

# ADAPTIVE ESTIMATES FOR AUTOREGRESSIVE PROCESSES\*

RUDOLF BERAN

(Received Dec. 4, 1972; revised Jan. 6, 1975)

## Abstract

Let  $\{X_t: t=0, \pm 1, \pm 2, \dots\}$  be a stationary  $r$ th order autoregressive process whose generating disturbances are independent identically distributed random variables with marginal distribution function  $F$ . Adaptive estimates for the parameters of  $\{X_t\}$  are constructed from the observed portion of a sample path. The asymptotic efficiency of these estimates relative to the least squares estimates is greater than or equal to one for all regular  $F$ . The nature of the adaptive estimates encourages stable behavior for moderate sample sizes. A similar approach can be taken to estimation problems in the general linear model.

## 1. Introduction

A discrete stochastic process  $\{X_t: t=0, \pm 1, \pm 2, \dots\}$  may be called a stationary  $r$ th order autoregressive process if it has the following properties: the process is strictly stationary and satisfies a difference equation of the form

$$(1.1) \quad X_t - \mu = \sum_{j=1}^r \alpha_j (X_{t-j} - \mu) + E_t,$$

where  $\{E_t\}$  is a sequence of independent identically distributed random variables with mean 0 and finite variance  $\sigma^2$ , the parameters  $\mu$ ,  $\{\alpha_j\}$  are real-valued, and the roots of the polynomial equation

$$(1.2) \quad x^r = \sum_{j=1}^r \alpha_j x^{r-j}$$

all have modulus less than one. Under these assumptions, the difference equation (1.1) has a unique solution, expressible in the form

---

\* This research was partially supported by National Science Foundation Grant GP-31091X. American Mathematical Society 1970 subject classification. Primary 62N10; Secondary 62G35. Key words and phrases: autoregressive process, adaptive estimates, robust estimates.

$$(1.3) \quad X_t = \mu + \sum_{j=0}^{\infty} \beta_j E_{t-j},$$

the sum on the right converging in mean square (Anderson [1] or Mann and Wald [7]).

Stationary autoregressive processes serve as convenient parametric models for certain time series; their prediction theory is particularly simple. Estimation of the parameters from the observed portion  $(X_1, X_2, \dots, X_N)$  of a sample path is commonly carried out by writing (1.1) in the form

$$(1.4) \quad X_t = \alpha_0 + \sum_{j=1}^r \alpha_j X_{t-j} + E_t$$

and by applying the method of least squares. If  $X$  is the  $(N-r) \times (r+1)$  matrix whose  $i$ th row is  $(1, X_{r+i-1}, X_{r+i-2}, \dots, X_i)$  and  $Y$  is the  $(N-r) \times 1$  vector  $(X_{r+1}, X_{r+2}, \dots, X_N)^T$ , the least squares estimate (LSE) of  $\rho = (\alpha_0, \alpha_1, \dots, \alpha_r)^T$  based upon  $(X_1, X_2, \dots, X_N)$ ,  $N \geq 2r+1$ , is

$$(1.5) \quad \hat{\rho}_N = (X^T X)^{-1} X^T Y.$$

The asymptotic distribution of  $N^{1/2}(\hat{\rho}_N - \rho)$  as  $N \rightarrow \infty$  is normal  $(0, \sigma^2 \Gamma^{-1})$ , where  $\Gamma = (N-r)^{-1} E(X^T X)$  (Anderson [1] or Mann and Wald [7]).

Let  $F$  denote the marginal distribution function of  $E_t$  and let  $(\cdot, \cdot)$  and  $\|\cdot\|$  denote the inner product and norm in  $L_2(F)$ . Under regularity conditions (cf. Section 2), there exists an estimate  $\hat{\rho}_N(\phi_F)$ , depending on

$$(1.6) \quad \phi_F(x) = -f'(x)/f(x),$$

$f$  being the density of  $F$ , such that the asymptotic distribution of  $N^{1/2}(\hat{\rho}_N(\phi_F) - \rho)$  is normal  $(0, \|\phi_F\|^{-2} \Gamma^{-1})$ . The asymptotic efficiency of the LSE  $\hat{\rho}_N$  relative to  $\hat{\rho}_N(\phi_F)$  is less than or equal to one, with equality if and only if  $F$  is normal. Thus, the problem arises of constructing a practical estimate of  $\rho$  whose asymptotic performance will dominate that of the LSE for all  $F$  of interest. This paper develops one possible answer, under mild regularity assumptions on  $F$ .

## 2. Linearized estimates

Let  $\rho_0$  denote the parameter vector of the autoregressive model under which the data is realized. For every  $\rho = (\alpha_0, \alpha_1, \dots, \alpha_r)^T \in R^{r+1}$ , define the residual process

$$(2.1) \quad R_t(\rho) = X_t - \alpha_0 - \sum_{j=1}^r \alpha_j X_{t-j}, \quad r+1 \leq t \leq N.$$

Corresponding to any scalar-valued function  $\phi$  defined on  $R^1$ , let  $\Psi(\rho)$

denote the  $(N-r) \times 1$  vector  $(\phi \circ R_{r+1}(\rho), \phi \circ R_{r+2}(\rho), \dots, \phi \circ R_N(\rho))^T$ , where  $\phi \circ R_i(\rho)$  denotes  $\phi(R_i(\rho))$ . Under the usual assumptions, a conditional maximum likelihood estimate for  $\rho_0$  must be a solution to the equation

$$(2.2) \quad X^T \Phi_F(\rho) = 0 ;$$

$\Phi_F$  is associated with  $\phi_F$  in the obvious way. Consequently an  $M$ -estimate for  $\rho_0$  may be defined (Huber [4]) as any estimate which satisfies an equation of the form

$$(2.3) \quad X^T \Psi(\rho) = 0$$

for some score function  $\phi$ .

