_1} \phi_{j_2} \phi_{j_3} \phi_{j_4} \\ &\quad \times E\left[\hat{\gamma}(k-j_1)\hat{\gamma}(k-j_2)\hat{\gamma}(k-j_3)\hat{\gamma}(k-j_4)\right], \end{aligned}$$

where  $\hat{\gamma}(k-j) = n^{-1} \sum_{t=1}^n X_{t-k} X_{t-j}$ . Applying repeatedly the Cauchy-Schwartz inequality and taking into account that  $(X_t)_{t \in \mathbb{Z}}$  is Gaussian and stationary with  $E|X_t|^2 = 1$ , we obtain

$$E[\hat{\gamma}(k-j_1)\hat{\gamma}(k-j_2)\hat{\gamma}(k-j_3)\hat{\gamma}(k-j_4)] \leq E|X_t|^8 \leq 105.$$

Therefore,

$$E|\mathcal{S}_{2,k}|^4 \leq E|X_t|^8 \left(\sum_{j=d+1}^{\infty} |\phi_j|\right)^4 \leq \frac{105}{\rho^{4d} l^4 (\rho-1)^4},$$

and thus

$$E\|\mathcal{S}_2\|^4 = E \sum_{k,l=1}^d |\mathcal{S}_{2,k}|^2 |\mathcal{S}_{2,l}|^2 \leq \frac{105d^2}{\rho^{4d} l^4 (\rho-1)^4}.$$

Now we derive an upper bound on  $E\|\mathcal{S}_3\|^4$ . Denote

$$\mathcal{S}_{3,k} := \frac{1}{n} \sum_{t=1}^n X_{t-k} \varepsilon_t = \frac{S_n}{n}, \quad k = 1, 2, \dots, d.$$

To bound  $E|\mathcal{S}_{3,k}|^4 = n^{-4}|S_n|^4$  from above we note that  $\{S_i, \mathcal{F}_{-\infty}^i, 1 \leq i \leq n\}$  is a martingale ( $\mathcal{F}_{-\infty}^i = \sigma(\varepsilon_i, \varepsilon_{-1}, \dots)$ ). Therefore due to Burkholder's inequality [see, e.g., Hall and Heyde (1980), page 23] we have

$$E|S_n|^4 = E\left|\sum_{t=1}^n X_{t-k}\varepsilon_t\right|^4 \leq K_1 E\left|\sum_{t=1}^n (X_{t-k}\varepsilon_t)^2\right|^2,$$

where  $K_1$  is an absolute constant. Thus,

$$E|\mathcal{S}_{3,k}|^4 \leq \frac{K_1}{n^4} \sum_{t,\tau=1}^n E[X_{t-k}^2 X_{\tau-k}^2 \varepsilon_t^2 \varepsilon_\tau^2] \leq \frac{K_1}{n^2} E|X_t|^4 E|\varepsilon_t|^4 \leq K_2 \frac{\sigma_\varepsilon^4}{n^2},$$

and finally

$$E\|\mathcal{S}_3\|^4 \leq K_2 \frac{d^2 \sigma_\varepsilon^4}{n^2} \leq K_2 \frac{d^2}{l^4 n^2},$$

where the last inequality follows from (17). This completes the proof.  $\square$

