



## An Adaptive Regression Model

Thomas F. Cooley; Edward C. Prescott

*International Economic Review*, Vol. 14, No. 2. (Jun., 1973), pp. 364-371.

Stable URL:

<http://links.jstor.org/sici?sici=0020-6598%28197306%2914%3A2%3C364%3AAARM%3E2.0.CO%3B2-%23>

*International Economic Review* is currently published by Economics Department of the University of Pennsylvania.

---

Your use of the JSTOR archive indicates your acceptance of JSTOR's Terms and Conditions of Use, available at <http://www.jstor.org/about/terms.html>. JSTOR's Terms and Conditions of Use provides, in part, that unless you have obtained prior permission, you may not download an entire issue of a journal or multiple copies of articles, and you may use content in the JSTOR archive only for your personal, non-commercial use.

Please contact the publisher regarding any further use of this work. Publisher contact information may be obtained at [http://www.jstor.org/journals/ier\\_pub.html](http://www.jstor.org/journals/ier_pub.html).

Each copy of any part of a JSTOR transmission must contain the same copyright notice that appears on the screen or printed page of such transmission.

---

JSTOR is an independent not-for-profit organization dedicated to and preserving a digital archive of scholarly journals. For more information regarding JSTOR, please contact [support@jstor.org](mailto:support@jstor.org).

## AN ADAPTIVE REGRESSION MODEL\*

BY THOMAS F. COOLEY AND EDWARD C. PRESCOTT<sup>1</sup>

ECONOMETRICIANS frequently approximate complex behavioral and technological relationships using equations that are linear in a small number of unknown parameters. The effect of omitted variables, aggregation errors, and other errors in specification are included in the additive disturbance which is assumed, among other things, to be temporally uncorrelated. Utilizing time series data, linear regression analysis is then used to estimate parameters. The adaptive regression model developed in this paper would be used in the same manner, but it does not assume the disturbances are independent. Instead, it assumes the disturbances are the sum of not only a transitory element that has effect in the current period but also a permanent component whose effect persists into the future. If for example, omitted variables are subject to permanent and transitory changes, as is sometimes assumed in economic theory [4] and by the widely used adaptive forecasting model [6], these disturbances will have both permanent and transitory components. In the adaptive regression model the transitory disturbance can be thought of as the usual additive error term, while the permanent component causes random changes in the intercept value.

It is common practice in econometric research to test for serial correlation in the residuals. If the test indicates serial correlation is present it is typically assumed that the disturbances are subject to a first order auto-regressive process. In fact, such processes are likely to describe the true distribution of the disturbances only in rare instances. An auto-regressive error process implies that the effects of omitted factors *all* decay exponentially with time and at the *same* rate. This is an unreasonable assumption for most economic applications. Some omitted factors, such as labor union strikes or the vagaries of the weather, will have only transitory effects while other factors, like changes in tastes or technological developments, will have effects which persist into the future without decay. The auto-regressive assumption is often justified by the argument that omitted variables are subject to an auto-regressive process. This argument holds, however, only if *all* omitted factors contributing to the additive disturbance are subject to auto-regressive processes with the *same* parameter. The widespread use of the auto-regressive correction in econometrics is explained by the fact that it accounts for serial correlation and is computationally efficient. The adaptive regression also explains serial correlation, is computationally efficient, and assumes an error structure which, in many situations, provides a better approximation of reality.

Of particular methodological interest is the proof of consistency of the maximum likelihood estimator for a subset of the unknown parameters when the

\* Manuscript received August 16, 1972; revised January 5, 1973.