Application of Newton's method to (2.3), with the LSE  $\hat{\rho}_N$  as starting point, suggests the linearized estimate

$$(2.4) \quad \hat{\rho}_N(\phi) = \hat{\rho}_N + [\hat{D}(\phi, F)]^{-1} (X^T X)^{-1} X^T \Psi(\hat{\rho}_N) ,$$

where  $\hat{D}(\phi, F)$  is a consistent estimate of a functional  $D(\phi, F)$  defined below ; one possibility is  $D(\phi, F) = (\phi', 1)$ . This technique for adjusting estimates is well-known in the literature if  $\phi = \phi_F$  (for example, see LeCam [6]). We will study the asymptotic behavior of  $\hat{\rho}_N(\phi)$  as a preliminary to the construction of adaptive estimates.

A basic result which underlies the asymptotic theory developed in this paper is

LEMMA 2.1. *If  $\{X_t\}$  is a stationary  $r$ th order autoregressive process, then*

$$N^{-1/2} \max_{1 \leq t \leq N} |X_t| \xrightarrow{p} 0$$

$$N^{-1} X^T X \xrightarrow{p} \Gamma \quad \text{nonsingular}$$

as  $N \rightarrow \infty$ .

The first property follows from  $E X_t^2 < \infty$ , which is implied by our definition of  $X_t$ ; the second is proved in Anderson [1], for example.

The following assumptions on  $\phi$  and  $F$  will be used :

- A1.  $\phi = \phi_+ - \phi_-$ , where  $\phi_{\pm} \in L_2(F)$  and is monotone nondecreasing ;  $(\phi, 1) = 0$ .
- A2.  $\lim_{h \rightarrow 0} \|\phi_{\pm}(x+h) - \phi_{\pm}(x-h)\|^2 = 0$ , and for some  $\epsilon > 0$ ,  
 $\sup_{|a| \leq \epsilon, |h| \leq \epsilon} |h|^{-1} [(\phi_{\pm}(x+a+h) - \phi_{\pm}(x+a)), 1] < \infty$ .
- A3.  $\lim_{h \rightarrow 0} (2h)^{-1} [(\phi_{\pm}(x+h) - \phi_{\pm}(x-h)), 1] = D(\phi_{\pm}, F) < \infty$ .

In applications of our results,  $\phi$  will belong to a known, relatively small family of functions. With this family specified in advance, assumptions A1, A2, A3 amount to regularity assumptions on  $F$ . For the sake of

brevity in exposition, any  $F$  satisfying A1, A2, A3 for a given family of functions  $\phi$  will be called *regular*.

For every  $z \in R^{r+1}$ , let  $|z|$  denote  $\max_{1 \leq i \leq r+1} |z_i|$ . Let  $D(\phi, F) = D(\phi_+, F) - D(\phi_-, F)$  and let  $\hat{D}(\phi, F)$  denote a consistent estimate of  $D(\phi, F)$ .

**THEOREM 2.1.** *If A1, A2, A3 are satisfied, then for every  $C > 0$*

$$(2.6) \quad N^{-1/2} \sup_{|\rho - \rho_0| \leq CN^{-1/2}} |X^T \Psi(\rho) - X^T \Psi(\rho_0) + (X^T X)(\rho - \rho_0) \hat{D}(\phi, F)| \xrightarrow{p} 0$$

*in  $\rho_0$ -probability as  $N \rightarrow \infty$ .*

The proof of this theorem is deferred to Section 5. A similar result for the general linear model has been established by Bickel [2].

**THEOREM 2.2.** *If A1, A2, A3 are satisfied, the asymptotic distribution of  $N^{1/2}(\hat{\rho}_N(\phi) - \rho_0)$  as  $N \rightarrow \infty$  is normal  $(0, V(\phi, F)\Gamma^{-1})$ , where*

$$(2.7) \quad V(\phi, F) = \|\phi\|^2 / D^2(\phi, F).$$

**PROOF.** Since  $N^{1/2}(\hat{\rho}_N - \rho_0)$  is bounded in probability asymptotically, it follows from (2.4) and Theorem 2.1 that

$$(2.8) \quad N^{1/2}(\hat{\rho}_N(\phi) - \rho_0) = N^{1/2}[\hat{D}(\phi, F)]^{-1}(X^T X)^{-1}X^T \Psi(\rho_0) + o_p(1).$$

Let  $\mathcal{A}_t$  be the  $\sigma$ -algebra generated by  $(X_1, X_2, \dots, X_t)$ ,  $t \geq 1$ . For every  $c \in R^{r+1}$ , let

$$(2.9) \quad S_N(c) = c^T X^T \Psi(\rho_0).$$

Because of A1 and (1.3),  $\{S_N(c), \mathcal{A}_N; N \geq 2r+1\}$  is a martingale. Application of a suitable central limit theorem (Brown [3], p. 60) shows that the asymptotic distribution of  $N^{-1/2}S_N(c)$  is normal  $(0, \|\phi\|^2 c^T \Gamma c)$ . The theorem follows.

