

## Chapter 3

---

### ARIMA Models

Classical regression is often insufficient for explaining all of the interesting dynamics of a time series. For example, the ACF of the residuals of the simple linear regression fit to the price of chicken data (see [Example 2.4](#)) reveals additional structure in the data that regression did not capture. Instead, the introduction of correlation that may be generated through lagged linear relations leads to proposing the *autoregressive (AR)* and *autoregressive moving average (ARMA)* models that were presented in Whittle (1951). Adding nonstationary models to the mix leads to the *autoregressive integrated moving average (ARIMA)* model popularized in the landmark work by Box and Jenkins (1970). The *Box–Jenkins method* for identifying ARIMA models is given in this chapter along with techniques for *parameter estimation* and *forecasting* for these models. A partial theoretical justification of the use of ARMA models is discussed in [Section B.4](#).

### 3.1 Autoregressive Moving Average Models

The classical regression model of [Chapter 2](#) was developed for the static case, namely, we only allow the dependent variable to be influenced by current values of the independent variables. In the time series case, it is desirable to allow the dependent variable to be influenced by the past values of the independent variables and possibly by its own past values. If the present can be plausibly modeled in terms of only the past values of the independent inputs, we have the enticing prospect that forecasting will be possible.

#### INTRODUCTION TO AUTOREGRESSIVE MODELS

Autoregressive models are based on the idea that the current value of the series,  $x_t$ , can be explained as a function of  $p$  past values,  $x_{t-1}, x_{t-2}, \dots, x_{t-p}$ , where  $p$  determines the number of steps into the past needed to forecast the current value. As a typical case, recall [Example 1.10](#) in which data were generated using the model

$$x_t = x_{t-1} - .90x_{t-2} + w_t,$$

where  $w_t$  is white Gaussian noise with  $\sigma_w^2 = 1$ . We have now assumed the current value is a particular *linear* function of past values. The regularity that persists in [Figure 1.9](#) gives an indication that forecasting for such a model might be a distinct possibility, say, through some version such as

$$x_{n+1}^n = x_n - .90x_{n-1},$$

where the quantity on the left-hand side denotes the forecast at the next period  $n + 1$  based on the observed data,  $x_1, x_2, \dots, x_n$ . We will make this notion more precise in our discussion of forecasting ([Section 3.4](#)).

The extent to which it might be possible to forecast a real data series from its own past values can be assessed by looking at the autocorrelation function and the lagged scatterplot matrices discussed in [Chapter 2](#). For example, the lagged scatterplot matrix for the Southern Oscillation Index (SOI), shown in [Figure 2.8](#), gives a distinct indication that lags 1 and 2, for example, are linearly associated with the current value. The ACF shown in [Figure 1.16](#) shows relatively large positive values at lags 1, 2, 12, 24, and 36 and large negative values at 18, 30, and 42. We note also the possible relation between the SOI and Recruitment series indicated in the scatterplot matrix shown in [Figure 2.9](#). We will indicate in later sections on transfer function and vector AR modeling how to handle the dependence on values taken by other series.

The preceding discussion motivates the following definition.

**Definition 3.1** An autoregressive model of order  $p$ , abbreviated  $AR(p)$ , is of the form

$$x_t = \phi_1 x_{t-1} + \phi_2 x_{t-2} + \dots + \phi_p x_{t-p} + w_t, \quad (3.1)$$

where  $x_t$  is stationary,  $w_t \sim wn(0, \sigma_w^2)$ , and  $\phi_1, \phi_2, \dots, \phi_p$  are constants ( $\phi_p \neq 0$ ). The mean of  $x_t$  in (3.1) is zero. If the mean,  $\mu$ , of  $x_t$  is not zero, replace  $x_t$  by  $x_t - \mu$  in (3.1),

$$x_t - \mu = \phi_1(x_{t-1} - \mu) + \phi_2(x_{t-2} - \mu) + \dots + \phi_p(x_{t-p} - \mu) + w_t,$$

or write

$$x_t = \alpha + \phi_1 x_{t-1} + \phi_2 x_{t-2} + \dots + \phi_p x_{t-p} + w_t, \quad (3.2)$$

where  $\alpha = \mu(1 - \phi_1 - \dots - \phi_p)$ .

We note that (3.2) is similar to the regression model of [Section 2.1](#), and hence the term auto (or self) regression. Some technical difficulties, however, develop from applying that model because the regressors,  $x_{t-1}, \dots, x_{t-p}$ , are random components, whereas  $z_t$  was assumed to be fixed. A useful form follows by using the backshift operator (2.29) to write the  $AR(p)$  model, (3.1), as

$$(1 - \phi_1 B - \phi_2 B^2 - \dots - \phi_p B^p)x_t = w_t, \quad (3.3)$$

or even more concisely as

$$\phi(B)x_t = w_t. \quad (3.4)$$

The properties of  $\phi(B)$  are important in solving (3.4) for  $x_t$ . This leads to the following definition.

**Definition 3.2** The autoregressive operator is defined to be

$$\phi(B) = 1 - \phi_1 B - \phi_2 B^2 - \cdots - \phi_p B^p. \quad (3.5)$$

**Example 3.1 The AR(1) Model**

We initiate the investigation of AR models by considering the first-order model, AR(1), given by  $x_t = \phi x_{t-1} + w_t$ . Iterating backwards  $k$  times, we get

$$\begin{aligned} x_t &= \phi x_{t-1} + w_t = \phi(\phi x_{t-2} + w_{t-1}) + w_t \\ &= \phi^2 x_{t-2} + \phi w_{t-1} + w_t \\ &\vdots \\ &= \phi^k x_{t-k} + \sum_{j=0}^{k-1} \phi^j w_{t-j}. \end{aligned}$$

This method suggests that, by continuing to iterate backward, and provided that  $|\phi| < 1$  and  $\sup_t \text{var}(x_t) < \infty$ , we can represent an AR(1) model as a linear process given by<sup>3.1</sup>

$$x_t = \sum_{j=0}^{\infty} \phi^j w_{t-j}. \quad (3.6)$$

Representation (3.6) is called the stationary solution of the model. In fact, by simple substitution,

$$\underbrace{\sum_{j=0}^{\infty} \phi^j w_{t-j}}_{x_t} = \phi \underbrace{\left( \sum_{k=0}^{\infty} \phi^k w_{t-1-k} \right)}_{x_{t-1}} + w_t.$$

The AR(1) process defined by (3.6) is stationary with mean

$$E(x_t) = \sum_{j=0}^{\infty} \phi^j E(w_{t-j}) = 0,$$

and autocovariance function,

$$\begin{aligned} \gamma(h) &= \text{cov}(x_{t+h}, x_t) = E \left[ \left( \sum_{j=0}^{\infty} \phi^j w_{t+h-j} \right) \left( \sum_{k=0}^{\infty} \phi^k w_{t-k} \right) \right] \\ &= E \left[ \left( w_{t+h} + \cdots + \phi^h w_t + \phi^{h+1} w_{t-1} + \cdots \right) (w_t + \phi w_{t-1} + \cdots) \right] \\ &= \sigma_w^2 \sum_{j=0}^{\infty} \phi^{h+j} \phi^j = \sigma_w^2 \phi^h \sum_{j=0}^{\infty} \phi^{2j} = \frac{\sigma_w^2 \phi^h}{1 - \phi^2}, \quad h \geq 0. \end{aligned} \quad (3.7)$$

<sup>3.1</sup> Note that  $\lim_{k \rightarrow \infty} E \left( x_t - \sum_{j=0}^{k-1} \phi^j w_{t-j} \right)^2 = \lim_{k \rightarrow \infty} \phi^{2k} E \left( x_{t-k}^2 \right) = 0$ , so (3.6) exists in the mean square sense (see [Appendix A](#) for a definition).

Recall that  $\gamma(h) = \gamma(-h)$ , so we will only exhibit the autocovariance function for  $h \geq 0$ . From (3.7), the ACF of an AR(1) is

$$\rho(h) = \frac{\gamma(h)}{\gamma(0)} = \phi^h, \quad h \geq 0, \quad (3.8)$$

and  $\rho(h)$  satisfies the recursion

$$\rho(h) = \phi \rho(h-1), \quad h = 1, 2, \dots \quad (3.9)$$

We will discuss the ACF of a general AR( $p$ ) model in Section 3.3.

### Example 3.2 The Sample Path of an AR(1) Process

Figure 3.1 shows a time plot of two AR(1) processes, one with  $\phi = .9$  and one with  $\phi = -.9$ ; in both cases,  $\sigma_w^2 = 1$ . In the first case,  $\rho(h) = .9^h$ , for  $h \geq 0$ , so observations close together in time are positively correlated with each other. This result means that observations at contiguous time points will tend to be close in value to each other; this fact shows up in the top of Figure 3.1 as a very smooth sample path for  $x_t$ . Now, contrast this with the case in which  $\phi = -.9$ , so that  $\rho(h) = (-.9)^h$ , for  $h \geq 0$ . This result means that observations at contiguous time points are negatively correlated but observations two time points apart are positively correlated. This fact shows up in the bottom of Figure 3.1, where, for example, if an observation,  $x_t$ , is positive, the next observation,  $x_{t+1}$ , is typically negative, and the next observation,  $x_{t+2}$ , is typically positive. Thus, in this case, the sample path is very choppy.