<sup>1</sup> The authors acknowledge helpful comments of Professors F. Gerard Adams, Michael D. McCarthy and Melvin Hinich.

observations are not identically and independently distributed and a consistent estimator does not exist for the entire parameter set.<sup>2</sup>

## 2. THE MODEL

The assumed structure is

$$(1.1) \quad y = \beta_{0t} + X'_t \beta^* + u_t$$

where  $y_t$  is the  $t$ -th observation of the dependent variable,  $\beta_{0t}$  the random intercept parameter for period  $t$ ,  $x_t$  a  $(k - 1)$  component vector of predetermined explanatory variables,  $\beta^*$  a vector of unknown slope coefficients, and  $u_t$  the additive transitory disturbance. The elements  $\beta_{0t}$  are subject to permanent changes  $v_t$ :

$$(1.2) \quad \beta_{0,t+1} = \beta_{0,t} + v_t.$$

The  $u_t$  and  $v_t$  are all independent normal variates with mean 0 and variances

$$(1.3) \quad \text{var}(u_t) = (1 - \gamma)\sigma^2 \text{ and } \text{var}(v_t) = \gamma\sigma^2,$$

with  $0 \leq \gamma \leq 1$ . The unknown parameter  $\gamma$  measures the relative importance of the permanent component, the larger its value the greater the importance of permanent change.

If  $\gamma = 0$ , the above structure is the multiple regression model with an unchanging intercept. On the other hand, if  $\beta^*$  equals the zero vector, this system reduces to the adaptive forecasting structure of [6]. Thus, the model is a generalization of both regression analysis and adaptive forecasting.<sup>3</sup>

The process generating the intercepts is not stationary and writing down the likelihood function is impossible. The likelihood function conditional on the value of the process at some point in time, however, is well defined for any finite set of these elements. One approach is to treat  $\beta_{0,0}$  as an unknown parameter and the other  $\beta_{0,t}$  as realizations of the random process, but, a more convenient selection for forecasting is the value of the intercept one period subsequent to the sample,  $\beta_{0,T+1}$ . Defining  $\beta_0$  to be  $\beta_{0,T+1}$  and using (1.3)

$$(1.4) \quad \beta_{0,t} = \beta_0 - \sum_{s=t}^T v_s.$$

Substituting for  $\beta_{0,t}$  in (1.1) yields

$$(1.5) \quad y_t = \beta_0 + x'_t \beta^* + u_t - \sum_{s=t}^T v_s.$$

Except for the the dependence of the error terms, this is the usual linear regres-

<sup>2</sup> In a subsequent paper [2], we compared the predictive and estimation efficiencies of adaptive and conventional regression analysis with and without the auto-regressive correction. A Fortran Program for adaptive regression is available upon request.

<sup>3</sup> The structure is related to those considered by [1] and [7] from a Bayesian point of view. Our estimation procedure is computationally more efficient and we are able to develop asymptotic properties of the estimators.

sion structure.

To simplify notation let  $y$  be the  $T$  component vector of the  $y_t$ ,  $\beta$  the  $k$  component vector

$$(1.6) \quad \beta' = [\beta_0, \beta_1^*, \beta_2^*, \dots, \beta_{k-1}^*],$$

and  $X$  the  $T \times k$  matrix

$$(1.7) \quad X = \begin{bmatrix} 1 & x_{11} & \cdot & \cdot & \cdot & x_{k-1,1} \\ 1 & x_{12} & \cdot & \cdot & \cdot & x_{k-1,2} \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ 1 & x_{1T} & \cdot & \cdot & \cdot & x_{k-1,T} \end{bmatrix}.$$

Define the  $T \times T$  matrix  $Q_\gamma$  to be

$$(1.8) \quad Q_\gamma = (1 - \gamma)I + \gamma R,$$

where the  $T \times T$  matrix  $R$  has  $(i, j)$ -th element

$$(1.9) \quad r_{ij} = \min [T - i + 1, T - j + 1].$$

With this notation along with (1.5), it is easily verified that

$$(1.10) \quad y \sim N[X\beta, \sigma^2 Q_\gamma].$$

If  $\gamma$  were known,  $Q_\gamma$  would be known and estimation would be a simple application of Aitken's generalized least squares analysis. The maximum likelihood estimator of  $\beta$  would be

$$(1.11) \quad B(\gamma) = (X'Q_\gamma^{-1}X)^{-1}XQ_\gamma^{-1}y,$$

and of  $\sigma^2$

$$(1.12) \quad s^2(\gamma) = T^{-1}[y - XB(\gamma)]'Q_\gamma^{-1}[y - XB(\gamma)].$$