Suppose that  $F$  has an absolutely continuous density  $f$  and  $\phi_F$  satisfies A1, A2, A3 with  $D(\phi_F, F) = (\phi'_F, 1) = \|\phi_F\|^2$ . Then the estimate  $\hat{\rho}_N(\phi_F)$ , defined according to (2.4), achieves the minimal asymptotic covariance matrix  $\|\phi_F\|^{-2}\Gamma^{-1}$  mentioned earlier.

### 3. Adaptive estimates

Since  $F$  is usually not known, the estimate  $\hat{\rho}_N(\phi_F)$  cannot be found in practice. A natural idea is to estimate  $\phi_F$  from the data and use this estimate in place of  $\phi_F$  in (2.4). While theoretically possible, this approach encounters the difficulty that consistent estimates of  $\phi_F$  may converge very slowly as the sample size increases.

In estimating  $\phi_F$ , as in estimating densities, there is a trade-off

between asymptotic variance and asymptotic bias. This suggests a more modest approach that separates these two considerations: Replace  $\phi_F$  with an approximation  $\phi_{F,H}$  which is easier to estimate from the data but keeps  $V(\phi_{F,H}, F)$  close to  $V(\phi_F, F)$  for a range of interesting  $F$ . Estimate  $\phi_{F,H}$  consistently by  $\hat{\phi}_{F,H}$  and estimate  $\rho_0$  by  $\hat{\rho}_N(\hat{\phi}_{F,H})$ .

This program motivates the developments in this section. However, the results will be formulated in a manner that does not involve  $\phi_F$ .

A real-valued function  $g$  defined on  $R^1$  will be said to satisfy condition C if  $g = g_+ - g_-$ , where  $g_{\pm}$  is monotone nondecreasing, and  $\lim_{h \rightarrow 0} ([g_{\pm}(x+h) - g_{\pm}(x-h)], 1) = 0$ . Let  $\{\phi_i: 1 \leq i \leq k\}$  be a family of score functions which fulfill some or all of the following assumptions, as required.

- B1. Each  $\phi_i$  satisfies A1 and A2.
- B2. Each  $\phi_{i\pm}$  is absolutely continuous and  $\phi'_{i\pm}$  satisfies condition C.
- B3. Each  $\phi_i \phi_j$  satisfies condition C.
- B4. If  $\left\| \sum_{j=1}^k c_j \phi_j \right\| = 0$  for some constants  $\{c_j\}$ , then  $c_j = 0, 1 \leq j \leq k$ .

Note that B2 implies that A3 holds for each  $\phi_i$ , with  $D(\phi_{i\pm}, F) = (\phi'_{i\pm}, 1)$ . Let  $H$  be the subspace of  $L_2(F)$  spanned by all linear combinations of the  $\{\phi_i\}$ . Let  $W$  denote the  $k \times k$  matrix whose  $(i, j)$ th element is  $(\phi_i, \phi_j)$  and let  $v = ((\phi'_1, 1), (\phi'_2, 1), \dots, (\phi'_k, 1))^T$ . Assumption B4 ensures that  $W$  is nonsingular. Define the vector  $a = (a_1, a_2, \dots, a_k)^T$  by  $a = W^{-1}v$  and let

$$(3.1) \quad \phi_{F,H} = \sum_{i=1}^k a_i \phi_i .$$

LEMMA 3.1. *If each  $\phi_i \in L_2(F)$  and B2, B4 are satisfied, then*

$$(3.2) \quad V(\phi_{F,H}, F) = \min_{\phi \in H} V(\phi, F) .$$

PROOF. If  $\phi \in H$ , there exist constants  $c = (c_1, c_2, \dots, c_k)^T$  such that  $\phi = \sum_{i=1}^k c_i \phi_i$ . Since  $v = Wa$ ,

$$(3.3) \quad D(\phi, F) = (\phi', 1) = c^T v = c^T W a = (\phi, \phi_{F,H}) .$$

Hence

$$(3.4) \quad V(\phi, F) = \frac{\|\phi\|^2}{(\phi, \phi_{F,H})^2} \geq \frac{1}{\|\phi_{F,H}\|^2} = V(\phi_{F,H}, F) ,$$

with equality if and only if  $\phi$  is proportional to  $\phi_{F,H}$ .

An interesting interpretation can be given to this lemma when  $F'$  has an absolutely continuous density and  $\phi_F \in L_2(F')$ . In this case,  $(\phi', 1) = (\phi, \phi_F)$  for every  $\phi \in L_2(F')$  and  $\phi_{F,H}$  is simply the projection of  $\phi_F$  into

H. Note also that any multiple of  $\phi_{F,H}$  will retain the minimizing property (3.2).

LEMMA 3.2. *If B1, B2, B3 are satisfied, then for every  $C > 0$ ,  $1 \leq i, j \leq K$ ,*

$$(3.5) \quad (N-r)^{-1} \sup_{|\rho-\rho_0| \leq CN^{-1/2}} \left| \sum_{t=r+1}^N \phi'_i \circ R_t(\rho) - \sum_{t=r+1}^N \phi'_i(E_t) \right|^2 \rightarrow 0$$

$$(N-r)^{-1} \sup_{|\rho-\rho_0| \leq CN^{-1/2}} \left| \sum_{t=r+1}^N [\phi_i \circ R_t(\rho)][\phi_j \circ R_t(\rho)] - \sum_{t=r+1}^N \phi_i(E_t)\phi_j(E_t) \right|^2 \rightarrow 0$$

in  $\rho_0$ -probability as  $N \rightarrow \infty$ .