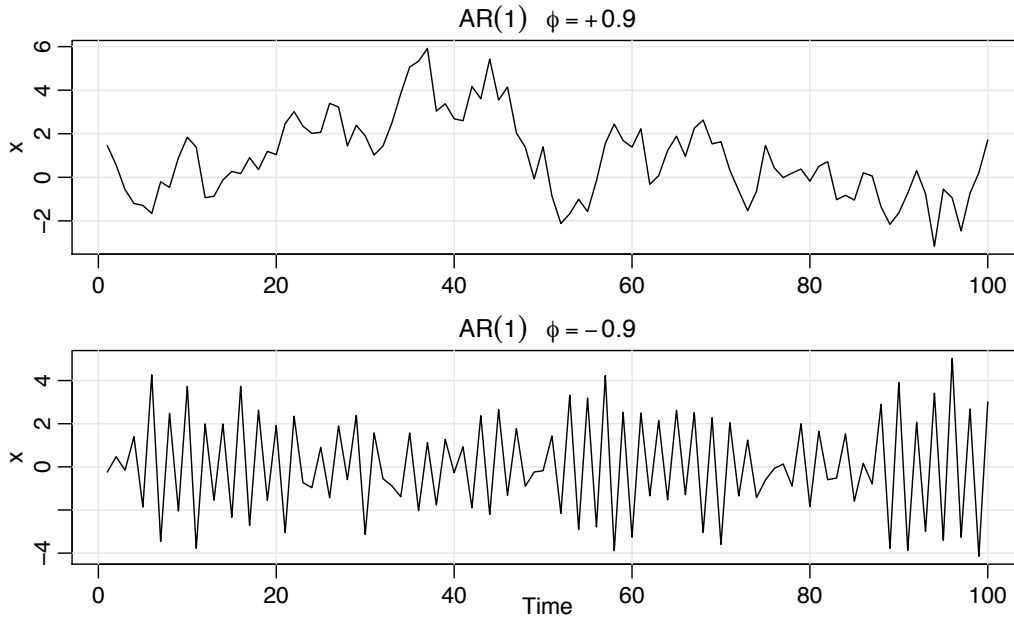
The following R code can be used to obtain a figure similar to Figure 3.1:

```
par(mfrow=c(2,1))
plot(arima.sim(list(order=c(1,0,0), ar=.9), n=100), ylab="x",
     main=expression(AR(1)~~~phi==+.9))
plot(arima.sim(list(order=c(1,0,0), ar=-.9), n=100), ylab="x",
     main=expression(AR(1)~~~phi==-.9))
```

### Example 3.3 Explosive AR Models and Causality

In Example 1.18, it was discovered that the random walk  $x_t = x_{t-1} + w_t$  is not stationary. We might wonder whether there is a stationary AR(1) process with  $|\phi| > 1$ . Such processes are called explosive because the values of the time series quickly become large in magnitude. Clearly, because  $|\phi|^j$  increases without bound as  $j \rightarrow \infty$ ,  $\sum_{j=0}^{k-1} \phi^j w_{t-j}$  will not converge (in mean square) as  $k \rightarrow \infty$ , so the intuition used to get (3.6) will not work directly. We can, however, modify that argument to obtain a stationary model as follows. Write  $x_{t+1} = \phi x_t + w_{t+1}$ , in which case,

$$\begin{aligned} x_t &= \phi^{-1} x_{t+1} - \phi^{-1} w_{t+1} = \phi^{-1} \left( \phi^{-1} x_{t+2} - \phi^{-1} w_{t+2} \right) - \phi^{-1} w_{t+1} \\ &\vdots \\ &= \phi^{-k} x_{t+k} - \sum_{j=1}^{k-1} \phi^{-j} w_{t+j}, \end{aligned} \quad (3.10)$$



**Fig. 3.1.** Simulated AR(1) models:  $\phi = .9$  (top);  $\phi = -.9$  (bottom).

by iterating forward  $k$  steps. Because  $|\phi|^{-1} < 1$ , this result suggests the stationary future dependent AR(1) model

$$x_t = - \sum_{j=1}^{\infty} \phi^{-j} w_{t+j}. \quad (3.11)$$

The reader can verify that this is stationary and of the AR(1) form  $x_t = \phi x_{t-1} + w_t$ . Unfortunately, this model is useless because it requires us to know the future to be able to predict the future. When a process does not depend on the future, such as the AR(1) when  $|\phi| < 1$ , we will say the process is *causal*. In the explosive case of this example, the process is stationary, but it is also future dependent, and not causal.

### Example 3.4 Every Explosion Has a Cause

Excluding explosive models from consideration is not a problem because the models have causal counterparts. For example, if

$$x_t = \phi x_{t-1} + w_t \quad \text{with} \quad |\phi| > 1$$

and  $w_t \sim \text{iid } N(0, \sigma_w^2)$ , then using (3.11),  $\{x_t\}$  is a non-causal stationary Gaussian process with  $E(x_t) = 0$  and

$$\begin{aligned} \gamma_x(h) &= \text{cov}(x_{t+h}, x_t) = \text{cov}\left(- \sum_{j=1}^{\infty} \phi^{-j} w_{t+h+j}, - \sum_{k=1}^{\infty} \phi^{-k} w_{t+k}\right) \\ &= \sigma_w^2 \phi^{-2} \phi^{-h} / (1 - \phi^{-2}). \end{aligned}$$

Thus, using (3.7), the causal process defined by

$$y_t = \phi^{-1} y_{t-1} + v_t$$

where  $v_t \sim \text{iid } N(0, \sigma_w^2 \phi^{-2})$  is stochastically equal to the  $x_t$  process (i.e., all finite distributions of the processes are the same). For example, if  $x_t = 2x_{t-1} + w_t$  with  $\sigma_w^2 = 1$ , then  $y_t = \frac{1}{2}y_{t-1} + v_t$  with  $\sigma_v^2 = 1/4$  is an equivalent causal process (see Problem 3.3). This concept generalizes to higher orders, but it is easier to show using Chapter 4 techniques; see Example 4.8.

The technique of iterating backward to get an idea of the stationary solution of AR models works well when  $p = 1$ , but not for larger orders. A general technique is that of matching coefficients. Consider the AR(1) model in operator form

$$\phi(B)x_t = w_t, \quad (3.12)$$

where  $\phi(B) = 1 - \phi B$ , and  $|\phi| < 1$ . Also, write the model in equation (3.6) using operator form as

$$x_t = \sum_{j=0}^{\infty} \psi_j w_{t-j} = \psi(B)w_t, \quad (3.13)$$

where  $\psi(B) = \sum_{j=0}^{\infty} \psi_j B^j$  and  $\psi_j = \phi^j$ . Suppose we did not know that  $\psi_j = \phi^j$ . We could substitute  $\psi(B)w_t$  from (3.13) for  $x_t$  in (3.12) to obtain

$$\phi(B)\psi(B)w_t = w_t. \quad (3.14)$$

The coefficients of  $B$  on the left-hand side of (3.14) must be equal to those on right-hand side of (3.14), which means

$$(1 - \phi B)(1 + \psi_1 B + \psi_2 B^2 + \cdots + \psi_j B^j + \cdots) = 1. \quad (3.15)$$

Reorganizing the coefficients in (3.15),

$$1 + (\psi_1 - \phi)B + (\psi_2 - \psi_1\phi)B^2 + \cdots + (\psi_j - \psi_{j-1}\phi)B^j + \cdots = 1,$$

we see that for each  $j = 1, 2, \dots$ , the coefficient of  $B^j$  on the left must be zero because it is zero on the right. The coefficient of  $B$  on the left is  $(\psi_1 - \phi)$ , and equating this to zero,  $\psi_1 - \phi = 0$ , leads to  $\psi_1 = \phi$ . Continuing, the coefficient of  $B^2$  is  $(\psi_2 - \psi_1\phi)$ , so  $\psi_2 = \phi^2$ . In general,

$$\psi_j = \psi_{j-1}\phi,$$

with  $\psi_0 = 1$ , which leads to the solution  $\psi_j = \phi^j$ .

Another way to think about the operations we just performed is to consider the AR(1) model in operator form,  $\phi(B)x_t = w_t$ . Now multiply both sides by  $\phi^{-1}(B)$  (assuming the inverse operator exists) to get

$$\phi^{-1}(B)\phi(B)x_t = \phi^{-1}(B)w_t,$$

or

$$x_t = \phi^{-1}(B)w_t.$$

We know already that

$$\phi^{-1}(B) = 1 + \phi B + \phi^2 B^2 + \cdots + \phi^j B^j + \cdots,$$

that is,  $\phi^{-1}(B)$  is  $\psi(B)$  in (3.13). Thus, we notice that working with operators is like working with polynomials. That is, consider the polynomial  $\phi(z) = 1 - \phi z$ , where  $z$  is a complex number and  $|\phi| < 1$ . Then,

$$\phi^{-1}(z) = \frac{1}{(1 - \phi z)} = 1 + \phi z + \phi^2 z^2 + \cdots + \phi^j z^j + \cdots, \quad |z| \leq 1,$$

and the coefficients of  $B^j$  in  $\phi^{-1}(B)$  are the same as the coefficients of  $z^j$  in  $\phi^{-1}(z)$ . In other words, we may treat the backshift operator,  $B$ , as a complex number,  $z$ . These results will be generalized in our discussion of ARMA models. We will find the polynomials corresponding to the operators useful in exploring the general properties of ARMA models.