Since  $\gamma$  is not known, search techniques must be utilized to determine the parameter set with maximum likelihood. The log likelihood function of the observations is (except for a constant)

$$(1.13) \quad L(y; \beta, \sigma^2, \gamma, X) = -\frac{T}{2} \ln \sigma^2 - \frac{1}{2} \ln |Q_\gamma| \\ - \frac{1}{2\sigma^2} (y - X\beta)' Q_\gamma^{-1} (y - X\beta).$$

Inserting the conditional maximum likelihood estimators of  $\beta$  and  $\sigma^2$  yields the concentrated likelihood function (except for a constant)

$$(1.14) \quad L_c(y; \gamma) = -\frac{T}{2} \ln s^2(\gamma) - \frac{1}{2} \ln |Q_\gamma|.$$

The estimation strategy is to search over the interval  $0 \leq \gamma \leq 1$  and choose as the estimator of  $\gamma$  that value, say  $g$ , such that

$$(1.15) \quad L_c(y; g) \geq L_c(y; \gamma) \text{ all } \gamma \in [0, 1].$$

The corresponding maximum likelihood estimates of  $\beta$  and  $\sigma^2$  are then  $B(g)$  and  $s^2(g)$  respectively.

The above procedure is straight-forward but excessively expensive, given current computer technology, because the  $T \times T$  matrix  $Q_\gamma$  must be inverted for each value of  $\gamma$  that is searched in the interval  $[0, 1]$ . Fortunately, the variables can be transformed so the covariance matrix is diagonal and the transformation does *not* depend upon  $\gamma$ .

Let  $P$  be the matrix whose rows are the set of orthonormal eigenvectors for  $R$ . Then  $Py \sim N[PX\beta, \sigma^2 P[(1 - \gamma)I + \gamma R]P']$  or

$$(1.16) \quad Py \sim N[PX\beta, \sigma^2 D(\gamma)]$$

where  $D(\gamma)$  is diagonal. Letting  $r_i$  be the eigenvalue corresponding to the  $i$ -th row of  $P$  and  $d_i$  the  $(i, i)$ -th element of  $D(\gamma)$ ,

$$(1.17) \quad d_i(\gamma) = (1 - \gamma) + \gamma r_i.$$

The analytic expressions for the  $r_i$  and  $p_{ij}$  are given by the following:

*Result*

$$(1.18) \quad r_i = [2 + 2 \cos \{2\pi(T - i + 1)/(2T + 1)\}]^{-1}$$

and

$$(1.19) \quad p_{ij} = (-1)^j 2(2T + 1)^{-1/2} \sin [2\pi(T - i + 1)(T - j + 1)/(2T + 1)].$$

PROOF. The inverse of  $R$  is a tridiagonal matrix with 2's down the main diagonal except for a 1 in the first position, and  $-1$  for the elements one position off the main diagonal. It is readily verified that the  $i$ -th row of  $P$  is an eigenvector corresponding to  $r_i^{-1}$  of the matrix  $R^{-1}$ . But, the eigenvectors of a symmetric matrix are the same as those of the inverse and the eigenvalues are the reciprocals. Thus, the rows of  $P$  are a set of orthonormal eigenvectors of  $R$  and the eigenvalue corresponding to row  $i$  is  $r_i$ .

Thus, the elements of  $D(\gamma)$  are defined by

$$(1.20) \quad d_i(\gamma) = (1 - \gamma) + \gamma \left[ 2 + 2 \cos \left( \frac{2\pi(T - i + 1)}{2T + 1} \right) \right]^{-1}.$$

No  $T \times T$  matrix need be inverted, and estimation costs are comparable those of ordinary least squares with auto-regressive correction.