PROOF. Let  $\rho_0 - \rho = (A_0, A_1, \dots, A_r)^T$ . For each  $\delta > 0$ , define  $U_i(\delta)$  by

$$(3.6) \quad U_i(\delta) = \begin{cases} X_i & \text{if } |X_i| \leq \delta N^{1/2} \\ 0 & \text{otherwise.} \end{cases}$$

Without loss of generality, assume  $\phi'$  is monotone nondecreasing. Under B2,

$$(3.7) \quad (N-r)^{-1} \mathbb{E} \left[ \sup_{|\rho-\rho_0| \leq CN^{-1/2}} \left| \sum_{t=r+1}^N \phi'_i \left( E_t + A_0 + \sum_{j=1}^r A_j U_{t-j}(\delta) \right) - \sum_{t=r+1}^N \phi'_i(E_t) \right| \right]$$

$$\leq \mathbb{E} [\phi'_i(E_t + CN^{-1/2} + rC\delta) - \phi'_i(E_t - CN^{-1/2} - rC\delta)] \rightarrow 0$$

as  $N \rightarrow \infty$  and  $\delta \rightarrow 0$ . Moreover, because of Lemma 2.1,

$$(3.8) \quad \mathbb{P}[(U_1(\delta), U_2(\delta), \dots, U_N(\delta)) \neq (X_1, X_2, \dots, X_N)]$$

$$\leq \mathbb{P}[\max_{1 \leq i \leq N} |X_i| > \delta N^{1/2}] \rightarrow 0$$

for every  $\delta > 0$  as  $N \rightarrow \infty$ . The first line in (3.5) is implied by (3.7) and (3.8); the second is proved analogously using B3.

Let  $\hat{v}$  be the  $k \times 1$  vector whose  $i$ th component is  $(N-r)^{-1} \sum_{t=r+1}^N \phi'_i \circ R_t(\hat{\rho}_N)$  and let  $\hat{W}$  be the  $k \times k$  matrix whose  $(i, j)$ th element is  $(N-r)^{-1} \sum_{t=r+1}^N [\phi_i \circ R_t(\hat{\rho}_N)][\phi_j \circ R_t(\hat{\rho}_N)]$ . From the lemma above, it follows that  $\hat{v} \xrightarrow{p} v$  and  $\hat{W} \xrightarrow{p} W$  as  $N \rightarrow \infty$ . Define  $\hat{W}^{-1}$  as the inverse of  $\hat{W}$  when possible and arbitrarily otherwise. Since  $W$  is nonsingular under B4,  $\hat{W}^{-1} \xrightarrow{p} W^{-1}$  and  $\hat{a} = \hat{W}^{-1} \hat{v} \xrightarrow{p} a$  as  $N \rightarrow \infty$ . The implied estimate of  $\phi_{F,H}$  is

$$(3.9) \quad \hat{\phi}_{F,H} = \sum_{i=1}^k \hat{a}_i \phi_i.$$

By setting  $\phi = \hat{\phi}_{F,H}$  in the linearized estimate  $\hat{\rho}_N(\phi)$  and noting that  $\hat{a}^T \hat{W} \hat{a}$  is a consistent estimate of  $D(\phi_{F,H}, F) = \|\phi_{F,H}\|^2 = a^T W a$ , we arrive at the adaptive estimate

$$(3.10) \quad \hat{\rho}_N(H) = \hat{\rho}_N + (\hat{a}^T \hat{W} \hat{a})^{-1} (X^T X)^{-1} X^T \hat{\phi}_{F,H}(\hat{\rho}),$$

where  $\hat{\phi}_{F,H}(\hat{\rho}_N)$  is the  $(N-r) \times 1$  vector of scored residuals  $(\hat{\phi}_{F,H} \circ R_{r+1}(\hat{\rho}_N), \dots, \hat{\phi}_{F,H} \circ R_N(\hat{\rho}_N))^T$ .

**THEOREM 3.1.** *If B1, B2, B3, B4 are satisfied, then as  $N \rightarrow \infty$ ,  $N^{1/2}(\hat{\rho}_N(H) - \hat{\rho}_N(\phi_{F,H})) \xrightarrow{p} 0$  in  $\rho_0$ -probability and the asymptotic distribution of  $N^{1/2}(\hat{\rho}_N(H) - \rho_0)$  is normal  $(0, \|\phi_{F,H}\|^{-2} \Gamma^{-1})$ .*

**PROOF.** As in (2.2), let  $\Phi_{F,H}$ ,  $\Psi_i$  denote  $(N-r) \times 1$  vectors of scored residuals, the score functions being  $\phi_{F,H}$ ,  $\psi_i$  respectively. Since

$$(3.11) \quad \Phi_{F,H}(\rho) = \sum_{i=1}^k a_i \Psi_i(\rho) \quad \hat{\Phi}_{F,H}(\rho) = \sum_{i=1}^k \hat{a}_i \Psi_i(\rho),$$

we may write

$$(3.12) \quad N^{-1/2} X^T \hat{\Phi}_{F,H}(\hat{\rho}_N) - N^{-1/2} X^T \Phi_{F,H}(\hat{\rho}_N) = \sum_{i=1}^k (\hat{a}_i - a_i) (N^{-1/2} X^T \Psi_i(\hat{\rho}_N)).$$

Theorem 2.1 implies that  $N^{-1/2} X^T \Psi_i(\hat{\rho}_N)$  is bounded in probability asymptotically. Since also  $\hat{a} \xrightarrow{p} a$ , the difference (3.12) converges in probability to zero as  $N \rightarrow \infty$ . The theorem follows with the help of Theorem 2.2.