#### INTRODUCTION TO MOVING AVERAGE MODELS

As an alternative to the autoregressive representation in which the  $x_t$  on the left-hand side of the equation are assumed to be combined linearly, the moving average model of order  $q$ , abbreviated as  $MA(q)$ , assumes the white noise  $w_t$  on the right-hand side of the defining equation are combined linearly to form the observed data.

**Definition 3.3** *The moving average model of order  $q$ , or  $MA(q)$  model, is defined to be*

$$x_t = w_t + \theta_1 w_{t-1} + \theta_2 w_{t-2} + \cdots + \theta_q w_{t-q}, \quad (3.16)$$

where  $w_t \sim wn(0, \sigma_w^2)$ , and  $\theta_1, \theta_2, \dots, \theta_q$  ( $\theta_q \neq 0$ ) are parameters.<sup>3.2</sup>

The system is the same as the infinite moving average defined as the linear process (3.13), where  $\psi_0 = 1$ ,  $\psi_j = \theta_j$ , for  $j = 1, \dots, q$ , and  $\psi_j = 0$  for other values. We may also write the  $MA(q)$  process in the equivalent form

$$x_t = \theta(B)w_t, \quad (3.17)$$

using the following definition.

**Definition 3.4** *The moving average operator is*

$$\theta(B) = 1 + \theta_1 B + \theta_2 B^2 + \cdots + \theta_q B^q. \quad (3.18)$$

Unlike the autoregressive process, the moving average process is stationary for any values of the parameters  $\theta_1, \dots, \theta_q$ ; details of this result are provided in Section 3.3.

<sup>3.2</sup> Some texts and software packages write the MA model with negative coefficients; that is,  $x_t = w_t - \theta_1 w_{t-1} - \theta_2 w_{t-2} - \cdots - \theta_q w_{t-q}$ .

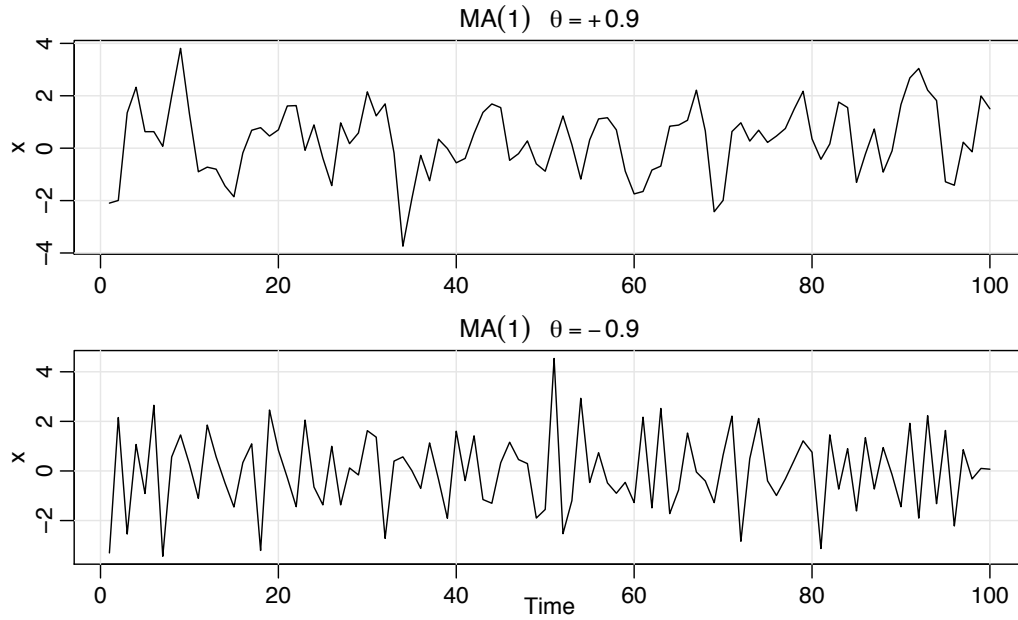


Fig. 3.2. Simulated MA(1) models:  $\theta = .9$  (top);  $\theta = -.9$  (bottom).

### Example 3.5 The MA(1) Process

Consider the MA(1) model  $x_t = w_t + \theta w_{t-1}$ . Then,  $E(x_t) = 0$ ,

$$\gamma(h) = \begin{cases} (1 + \theta^2)\sigma_w^2 & h = 0, \\ \theta\sigma_w^2 & h = 1, \\ 0 & h > 1, \end{cases}$$

and the ACF is

$$\rho(h) = \begin{cases} \frac{\theta}{(1+\theta^2)} & h = 1, \\ 0 & h > 1. \end{cases}$$

Note  $|\rho(1)| \leq 1/2$  for all values of  $\theta$  (Problem 3.1). Also,  $x_t$  is correlated with  $x_{t-1}$ , but not with  $x_{t-2}, x_{t-3}, \dots$ . Contrast this with the case of the AR(1) model in which the correlation between  $x_t$  and  $x_{t-k}$  is never zero. When  $\theta = .9$ , for example,  $x_t$  and  $x_{t-1}$  are positively correlated, and  $\rho(1) = .497$ . When  $\theta = -.9$ ,  $x_t$  and  $x_{t-1}$  are negatively correlated,  $\rho(1) = -.497$ . Figure 3.2 shows a time plot of these two processes with  $\sigma_w^2 = 1$ . The series for which  $\theta = .9$  is smoother than the series for which  $\theta = -.9$ .

A figure similar to Figure 3.2 can be created in R as follows:

```
par(mfrow = c(2,1))
plot(arima.sim(list(order=c(0,0,1), ma=.9), n=100), ylab="x",
     main=(expression(MA(1)~~~theta==+.5)))
plot(arima.sim(list(order=c(0,0,1), ma=-.9), n=100), ylab="x",
     main=(expression(MA(1)~~~theta==-.5)))
```



**Example 3.6 Non-uniqueness of MA Models and Invertibility**

Using Example 3.5, we note that for an MA(1) model,  $\rho(h)$  is the same for  $\theta$  and  $\frac{1}{\theta}$ ; try 5 and  $\frac{1}{5}$ , for example. In addition, the pair  $\sigma_w^2 = 1$  and  $\theta = 5$  yield the same autocovariance function as the pair  $\sigma_w^2 = 25$  and  $\theta = 1/5$ , namely,

$$\gamma(h) = \begin{cases} 26 & h = 0, \\ 5 & h = 1, \\ 0 & h > 1. \end{cases}$$

Thus, the MA(1) processes

$$x_t = w_t + \frac{1}{5}w_{t-1}, \quad w_t \sim \text{iid } N(0, 25)$$

and

$$y_t = v_t + 5v_{t-1}, \quad v_t \sim \text{iid } N(0, 1)$$

are the same because of normality (i.e., all finite distributions are the same). We can only observe the time series,  $x_t$  or  $y_t$ , and not the noise,  $w_t$  or  $v_t$ , so we cannot distinguish between the models. Hence, we will have to choose only one of them. For convenience, by mimicking the criterion of causality for AR models, we will choose the model with an infinite AR representation. Such a process is called an *invertible* process.

To discover which model is the invertible model, we can reverse the roles of  $x_t$  and  $w_t$  (because we are mimicking the AR case) and write the MA(1) model as  $w_t = -\theta w_{t-1} + x_t$ . Following the steps that led to (3.6), if  $|\theta| < 1$ , then  $w_t = \sum_{j=0}^{\infty} (-\theta)^j x_{t-j}$ , which is the desired infinite AR representation of the model. Hence, given a choice, we will choose the model with  $\sigma_w^2 = 25$  and  $\theta = 1/5$  because it is invertible.

As in the AR case, the polynomial,  $\theta(z)$ , corresponding to the moving average operators,  $\theta(B)$ , will be useful in exploring general properties of MA processes. For example, following the steps of equations (3.12)–(3.15), we can write the MA(1) model as  $x_t = \theta(B)w_t$ , where  $\theta(B) = 1 + \theta B$ . If  $|\theta| < 1$ , then we can write the model as  $\pi(B)x_t = w_t$ , where  $\pi(B) = \theta^{-1}(B)$ . Let  $\theta(z) = 1 + \theta z$ , for  $|z| \leq 1$ , then  $\pi(z) = \theta^{-1}(z) = 1/(1 + \theta z) = \sum_{j=0}^{\infty} (-\theta)^j z^j$ , and we determine that  $\pi(B) = \sum_{j=0}^{\infty} (-\theta)^j B^j$ .

