# A DUALITY BETWEEN AUTOREGRESSIVE AND MOVING AVERAGE PROCESSES CONCERNING THEIR LEAST SQUARES PARAMETER ESTIMATES[1]

By David A. Pierce

*University of Missouri*

**1. Introduction.** The methods employed in least squares parameter estimation for moving average (MA) processes differ from those appropriate for autoregressive (AR) processes, as only the latter are linear in the parameters. There is nevertheless an interesting duality between these two classes of time series models: if AR and MA series, each of the same order and with the same parameter values, are generated from the same sequence of errors, then to a close approximation the least squares estimates calculated from the MA series will *under*estimate the true parameter values by the same amount that those determined from the AR series will *over*-estimate them. This relation is established in Section 3 via a linear approximation of the moving average errors (considered as functions of the parameters for a given series) in a neighborhood of the true parameter values. In Section 4 the large-sample distribution of the estimates in any MA process is then obtained as a direct consequence of this relation and of known results for AR processes, and some properties of the least squares estimates in both classes of processes are examined.

Let us introduce notation and terminology to be employed. A sequence $\{x_t\}$ follows an AR *process of order p* if it is governed by the relation

$$(1.1) \qquad x_t = \sum_{j=1}^{p} \theta_j x_{t-j} + a_t$$

where the $\{a_t\}$ are $N(0, \sigma^2)$ and independent, and the parameters $\theta = (\theta_1, \cdots, \theta_p)'$ are such that the roots of the auxiliary equation

$$(1.2) \qquad \theta(u) = 1 - \theta_1 u - \cdots - \theta_p u^p = 0$$

lie outside the unit circle (the set of all $\theta$ satisfying this condition will be referred to as the *admissible parameter space*). Following [3] we may define a *backward shift operator B* by the relation $Bw_t = w_{t-1}$ for any sequence $\{w_t\}$, and equation (1.1) may then be written

$$(1.3) \qquad \theta(B)x_t = a_t$$

where $\theta(B) = (1 - \sum_{j=1}^{p} \theta_j B^j)$ analogous to (1.2).

In the MA process the recursive relationship is in terms of the deviates $\{a_t\}$ rather than the observations themselves; a sequence $\{y_t\}$ is a moving average of order $p$ if

$$(1.4) \qquad y_t = -\sum_{j=1}^{p} \theta_j a_{t-j} + a_t = \theta(B)a_t$$

with $\{a_t\}$ and $\theta = (\theta_1, \cdots, \theta_p)'$ as before.

---

422

**2. Least squares estimation.** The equations (1.3) and (1.4) are well defined for any $t = \cdots, -1, 0, 1, \cdots$; however, in practice we will have available a sample series [say $x_1, x_2, \cdots, x_n$ or $y_1, y_2, \cdots, y_n$], and our problem is then to estimate the unknown parameters $\theta$. The method of least squares leads us to the minimization of the sum of squares

$$(2.1) \qquad S(\dot{\theta}) = \sum_t [a_t(\dot{\theta})]^2$$

as a function of the indeterminate $\dot{\theta}$. In (2.1), the quantity $a_t(\dot{\theta})$ is obtained by solving either (1.3) or (1.4) for $a_t = a_t(\theta)$ which then becomes (given $\{x_t\}$ or $\{y_t\}$) a function of the point $\theta = \dot{\theta}$. There is of course a need to determine or have available a set of "starting values" $[x_0, x_{-1}, \cdots, x_{-p+1}$ for (1.3) or $a_0, a_{-1}, \cdots, a_{-p+1}$ for (1.4)]; however, the effect of any particular choice of these values becomes negligible asymptotically, and we can thus ignore considerations pertaining to the initialization of the process in the following discussion.