A desirable property possessed by the LSE  $\hat{\rho}_N$  is invariance under rescaling of the observations in the sense that the mapping  $X_t \rightarrow cX_t$ ,  $c > 0$ , induces the mappings  $\hat{a}_0 \rightarrow c\hat{a}_0$  and  $\hat{a}_i \rightarrow \hat{a}_i$  for  $1 \leq i \leq r$ . For suitable  $H$ , the adaptive estimate  $\hat{\rho}_N(H)$  is also scale invariant.

**DEFINITION.** A subspace  $H$  of  $L_2(F)$  is said to be closed under scaling if  $\phi(\cdot) \in H$  implies that  $\phi(c \cdot) \in H$  for every scalar  $c > 0$ .

**THEOREM 3.2.** *If  $H$  is closed under scaling and  $\hat{W}$  is nonsingular, the adaptive estimate  $\hat{\rho}_N(H)$  is scale invariant.*

**PROOF.** Let  $F_N$  denote the empirical distribution function of the residuals  $\{R_t(\hat{\rho}_N): r+1 \leq t \leq N\}$  and let

$$(3.13) \quad \hat{V}(\phi, F_N) = \frac{\int \phi^2(t) dF_N(t)}{\left[ \int \phi'(t) dF_N(t) \right]^2}.$$

An argument like that for Lemma 3.1 shows that  $\hat{\phi}_{F,H}$  is characterized, up to a constant of proportionality, by the property

$$(3.14) \quad \hat{V}(\hat{\phi}_{F,H}, F_N) = \min_{\phi \in H} \hat{V}(\phi, F_N).$$

The scaling  $X_t \rightarrow cX_t$ ,  $1 \leq t \leq N$ ,  $c > 0$ , induces the following mappings:  $R_t(\hat{\rho}_N) \rightarrow cR_t(\hat{\rho}_N)$  and  $F_N(\cdot) \rightarrow F_N(\cdot/c)$ . Consideration of (3.13), (3.14) and

the closure of  $H$  yields

$$(3.15) \quad \hat{V}(\hat{\phi}_{F,H}(\cdot/c), F_N(\cdot/c)) = c^2 \hat{V}(\hat{\phi}_{F,H}, F_N) = c^2 \min_{\phi \in H} \hat{V}(\phi, F_N) \\ = \min_{\phi \in H} \hat{V}(\phi, F_N(\cdot/c)).$$

Consequently, the scaling  $X_i \rightarrow cX_i$  must map  $\hat{\phi}_{F,H}(\cdot)$  into a multiple of  $\hat{\phi}_{F,H}(\cdot/c)$ . The theorem follows from this fact and the invariance of  $\hat{\rho}_N$ .

Under the assumptions of Lemma 3.2,  $\hat{V}(\phi, F_N)$  is a consistent estimate of  $V(\phi, F)$ . Therefore,  $\hat{\phi}_{F,H}$  is an element of  $H$  that minimizes, in the obvious sense,  $N(X^T X)^{-1} \hat{V}(\phi, F_N)$ , the estimated asymptotic covariance matrix of  $\hat{\rho}_N(\phi)$ . This fact makes  $\hat{\rho}_N(H)$  an extended analogue of the adaptive  $L$ -estimates for location studied by Jaeckel [5].

#### 4. Applications

Scale invariant adaptive estimates of  $\rho$  can be constructed as follows. Assume  $F$  is symmetric about the origin and take as a basis for  $H$  the set of functions  $\phi_i(x) = |x|^{r_i} \text{sign}(x)$ ,  $r_i > 0$ ,  $1 \leq i \leq k$ . In this case,  $H$  is closed under scaling and the other assumptions required by Theorem 3.1 can be checked readily. Indeed, let  $\phi(x) = |x|^r \text{sign}(x)$ ,  $r > 0$ . If  $\int |x|^{2r} dF(x) < \infty$ , assumption B1 holds for  $\phi$ , and if also  $r \geq 1$ , so does B2. On the other hand, if  $0 < r < 1$ , fulfillment of B2 is assured whenever  $F$  has a bounded density which is uniformly continuous in a neighborhood of the origin. Assumption B3 holds under a moment condition similar to that for B1, while B4 is satisfied if  $F$  is absolutely continuous and the exponents  $\{r_i\}$  are distinct.

Whenever the particular score function  $\phi(x) = x$  belongs to  $H$ , the asymptotic efficiency of  $\hat{\rho}_N(H)$  relative to the LSE  $\hat{\rho}_N$  is greater than or equal to one for all regular  $F$ ; this is a consequence of Theorem 3.1 and Lemma 3.1.