**AUTOREGRESSIVE MOVING AVERAGE MODELS**

We now proceed with the general development of autoregressive, moving average, and mixed *autoregressive moving average* (ARMA), models for stationary time series.

**Definition 3.5** A time series  $\{x_t; t = 0, \pm 1, \pm 2, \dots\}$  is **ARMA**( $p, q$ ) if it is stationary and

$$x_t = \phi_1 x_{t-1} + \dots + \phi_p x_{t-p} + w_t + \theta_1 w_{t-1} + \dots + \theta_q w_{t-q}, \quad (3.19)$$

with  $\phi_p \neq 0$ ,  $\theta_q \neq 0$ , and  $\sigma_w^2 > 0$ . The parameters  $p$  and  $q$  are called the *autoregressive* and the *moving average orders*, respectively. If  $x_t$  has a nonzero mean  $\mu$ , we set  $\alpha = \mu(1 - \phi_1 - \dots - \phi_p)$  and write the model as

$$x_t = \alpha + \phi_1 x_{t-1} + \cdots + \phi_p x_{t-p} + w_t + \theta_1 w_{t-1} + \cdots + \theta_q w_{t-q}, \quad (3.20)$$

where  $w_t \sim wn(0, \sigma_w^2)$ .

As previously noted, when  $q = 0$ , the model is called an autoregressive model of order  $p$ ,  $AR(p)$ , and when  $p = 0$ , the model is called a moving average model of order  $q$ ,  $MA(q)$ . To aid in the investigation of ARMA models, it will be useful to write them using the AR operator, (3.5), and the MA operator, (3.18). In particular, the  $ARMA(p, q)$  model in (3.19) can then be written in concise form as

$$\phi(B)x_t = \theta(B)w_t. \quad (3.21)$$

The concise form of the model points to a potential problem in that we can unnecessarily complicate the model by multiplying both sides by another operator, say

$$\eta(B)\phi(B)x_t = \eta(B)\theta(B)w_t,$$

without changing the dynamics. Consider the following example.

### Example 3.7 Parameter Redundancy

Consider a white noise process  $x_t = w_t$ . If we multiply both sides of the equation by  $\eta(B) = 1 - .5B$ , then the model becomes  $(1 - .5B)x_t = (1 - .5B)w_t$ , or

$$x_t = .5x_{t-1} - .5w_{t-1} + w_t, \quad (3.22)$$

which looks like an  $ARMA(1, 1)$  model. Of course,  $x_t$  is still white noise; nothing has changed in this regard [i.e.,  $x_t = w_t$  is the solution to (3.22)], but we have hidden the fact that  $x_t$  is white noise because of the *parameter redundancy* or *over-parameterization*.

The consideration of parameter redundancy will be crucial when we discuss estimation for general ARMA models. As this example points out, we might fit an  $ARMA(1, 1)$  model to white noise data and find that the parameter estimates are significant. If we were unaware of parameter redundancy, we might claim the data are correlated when in fact they are not (Problem 3.20). Although we have not yet discussed estimation, we present the following demonstration of the problem. We generated 150 iid normals and then fit an  $ARMA(1, 1)$  to the data. Note that  $\hat{\phi} = -.96$  and  $\hat{\theta} = .95$ , and both are significant. Below is the R code (note that the estimate called ‘intercept’ is really the estimate of the mean).

```
set.seed(8675309)      # Jenny, I got your number
x = rnorm(150, mean=5)  # generate iid N(5,1)s
arima(x, order=c(1,0,1)) # estimation
Coefficients:
      ar1      ma1 intercept<= misnomer
    -0.9595  0.9527      5.0462
s.e.    0.1688  0.1750      0.0727
```

Thus, forgetting the mean estimate, the fitted model looks like

$$(1 + .96B)x_t = (1 + .95B)w_t,$$

which we should recognize as an over-parametrized model.

Example 3.3, Example 3.6, and Example 3.7 point to a number of problems with the general definition of ARMA( $p, q$ ) models, as given by (3.19), or, equivalently, by (3.21). To summarize, we have seen the following problems:

- (i) parameter redundant models,
- (ii) stationary AR models that depend on the future, and
- (iii) MA models that are not unique.

To overcome these problems, we will require some additional restrictions on the model parameters. First, we make the following definitions.

**Definition 3.6** *The AR and MA polynomials are defined as*

$$\phi(z) = 1 - \phi_1 z - \cdots - \phi_p z^p, \quad \phi_p \neq 0, \quad (3.23)$$

and

$$\theta(z) = 1 + \theta_1 z + \cdots + \theta_q z^q, \quad \theta_q \neq 0, \quad (3.24)$$

respectively, where  $z$  is a complex number.

To address the first problem, we will henceforth refer to an ARMA( $p, q$ ) model to mean that it is in its simplest form. That is, in addition to the original definition given in equation (3.19), we will also require that  $\phi(z)$  and  $\theta(z)$  have no common factors. So, the process,  $x_t = .5x_{t-1} - .5w_{t-1} + w_t$ , discussed in Example 3.7 is not referred to as an ARMA(1, 1) process because, in its reduced form,  $x_t$  is white noise.

To address the problem of future-dependent models, we formally introduce the concept of *causality*.

**Definition 3.7** *An ARMA( $p, q$ ) model is said to be **causal**, if the time series  $\{x_t; t = 0, \pm 1, \pm 2, \dots\}$  can be written as a one-sided linear process:*

$$x_t = \sum_{j=0}^{\infty} \psi_j w_{t-j} = \psi(B)w_t, \quad (3.25)$$

where  $\psi(B) = \sum_{j=0}^{\infty} \psi_j B^j$ , and  $\sum_{j=0}^{\infty} |\psi_j| < \infty$ ; we set  $\psi_0 = 1$ .

In Example 3.3, the AR(1) process,  $x_t = \phi x_{t-1} + w_t$ , is causal only when  $|\phi| < 1$ . Equivalently, the process is causal only when the root of  $\phi(z) = 1 - \phi z$  is bigger than one in absolute value. That is, the root, say,  $z_0$ , of  $\phi(z)$  is  $z_0 = 1/\phi$  (because  $\phi(z_0) = 0$ ) and  $|z_0| > 1$  because  $|\phi| < 1$ . In general, we have the following property.

**Property 3.1 Causality of an ARMA( $p, q$ ) Process**

*An ARMA( $p, q$ ) model is causal if and only if  $\phi(z) \neq 0$  for  $|z| \leq 1$ . The coefficients of the linear process given in (3.25) can be determined by solving*

$$\psi(z) = \sum_{j=0}^{\infty} \psi_j z^j = \frac{\theta(z)}{\phi(z)}, \quad |z| \leq 1.$$

Another way to phrase **Property 3.1** is that *an ARMA process is causal only when the roots of  $\phi(z)$  lie outside the unit circle*; that is,  $\phi(z) = 0$  only when  $|z| > 1$ . Finally, to address the problem of uniqueness discussed in **Example 3.6**, we choose the model that allows an infinite autoregressive representation.

**Definition 3.8** An ARMA( $p, q$ ) model is said to be **invertible**, if the time series  $\{x_t; t = 0, \pm 1, \pm 2, \dots\}$  can be written as

$$\pi(B)x_t = \sum_{j=0}^{\infty} \pi_j x_{t-j} = w_t, \quad (3.26)$$

where  $\pi(B) = \sum_{j=0}^{\infty} \pi_j B^j$ , and  $\sum_{j=0}^{\infty} |\pi_j| < \infty$ ; we set  $\pi_0 = 1$ .

Analogous to **Property 3.1**, we have the following property.

**Property 3.2 Invertibility of an ARMA( $p, q$ ) Process**

*An ARMA( $p, q$ ) model is invertible if and only if  $\theta(z) \neq 0$  for  $|z| \leq 1$ . The coefficients  $\pi_j$  of  $\pi(B)$  given in (3.26) can be determined by solving*

$$\pi(z) = \sum_{j=0}^{\infty} \pi_j z^j = \frac{\phi(z)}{\theta(z)}, \quad |z| \leq 1.$$

Another way to phrase **Property 3.2** is that *an ARMA process is invertible only when the roots of  $\theta(z)$  lie outside the unit circle*; that is,  $\theta(z) = 0$  only when  $|z| > 1$ . The proof of **Property 3.1** is given in **Section B.2** (the proof of **Property 3.2** is similar). The following examples illustrate these concepts.

**Example 3.8 Parameter Redundancy, Causality, Invertibility**

Consider the process

$$x_t = .4x_{t-1} + .45x_{t-2} + w_t + w_{t-1} + .25w_{t-2},$$

or, in operator form,

$$(1 - .4B - .45B^2)x_t = (1 + B + .25B^2)w_t.$$

At first,  $x_t$  appears to be an ARMA(2, 2) process. But notice that

$$\phi(B) = 1 - .4B - .45B^2 = (1 + .5B)(1 - .9B)$$

and

$$\theta(B) = (1 + B + .25B^2) = (1 + .5B)^2$$

have a common factor that can be canceled. After cancellation, the operators are  $\phi(B) = (1 - .9B)$  and  $\theta(B) = (1 + .5B)$ , so the model is an ARMA(1, 1) model,  $(1 - .9B)x_t = (1 + .5B)w_t$ , or

$$x_t = .9x_{t-1} + .5w_{t-1} + w_t. \quad (3.27)$$

The model is causal because  $\phi(z) = (1 - .9z) = 0$  when  $z = 10/9$ , which is outside the unit circle. The model is also invertible because the root of  $\theta(z) = (1 + .5z)$  is  $z = -2$ , which is outside the unit circle.