The actual determination of the quantity $\hat{\theta}$ which minimizes the sum of squares (2.1) is facilitated in AR processes by the fact that Equation (1.3) is linear in the parameters $[\partial a_t(\dot{\theta})/\partial \dot{\theta}_j = -x_{t-j}$ independently of $\dot{\theta}]$, whereas (1.4) is not. Thus (1.3) can be written in matrix form as

$$(2.2) \qquad \mathbf{x} = X\theta + \mathbf{a}$$

where $\mathbf{a} = (a_1, \cdots, a_n)'$, $\mathbf{x} = (x_1, \cdots, x_n)'$, and $X$ is an $(n \times p)$ matrix whose $(tj)$th element is the lagged observation $x_{t-j}$. The least squares estimates $\hat{\theta}^{(1)}$ of $\theta$ for the AR series are then

$$(2.3) \qquad \hat{\theta}^{(1)} = (X'X)^{-1}X'\mathbf{x}$$

with a large-sample distribution very similar to that of the estimates in the fixed-$x$ linear normal regression model [7, Section 4].

The situation is fundamentally different for MA processes, however; letting

$$(2.4) \qquad \pi(u) = 1 - \sum_{j=1}^{\infty} \pi_j u^j$$

be the solution of $\pi(u)\theta(u) = 1$ [whose existence is guaranteed by the admissibility of $(\theta_1, \cdots, \theta_p)$], equation (1.3) becomes

$$[\theta(B)]^{-1}y_t = \pi(B)y_t = a_t,$$

and thus $a_t(\dot{\theta}) = [\dot{\theta}(B)]^{-1}y_t$, whose derivative with respect to $\dot{\theta}_j$ is a function of $\dot{\theta}$. Thus nonlinear estimation methods are employed in practice to obtain the least squares estimates, say $\hat{\theta}^{(2)}$, of $\theta$ in the MA process (1.4).

Despite this basic difference concerning linearity in these two classes of processes, however, there exists an important similarity in the behavior (asymptotically) of their least squares estimates, as summarized in Section 1. This duality is obtained in the following section by considering AR and MA processes each involving the same parameters $\theta$ and generated from the same deviates $\{a_t\}$.

**3. The correspondence between $\hat{\theta}^{(1)}$ and $\hat{\theta}^{(2)}$.** For a given sequence $\{a_t\}$ of independent $N(0, \sigma^2)$ deviates and a given vector of admissible parameters

$\theta = (\theta_1, \cdots, \theta_p)'$, suppose we have AR and MA series, $\{x_t\}$ and $\{y_t\}$, generated by the relations (1.3) and (1.4) respectively. In particular, then,

$$(3.1) \qquad\qquad y_t = [\theta(B)]^2 x_t.$$

Suppose $\dot{\theta}$ is a point in the admissible parameter space which is "close" (in a sense to become evident later) to the true parameter values $\theta$. Then in terms of $\dot{\theta}$ we can define, corresponding to (2.1),

(i) an AR error

$$(3.2) \qquad\qquad \dot{a}_t^{(1)} = a_t^{(1)}(\dot{\theta}) = \dot{\theta}(B)x_t, \qquad\qquad \text{and}$$

(ii) an MA error

$$(3.3) \qquad\qquad \dot{a}_t^{(2)} = a_t^{(2)}(\dot{\theta}) = [\dot{\theta}(B)]^{-1} y_t.$$

Now since $\dot{a}_t^{(1)}$ is linear in $\dot{\theta}$, we have for all $\dot{\theta}$ the exact expression [noting that at $\dot{\theta} = \theta$, $\dot{a}_t^{(1)} = \dot{a}_t^{(2)} = a_t$]

$$(3.4) \qquad\qquad \dot{a}_t^{(1)} = \sum_{j=1}^{p} (\dot{\theta}_j - \theta_j)\frac{\partial \dot{a}_t^{(1)}}{\partial \dot{\theta}_j} + a_t$$

$$= \sum_{j=1}^{p} (\dot{\theta}_j - \theta_j)x_{t-j} + a_t$$

or in matrix form

$$(3.5) \qquad\qquad \dot{\mathbf{a}}^{(1)} = X(\theta - \dot{\theta}) + \mathbf{a},$$

analogous to (2.2). The least squares estimation may then be performed directly on the "observations" $\dot{a}_1^{(1)}, \cdots, \dot{a}_n^{(1)}$ in (3.5) to give

$$(3.6) \qquad\qquad \hat{\theta}^{(1)} - \dot{\theta} = (X'X)^{-1}X'\dot{\mathbf{a}}^{(1)}.$$