The adaptive estimate  $\hat{\rho}_N(H)$  can be applied to hypothesis testing. Let  $C$  be a  $q \times (r+1)$  matrix constant of rank  $q$ , let  $\tau = C\rho$  and let  $\hat{\tau}_N = C\hat{\rho}_N(H)$ . To test the hypothesis  $H: \sum_{i=1}^q \tau_i^2 = 0$  versus  $K: \sum_{i=1}^q \tau_i^2 > 0$ , calculate

$$(4.1) \quad T_N = [\hat{\tau}_N^T (C\hat{I}_N^{-1}C')^{-1} \hat{\tau}_N] [\hat{a}^T \hat{W} \hat{a}],$$

where  $\hat{I}_N = N^{-1}X^T X$ , and reject  $H$  for values of  $T_N$  that are large relative to the asymptotic  $\chi_q^2$  distribution implied by Theorem 3.1. Under a sequence of alternatives  $K_N: \tau = N^{-1/2}\Delta$ ,  $\Delta$  a non-null  $q \times 1$  vector, the asymptotic efficiency of this test relative to the corresponding test based on  $\hat{\rho}_N$  is the same as in the estimation problem.



The central ideas underlying the definition and properties of  $\hat{\rho}_N(H)$  carry over unchanged to the general linear model. For counterparts of Theorems 2.1 and 2.2 in that case, see Bickel [2]. The analogues of Theorems 3.1 and 3.2 will be evident to the reader.

*Practical aspects.* The following practical suggestions are made on partly heuristic grounds and need further investigation.

1. For samples from a moderately contaminated normal distribution, try  $\phi_1(x)=|x|^{1/2} \text{sign}(x)$ ,  $\phi_2(x)=x$  as a basis for  $H$ . If  $F$  is normal,  $\hat{\rho}_N(H)$  is still fully efficient asymptotically. If  $F$  is actually double exponential,  $\hat{\phi}_{F,H}(x)$  will converge in probability to  $\phi_{F,H}(x) \doteq 1.90\phi_1(x) - .76\phi_2(x)$ , and therefore  $\hat{\rho}_N(H)$  will discount outlying residuals in large samples. Note that  $\phi_{F,H}$  becomes negative-valued only far out in the tails of the double exponential distribution. With  $F$  double exponential, the asymptotic efficiency of  $\hat{\rho}_N(H)$  relative to the best estimate is  $(1/2)[1+\pi/(32-9\pi)] \doteq .92$ . The efficiency of  $\hat{\rho}_N$  in this case is only .50 and the efficiency of  $\hat{\rho}_N(\phi_1)$  is .79.

2. If more serious departures from normality are anticipated, bring into  $H$  selected functions of the form  $\phi(x)=|x|^r \text{sign}(x)$ , with  $0 \leq r < 1$  for heavier tailed  $F$  and  $r > 1$  for lighter tailed  $F$ . In some cases it will be necessary to replace  $\hat{\rho}_N$  with a more robust estimate, such as the  $M$ -estimate corresponding to the score function  $\phi(x)=\text{sign}(x)$ .

3. If the sample size  $N$  is large, some experimentation with the choice of  $H$  may be worthwhile. Plot  $\hat{\phi}_{F,H}(x)$  and note changes that occur as functions are added to or removed from the basis of  $H$ . The aim is to discover, at least qualitatively, the shape of  $\phi_F$ . Keep  $K$  small relative to  $N$ .

*Numerical example.* To check the numerical practicality of the adaptive estimator, a pseudo-random sample of size 50 was generated from the autoregressive process  $X_t = .5 + .5X_{t-1} + E_t$ , where  $E_t$  has a double-exponential distribution with scale parameter .75. The first column of Table 1 records the sample values. For this data, the LSE of  $\rho = (\alpha_0, \alpha_1)^T = (.5, .5)^T$  is  $\hat{\rho}_N = (.457692, .534041)^T$  and the least squares residuals  $R_t(\hat{\rho}_N) = X_t - \hat{\alpha}_0 - \hat{\alpha}_1 X_{t-1}$  are given by the second column of Table 1.

The basis for the subspace  $H$  consists of two functions:  $\phi_1(x) = |x|^{1/2} \text{sgn}(x)$  and  $\phi_2(x) = x$ . The estimates of  $W$  and  $a$  are

$$(4.2) \quad W = \begin{pmatrix} .714230 & .784176 \\ .784176 & .959636 \end{pmatrix}, \quad \hat{a} = \begin{pmatrix} 1.51483 \\ -.19579 \end{pmatrix}.$$

The actual values of  $W$  and  $a$  under the double-exponential model that generated the sample are

Table 1

Autoregressive series	Residuals	Scored residuals
1.13639	.524676	.99453
1.58924	.866306	1.24032
2.17272	— .375636	— .854877
1.24238	— .403473	— .883213
.7177	1.28047	1.46344
2.12145	.71641	1.1419
2.30704	.416844	.896409
2.10659	—2.3101	—1.85008
— .727401	.239953	.695056
.309182	— .130794	— .522236
.492013	2.8828	2.00756
3.60325	— .437629	— .916427
1.94435	1.92483	1.72477
3.42088	2.42519	1.88421
4.70978	— .70798	—1.13598
2.26493	.65645E-2	.121448
1.67382	—1.1386	—1.39347
.212983	1.49821	1.56083
2.06964	.815332E-1	.41658
1.6445	—1.77403	—1.67029
— .438107	— .731375	—1.15229
— .507651	.261263	.723134
.447848	— .955391	—1.29359
— .25853	— .473251	— .949439
— .153625	— .396044	— .875769
— .020395	— .805794E-1	— .41423
.36622	—1.6481	—1.62202
— .994836	.820102E-1	.41775
.841844E-2	.109264	.479334
.571451	.464126	.94113
1.227	— .811606E-1	— .415664
1.0318	— .722722	—1.1463
.285992	—1.18925	—1.41911
— .578829	1.18348	1.41623
1.33206	— .904759E-1	— .437934
1.07859	.337221	.813645
1.37092	—1.50854	—1.56519
— .318721	— .340982	— .817801
— .535011E-1	— .773312	—1.1807
— .344192	— .217533	— .663931
.563454E-1	.163325	.580217
.651107	— .313339	— .786601
.49207	.524206	.99413
1.24468	.116964	.495169
1.23937	.51387	.985287
1.63344	— .30782	— .78018
1.02219	— .39051	— .870168
.613075	.472668	.948911
1.25777	.406439	.886163
1.53583		