**Acknowledgment.** We thank the referees for some helpful comments on a previous version of this manuscript.

## REFERENCES

- AKAIKE, H. (1970). Statistical predictor identification. *Ann. Inst. Statist. Math.* **21** 243–247.
- AKAIKE, H. (1974). A new look at the statistical model identification. *IEEE Trans. Automat. Control* **AC-19** 716–723.
- AN, H.-Z., CHEN, Z.-G. C. and HANNAN, E. J. (1982). Autocorrelation, autoregression and autoregressive approximation. *Ann. Statist.* **10** 926–936.
- BERK, K. N. (1974). Consistent autoregressive spectral estimates. *Ann. Statist.* **2** 489–502.
- BHANSALI, R. J. (1981). Effects of not knowing the order of an autoregressive process on the mean-squared error of prediction. *J. Amer. Statist. Assoc.* **76** 588–597.
- BHANSALI, R. J. (1986). Asymptotically efficient selection of the order by the criterion autoregressive transfer function. *Ann. Statist.* **14** 315–325.
- BOSQ, D. (1996). *Nonparametric Statistics for Stochastic Processes*. Springer, New York.
- BRADLEY, R. C. (1986). Basic properties of strong mixing conditions. In *Dependence in Probability and Statistics* (E. Eberlein and M. Taqqu, eds.) 165–192. Birkhäuser, Boston.
- BRILLINGER, D. R. (1975). *Times Series: Data Analysis and Theory*. Holt, Rinehart and Winston, New York.
- BÜHLMANN, P. (1995). Moving average representations of autoregressive approximations. *Stochastic Process. Appl.* **60** 331–342.
- DOUKHAN, P. (1994). *Mixing: Properties and Examples*. Springer, New York.
- EFROMOVICH, S. (1998). Data-driven efficient estimation of the spectral density. *J. Amer. Statist. Assoc.* **92** 762–769.
- GOLDENSHLUGER, A. (1998). Nonparametric estimation of transfer functions: rates of convergence and adaptation. *IEEE Trans. Inform. Theory* **44** 644–658.
- GOLUBEV, G. and LEVIT, B. (1996). Asymptotically efficient estimation for analytic distributions. *Math. Methods Statist.* **5** 357–368.
- GOLUBEV, G. K., LEVIT, B. and TSYBAKOV A. (1996). Asymptotically efficient estimation of analytic functions in Gaussian noise. *Bernoulli* **2** 167–181.
- GERENCSÉR, L. (1992). AR( $\infty$ ) estimation and non-parametric stochastic complexity. *IEEE Trans. Inform Theory* **38** 1768–1778.

- GRENANDER, U. and SZEGÖ, G. (1984). *Toeplitz Forms and Its Applications*, 2nd ed. Chelsea, New York.
- HALL, P. and HEYDE, C. C. (1980). *Martingale Limit Theory and Its Application*. Academic Press, New York.
- HANNAN, E. J. and KAVALIERIS, L. (1986). Regression, autoregression models. *J. Time Ser. Anal.* **7** 27–49.
- IBRAGIMOV, I. A. and HAS'MINSKII, R. Z. (1981). *Statistical Estimation*. Springer, New York.
- IBRAGIMOV, I. A. and ROZANOV, Y. A. (1978). *Gaussian Random Processes*. Springer, New York.
- KOROSTELEV, A. P. and TSYBAKOV, A. B. (1993). *Minimax Theory of Image Reconstruction*. Springer, New York.
- PARZEN, E. (1974). Some recent advances in time series modelling. *IEEE Trans. Automat. Control* **AC-19** 723–730.
- PARZEN, E. (1983). Autoregressive spectral estimation. In *Time Series in the Frequency Domain* 221–247. North-Holland, Amsterdam.
- RISSANEN, J. (1983). Universal coding, information, prediction, and estimation. *IEEE Trans. Inform. Theory* **30** 629–636.
- RUDIN, W. (1974). *Real and Complex Analysis*, 2nd ed. McGraw-Hill, New York.
- SAULIS, L. and STATULEVIČUS, V. A. (1991). *Limit Theorems for Large Deviations*. Kluwer, Dordrecht.
- SCHWARZ, G. (1978). Estimating the dimension of a model. *Ann. Statist.* **6** 461–464.
- SHIBATA, R. (1980). Asymptotically efficient selection of the order of the model for estimating parameters of a linear process. *Ann. Statist.* **8** 147–164.
- SHIBATA, R. (1981). An optimal autoregressive spectral estimate. *Ann. Statist.* **9** 300–306.

DEPARTMENT OF STATISTICS  
UNIVERSITY OF HAIFA  
HAIFA  
ISRAEL 31905  
E-MAIL: goldensh@rstat.haifa.ac.il

INFORMATION SYSTEMS LAB  
STANFORD UNIVERSITY  
STANFORD, CALIFORNIA 94305-9510  
E-MAIL: assaf@isl.stanford.edu