

Univariate Time Series Models

Econ 241b, Fall 2005

T. Rothenberg

Contents

1	Discrete-Time Stationary Processes	2
1.1	Introduction	2
1.2	Autoregressive Models	3
1.3	Beveridge-Nelson Decompositions	5
1.4	Stationary ARMA Models	6
1.5	Nonstationary ARMA Models	7
1.6	State-Space Models	7
1.7	Long Memory Processes	8
2	Prediction	8
3	Estimating ARMA Models	10
4	ARMAX Models	12
4.1	Distributed Lags	12
4.2	Problems in Interpreting Dynamic Regression Models	13
5	ARCH and GARCH Models	14

1 Discrete-Time Stationary Processes

1.1 Introduction

A discrete-time random process is a doubly infinite sequence of random variables $\{y_t\}$ defined for all integer t . Such a process is said to be *strictly stationary* if, for any integer r , the joint distribution of any finite subset, say y_a, y_b, \dots, y_k , is the same as that of $y_{a+r}, y_{b+r}, \dots, y_{k+r}$. This is often expressed loosely by the phrase "calendar time does not matter." The sequence is said to be *weakly* (or covariance) stationary if each y_t has the same mean and if the covariance between each pair y_t and y_s depends only on the time difference $|t - s|$. Strict stationarity implies weak stationarity as long as the second moments exist. Unless otherwise stated, "stationarity" will be used in these notes to mean weak stationarity.

Stationarity puts a great deal of structure on a time series. To characterize the second-order moments of an arbitrary finite random sequence y_1, \dots, y_T , we need to specify T means, T variances, and $T(T-1)/2$ covariances. Stationarity reduces this to one mean, one variance, and $T-1$ covariances. More generally, the second-order moment properties of a weakly stationary process are completely described by its mean and its autocovariances $\gamma_r = E(y_t y_{t+r})$ for all nonnegative integer r . (Note $\gamma_{-r} = \gamma_r$.) It is often more convenient to present this information in the form of the mean, the variance γ_0 , and the autocorrelations $\rho_r = \gamma_r/\gamma_0$ for $r > 0$. To simplify the notation, we shall typically assume that the mean is zero.

The weakly stationary process $\{\varepsilon_t\}$ is said to be *white noise* if the ε_t are uncorrelated with mean zero and finite variance σ^2 . An infinite sequence $\{y_t\}$ is said to be a *time invariant linear process* if it can be expressed in the form

$$y_t = \sum_{j=-\infty}^{\infty} c_j \varepsilon_{t-j}$$

for some white noise process $\{\varepsilon_t\}$ and some infinite sequence of real numbers $\{c_j\}$ that are square summable; that is $\sum c_j^2 < \infty$. Then y_t is well defined in the sense of mean square convergence of the infinite sum. Usually we will make the stronger assumption that the c 's are absolutely summable: $\sum |c_j| < \infty$. Using the lag operator

$$C(L) = \sum_{j=-\infty}^{\infty} c_j L^j$$

where L is defined by $Lx_t \equiv x_{t-1}$, $L^{-1}x_t \equiv x_{t+1}$, and $L^0x_t = x_t$, we can write more succinctly $y_t = C(L)\varepsilon_t$. $C(L)$ is sometimes called a *linear filter*. If $c_j = 0$ when $j < 0$ so $C(L)$ is a polynomial in L , $C(L)$ is said to be a *one-sided backwards filter* and y_t is called a *moving average process*. We will usually set $c_0 = 1$. If $C_q(L)$ is a polynomial of finite order q , the moving average process $y_t = C_q(L)\varepsilon_t$ is said to be an MA(q) process. Its autocovariances are zero after lag q . Infinite order moving average processes with absolutely summable coefficients are always weakly stationary with autocovariances given by

$$\gamma_r = \sigma^2 \sum_{j=0}^{\infty} c_j c_{j+|r|}.$$

These autocovariances will also be absolutely summable.

Not every sequence $\{\gamma_r\}$ can be an autocovariance sequence. Furthermore, there are generally many different sequences $\{c_j\}$ that can generate the same valid autocovariance

sequence. For example, it is easily verified that the two MA(1) models

$$\begin{aligned} y_t &= (1 + \theta L)\varepsilon_t & \{\varepsilon_t\} \text{ white noise with variance } \sigma^2 \\ y_t &= (1 + \frac{1}{\theta}L)\eta_t & \{\eta_t\} \text{ white noise with variance } \sigma^2\theta^2 \end{aligned}$$

have the same autocovariances, so the MA representation is not unique when $\theta \neq 1$. In general we have the following result: The polynomial $C_q(L)$ can always be factored as

$$C_q(L) = (\lambda_1 - L)(\lambda_2 - L) \cdots (\lambda_q - L) \frac{c_0}{\lambda_1 \cdots \lambda_q}$$

where the λ 's are solutions of the equation $C_q(\lambda) = 0$. Let $C_q^*(L)$ be the polynomial obtained by replacing one or more of the λ 's by its reciprocal. Then there exists a positive scalar ω such that the models

$$\begin{aligned} y_t &= C_q(L)\varepsilon_t & \{\varepsilon_t\} \text{ white noise with variance } \sigma^2 \\ y_t &= C_q^*(L)\eta_t & \{\eta_t\} \text{ white noise with variance } \omega\sigma^2 \end{aligned}$$

have the same autocovariances.