$$(4.3) \quad W = \begin{pmatrix} .750000 & .863432 \\ .863432 & 1.125000 \end{pmatrix}, \quad a = \begin{pmatrix} 2.92981 \\ -1.35972 \end{pmatrix}.$$

Although  $\hat{a}$  differs markedly from  $a$ , the components of  $\hat{a}$  have the correct signs. Hence, scoring the residuals according to  $a$  has the desired effect of discounting the larger residual values. The scored residual vector  $\hat{\phi}_{F,H}(\hat{\rho}_N)$  is listed in the third column of Table 1.

The adaptive estimate  $\hat{\rho}_N(H)$  is  $(.451717, .487241)^T$ . Comparison with the LSE and the actual parameter values shows that the adaptive estimate of  $\alpha$  is slightly worse than the LSE of  $\alpha_0$ , but the situation is reversed for  $\alpha_1$  and the gain in accuracy outweighs the loss.

5. Proof of Theorem 2.1

The random vector  $X^t\psi(\rho)$  may be written out as

$$\left( \sum_{t=r+1}^N \phi \circ R_t(\rho), \sum_{t=r+1}^N X_{t-1}\phi \circ R_t(\rho), \dots, \sum_{t=r+1}^N X_{t-r}\phi \circ R_t(\rho) \right)^T.$$

Corresponding to these components, define

$$(5.1) \quad \begin{aligned} T_{N0}(\rho) &= N^{-1/2} \sum_{t=r+1}^N [\phi \circ R_t(\rho) - E(\phi \circ R_t(\rho) | A_{t-1})] \\ T_{Ni}(\rho) &= N^{-1/2} \sum_{t=r+1}^N X_{t-i}[\phi \circ R_t(\rho) - E(\phi \circ R_t(\rho) | A_{t-1})], \quad 1 \leq i \leq r. \end{aligned}$$

The method of expansion adopted in this section uses ideas from Bickel [2].

LEMMA 5.1. *If A1, A2 are satisfied, then for every  $C > 0$*

$$(5.2) \quad \sup_{|\rho - \rho_0| \leq CN^{-1/2}} |T_{Ni}(\rho) - T_{Ni}(\rho_0)| \xrightarrow{p} 0, \quad 0 \leq i \leq r$$

*in  $\rho_0$ -probability as  $N \rightarrow \infty$ .*

PROOF. Without loss of generality, assume that  $i \geq 1$  and  $\phi$  is monotone nondecreasing. We begin by showing that if  $|\rho - \rho_0| \leq CN^{-1/2}$ ,

$$(5.3) \quad T_{Ni}(\rho) - T_{Ni}(\rho_0) \xrightarrow{p} 0$$

as  $N \rightarrow \infty$ . Indeed, let  $\rho_0 - \rho = (A_0, A_1, \dots, A_r)^T$  and let

$$(5.4) \quad \begin{aligned} R_t^i(\rho) &= E_t + A_0 + \sum_{j=1}^r A_j U_{t-j}(\delta) \\ T_{Ni}^i(\rho) &= N^{-1/2} \sum_{t=r+1}^N X_{t-i}[\phi \circ R_t^i(\rho) - E(\phi \circ R_t^i(\rho) | A_{t-1})], \end{aligned}$$

where the  $\{U_j(\delta)\}$  are defined as in (3.6). Because of (3.8),

$$(5.5) \quad P [T_{Ni}(\rho) \neq T_{Ni}^s(\rho)] \rightarrow 0$$

as  $N \rightarrow \infty$ .

$$(5.6) \quad \begin{aligned} & E [T_{Ni}^s(\rho) - T_{Ni}(\rho_0)]^2 \\ &= N^{-1} \sum_{t=r+1}^N E \{ X_{t-i}^2 [\phi \circ R_t^s(\rho) - \phi \circ R_t(\rho_0) \\ &\quad - E(\phi \circ R_t^s(\rho) - \phi \circ R_t(\rho_0) | A_{t-1})]^2 \} \\ &\leq N^{-1} E \sum_{t=r+1}^N X_{t-i}^2 \int \left[ \phi \left( x + \Delta_0 + \sum_{j=1}^r \Delta_j U_{t-j}(\delta) \right) - \phi(x) \right]^2 dF(x) \\ &\leq N^{-1} E \sum_{t=r+1}^N X_{t-i}^2 \int [\phi(x+h) - \phi(x-h)]^2 dF(x), \end{aligned}$$

where  $h = CN^{-1/2} + rC\delta$ . Now (5.3) follows because of A2 and (5.5).