## 2. LARGE SAMPLE ANALYSIS

The intercept cannot be estimated consistently because it is continually subject to random changes. But, if  $\gamma$  were known, the maximum likelihood estimator of  $\beta$  would be efficient in the sense that it would have minimum variance in the

class of unbiased estimators. In this section we prove that  $g$ , the m.l.e. of  $\gamma$ , is consistent implying  $B(g)$ , the m.l.e. of  $\beta$ , is asymptotically efficient. Further, the asymptotic distribution of  $B(g)$  is normal with mean  $\beta$  and covariance  $\sigma^2(X'Q_\gamma^{-1}X)^{-1}$ .

Subsequently  $(\gamma_0, \sigma_0^2)$  will denote the true value of  $(\gamma, \sigma^2)$  and the  $T$  subscript will be implied for those elements depending upon the sample size. Letting  $S$  be the generalized sum of squared residuals condition on  $\gamma$ , the concentrated log likelihood function, (1.15), divided by  $T/2$  is

$$(2.1) \quad L(\gamma; T) = -\ln S/T - T^{-1} \ln |D_\gamma|.$$

LEMMA A. For  $\gamma \in [0, 1]$

$$(2.2) \quad \text{plim } L(\gamma; T) = -\ln [\sigma_0^2 \sum d_t(\gamma_0)/d_t(\gamma)] - T^{-1} \ln |D_\gamma| \equiv f(\gamma; T)$$

where here and subsequently summations are from 1 to  $T$ .

PROOF. Assuming the variables have been transformed via transformation  $P$ , the generalized sum of squared residuals is

$$(2.3) \quad S = w'[I - M_\gamma]w$$

where  $w \sim N(0, \sigma_0^2 D_\gamma^{-1} D_0)$  and  $M_\gamma = D_\gamma^{-1/2} X(X' D_\gamma^{-1} X)^{-1} X' D_\gamma^{-1/2}$ . But from (1.20)

$$(2.4) \quad 0 < d_t(\gamma_0)/d_t(\gamma) \leq 4 + \gamma_0/\gamma$$

which uniformly bounds the variances and fourth moments of  $w_t$  for  $\gamma > 0$ . This along with the fact  $M_\gamma$  is idempotent of rank  $k$  implies  $w' M_\gamma w/T$  converges in mean to 0 and  $w'w/T$  to  $\sigma_0^2 \sum d_t(\gamma_0)/d_t(\gamma)$ . As convergence in mean implies convergence in probability the result is proven for  $\gamma > 0$ .

If  $\gamma = 0$  and  $\gamma_0 > 0$

$$L(\gamma; T) - f(\gamma; T) = -\ln [\sum w_i^2/\sigma_0^2 \sum d_t(\gamma_0)].$$

The variance of the term in brackets is

$$3\sigma_0^{-2} \sum d_t(\gamma_0)^2 / [\sum d_t(\gamma_0)^2].$$

Using well-known properties of eigenvalues [5, (273)], the numerator is of order  $T^3$  as it equals the sum of all the elements of  $Q_\gamma$  squared while the denominator is of order  $T^4$  as it is the square of the sum of the diagonal elements of  $Q_\gamma$ . Thus, the variance is of order  $T^{-1}$  implying convergence of the random variable to its mean in probability. If  $\gamma = 0$  and  $\gamma_0 = 0$ , the result is trivial. Thus,

$$\text{plim } [L(\gamma; T) - f(\gamma; T)] = -\ln [\sum Ew_i^2/\sigma_0^2 \sum d_t(\gamma_0)] = 0,$$

completing the proof.

REMARK. Pointwise convergence in probability does not imply convergence uniform in  $\gamma$  in probability. Uniform convergence (that is convergence in the  $L_\infty$ -norm) is needed to conclude  $L(\gamma; T)$  and  $f(\gamma, T)$  have the same maximum in the limit. Lemma B establishes a continuity condition for the  $L(\gamma; T)$  and

$f(\gamma; T)$  which implies uniform convergence given pointwise convergence. The final step in the proof is to show that  $\lim f(\gamma; T)$  is continuous and has a unique maxima at  $\gamma_0$ . Then one can conclude that  $g$  converges in probability to  $\gamma$ .