1.2 Autoregressive Models

A common approach to modelling assumes that the observable time series $\{y_t\}$ is a stationary solution of a stochastic difference equation. The AR(1) model where, for all integer t , the y 's satisfy $y_t = \alpha y_{t-1} + \varepsilon_t$ for white noise $\{\varepsilon_t\}$ is a simple example. Sometimes the difference equation describes a behavioral or technological relation or has been obtained by algebraic manipulation from such a relation; more commonly, the equation is simply a convenient way to describe a series with a particular autocorrelation pattern. These notes discuss some statistical issues for a general class of difference equation models.

If $|\alpha| < 1$ in the above AR(1) model, we obtain by recursive substitution the moving average

$$y_t = \sum_{j=0}^{\infty} \alpha^j \varepsilon_{t-j}. \tag{1}$$

Since the moving average coefficients are square summable, this is a well-defined stationary solution to the difference equation. This representation can also be obtained using lag operators. Writing the difference equation as $(1 - \alpha L)y_t = \varepsilon_t$, multiplying both sides by $1/(1 - \alpha L)$, and using its power series expansion, we find

$$y_t = \frac{1}{1 - \alpha L} \varepsilon_t = (1 + \alpha L + \alpha^2 L^2 + \cdots) \varepsilon_t = \sum_{j=0}^{\infty} \alpha^j \varepsilon_{t-j}.$$

From difference-equation theory, we know that any time series $\{y_t\}$ of the form

$$y_t = \sum_{j=0}^{\infty} \alpha^j \varepsilon_{t-j} + A\alpha^t$$

for some constant A is a solution to $y_t = \alpha y_{t-1} + \varepsilon_t$. This is a well-defined stationary process only if $A = 0$ and $|\alpha| < 1$.

If $|\alpha| > 1$, the infinite sum (1) is not the mean-square limit of a sequence of random variables and does not constitute a meaningful solution to the difference equation. But, setting $\beta = 1/\alpha$, we find that the difference equation has the well-defined stationary solution

$$y_t = \frac{1}{1 - \alpha L} \varepsilon_t = -\frac{1}{L} \frac{\beta}{1 - \beta L^{-1}} \varepsilon_t = -\frac{\beta}{L} (1 + \beta L^{-1} + \beta^2 L^{-2} + \dots) \varepsilon_t = -\sum_{j=1}^{\infty} \beta^j \varepsilon_{t+j},$$

a convergent *forward* filter of $\{\varepsilon_t\}$. (This solution could also be obtained by recursive substitution using the implied forward difference equation $y_t = \beta y_{t+1} - \beta \varepsilon_{t+1}$.) Note that now y_{t-1} is correlated with ε_t so the equation $y_t = \alpha y_{t-1} + \varepsilon_t$ does not represent a linear regression.

Thus, as long as $|\alpha| \neq 1$, there exists a unique stationary solution to the difference equation $y_t = \alpha y_{t-1} + \varepsilon_t$. The autocovariances for this stationary AR(1) process are $\gamma_r = \alpha^{|r|} \sigma_\varepsilon^2 / (1 - \alpha^2)$ if $|\alpha| < 1$ and $\gamma_r = \alpha^{-|r|-2} \sigma_\varepsilon^2 / (1 - \alpha^{-2})$ if $|\alpha| > 1$. This implies that, if $|\alpha|$ is neither zero nor one, the two models

$$\begin{aligned} y_t &= \alpha y_{t-1} + \varepsilon_t, & \{\varepsilon_t\} &\text{white noise with variance } \sigma^2 \\ y_t &= \alpha^{-1} y_{t-1} + \eta_t, & \{\eta_t\} &\text{white noise with variance } \sigma^2 / \alpha^2 \end{aligned}$$

have stationary solutions with the same second-order moments for the observable data and are in this sense equivalent models. We shall normally work with the model whose autoregressive coefficient is less than one in absolute value so a backward moving average representation exists and the "error" will be uncorrelated with the "regressor." Indeed, many authors reserve the words "stationary AR(1) model" for this invertible case.

Let

$$A_p(L) = 1 - a_1 L - \dots - a_p L^p$$

be a polynomial of order p in the lag operator L so that

$$A_p(L)y_t = y_t - a_1 y_{t-1} - \dots - a_p y_{t-p}.$$

Then, a p th order autoregressive model can be written concisely as $A_p(L)y_t = \varepsilon_t$. If there exists a lag polynomial $C(L)$ (possibly of infinite order as long as it contains only nonnegative powers of L and the c_j are absolutely summable) such that $A_p(L)C(L) = 1$, then $C(L) \equiv A_p(L)^{-1}$ is called the inverse of $A_p(L)$. Since $A_p(L)$ can be factored as

$$A_p(L) = (\lambda_1 - L)(\lambda_2 - L) \dots (\lambda_p - L) \frac{1}{\lambda_1 \dots \lambda_p}$$

where the λ 's are the p (possibly complex) roots of the polynomial equation $A_p(\lambda) = 0$, $A_p(L)$ has an inverse if and only if all these roots lie outside the unit circle in the complex plane. If $A_p(L)$ is invertible, the difference equation $A_p(L)y_t = \varepsilon_t$ can be solved to express y_t as an infinite moving average $y_t = A_p^{-1}(L)\varepsilon_t$, where

$$A^{-1}(L_p) = \prod_{j=1}^p (1 - \lambda_j^{-1} L)^{-1} = \prod_{j=1}^p (1 + \lambda_j^{-1} L + \lambda_j^{-2} L^2 + \dots) = 1 + c_1 L + c_2 L^2 + \dots$$

This is the unique stationary solution to the difference equation. (Note that if any of the λ 's are inside the unit circle, replacing them with their reciprocals in the factorization yields an invertible polynomial with the same autocorrelation function as before. Since the moving average representation is useful for developing forecasts and estimates, we will always use the

invertible version of $A_p(L)$. Again, many authors call this version "the" stationary AR(p) model. Others call it a "causal" AR model.)