Recalling that the MA error $\dot{a}_t^{(2)}$ is not linear in $\dot{\theta}$, the analogue of (3.4) for moving average series is an approximation via a first order Taylor expansion about $\dot{\theta} = \theta$:

$$(3.7) \qquad\qquad \dot{a}_t^{(2)} = \sum_{j=1}^{p} (\dot{\theta}_j - \theta_j)\frac{\partial \dot{a}_t^{(2)}}{\partial \dot{\theta}_j}\bigg|_{\dot{\theta}=\theta} + a_t$$

where the error thereby introduced is $O(|\dot{\theta} - \theta|^2)$, with $|\dot{\theta} - \theta| = [\sum(\dot{\theta}_j - \theta_j)^2]^{\frac{1}{2}}$ denoting the Euclidean distance in $p$-space between $\dot{\theta}$ and $\theta$. The derivatives in (3.7) are

$$\frac{\partial \dot{a}_t^{(2)}}{\partial \dot{\theta}_j}\bigg|_{\dot{\theta}=\theta} = \frac{\partial}{\partial \dot{\theta}_j}[\dot{\theta}(B)]^{-1}y_t\bigg|_{\dot{\theta}=\theta}$$

$$(3.8) \qquad\qquad = \dot{\theta}^{-2}(B)B^j y_t\big|_{\dot{\theta}=\theta}$$

$$= \theta^{-2}(B)y_{t-j}$$

$$= x_{t-j},$$

the last equality following from (3.1). Thus (3.7) becomes $\dot{a}_t^{(2)} = \sum (\dot{\theta}_j - \theta_j) x_{t-j} + a_t$ or

$$(3.9) \qquad \dot{\mathbf{a}}^{(2)} = X(\dot{\theta} - \theta) + \mathbf{a}$$

where $X$ is the *same* as in (3.5) and (2.2). Analogous to (3.6) we now have approximately

$$(3.10) \qquad \dot{\theta} - \hat{\theta}^{(2)} = (X'X)^{-1} X' \dot{\mathbf{a}}^{(2)}.$$

The equations (3.6) and (3.10) both involve the least squares estimates $\hat{\theta}^{(i)}$, $i = 1, 2$ and the fixed but arbitrary values $\dot{\theta}$. If we now set $\dot{\theta} = \theta$, these equations become

$$(3.11) \qquad \hat{\theta}^{(1)} - \theta = (X'X)^{-1} X' \mathbf{a} \qquad\qquad \text{and}$$

$$(3.12) \qquad \theta - \hat{\theta}^{(2)} = (X'X)^{-1} X' \mathbf{a}, \qquad\qquad \text{whence}$$

$$(3.13) \qquad (\theta - \hat{\theta}^{(2)}) = -(\theta - \hat{\theta}^{(1)})$$

where the error of approximation in (3.13) is of order $n^{-1}$ in probability and becomes asymptotically negligible.

Equation (3.13) shows that for large samples $\hat{\theta}^{(2)}$ from the MA series *under-*estimates the true value $\theta$ by the same amount that $\hat{\theta}^{(1)}$ from the AR series *over-*estimates $\theta$, this being for the particular instance when both series are generated from the same white noise. It follows from this that the asymptotic distribution of $\hat{\theta}^{(2)}$ for any MA process (1.4) is the same as that of the dual AR process (1.3), since the latter distribution is known to be normal and hence symmetric. (The large-sample distribution of $\hat{\theta}^{(2)}$ has been established via other methods by Whittle [8] and Box and Jenkins [3]; Box and Jenkins also obtain the joint distribution of the estimates in a mixed AR–MA process.) We briefly discuss this common distribution and the properties of $\hat{\theta}^{(i)}$ in the following section.

**4. Distribution and properties of the least squares estimates; effect of nonnormality of $\{a_t\}$.** The large-sample distribution of the least squares estimates in an AR process was obtained by Mann and Wald [7] who showed (i) that asymptotically,

$$(4.1) \qquad n^{\frac{1}{2}}(\hat{\theta}^{(1)} - \theta) \sim N(0, \sigma^2 \Gamma^{-1})$$

where $\Gamma = (\gamma_{|i-j|})$ is a $(p \times p)$ matrix of the *autocovariances* $\gamma_k = E(x_t x_{t+k})$ of the AR process (1.3); and (ii) that

$$(4.2) \qquad \Gamma = \operatorname{plim} n^{-1} X'X$$

with $X$ as in (2.2). It follows from (3.13) that equation (4.1) also gives the asymptotic distribution of the estimates $\hat{\theta}^{(2)}$ of the MA process (1.4), with the stipulation, of course, that $\Gamma$ is the autocovariance matrix not of this process but of the dual AR process (1.3).