To write the model as a linear process, we can obtain the  $\psi$ -weights using Property 3.1,  $\phi(z)\psi(z) = \theta(z)$ , or

$$(1 - .9z)(1 + \psi_1 z + \psi_2 z^2 + \cdots + \psi_j z^j + \cdots) = 1 + .5z.$$

Rearranging, we get

$$1 + (\psi_1 - .9)z + (\psi_2 - .9\psi_1)z^2 + \cdots + (\psi_j - .9\psi_{j-1})z^j + \cdots = 1 + .5z.$$

Matching the coefficients of  $z$  on the left and right sides we get  $\psi_1 - .9 = .5$  and  $\psi_j - .9\psi_{j-1} = 0$  for  $j > 1$ . Thus,  $\psi_j = 1.4(.9)^{j-1}$  for  $j \geq 1$  and (3.27) can be written as

$$x_t = w_t + 1.4 \sum_{j=1}^{\infty} .9^{j-1} w_{t-j}.$$

The values of  $\psi_j$  may be calculated in R as follows:

```
ARMAtoMA(ar = .9, ma = .5, 10) # first 10 psi-weights
[1] 1.40 1.26 1.13 1.02 0.92 0.83 0.74 0.67 0.60 0.54
```

The invertible representation using Property 3.1 is obtained by matching coefficients in  $\theta(z)\pi(z) = \phi(z)$ ,

$$(1 + .5z)(1 + \pi_1 z + \pi_2 z^2 + \pi_3 z^3 + \cdots) = 1 - .9z.$$

In this case, the  $\pi$ -weights are given by  $\pi_j = (-1)^j 1.4(.5)^{j-1}$ , for  $j \geq 1$ , and hence, because  $w_t = \sum_{j=0}^{\infty} \pi_j x_{t-j}$ , we can also write (3.27) as

$$x_t = 1.4 \sum_{j=1}^{\infty} (-.5)^{j-1} x_{t-j} + w_t.$$

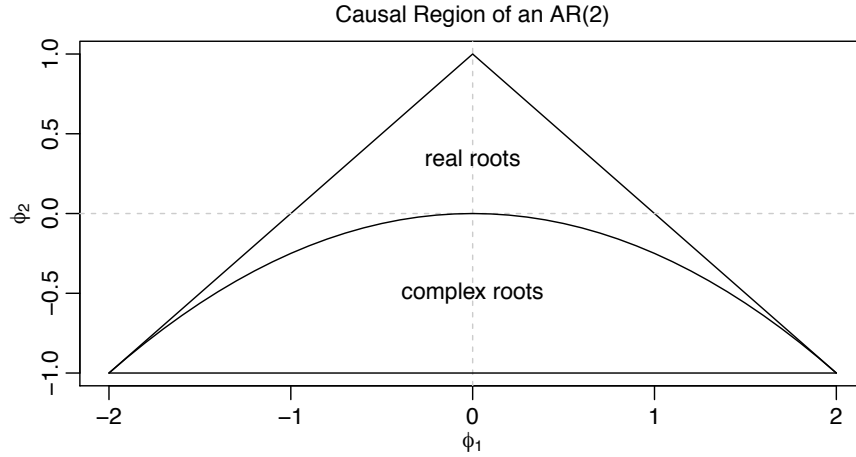
The values of  $\pi_j$  may be calculated in R as follows by reversing the roles of  $w_t$  and  $x_t$ ; i.e., write the model as  $w_t = -.5w_{t-1} + x_t - .9x_{t-1}$ :

```
ARMAtoMA(ar = -.5, ma = -.9, 10) # first 10 pi-weights
[1] -1.400 .700 -.350 .175 -.087 .044 -.022 .011 -.006 .003
```

### Example 3.9 Causal Conditions for an AR(2) Process

For an AR(1) model,  $(1 - \phi B)x_t = w_t$ , to be causal, the root of  $\phi(z) = 1 - \phi z$  must lie outside of the unit circle. In this case,  $\phi(z) = 0$  when  $z = 1/\phi$ , so it is easy to go from the causal requirement on the root,  $|1/\phi| > 1$ , to a requirement on the parameter,  $|\phi| < 1$ . It is not so easy to establish this relationship for higher order models.

For example, the AR(2) model,  $(1 - \phi_1 B - \phi_2 B^2)x_t = w_t$ , is causal when the two roots of  $\phi(z) = 1 - \phi_1 z - \phi_2 z^2$  lie outside of the unit circle. Using the quadratic formula, this requirement can be written as



**Fig. 3.3.** Causal region for an AR(2) in terms of the parameters.

$$\left| \frac{\phi_1 \pm \sqrt{\phi_1^2 + 4\phi_2}}{-2\phi_2} \right| > 1.$$

The roots of  $\phi(z)$  may be real and distinct, real and equal, or a complex conjugate pair. If we denote those roots by  $z_1$  and  $z_2$ , we can write  $\phi(z) = (1 - z_1^{-1}z)(1 - z_2^{-1}z)$ ; note that  $\phi(z_1) = \phi(z_2) = 0$ . The model can be written in operator form as  $(1 - z_1^{-1}B)(1 - z_2^{-1}B)x_t = w_t$ . From this representation, it follows that  $\phi_1 = (z_1^{-1} + z_2^{-1})$  and  $\phi_2 = -(z_1 z_2)^{-1}$ . This relationship and the fact that  $|z_1| > 1$  and  $|z_2| > 1$  can be used to establish the following equivalent condition for causality:

$$\phi_1 + \phi_2 < 1, \quad \phi_2 - \phi_1 < 1, \quad \text{and} \quad |\phi_2| < 1. \quad (3.28)$$

This causality condition specifies a triangular region in the parameter space; see [Figure 3.3](#). We leave the details of the equivalence to the reader ([Problem 3.5](#)).

## 3.2 Difference Equations

The study of the behavior of ARMA processes and their ACFs is greatly enhanced by a basic knowledge of difference equations, simply because they are difference equations. We will give a brief and heuristic account of the topic along with some examples of the usefulness of the theory. For details, the reader is referred to Mickens (1990).

Suppose we have a sequence of numbers  $u_0, u_1, u_2, \dots$  such that

$$u_n - \alpha u_{n-1} = 0, \quad \alpha \neq 0, \quad n = 1, 2, \dots \quad (3.29)$$

For example, recall [\(3.9\)](#) in which we showed that the ACF of an AR(1) process is a sequence,  $\rho(h)$ , satisfying

$$\rho(h) - \phi\rho(h-1) = 0, \quad h = 1, 2, \dots$$

Equation (3.29) represents a *homogeneous difference equation of order 1*. To solve the equation, we write:

$$\begin{aligned} u_1 &= \alpha u_0 \\ u_2 &= \alpha u_1 = \alpha^2 u_0 \\ &\vdots \\ u_n &= \alpha u_{n-1} = \alpha^n u_0. \end{aligned}$$

Given an initial condition  $u_0 = c$ , we may solve (3.29), namely,  $u_n = \alpha^n c$ .

In operator notation, (3.29) can be written as  $(1 - \alpha B)u_n = 0$ . The polynomial associated with (3.29) is  $\alpha(z) = 1 - \alpha z$ , and the root, say,  $z_0$ , of this polynomial is  $z_0 = 1/\alpha$ ; that is  $\alpha(z_0) = 0$ . We know a solution (in fact, *the* solution) to (3.29), with initial condition  $u_0 = c$ , is

$$u_n = \alpha^n c = \left(z_0^{-1}\right)^n c. \quad (3.30)$$

That is, the solution to the difference equation (3.29) depends only on the initial condition and the inverse of the root to the associated polynomial  $\alpha(z)$ .