The autocorrelation function for an AR(p) process with invertible lag polynomial can be found by solving a difference equation. Suppose

$$y_t = \alpha_1 y_{t-1} + \cdots + \alpha_p y_{t-p} + \varepsilon_t. \quad (2)$$

Since we are assuming invertibility, y_t has a moving average representation and $E y_t \varepsilon_s = 0$ when $t < s$. Then, for positive integer r , multiplying both sides of (2) by y_{t-r} , taking expectations, and dividing by γ_0 , we find the *Yule – Walker* equations

$$\rho_r = \alpha_1 \rho_{r-1} + \alpha_2 \rho_{r-2} + \cdots + \alpha_p \rho_{r-p}.$$

Thus the autocorrelations satisfy the homogeneous p th order difference equation that generates the y_t . The high-order autocorrelations ultimately die off exponentially to zero.

For any weakly stationary process $\{y_t\}$, let V_r be the $r \times r$ covariance matrix for (y_1, \dots, y_r) and let C_r be the r -dimensional column vector $(\gamma_1, \dots, \gamma_r)'$. Then the r th *partial autocorrelation* is defined as the r th element of $V_r^{-1} C_r$. The first partial autocorrelation coefficient is just the first ordinary autocorrelation coefficient ρ_1 . For the AR(p) process (2), the p th partial autocorrelation coefficient is α_p and the higher order coefficients are zero.

1.3 Beveridge-Nelson Decompositions

Any equation relating variables at different dates can be rewritten so that some level variables are replaced by first-differenced variables. Suppose, for example, $\{y_t\}$ is a stationary MA(2) process so $y_t = \varepsilon_t + c_1 \varepsilon_{t-1} + c_2 \varepsilon_{t-2}$. By rearranging terms, it is easily verified that this equation can be rewritten as

$$y_t = (1 + c_1 + c_2)\varepsilon_t - (c_1 + c_2)\Delta\varepsilon_t - c_2\Delta\varepsilon_{t-1}.$$

More generally, suppose $y_t = C(L)\varepsilon_t$ where $C(L) = \sum_{i=0}^{\infty} c_i L^i$. Then the MA representation can be rewritten as

$$y_t = C(1)\varepsilon_t + C^*(L)\Delta\varepsilon_t$$

where $C^*(L) = \sum_{i=0}^{\infty} c_i^* L^i$ and $c_i^* = -\sum_{j=i+1}^{\infty} c_j$. The coefficients c_i^* are absolutely summable as long as $\sum_{i=0}^{\infty} i|c_i| < \infty$. Defining $v_t = C^*(L)\varepsilon_t$, this implies that y_t has the representation

$$y_t = c\varepsilon_t + \Delta v_t$$

where $c = C(1)$ and v_t is a stationary process. This representation is often called the Beveridge-Nelson decomposition of the MA process and is very useful when deriving asymptotic properties of sample moments. For example, suppose one wanted to show that the standardized sample average of the dependent y data is asymptotically normal. Since

$$T^{-1/2} \sum_{t=1}^T y_t = cT^{-1/2} \sum_{t=1}^T \varepsilon_t + T^{-1/2}(v_T - v_0)$$

and the final term converges in probability to zero as $T \rightarrow \infty$, one only need verify that the standardized average of the ε 's is asymptotically normal. If the ε 's are i.i.d. with finite variance, this follows from a standard central limit theorem.

A similar decomposition holds for AR processes. Suppose $y_t = \alpha_1 y_{t-1} + \alpha_2 y_{t-2} + \varepsilon_t$. It is easy to verify that this equation can be rewritten as $\Delta y_t = (\alpha_1 + \alpha_2 - 1)y_{t-1} - \alpha_2 \Delta y_{t-1} + \varepsilon_t$.

More generally, if $A(L)y_t = \varepsilon_t$ for some (possible infinite order) lag polynomial $A(L) = \sum a_i L^i$ where $\sum i|a_i| < \infty$, then y_t has the alternative representation

$$\Delta y_t = -A(1)y_{t-1} + A^*(L)\Delta y_{t-1} + \varepsilon_t$$

where $A^*(L)$ is a lag polynomial with absolutely summable coefficients. The usefulness of this representation will become apparent later when we consider unit root models.

1.4 Stationary ARMA Models

Let $A_p(L)$ and $B_q(L)$ be polynomials in L of orders p and q , having no common factors. Let $\{\varepsilon_t\}$ be white noise. A time series $\{y_t\}$ satisfying the difference equation

$$A_p(L)y_t = B_q(L)\varepsilon_t$$

is called a zero-mean ARMA(p,q) process. A unique stationary solution exists as long as none of the roots of the equation $A(z) = 0$ lie on the unit circle. If all the roots lie strictly outside the unit circle, then A has an inverse and the series has a moving average representation $y_t = A(L)^{-1}B(L)\varepsilon_t$. If $B(L)$ is invertible, y_t also has a (possibly infinite) autoregressive representation $B(L)^{-1}A(L)y_t = \varepsilon_t$. If some of the roots of either $A(L) = 0$ or $B(z) = 0$ are inside the unit circle, they can be replaced by their reciprocals in the factorization of $A(L)$ or $B(L)$ without changing the autocorrelations of the $\{y_t\}$ process. Thus, as long as neither polynomial has a unit root, we will choose the invertible versions; the unique stationary solution to the difference equation will then have both an autoregressive and moving average representation. Indeed, ARMA models can be motivated by the assumption that $\{y_t\}$ has infinite AR and MA representations

$$y_t = C(L)\varepsilon_t \quad \text{and} \quad D(L)y_t = \varepsilon_t$$

where $C(L)$ can be well approximated by $B(L)/A(L)$ and $D(L)$ can be well approximated by $A(L)/B(L)$.