Decompose the  $r+1$  dimensional cube  $B = \{\rho: |\rho - \rho_0| \leq CN^{-1/2}\}$  into sub-cubes whose vertices are at the points  $\rho_0 + (j_1 \varepsilon N^{-1/2}, j_2 \varepsilon N^{-1/2}, \dots, j_{r+1} \varepsilon N^{-1/2})$ , where  $\varepsilon > 0$  is chosen to give an even division of  $B$  into sub-cubes and  $j_i = 0, \pm 1, \dots, \pm M(\varepsilon)$ . For each  $\rho \in B$ , let  $V(\rho)$  denote the vertex nearest  $\rho_0$  of the sub-cube containing  $\rho$  (or one of them, in case of ties). Then for each  $\varepsilon > 0$ , from (5.3)

$$(5.7) \quad \sup_{|\rho - \rho_0| \leq CN^{-1/2}} |T_{Ni}(\rho_0) - T_{Ni} \circ V(\rho)| \xrightarrow{P} 0$$

as  $N \rightarrow \infty$ .

Suppose that  $\rho \in B^*$ , a particular sub-cube of  $B$ , and let  $\rho^* = V(\rho)$ . Then

$$(5.8) \quad \begin{aligned} & \sup_{\rho \in B^*} |T_{Ni}(\rho) - T_{Ni} \circ V(\rho)| \\ & \leq \sup_{\rho \in B^*} \left| N^{-1/2} \sum_{t=r+1}^N X_{t-i} (\phi \circ R_t(\rho) - \phi \circ R_t(\rho^*)) \right| \\ & \quad + \sup_{\rho \in B^*} \left| N^{-1/2} \sum_{t=r+1}^N X_{t-i} E(\phi \circ R_t(\rho) - \phi \circ R_t(\rho^*) | A_{t-1}) \right| \\ & = D_1 + D_2. \end{aligned}$$

By an argument like that for (5.3),

$$(5.9) \quad D_1 \leq N^{-1/2} \sum_{t=r+1}^N |X_{t-i}| E[\phi(R_t(\rho^*) + \varepsilon S_t) - \phi(R_t(\rho^*) - \varepsilon S_t) | A_{t-1}] + o_p(1),$$

where  $S_t = N^{-1/2} \left( 1 + \sum_{j=1}^r |X_{t-j}| \right)$ . Hence

$$(5.10) \quad \begin{aligned} & \sup_{\rho \in B^*} |T_{Ni}(\rho) - T_{Ni} \circ V(\rho)| \\ & \leq 2N^{-1/2} \sum_{t=r+1}^N |X_{t-i}| E[\phi(R_t(\rho^*) + \varepsilon S_t) - \phi(R_t(\rho^*) - \varepsilon S_t) | A_{t-1}] \\ & \quad + o_p(1) \end{aligned}$$

$$\leq 2N^{-1/2} \sum_{t=r+1}^N |X_{t-i}| \int [\phi(x+U_t+\varepsilon S_t) - \phi(x+U_t-\varepsilon S_t)] dF(x) \\ + o_p(1),$$

where  $|U_t| \leq CS_t$ . In view of this, A2, and Lemma 2.1,

$$(5.11) \quad \sup_{\rho \in B^*} |T_{Nt}(\rho) - T_{Nt} \circ V(\rho)| = O_p(\varepsilon)$$

as  $N \rightarrow \infty$ . The lemma follows from (5.11) and (5.7).

PROOF OF THEOREM 2.1. Note that

$$(5.12) \quad N^{-1/2} \sum_{t=r+1}^N X_{t-i} E(\phi \circ R_t(\rho) - \phi \circ R_t(\rho_0) | A_{t-1}) \\ = N^{-1/2} \sum_{t=r+1}^N X_{t-i} V_t \cdot V_t^{-1} \int [\phi(x+V_t) - \phi(x)] dF(x),$$

where  $V_t = \Delta_0 + \sum_{j=1}^N \Delta_j X_{t-j}$ . As  $N \rightarrow \infty$ , the difference between (5.12) and  $N^{-1/2} \left( \sum_{t=r+1}^N X_{t-1} V_t \right) D(\phi, F)$  converges in probability to zero because of A3 and Lemma 2.1. The theorem follows from Lemma 5.1.

UNIVERSITY OF CALIFORNIA, BERKELEY

#### REFERENCES

- [1] Anderson, T. W. (1971). *The Statistical Analysis of Time Series*, John Wiley & Sons.
- [2] Bickel, P. J. (1975). One-step Huber estimates in the linear model, *J. Amer. Statist. Ass.*, **70**, 428-434.
- [3] Brown, B. M. (1971). Martingale central limit theorems, *Ann. Math. Statist.*, **42**, 59-66.
- [4] Huber, P. J. (1964). Robust estimation of a location parameter, *Ann. Math. Statist.*, **35**, 73-101.
- [5] Jaeckel, L. A. (1971). Some flexible estimates of location, *Ann. Math. Statist.*, **42**, 1540-1552.
- [6] LeCam, L. (1956). On the asymptotic theory of estimation and testing hypotheses, *Proc. 3rd Berkeley Symp. Math. Statist. Prob.*, **1**, 129-156.
- [7] Mann, H. B. and Wald, A. (1943). On the statistical treatment of linear stochastic difference equations, *Econometrica*, **11**, 217-226.