Now suppose that the sequence satisfies

$$u_n - \alpha_1 u_{n-1} - \alpha_2 u_{n-2} = 0, \quad \alpha_2 \neq 0, \quad n = 2, 3, \dots \quad (3.31)$$

This equation is a *homogeneous difference equation of order 2*. The corresponding polynomial is

$$\alpha(z) = 1 - \alpha_1 z - \alpha_2 z^2,$$

which has two roots, say,  $z_1$  and  $z_2$ ; that is,  $\alpha(z_1) = \alpha(z_2) = 0$ . We will consider two cases. First suppose  $z_1 \neq z_2$ . Then the general solution to (3.31) is

$$u_n = c_1 z_1^{-n} + c_2 z_2^{-n}, \quad (3.32)$$

where  $c_1$  and  $c_2$  depend on the initial conditions. The claim it is a solution can be verified by direct substitution of (3.32) into (3.31):

$$\begin{aligned} &\underbrace{(c_1 z_1^{-n} + c_2 z_2^{-n})}_{u_n} - \alpha_1 \underbrace{(c_1 z_1^{-(n-1)} + c_2 z_2^{-(n-1)})}_{u_{n-1}} - \alpha_2 \underbrace{(c_1 z_1^{-(n-2)} + c_2 z_2^{-(n-2)})}_{u_{n-2}} \\ &= c_1 z_1^{-n} (1 - \alpha_1 z_1 - \alpha_2 z_1^2) + c_2 z_2^{-n} (1 - \alpha_1 z_2 - \alpha_2 z_2^2) \\ &= c_1 z_1^{-n} \alpha(z_1) + c_2 z_2^{-n} \alpha(z_2) = 0. \end{aligned}$$

Given two initial conditions  $u_0$  and  $u_1$ , we may solve for  $c_1$  and  $c_2$ :

$$u_0 = c_1 + c_2 \quad \text{and} \quad u_1 = c_1 z_1^{-1} + c_2 z_2^{-1},$$

where  $z_1$  and  $z_2$  can be solved for in terms of  $\alpha_1$  and  $\alpha_2$  using the quadratic formula, for example.

When the roots are equal,  $z_1 = z_2 (= z_0)$ , a general solution to (3.31) is

$$u_n = z_0^{-n}(c_1 + c_2 n). \quad (3.33)$$

This claim can also be verified by direct substitution of (3.33) into (3.31):

$$\begin{aligned} & \underbrace{z_0^{-n}(c_1 + c_2 n)}_{u_n} - \alpha_1 \underbrace{(z_0^{-(n-1)}[c_1 + c_2(n-1)])}_{u_{n-1}} - \alpha_2 \underbrace{(z_0^{-(n-2)}[c_1 + c_2(n-2)])}_{u_{n-2}} \\ &= z_0^{-n}(c_1 + c_2 n) \left(1 - \alpha_1 z_0 - \alpha_2 z_0^2\right) + c_2 z_0^{-n+1} (\alpha_1 + 2\alpha_2 z_0) \\ &= c_2 z_0^{-n+1} (\alpha_1 + 2\alpha_2 z_0). \end{aligned}$$

To show that  $(\alpha_1 + 2\alpha_2 z_0) = 0$ , write  $1 - \alpha_1 z - \alpha_2 z^2 = (1 - z_0^{-1} z)^2$ , and take derivatives with respect to  $z$  on both sides of the equation to obtain  $(\alpha_1 + 2\alpha_2 z) = 2z_0^{-1}(1 - z_0^{-1} z)$ . Thus,  $(\alpha_1 + 2\alpha_2 z_0) = 2z_0^{-1}(1 - z_0^{-1} z_0) = 0$ , as was to be shown. Finally, given two initial conditions,  $u_0$  and  $u_1$ , we can solve for  $c_1$  and  $c_2$ :

$$u_0 = c_1 \quad \text{and} \quad u_1 = (c_1 + c_2)z_0^{-1}.$$

It can also be shown that these solutions are unique.

To summarize these results, in the case of distinct roots, the solution to the homogeneous difference equation of degree two was

$$\begin{aligned} u_n &= z_1^{-n} \times (\text{a polynomial in } n \text{ of degree } m_1 - 1) \\ &+ z_2^{-n} \times (\text{a polynomial in } n \text{ of degree } m_2 - 1), \end{aligned} \quad (3.34)$$

where  $m_1$  is the multiplicity of the root  $z_1$  and  $m_2$  is the multiplicity of the root  $z_2$ . In this example, of course,  $m_1 = m_2 = 1$ , and we called the polynomials of degree zero  $c_1$  and  $c_2$ , respectively. In the case of the repeated root, the solution was

$$u_n = z_0^{-n} \times (\text{a polynomial in } n \text{ of degree } m_0 - 1), \quad (3.35)$$

where  $m_0$  is the multiplicity of the root  $z_0$ ; that is,  $m_0 = 2$ . In this case, we wrote the polynomial of degree one as  $c_1 + c_2 n$ . In both cases, we solved for  $c_1$  and  $c_2$  given two initial conditions,  $u_0$  and  $u_1$ .

These results generalize to the homogeneous difference equation of order  $p$ :

$$u_n - \alpha_1 u_{n-1} - \cdots - \alpha_p u_{n-p} = 0, \quad \alpha_p \neq 0, \quad n = p, p+1, \dots \quad (3.36)$$

The associated polynomial is  $\alpha(z) = 1 - \alpha_1 z - \cdots - \alpha_p z^p$ . Suppose  $\alpha(z)$  has  $r$  distinct roots,  $z_1$  with multiplicity  $m_1$ ,  $z_2$  with multiplicity  $m_2$ ,  $\dots$ , and  $z_r$  with multiplicity  $m_r$ , such that  $m_1 + m_2 + \cdots + m_r = p$ . The general solution to the difference equation (3.36) is

$$u_n = z_1^{-n} P_1(n) + z_2^{-n} P_2(n) + \cdots + z_r^{-n} P_r(n), \quad (3.37)$$

where  $P_j(n)$ , for  $j = 1, 2, \dots, r$ , is a polynomial in  $n$ , of degree  $m_j - 1$ . Given  $p$  initial conditions  $u_0, \dots, u_{p-1}$ , we can solve for the  $P_j(n)$  explicitly.



**Example 3.10 The ACF of an AR(2) Process**

Suppose  $x_t = \phi_1 x_{t-1} + \phi_2 x_{t-2} + w_t$  is a causal AR(2) process. Multiply each side of the model by  $x_{t-h}$  for  $h > 0$ , and take expectation:

$$E(x_t x_{t-h}) = \phi_1 E(x_{t-1} x_{t-h}) + \phi_2 E(x_{t-2} x_{t-h}) + E(w_t x_{t-h}).$$

The result is

$$\gamma(h) = \phi_1 \gamma(h-1) + \phi_2 \gamma(h-2), \quad h = 1, 2, \dots \quad (3.38)$$

In (3.38), we used the fact that  $E(x_t) = 0$  and for  $h > 0$ ,

$$E(w_t x_{t-h}) = E\left(w_t \sum_{j=0}^{\infty} \psi_j w_{t-h-j}\right) = 0.$$

Divide (3.38) through by  $\gamma(0)$  to obtain the difference equation for the ACF of the process:

$$\rho(h) - \phi_1 \rho(h-1) - \phi_2 \rho(h-2) = 0, \quad h = 1, 2, \dots \quad (3.39)$$

The initial conditions are  $\rho(0) = 1$  and  $\rho(-1) = \phi_1/(1 - \phi_2)$ , which is obtained by evaluating (3.39) for  $h = 1$  and noting that  $\rho(1) = \rho(-1)$ .

Using the results for the homogeneous difference equation of order two, let  $z_1$  and  $z_2$  be the roots of the associated polynomial,  $\phi(z) = 1 - \phi_1 z - \phi_2 z^2$ . Because the model is causal, we know the roots are outside the unit circle:  $|z_1| > 1$  and  $|z_2| > 1$ . Now, consider the solution for three cases:

(i) When  $z_1$  and  $z_2$  are real and distinct, then

$$\rho(h) = c_1 z_1^{-h} + c_2 z_2^{-h},$$

so  $\rho(h) \rightarrow 0$  exponentially fast as  $h \rightarrow \infty$ .

(ii) When  $z_1 = z_2 (= z_0)$  are real and equal, then

$$\rho(h) = z_0^{-h} (c_1 + c_2 h),$$

so  $\rho(h) \rightarrow 0$  exponentially fast as  $h \rightarrow \infty$ .

(iii) When  $z_1 = \bar{z}_2$  are a complex conjugate pair, then  $c_2 = \bar{c}_1$  (because  $\rho(h)$  is real), and

$$\rho(h) = c_1 z_1^{-h} + \bar{c}_1 \bar{z}_1^{-h}.$$

Write  $c_1$  and  $z_1$  in polar coordinates, for example,  $z_1 = |z_1|e^{i\theta}$ , where  $\theta$  is the angle whose tangent is the ratio of the imaginary part and the real part of  $z_1$  (sometimes called  $\arg(z_1)$ ; the range of  $\theta$  is  $[-\pi, \pi]$ ). Then, using the fact that  $e^{i\alpha} + e^{-i\alpha} = 2 \cos(\alpha)$ , the solution has the form

$$\rho(h) = a|z_1|^{-h} \cos(h\theta + b),$$

where  $a$  and  $b$  are determined by the initial conditions. Again,  $\rho(h)$  dampens to zero exponentially fast as  $h \rightarrow \infty$ , but it does so in a sinusoidal fashion. The implication of this result is shown in the next example.