If one or more roots of $A(L) = 0$ lie on the unit circle, then there is no stationary solution to the difference equation. If those roots are real and positive, then taking differences will produce a stationary series. Unit roots of $B(L) = 0$ do not affect the existence of a stationary solution to the difference equation, but typically indicate that the time series $\{y_t\}$ has been overdifferenced. In an influential book, Box and Jenkins argue that, after perhaps taking first or second differences, many real-world time series can be well approximated as stationary invertible ARMA(p,q) processes with small values for p and q . If d differences are needed to produce stationarity, the process is denoted ARIMA(p,d,q). If $\{y_t\}$ is ARIMA (p,d,q), then $\{\Delta^d y_t\}$ is ARMA(p,q); hence, properties of ARIMA processes can easily be derived from the properties of ARMA processes.

The autocovariances for an ARMA(p,q) process with invertible lag polynomials can be found as follows. Since $Ey_t\varepsilon_s = 0$ when $t < s$, if we multiply both sides of

$$y_t = \alpha_1 y_{t-1} + \cdots + \alpha_p y_{t-p} + \varepsilon_t + \beta_1 \varepsilon_{t-1} + \cdots + \beta_q \varepsilon_{t-q} \quad (3)$$

by the term y_{t-r} and take expectations, we find

$$\gamma_r = \alpha_1 \gamma_{r-1} + \alpha_2 \gamma_{r-2} + \cdots + \alpha_p \gamma_{r-p} \quad r > q \quad (4)$$

$$\gamma_r = \alpha_1 \gamma_{r-1} + \alpha_2 \gamma_{r-2} + \cdots + \alpha_p \gamma_{r-p} + \beta_r E y_{t-r} \varepsilon_{t-r} + \cdots + \beta_q E y_{t-r} \varepsilon_{t-q} \quad 0 \leq r \leq q$$

Thus the autocovariances of order $q + 1$ and higher satisfy the same difference equation as the autocovariances of an AR(p) process. The high-order autocovariances ultimately die off exponentially to zero. The low-order autocovariances can be found by multiplying both sides of equation (3) by ε_{t-r} for $r = 0, \dots, q$ and taking expectations. The resulting equations can be used with (4) above to solve for the γ 's.

ARMA models arise quite naturally from the aggregation of AR models. Suppose $A_p(L)x_t = u_t$ and $B_q(L)z_t = v_t$ where u_t and v_t are white noise. Then, if the lag polynomials are invertible,

$$y_t \equiv x_t + z_t = \frac{u_t}{A_p(L)} + \frac{v_t}{B_q(L)} = \frac{B_q(L)u_t + A_p(L)v_t}{A_p(L)B_q(L)}.$$

But $B_q(L)u_t + A_p(L)v_t$ is a moving average process of order $r = \max(p, q)$ and $A_p(L)B_q(L)$ is a lag polynomial of order $s = p + q$. Hence y_t is an ARMA(s, r) process. More generally, a linear combination of AR processes will be an ARMA process.

1.5 Nonstationary ARMA Models

The discussion in the previous paragraphs assumes the equation $A_p(L)y_t = B_q(L)\varepsilon_t$ holds for all integer t . This is unnecessarily strong given that we often only need to model an observed finite-length series y_1, \dots, y_T . An alternative is to assume that the difference equation holds only for $t = p + 1, \dots, T$ and to make specific assumptions on the initial random variables y_1, \dots, y_p . It can be shown that, as long as $A(L)$ has an inverse, there always exist initial conditions such that the observed series is stationary and has the same autocovariances (up to order $T - 1$) as that generated by the difference equation with infinite history.

Of course, we do not have to assume stationarity. Sometimes it is convenient to assume special initial conditions that may introduce a little nonstationarity but simplify forecasting and estimation. Or, if the difference equation arises from some physical experiment as in many engineering applications, the initial conditions may describe the actual state of the system when the experiment started. With arbitrary initial conditions, the nonstationary ARMA process is well defined even if the autoregressive lag polynomial has unit roots.

1.6 State-Space Models

Suppose the scalar process $\{y_t\}$ is generated as

$$y_t = b' \alpha_t + u_t \quad t = 1, 2, \dots$$

where $\{\alpha_t\}$ is an *unobserved* p -dimensional vector process generated by the first-order difference equation

$$\alpha_t = A\alpha_{t-1} + v_t \quad t = 1, 2, \dots$$

A is a nonrandom square matrix, b is a nonrandom vector, $\{v_t\}$ is a vector white noise process, and $\{u_t\}$ is a scalar white noise process independent of $\{v_t\}$ and the initial value α_0 . The α_t are called state variables and the process $\{y_t\}$ is said to have a state-space representation. It turns out that every ARMA process has a state-space representation. Thus the class of state space models includes the class of ARMA models. As we shall show later, some results about ARMA models can be derived most easily as special cases of results for state-space models.