Furthermore, under the normality assumption the minimization of (2.1) is equivalent (asymptotically) to the maximization of the likelihood function, so that $\hat{\theta}^{(i)}$ is also the ML estimate of $\theta$, $i = 1, 2$. Thus for large samples, least squares estimates in both AR and MA processes are efficient, and the information matrices

for either are given by $(1/\sigma^2)\Gamma$. In particular, the problem of efficient estimation in MA processes ([5], [6, Chapter 2], [8, page 431]) is seen to be equivalent to the problem of minimizing (2.1), which is readily accomplished with the aid of a computer and employing standard nonlinear estimation techniques [2], [4].

While the preceding discussion has been given under the assumption of a normal distribution for $\{a_t\}$, many large-sample properties of the estimates $\hat{\theta}^{(i)}$ hold under wider conditions. Mann and Wald in their investigation obtained the distribution (4.1) for $\hat{\theta}^{(1)}$ assuming only that the $\{a_t\}$ possessed moments of all orders, and T. W. Anderson [1] showed that the existence of moments of order two is sufficient. By (3.13) these remarks can also be made for least squares parameter estimates in MA processes. Clearly also the large sample covariance matrices of $\hat{\theta}^{(1)}$ and $\hat{\theta}^{(2)}$ will be the same (and given by (4.1)), though in general these will no longer be related to the information matrices. What is lost when the normality assumption is not satisfied is efficiency, as the ML estimates will no longer necessarily coincide with the least squares estimates. (The latter will however still have minimum asymptotic variance within a class of estimates limited in a way analogous to the restriction of linearity and unbiasedness under which the Gauss–Markov theorem holds in the non-normal regression model; Whittle [8] indicates these restrictions in terms of the spectrum.)

Let us conclude by briefly illustrating these considerations with the first order AR and MA processes,

$$(4.3) \qquad x_t = \theta x_{t-1} + a_t \qquad \text{and}$$

$$(4.4) \qquad y_t = -\theta a_{t-1} + a_t,$$

respectively. The admissibility restriction on the parameter is that $|\theta| < 1$. If $E(a_t^2) < \infty$, then the least squares estimate of $\theta$ in either process is asymptotically normal with mean $\theta$ and variance $(1-\theta^2)/n$ [the $(1 \times 1)$ covariance matrix is $\Gamma = \text{Var}(x_t) = \sigma^2/(1-\theta^2)$]. If in addition the $\{a_t\}$ are normally distributed then for large samples the information is $I(\theta) = \sum_0^\infty \theta^{2j} = (1-\theta^2)^{-1}$, and $\hat{\theta}^{(1)}$, $\hat{\theta}^{(2)}$ are efficient.

## REFERENCES

[1] ANDERSON, T. W. (1959). On asymptotic distributions of estimates of parameters of stochastic difference equations. *Ann. Math. Statist.* **30** 676–687.

[2] BOX, G. E. P. and JENKINS, G. M. (1970). *Time Series Analysis, Forecasting and Control.* Holden–Day, San Francisco. Chapter 7.

[3] BOX, G. E. P., JENKINS, G. M. and BACON, D. W. (1967). Models for forecasting seasonal and nonseasonal time series. *Spectral Analysis of Time Series*, ed. B. Harris. Wiley, New York, 271–311.

[4] DRAPER, N. R. and SMITH, H. (1966). *Applied Regression Analysis.* Wiley, New York, Chapter 10.

[5] DURBIN, J. (1959). Efficient estimation of parameters in moving average models. *Biometrika* **46** 306–316.

[6] HANNAN, E. J. (1960). *Time Series Analysis.* Methuen, London.

[7] MANN, H. B. and WALD, A. (1943). On the statistical treatment of linear stochastic difference equations. *Econometrica* **11** 173–220.

[8] WHITTLE, P. (1953). Estimation and information in stationary time series. *Ark. Mat.* **2** 423–434.