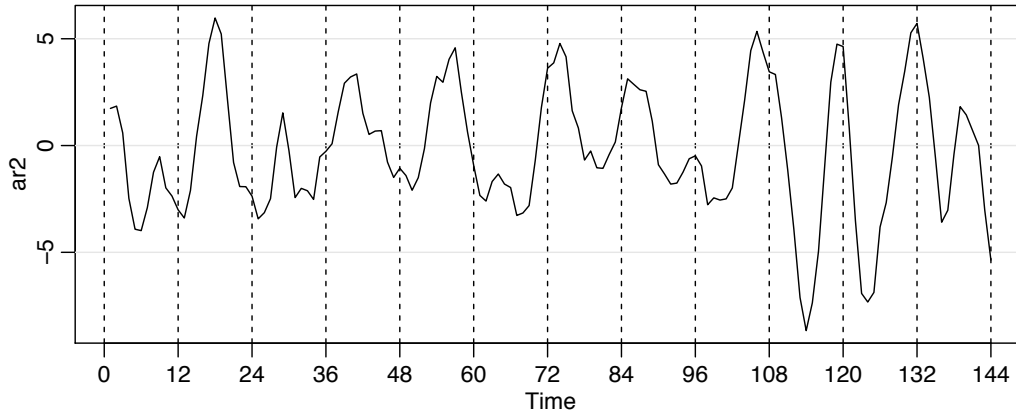


Fig. 3.4. Simulated AR(2) model,  $n = 144$  with  $\phi_1 = 1.5$  and  $\phi_2 = -.75$ .

### Example 3.11 An AR(2) with Complex Roots

Figure 3.4 shows  $n = 144$  observations from the AR(2) model

$$x_t = 1.5x_{t-1} - .75x_{t-2} + w_t,$$

with  $\sigma_w^2 = 1$ , and with complex roots chosen so the process exhibits pseudo-cyclic behavior at the rate of one cycle every 12 time points. The autoregressive polynomial for this model is  $\phi(z) = 1 - 1.5z + .75z^2$ . The roots of  $\phi(z)$  are  $1 \pm i/\sqrt{3}$ , and  $\theta = \tan^{-1}(1/\sqrt{3}) = 2\pi/12$  radians per unit time. To convert the angle to cycles per unit time, divide by  $2\pi$  to get  $1/12$  cycles per unit time. The ACF for this model is shown in left-hand-side of Figure 3.5.

To calculate the roots of the polynomial and solve for arg in R:

```
z = c(1,-1.5,.75)      # coefficients of the polynomial
(a = polyroot(z)[1])   # print one root = 1 + i/sqrt(3)
[1] 1+0.57735i
arg = Arg(a)/(2*pi)     # arg in cycles/pt
1/arg                   # the pseudo period
[1] 12
```

To reproduce Figure 3.4:

```
set.seed(8675309)
ar2 = arima.sim(list(order=c(2,0,0), ar=c(1.5,-.75)), n = 144)
plot(ar2, axes=FALSE, xlab="Time")
axis(2); axis(1, at=seq(0,144,by=12)); box()
abline(v=seq(0,144,by=12), lty=2)
```

To calculate and display the ACF for this model:

```
ACF = ARMAacf(ar=c(1.5,-.75), ma=0, 50)
plot(ACF, type="h", xlab="lag")
abline(h=0)
```

**Example 3.12 The  $\psi$ -weights for an ARMA Model**

For a causal ARMA( $p, q$ ) model,  $\phi(B)x_t = \theta(B)w_t$ , where the zeros of  $\phi(z)$  are outside the unit circle, recall that we may write

$$x_t = \sum_{j=0}^{\infty} \psi_j w_{t-j},$$

where the  $\psi$ -weights are determined using [Property 3.1](#).

For the pure MA( $q$ ) model,  $\psi_0 = 1$ ,  $\psi_j = \theta_j$ , for  $j = 1, \dots, q$ , and  $\psi_j = 0$ , otherwise. For the general case of ARMA( $p, q$ ) models, the task of solving for the  $\psi$ -weights is much more complicated, as was demonstrated in [Example 3.8](#). The use of the theory of homogeneous difference equations can help here. To solve for the  $\psi$ -weights in general, we must match the coefficients in  $\phi(z)\psi(z) = \theta(z)$ :

$$(1 - \phi_1 z - \phi_2 z^2 - \dots)(\psi_0 + \psi_1 z + \psi_2 z^2 + \dots) = (1 + \theta_1 z + \theta_2 z^2 + \dots).$$

The first few values are

$$\begin{aligned} \psi_0 &= 1 \\ \psi_1 - \phi_1 \psi_0 &= \theta_1 \\ \psi_2 - \phi_1 \psi_1 - \phi_2 \psi_0 &= \theta_2 \\ \psi_3 - \phi_1 \psi_2 - \phi_2 \psi_1 - \phi_3 \psi_0 &= \theta_3 \\ &\vdots \end{aligned}$$

where we would take  $\phi_j = 0$  for  $j > p$ , and  $\theta_j = 0$  for  $j > q$ . The  $\psi$ -weights satisfy the homogeneous difference equation given by

$$\psi_j - \sum_{k=1}^p \phi_k \psi_{j-k} = 0, \quad j \geq \max(p, q+1), \quad (3.40)$$

with initial conditions

$$\psi_j - \sum_{k=1}^j \phi_k \psi_{j-k} = \theta_j, \quad 0 \leq j < \max(p, q+1). \quad (3.41)$$

The general solution depends on the roots of the AR polynomial  $\phi(z) = 1 - \phi_1 z - \dots - \phi_p z^p$ , as seen from (3.40). The specific solution will, of course, depend on the initial conditions.

Consider the ARMA process given in (3.27),  $x_t = .9x_{t-1} + .5w_{t-1} + w_t$ . Because  $\max(p, q+1) = 2$ , using (3.41), we have  $\psi_0 = 1$  and  $\psi_1 = .9 + .5 = 1.4$ . By (3.40), for  $j = 2, 3, \dots$ , the  $\psi$ -weights satisfy  $\psi_j - .9\psi_{j-1} = 0$ . The general solution is  $\psi_j = c \cdot .9^j$ . To find the specific solution, use the initial condition  $\psi_1 = 1.4$ , so  $1.4 = .9c$  or  $c = 1.4/.9$ . Finally,  $\psi_j = 1.4(.9)^{j-1}$ , for  $j \geq 1$ , as we saw in [Example 3.8](#).

To view, for example, the first 50  $\psi$ -weights in R, use:

```
ARMAtoMA(ar=.9, ma=.5, 50)      # for a list
plot(ARMAtoMA(ar=.9, ma=.5, 50)) # for a graph
```

### 3.3 Autocorrelation and Partial Autocorrelation

We begin by exhibiting the ACF of an MA( $q$ ) process,  $x_t = \theta(B)w_t$ , where  $\theta(B) = 1 + \theta_1 B + \dots + \theta_q B^q$ . Because  $x_t$  is a finite linear combination of white noise terms, the process is stationary with mean

$$E(x_t) = \sum_{j=0}^q \theta_j E(w_{t-j}) = 0,$$

where we have written  $\theta_0 = 1$ , and with autocovariance function

$$\begin{aligned} \gamma(h) &= \text{cov}(x_{t+h}, x_t) = \text{cov}\left(\sum_{j=0}^q \theta_j w_{t+h-j}, \sum_{k=0}^q \theta_k w_{t-k}\right) \\ &= \begin{cases} \sigma_w^2 \sum_{j=0}^{q-h} \theta_j \theta_{j+h}, & 0 \leq h \leq q \\ 0 & h > q. \end{cases} \end{aligned} \quad (3.42)$$

Recall that  $\gamma(h) = \gamma(-h)$ , so we will only display the values for  $h \geq 0$ . Note that  $\gamma(q)$  cannot be zero because  $\theta_q \neq 0$ . The cutting off of  $\gamma(h)$  after  $q$  lags is the signature of the MA( $q$ ) model. Dividing (3.42) by  $\gamma(0)$  yields the **ACF of an MA( $q$ )**:

$$\rho(h) = \begin{cases} \frac{\sum_{j=0}^{q-h} \theta_j \theta_{j+h}}{1 + \theta_1^2 + \dots + \theta_q^2} & 1 \leq h \leq q \\ 0 & h > q. \end{cases} \quad (3.43)$$

For a causal ARMA( $p, q$ ) model,  $\phi(B)x_t = \theta(B)w_t$ , where the zeros of  $\phi(z)$  are outside the unit circle, write

$$x_t = \sum_{j=0}^{\infty} \psi_j w_{t-j}. \quad (3.44)$$

It follows immediately that  $E(x_t) = 0$  and the autocovariance function of  $x_t$  is

$$\gamma(h) = \text{cov}(x_{t+h}, x_t) = \sigma_w^2 \sum_{j=0}^{\infty} \psi_j \psi_{j+h}, \quad h \geq 0. \quad (3.45)$$