1.7 Long Memory Processes

In contrast to ARMA processes whose autocorrelations ultimately drop off to zero at a geometric rate, *long memory processes* have autocorrelations that drop off much more slowly. An example is the fractionally differenced process generated by the stochastic difference equation

$$(1 - L)^d y_t = \varepsilon_t \quad (5)$$

where $\{\varepsilon_t\}$ is white noise, $d \in (-0.5, 0.5)$, and the fractional difference operator $(1 - L)^d$ is defined by its power series expansion

$$(1 - L)^d = 1 - dL + \frac{1}{2}d(d-1)L^2 + \dots .$$

It can be shown that equation (5) always has a stationary solution with moving average representation of the form

$$y_t = (1 - L)^{-d} \varepsilon_t = \sum_{j=0}^{\infty} c_j \varepsilon_{t-j}$$

where, for large j and $d \neq 0$,

$$c_j \approx \frac{j^{d-1}}{\Gamma(d)} \quad \text{and} \quad \rho_j \approx \frac{\Gamma(1-d)}{\Gamma(d)} j^{2d-1} .$$

When $0 < d < 0.5$, the moving average coefficients $\{c_j\}$ are square summable, but not absolutely summable. The autocorrelation coefficients $\{\rho_j\}$ drop off slowly and are not even square summable when $0.25 \leq d < 0.5$.

A generalization of (5) is the fractionally differenced ARIMA(p,d,q) model

$$A_p(L)(1 - L)^d y_t = B_q(L)\varepsilon_t$$

where A_p and B_q are invertible lag polynomials of finite orders and $d \in (-0.5, 0.5)$. In contrast to the case where d is a positive integer, this difference equation has a stationary solution $\{y_t\}$, although with long memory properties.

2 Prediction

Let Y be a scalar random variable and let x be a (possibly infinite dimensional) vector random variable defined on the same probability space. The best predictor of Y given x is defined to be that function $g^*(x)$ such that $E[Y - g^*(x)]^2$ is minimized. Assuming that second moments exist, we find $g^*(x) = E(Y|x)$. The best *linear* predictor of Y given x (denoted $P(Y|x)$) is defined to be that scalar $a + b'x$ such that $E[Y - a - b'x]^2$ is minimized. Again assuming second moments exist, we find that the solution must satisfy

$$c_{xy} = V_{xx} b^* \quad \text{and} \quad a^* = EY - b^{*'} E x$$

where c_{xy} is the column vector of covariances between Y and x and V_{xx} is the covariance matrix for x . If V_{xx} is finite dimensional and nonsingular, we find

$$P(Y|x) = EY + c'_{xy} V_{xx}^{-1} (x - E x).$$

The mean square prediction error using this optimal linear predictor is

$$E(Y - P(Y|x))^2 = \text{var}(Y) - c'_{xy} V_{xx}^{-1} c_{xy} .$$

When x is a vector of high dimension, inverting V_{xx} is computationally demanding. Sometimes one can compute (or at least approximate) $P(Y|x)$ without explicitly finding the inverse. The best linear predictor has the following useful properties:

1. The prediction error $e \equiv Y - P(Y|x)$ has mean zero and is uncorrelated with every element of x .
2. For constants a , b , and c , $P(aY + bZ + c|x) = aP(Y|x) + bP(Z|x) + c$.
3. If Y is uncorrelated with each element of x , then $P(Y|x) = EY$.
4. If Y is an element of x or a linear combination of elements of x , then $P(Y|x) = Y$.
5. $P(Y|x) = P[P(Y|x, z)|x]$.

Convenient expressions for $P(Y|x)$ can be found in the special case where Y and x are realizations at different dates of the stationary ARMA(p,q) process $A(L)y_t = B(L)\varepsilon_t$. In particular, for any integer $s > 0$, let $P_T(y_{T+s})$ be the best linear predictor of y_{T+s} given the infinite history $y_T, y_{T-1}, y_{T-2}, \dots$. Suppose $A(L)$ and $B(L)$ are both invertible. Then y_{T+s} can be written as a moving average

$$y_{T+s} = \frac{B(L)}{A(L)}\varepsilon_{T+s} = \varepsilon_{T+s} + c_1\varepsilon_{T+s-1} + \dots + c_s\varepsilon_T + c_{s+1}\varepsilon_{T-1} + \dots$$

The best forecast of future ε 's is zero. Current and past ε 's can be perfectly predicted from the autoregressive representation. Thus

$$P_T(y_{T+s}) = c_s\varepsilon_T + c_{s+1}\varepsilon_{T-1} + \dots = \left(\frac{B(L)}{A(L)L^s} \right)_+ \varepsilon_T = \left(\frac{B(L)}{A(L)L^s} \right)_+ \frac{A(L)}{B(L)}y_T$$

where $D(L)_+$ is defined to be the lag polynomial $D(L)$ with terms containing negative powers of L dropped.

When $s = 1$, this expression simplifies since

$$\left(\frac{B(L)}{A(L)L} \right)_+ = \frac{1}{L} \left[\frac{B(L)}{A(L)} - 1 \right] = \frac{B(L) - A(L)}{A(L)L}.$$

We find

$$P_T(y_{T+1}) = \frac{B(L) - A(L)}{B(L)L}y_T$$

which implies $P_T(y_{T+1}) = y_{T+1} - \varepsilon_{T+1}$. That is, if $\{y_t\}$ is the stationary solution of an ARMA model with invertible lag polynomials, ε_t is the difference between y_t and the best linear predictor of y_t given its past history. The ε_t are often called *innovations* and interpreted as new information entering the system at time t .

Of course, in practice we can only make predictions based on a finite past history. If we truncate the infinite autoregression, it can serve as an approximation to the best linear predictor given a finite past history. As we shall see later, the Kalman filter provides an efficient computer algorithm for computing the best one-step-ahead linear predictor $P(y_{T+1}|y_1, \dots, y_T)$ exactly.