LEMMA B. *The convergence of  $L(\gamma; T)$  to  $f(\gamma; T)$  is uniform in probability; that is*

$$(2.5) \quad \text{plim sup}_{\gamma \in [0,1]} |L(\gamma; T) - f(\gamma; T)| = 0.$$

PROOF. Letting  $e(\gamma) = Py - PXB(\gamma)$ , then for  $\gamma_2 > \gamma_1$

$$\begin{aligned} S(\gamma_2) &\leq e(\gamma_1)'D_2^{-1}e(\gamma_1) \leq [1 + 3(\gamma_2 - \gamma_1)]e(\gamma_1)'D_1^{-1}e(\gamma_1) \\ &= [1 + 3(\gamma_2 - \gamma_1)]S(\gamma_1) \end{aligned}$$

because  $[1 + 3(\gamma_2 - \gamma_1)]D_1^{-1} - D_2^{-1}$  is positive definite, a result following from (1.20). Using this result

$$[S(\gamma_2) - S(\gamma_1)]/S(\gamma_1) \leq 3(\gamma_2 - \gamma_1)$$

which implies

$$(2.6) \quad \frac{d}{d\gamma}[L(\gamma; T) - \ln |D_\gamma|] \geq -3.$$

Similarly,

$$(2.7) \quad \frac{d}{d\gamma}[f(\gamma; T) - \ln |D_\gamma|] \geq -3.$$

If  $|L(\gamma_i; T) - f(\gamma_i; T)| \leq \varepsilon/3$  for  $i = 1, 2$ , then given (2.6) and (2.7)

$$(2.8) \quad |L(\gamma; T) - f(\gamma; T)| \leq 3(\gamma_1 - \gamma_2) + 2\varepsilon/3$$

for all  $\gamma_1 \leq \gamma \leq \gamma_2$ . By Lemma A for any  $\delta, \varepsilon > 0$  and  $N$ , there is a  $T^*$  such that for all  $T \geq T^*$

$$(2.9) \quad \Pr \{|f(i/N; T) - L(i/N; T)| < \varepsilon/3\} \geq 1 - \delta/(N + 1)$$

for  $i = 0, 1, \dots, N$ . This implies

$$\Pr \{\max_i |f(i/N; T) - L(i/N; T)| < \varepsilon/3\} \geq 1 - \delta.$$

But, by (2.8) and (2.9), for  $T \geq T^*$

$$\Pr \{\sup_{\gamma \in [0,1]} |f(\gamma; T) - L(\gamma; T)| \geq 2/3\varepsilon + 3/(N + 1)\} \leq 1 - \delta.$$

Selecting  $N > 8/\varepsilon$ , we find

$$\Pr \{\sup_{\gamma \in [0,1]} |f(\gamma; T) - L(\gamma; T)| \geq \varepsilon\} \leq 1 - \delta.$$

This proves the Lemma.

LEMMA C. *There is a  $K(\gamma; T) \geq K > 0$  such that*

$$f'(\gamma; T) = K(\gamma; T)(\gamma_0 - \gamma) \text{ for } \gamma \in (0, 1].$$

PROOF. Differentiating  $f(\gamma; T)$  yields

$$f'(\gamma; T) = \frac{1}{T} \sum \frac{d'_i(\gamma)d_i(\gamma_0)}{d_i(\gamma)^2} / \frac{1}{T} \sum \frac{d_i(\gamma_0)}{d_i(\gamma)} - \frac{1}{T} \sum \frac{d'_i(\gamma)}{d_i(\gamma)}.$$

Observe

$$\begin{aligned} f'(\gamma; T) &= \left[ \frac{1}{T} \sum \frac{d_i(\gamma_0)}{d_i(\gamma)} \right]^{-1} \frac{1}{T^2} \sum_{i,j} \frac{d'_i[d_i(\gamma_0)d_j(\gamma) - d_j(\gamma_0)d_i(\gamma)]}{d_i(\gamma)^2 d_j(\gamma)} \\ &= \left[ \frac{1}{T} \sum \frac{d_i(\gamma_0)}{d_i(\gamma)} \right]^{-1} \frac{\gamma_0 - \gamma}{2T^2 \gamma} \sum_{i,j} \frac{[d_i(\gamma) - d_j(\gamma)]^2}{d_i(\gamma)^2 d_j(\gamma)^2}. \end{aligned}$$