A stationary process $\{y_t\}$ is said to be *deterministic* if $P_T(y_{T+1}) = y_{T+1}$; that is, if it can be linearly predicted without error from its past history. Clearly, ARMA processes are never deterministic unless the innovation variance is zero. The *Wold decomposition theorem* asserts that every weakly stationary process can be written as the sum of a deterministic process and a moving average process. If the moving average coefficients drop off to zero rapidly after some finite lag, the moving average part can be approximated by an ARMA(p,q) process. Hence this theorem suggests that ARMA models should capture the second-order moment properties of many stationary time series after any perfectly predictable component has been removed. The exceptional cases where the moving average coefficients decline only slowly will be discussed later.

3 Estimating ARMA Models

The Box-Jenkins methodology for fitting a model to a time series $\{x_t\}$ consists of five steps:

1. Decide on the order of differencing d that is needed to produce a stationary series $y_t = (1 - L)^d x_t$ that can be approximated by an ARMA(p,q) model (with intercept if the mean of y_t is not zero).
2. By inspecting the sample autocorrelations and partial autocorrelations of $\{y_t\}$, determine tentative values for p and q .
3. Estimate the lag coefficients by (quasi) maximum likelihood assuming normally distributed innovations.
4. Compute approximate standard errors and confidence intervals for the unknown coefficients.
5. Using various diagnostics, check if the tentative model was indeed appropriate.

Here we shall briefly discuss a commonly used algorithm for computing parameter estimates. Alternative estimation procedures and asymptotic inference theory will be discussed later.

Suppose (y_1, \dots, y_T) were generated by an ARMA(p,q) model with normal innovations. Then, the y 's are normal with means and covariances depending on a parameter vector θ consisting of the unknown coefficients of the lag polynomials (and the intercept). If the lag orders p and q are known and small, θ can be estimated by maximum likelihood. However, this is computationally demanding when the sample size T is large since the calculation involves inverting the $T \times T$ covariance matrix. Later we shall show how a computational algorithm known as the Kalman filter can be employed to obtain exact maximum likelihood estimates. In practice, most practitioners maximize an approximation to the likelihood based on a slight change in the initial conditions. For example, in the $AR(p)$ model, if we condition on the first p values y_1, \dots, y_p and examine the process starting at $t = p + 1$, maximum likelihood is equivalent to a least squares regression of y_t on its p lagged values. We now show that a similar modification of initial conditions in the general ARMA case leads to a *nonlinear* least squares regression.

In every invertible ARMA(p,q) model, the observed y 's have a moving average representation expressing them linearly to current and lagged ε 's. Hence the likelihood function (the joint density of the observed y 's) can be derived from the density of the ε 's using familiar change-of-variable techniques. Unfortunately, the mapping y to ε is typically not one-to-one, so the calculation is nontrivial. However, conditioning on the observed values y_1, \dots, y_p and on $\varepsilon_p = \varepsilon_{p-1} = \dots = \varepsilon_{p-q} = 0$ necessarily leads to a one-to-one mapping. Let \mathbf{e} be the vector of white noise innovations $(\varepsilon_{p+1}, \dots, \varepsilon_T)'$ and let \mathbf{y} be the vector of observations $(y_{p+1}, \dots, y_T)'$. Then, by successive substitution, we can write $\mathbf{e} = \mathbf{D}\mathbf{y} + \mathbf{d}$ where \mathbf{D} is a triangular matrix with ones on the diagonal and \mathbf{d} is a vector of nonrandom quantities. Thus we can write $\mathbf{e} = \mathbf{e}(\mathbf{y}, \theta)$ where the Jacobian $|\partial\varepsilon/\partial y|$ does not depend on θ . If \mathbf{e}/σ is standard normal, the change-of-variable rule gives the density for \mathbf{y} having the form

$$f(\mathbf{y}) = (2\pi\sigma^2)^{-(T-p)/2} \exp\left\{-\frac{1}{2}\mathbf{e}(\mathbf{y}, \theta)' \mathbf{e}(\mathbf{y}, \theta)/\sigma^2\right\}$$

Maximum likelihood estimates of θ can be obtained by minimizing $\mathbf{e}'\mathbf{e}$. Although \mathbf{e} is linear in \mathbf{y} , it is nonlinear in θ as long as q (the number of moving average coefficients to be estimated) is greater than zero. This nonlinear least squares problem can be solved using the

Gauss-Newton algorithm. Let $m = p + q$ be the number of unknown parameter in θ and let $n = T - p$ be the number of observations after conditioning. (If there is an intercept, then $m = p + q + 1$.) Then, defining the $n \times m$ matrix $Z(\theta, y) \equiv \partial e / \partial \theta'$ and starting with some initial guess θ^0 , we use the recursion

$$\theta^{r+1} = \theta^r - (Z_r' Z_r)^{-1} Z_r' e^r$$

where $Z_r = Z(\theta^r, y)$ and $e^r = e(\theta^r, y)$. When $Z_r' \theta^r \simeq 0$, one stops.

For example, suppose $y_t = \alpha y_{t-1} + \varepsilon_t + \beta \varepsilon_{t-1}$ and we treat y_1 as fixed and assume $\varepsilon_1 = 0$. Then, using the equation $\varepsilon_t = y_t - \alpha y_{t-1} - \beta \varepsilon_{t-1}$ with $\varepsilon_1 = 0$, we find

$$\begin{aligned} \varepsilon_2 &= y_2 - \alpha y_1 \\ \varepsilon_3 &= y_3 - (\alpha + \beta) y_2 + \alpha \beta y_1 \\ \varepsilon_4 &= y_4 - (\alpha + \beta) y_3 + \beta(\alpha + \beta) y_2 - \alpha \beta^2 y_1 \\ &\vdots \end{aligned}$$

Clearly, the Jacobian matrix $\partial \varepsilon / \partial y$ is triangular with determinant equal to one. However, we do not have to explicitly solve for the ε 's in terms of the y 's. Under normality, the MLE which minimizes the sum of squared innovations can be computed using the following G-N algorithm:

1. For some initial $\theta^0 = (\alpha^0, \beta^0)'$ and starting with the initial condition

$$\varepsilon_1 = \partial \varepsilon_1 / \partial \alpha = \partial \varepsilon_1 / \partial \beta = 0,$$

build up the vector e^0 and the $(T - 1) \times 2$ matrix $Z_0 = [\partial e / \partial \alpha, \partial \varepsilon / \partial \beta]$ using the recursions

$$\begin{aligned} \varepsilon_t &= y_t - \alpha^0 y_{t-1} - \beta^0 \varepsilon_{t-1} \\ \partial \varepsilon_t / \partial \alpha &= -y_{t-1} - \beta^0 \partial \varepsilon_{t-1} / \partial \alpha \\ \partial \varepsilon_t / \partial \beta &= -\varepsilon_{t-1} - \beta^0 \partial \varepsilon_{t-1} / \partial \beta \end{aligned}$$

2. Regress e^0 on Z_0 and compute $\theta^1 = \theta^0 - (Z_0' Z_0)^{-1} Z_0' e^0$.
3. Redo steps one and two using θ^1 in place of θ^0 so $\theta^2 = \theta^1 - (Z_1' Z_1)^{-1} Z_1' e^1$.
4. Continue to convergence.

Of course, to insure that the algorithm converges to a global (and not just a local) minimum, alternative starting values should be used. Or at least the starting values should be chosen as some consistent method-of-moments estimator that has high probability of being near the true parameter value.

Finally, note that the G-N algorithm maximizes an approximation to the likelihood since ε_1 is not really 0 and y_1 is not really constant. However, the G-N estimate is usually close to the actual MLE as long as the parameters are far away from the noninvertibility boundaries $|\alpha| = 1$ and $|\beta| = 1$.

4 ARMAX Models

The ARMA model can be generalized to allow for the mean of y_t to depend on the exogenous variables $x_{1t}, x_{2t}, \dots, x_{Kt}$, possibly with some lag. Many dynamic economic relations can be expressed in the so-called ARMAX form

$$A(L)y_t = \mu + F_1(L)x_{1t} + \dots + F_k(L)x_{kt} + B(L)\varepsilon_t$$

where A , B , and the F 's are low-order lag polynomials and μ is a scalar intercept. Results stated above for estimation and hypothesis testing carry over with only slight modification.

4.1 Distributed Lags

ARMAX models are commonly employed in the context of distributed lag estimation. Suppose some economic variable y moves over time in response to some variable x , but perhaps with some lag. One might write

$$y_t = \alpha + \beta C(L)x_t + u_t, \quad D(L)u_t = \varepsilon_t \quad (6)$$

where the ε_t are white noise and D is an invertible lag polynomial. The first equation represents the (causal) distributed lag relation; the second equation models the autocorrelation structure of the errors (capturing omitted effects). The variable x is assumed to be exogenous; that is, $E[\varepsilon_t | \text{past } y\text{'s and all } x\text{'s}] = 0$. Three alternative estimation strategies are available:

1. Approximate $C(L)$ and $D(L)$ by ratios of low-order polynomials
2. Unconstrained OLS or GLS estimation of $C(L)$;
3. OLS or GLS estimation with polynomial constraints on the coefficients in $C(L)$.

A brief discussion of these methods is given below.

The *rational distributed lag* approach uses ratios of low-order polynomials to approximate the (possibly high order) polynomials $C(L)$ and $D(L)$. If

$$C(L) = \frac{P(L)}{Q(L)} \quad \text{and} \quad D(L) = \frac{A(L)}{B(L)},$$

then we obtain the ARMAX model

$$A(L)Q(L)y_t = \bar{\alpha} + \beta A(L)P(L)x_t + Q(L)B(L)\varepsilon_t. \quad (7)$$

The simplest version of the rational lag model (first proposed by Koyck) assumes no autocorrelation in the residuals and geometrically declining weights in $C(L)$ so $A(L) = B(L) = P(L) = 1$ and $Q(L) = 1 - \beta_1 L$. Then (7) becomes

$$y_t = \beta_0 + \beta_1 y_{t-1} + \beta_3 x_t + \varepsilon_t - \beta_1 \varepsilon_{t-1}$$

This is commonly estimated by the method of instrumental variables, using lagged x 's and y 's as instruments.

In general, models like (6) can be estimated more efficiently by nonlinear least squares. Unfortunately, consistent estimation of the ARMAX model by either instrumental variables or by Gauss-Newton requires that the lag orders are known. Box and Jenkins suggest possible ways to estimate these orders, but it is not clear that they are very successful.

In practice, we are often unwilling to pretend we know the orders of the lag polynomials. In that case, we might just run an OLS regression of y_t on $(x_t, x_{t-1}, \dots, x_{t-p})$ for some fairly large value of p . If the true order of the lag polynomial is less than p , this estimate of $C(L)$ will be unbiased even if the errors are autocorrelated. A look at the residuals might suggest a second-stage GLS estimator. Estimating a long distributed lag by OLS or GLS is likely to lead to collinearity problems. Although this will typically not seriously affect the precision of the estimate of the long-run multiplier, it will lead to large standard errors for individual lag parameters. To lessen collinearity, one can impose constraints on the lag coefficients. Forcing the c_j to lie on a polynomial in j leads to a linear regression with fewer parameters. For example, if $p = 4$ and

$$C(L) = 1 + c_1L + c_2L^2 + c_3L^3 + c_4L^4,$$

the linear constraint that $c_j = (4-j)/4$ yields a regression (with intercept) of y_t on the single explanatory variable

$$x_t^* = \frac{4x_t + 3x_{t-1} + 2x_{t-2} + x_{t-3}}{4}.$$

Since this polynomial distributed lag approach never introduces lagged y 's as regressors, standard linear model theory holds. The main cost is that the first p observations needed to construct x^* are lost; with 50 years of annual data and a lag length of 10, one would lose 20% of the observations.