From (1.20)  $d_i(\gamma) > 1/4$ , which implies the leading term exceeds  $1/4$ . Thus, with some additional algebra,

$$f'(\gamma; T) = K_1(\gamma; T) [T^{-1} \sum d_i(\gamma)^{-2} - (T^{-1} \sum d_i(\gamma)^{-1})^2](\gamma_0 - \gamma)$$

with  $K_1(\gamma; T) > 1/8$ . Using the definition of  $d_i(\gamma)$ , a positive lower bound can be obtained for the average squared deviation of the  $d_i(\gamma)^{-1}$  from their average which holds for all  $T$  sufficiently large.

Lemma C implies that the functions  $f(\gamma; T)$  have a unique maximum of  $\gamma_0$ . Furthermore, if  $|\gamma - \gamma_0| > \varepsilon$ , then  $|f(\gamma_0; T) - f(\gamma; T)| \geq \varepsilon/K$ . By Lemma B the probability that  $L(\gamma; T)$  will be arbitrarily close to  $f(\gamma; T)$  for all values of  $\gamma$  approaches one as  $T$  goes to infinity. Thus, the  $\gamma$  which maximize  $L(\gamma; T)$  converges in probability to  $\gamma_0$ . This discussion can be summarized by the following theorem:

**THEOREM.** *The maximum likelihood estimator  $g$  of  $\gamma$  converges in probability to  $\gamma_0$ , the true parameter value.*

Following the usual analysis, but using the log likelihood function concentrated on  $\beta$ , the asymptotic distribution of  $\theta' = (\gamma, \sigma^2)$  can be derived (c.f. [3]). The basic result is that  $\sqrt{T}(\theta - \theta_0)$  is asymptotically normal with mean 0 and covariance equal to the inverse of the information matrix,

$$I(\theta_0) = \frac{1}{2T} \begin{bmatrix} \sum \left( \frac{d'_i(\gamma_0)}{d_i(\gamma_0)} \right)^2 & \frac{1}{\sigma_0^2} \sum \frac{d'_i(\gamma_0)}{d_i(\gamma_0)} \\ \frac{1}{\sigma_0^2} \sum \frac{d'_i(\gamma_0)}{d_i(\gamma_0)} & \frac{T}{\sigma_0^4} \end{bmatrix}.$$

REMARK. It is of interest to note that the only use of normality in the consistency proof was to insure the existence of the fourth moment of the  $w_r$ . The normality assumption merely provided a convenient function to be maximized.

*Tufts University, U.S.A.,  
Carnegie-Mellon University, U.S.A.*



## REFERENCES

- [ 1 ] BOX, G. E. P. AND G. C. TIAO, "Changes in Level of a Non-stationary Time Series," *Biometrika*, LII (April, 1965), 181-192.
- [ 2 ] COOLEY, T. F. AND E. C. PRESCOTT, "Test of the Adaptive Regression Model," *Review of Economics and Statistics*, (forthcoming).
- [ 3 ] DHRYMES, P., *Distributed Lags: Problems of Formulation and Estimation* (San Francisco: Holden-Day, 1971).
- [ 4 ] FRIEDMAN, M., *A Theory of the Consumption Function* (Princeton: Princeton University Press, 1957).
- [ 5 ] HADLEY, G., *Linear Algebra* (Reading, Massachusetts: Addison-Wesley, 1961).
- [ 6 ] MUTH, J., "Optimal Properties of Exponentially Weighted Forecasts" *Journal of the American Statistical Society*, LV (June, 1960), 299-306.
- [ 7 ] TIAO, G. C. AND M. M. ALI, "Analysis of Correlated Random Effects: Linear Models with Two Random Component," *Biometrika*, LVIII (April, 1971), 37-51.