All three approaches to distributed lag estimation assume that the x 's are exogenous. The equation is usually given a causal interpretation as the response of some economic variable to changes in its determinants. How do we know that the direction of causality is not reversed? or that both variables are caused by a third variable? Sims has argued that, before attempting to fit a distributed lag, one should first regress y_t on some past, current and future x 's. If the future x 's appear with significant coefficients, the assumption that x is exogenous is suspect. We shall discuss such exogeneity tests in more detail when we consider multivariate dynamic models.

4.2 Problems in Interpreting Dynamic Regression Models

Even if we are sure that x is exogenous, there can be conceptual problems in interpreting dynamic economic models. Some of the problems can be illustrated by the following simple example. Suppose one theory is that people set y_t based on its past value y_{t-1} and on the past value of some exogenous variable x_{t-1} . Given y_{t-1} and x_{t-1} , the past is irrelevant. The model is

$$y_t = \alpha y_{t-1} + \beta x_{t-1} + \varepsilon_t$$

where ε_t represent other determinants of y_t and is assumed to be white noise. Another theory is that people set y_t as a geometrically declining weighted average of past x 's:

$$y_t = \frac{\beta}{1 - \alpha L} x_{t-1} + u_t$$

where u_t represents other determinants of y_t . If $(1 - \alpha L)u_t$ is white noise, the two models are equivalent: it is not possible to tell from the data which theory is correct. This does not matter if our goal is forecasting in a world where the process generating the data has not changed. In both cases, $E(y_t | \text{past data}) = \alpha y_{t-1} + \beta x_{t-1}$. But suppose the system were interrupted (say, due to war) and started over again. Then the two models would predict very different behavior in the second period. Prior knowledge about the form of autocorrelation in the errors can be used to distinguish between economically different models. But in the absence

of such knowledge, it is hard to distinguish among alternative behavioral models containing lagged endogenous variables. In the distributed lag model (6), the problem was assumed away by the specification that only lagged x 's enter the behavioral equation. In the absence of such assumptions, dynamic ARMAX models are open to many possible interpretations.

5 ARCH and GARCH Models

The models investigated so far are based on second-order moments and are designed primarily to help forecast future values of y_t from its past history. That is, the goal was to specify a simple model for the conditional mean of y_t given its past. Sometimes, however, we may want to base the analysis on higher-order moments. This is particularly the case when the goal is to model the conditional *variance* of y_t given the past. For example, when analyzing stock market returns for the purpose of pricing derivative securities, the Black-Scholes theory requires modelling both the conditional mean and the conditional variance of returns.

An AR(p) model with i.i.d. innovations specifies that the conditional expectation of y_t given past values is a linear function of the most recent p past values but the conditional variance of y_t given past values is constant. If we permit the conditional variance to vary with t we get a model with conditional heteroskedasticity. Suppose, for example,

$$C_p(L)y_t = u_t$$

where $E(u_t|\text{past } y\text{'s}) = 0$ and

$$\text{Var}(y_t|\text{past } y\text{'s}) = \beta + D_q(L)u_{t-1}^2$$

so that past deviations of y_t from its conditional mean help predict future variability of y_t . Then we have a model exhibiting *autoregressive conditional heteroskedasticity* (ARCH). The parameters of this model could be estimated simply (but generally inefficiently) by a two-step procedure. The u_t could be estimated by the residuals from a regression of y_t on p lagged y 's; then the remaining parameters could be estimated from a regression of \hat{u}_t^2 on q of its lagged values.

More efficient estimates can be obtained if we are willing to make the somewhat stronger assumption that the u_t 's can be expressed in terms of an underlying i.i.d. process. Specifically, if we assume that

$$u_t = \sqrt{h_t}\varepsilon_t$$

where $h_t = \text{Var}(y_t|\text{past } y\text{'s}) = \beta + D_q(L)u_{t-1}^2$ and the ε_t are i.i.d. with mean zero and variance one, then we get the ARCH process described in the previous paragraph. Moreover, if $f(\varepsilon)$ is the density function for each of the ε_t , the joint density of the observations y_{p+q+1}, \dots, y_T (conditional on observations up to period $p+q$) is

$$\prod_{t=p+q+1}^T f\left(\frac{C(L)y_t}{\sqrt{\beta + D(L)[C(L)y_t]^2}}\right) \frac{1}{\sqrt{\beta + D(L)[C(L)y_t]^2}}.$$

If the density f is known and we are willing to condition on the first $p+q$ values, the parameters of the model can be estimated by maximum likelihood.

Estimation in a high dimensional parameter space can be avoided by approximating $C(L)$ and $D(L)$ by ratios of low order polynomials. If, for example, $C(L) = A(L)/B(L)$ and $D(L) =$

$Q(L)/P(L)$, then we have a generalized ARCH (or GARCH) process for the conditional variance added on to an ARMA model for the mean:

$$\begin{aligned}A(L)y_t &= B(L)u_t \\u_t &= \sqrt{h_t}\varepsilon_t \\P(L)h_t &= \bar{\beta} + Q(L)u_{t-1}^2\end{aligned}$$

Again, under appropriate initial conditions, the parameters can be estimated by maximum likelihood if the density for ε_t is known.