We could then use (3.40) and (3.41) to solve for the  $\psi$ -weights. In turn, we could solve for  $\gamma(h)$ , and the ACF  $\rho(h) = \gamma(h)/\gamma(0)$ . As in **Example 3.10**, it is also possible to obtain a homogeneous difference equation directly in terms of  $\gamma(h)$ . First, we write

$$\begin{aligned} \gamma(h) &= \text{cov}(x_{t+h}, x_t) = \text{cov}\left(\sum_{j=1}^p \phi_j x_{t+h-j} + \sum_{j=0}^q \theta_j w_{t+h-j}, x_t\right) \\ &= \sum_{j=1}^p \phi_j \gamma(h-j) + \sigma_w^2 \sum_{j=h}^q \theta_j \psi_{j-h}, \quad h \geq 0, \end{aligned} \quad (3.46)$$

where we have used the fact that, for  $h \geq 0$ ,

$$\text{cov}(w_{t+h-j}, x_t) = \text{cov}\left(w_{t+h-j}, \sum_{k=0}^{\infty} \psi_k w_{t-k}\right) = \psi_{j-h} \sigma_w^2.$$

From (3.46), we can write a *general homogeneous equation for the ACF of a causal ARMA process*:

$$\gamma(h) - \phi_1 \gamma(h-1) - \cdots - \phi_p \gamma(h-p) = 0, \quad h \geq \max(p, q+1), \quad (3.47)$$

with initial conditions

$$\gamma(h) - \sum_{j=1}^p \phi_j \gamma(h-j) = \sigma_w^2 \sum_{j=h}^q \theta_j \psi_{j-h}, \quad 0 \leq h < \max(p, q+1). \quad (3.48)$$

Dividing (3.47) and (3.48) through by  $\gamma(0)$  will allow us to solve for the ACF,  $\rho(h) = \gamma(h)/\gamma(0)$ .

### Example 3.13 The ACF of an AR( $p$ )

In Example 3.10 we considered the case where  $p = 2$ . For the general case, it follows immediately from (3.47) that

$$\rho(h) - \phi_1 \rho(h-1) - \cdots - \phi_p \rho(h-p) = 0, \quad h \geq p. \quad (3.49)$$

Let  $z_1, \dots, z_r$  denote the roots of  $\phi(z)$ , each with multiplicity  $m_1, \dots, m_r$ , respectively, where  $m_1 + \cdots + m_r = p$ . Then, from (3.37), the general solution is

$$\rho(h) = z_1^{-h} P_1(h) + z_2^{-h} P_2(h) + \cdots + z_r^{-h} P_r(h), \quad h \geq p, \quad (3.50)$$

where  $P_j(h)$  is a polynomial in  $h$  of degree  $m_j - 1$ .

Recall that for a causal model, all of the roots are outside the unit circle,  $|z_i| > 1$ , for  $i = 1, \dots, r$ . If all the roots are real, then  $\rho(h)$  dampens exponentially fast to zero as  $h \rightarrow \infty$ . If some of the roots are complex, then they will be in conjugate pairs and  $\rho(h)$  will dampen, in a sinusoidal fashion, exponentially fast to zero as  $h \rightarrow \infty$ . In the case of complex roots, the time series will appear to be cyclic in nature. This, of course, is also true for ARMA models in which the AR part has complex roots.

### Example 3.14 The ACF of an ARMA(1, 1)

Consider the ARMA(1, 1) process  $x_t = \phi x_{t-1} + \theta w_{t-1} + w_t$ , where  $|\phi| < 1$ . Based on (3.47), the autocovariance function satisfies

$$\gamma(h) - \phi \gamma(h-1) = 0, \quad h = 2, 3, \dots,$$

and it follows from (3.29)–(3.30) that the general solution is

$$\gamma(h) = c \phi^h, \quad h = 1, 2, \dots \quad (3.51)$$

To obtain the initial conditions, we use (3.48):

$$\gamma(0) = \phi\gamma(1) + \sigma_w^2[1 + \theta\phi + \theta^2] \quad \text{and} \quad \gamma(1) = \phi\gamma(0) + \sigma_w^2\theta.$$

Solving for  $\gamma(0)$  and  $\gamma(1)$ , we obtain:

$$\gamma(0) = \sigma_w^2 \frac{1 + 2\theta\phi + \theta^2}{1 - \phi^2} \quad \text{and} \quad \gamma(1) = \sigma_w^2 \frac{(1 + \theta\phi)(\phi + \theta)}{1 - \phi^2}.$$

To solve for  $c$ , note that from (3.51),  $\gamma(1) = c\phi$  or  $c = \gamma(1)/\phi$ . Hence, the specific solution for  $h \geq 1$  is

$$\gamma(h) = \frac{\gamma(1)}{\phi} \phi^h = \sigma_w^2 \frac{(1 + \theta\phi)(\phi + \theta)}{1 - \phi^2} \phi^{h-1}.$$

Finally, dividing through by  $\gamma(0)$  yields the ACF

$$\rho(h) = \frac{(1 + \theta\phi)(\phi + \theta)}{1 + 2\theta\phi + \theta^2} \phi^{h-1}, \quad h \geq 1. \quad (3.52)$$

Notice that the general pattern of  $\rho(h)$  versus  $h$  in (3.52) is not different from that of an AR(1) given in (3.8). Hence, it is unlikely that we will be able to tell the difference between an ARMA(1,1) and an AR(1) based solely on an ACF estimated from a sample. This consideration will lead us to the partial autocorrelation function.

#### THE PARTIAL AUTOCORRELATION FUNCTION (PACF)

We have seen in (3.43), for MA( $q$ ) models, the ACF will be zero for lags greater than  $q$ . Moreover, because  $\theta_q \neq 0$ , the ACF will not be zero at lag  $q$ . Thus, the ACF provides a considerable amount of information about the order of the dependence when the process is a moving average process. If the process, however, is ARMA or AR, the ACF alone tells us little about the orders of dependence. Hence, it is worthwhile pursuing a function that will behave like the ACF of MA models, but for AR models, namely, the *partial autocorrelation function (PACF)*.

Recall that if  $X$ ,  $Y$ , and  $Z$  are random variables, then the partial correlation between  $X$  and  $Y$  given  $Z$  is obtained by regressing  $X$  on  $Z$  to obtain  $\hat{X}$ , regressing  $Y$  on  $Z$  to obtain  $\hat{Y}$ , and then calculating

$$\rho_{XY|Z} = \text{corr}\{X - \hat{X}, Y - \hat{Y}\}.$$

The idea is that  $\rho_{XY|Z}$  measures the correlation between  $X$  and  $Y$  with the linear effect of  $Z$  removed (or partialled out). If the variables are multivariate normal, then this definition coincides with  $\rho_{XY|Z} = \text{corr}(X, Y | Z)$ .

To motivate the idea for time series, consider a causal AR(1) model,  $x_t = \phi x_{t-1} + w_t$ . Then,

$$\begin{aligned} \gamma_x(2) &= \text{cov}(x_t, x_{t-2}) = \text{cov}(\phi x_{t-1} + w_t, x_{t-2}) \\ &= \text{cov}(\phi^2 x_{t-2} + \phi w_{t-1} + w_t, x_{t-2}) = \phi^2 \gamma_x(0). \end{aligned}$$

This result follows from causality because  $x_{t-2}$  involves  $\{w_{t-2}, w_{t-3}, \dots\}$ , which are all uncorrelated with  $w_t$  and  $w_{t-1}$ . The correlation between  $x_t$  and  $x_{t-2}$  is not zero, as it would be for an MA(1), because  $x_t$  is dependent on  $x_{t-2}$  through  $x_{t-1}$ . Suppose we break this chain of dependence by removing (or partial out) the effect  $x_{t-1}$ . That is, we consider the correlation between  $x_t - \phi x_{t-1}$  and  $x_{t-2} - \phi x_{t-1}$ , because it is the correlation between  $x_t$  and  $x_{t-2}$  with the linear dependence of each on  $x_{t-1}$  removed. In this way, we have broken the dependence chain between  $x_t$  and  $x_{t-2}$ . In fact,

$$\text{cov}(x_t - \phi x_{t-1}, x_{t-2} - \phi x_{t-1}) = \text{cov}(w_t, x_{t-2} - \phi x_{t-1}) = 0.$$

Hence, the tool we need is partial autocorrelation, which is the correlation between  $x_s$  and  $x_t$  with the linear effect of everything “in the middle” removed.

To formally define the PACF for mean-zero stationary time series, let  $\hat{x}_{t+h}$ , for  $h \geq 2$ , denote the regression<sup>3.3</sup> of  $x_{t+h}$  on  $\{x_{t+h-1}, x_{t+h-2}, \dots, x_{t+1}\}$ , which we write as

$$\hat{x}_{t+h} = \beta_1 x_{t+h-1} + \beta_2 x_{t+h-2} + \dots + \beta_{h-1} x_{t+1}. \quad (3.53)$$

No intercept term is needed in (3.53) because the mean of  $x_t$  is zero (otherwise, replace  $x_t$  by  $x_t - \mu_x$  in this discussion). In addition, let  $\hat{x}_t$  denote the regression of  $x_t$  on  $\{x_{t+1}, x_{t+2}, \dots, x_{t+h-1}\}$ , then

$$\hat{x}_t = \beta_1 x_{t+1} + \beta_2 x_{t+2} + \dots + \beta_{h-1} x_{t+h-1}. \quad (3.54)$$

Because of stationarity, the coefficients,  $\beta_1, \dots, \beta_{h-1}$  are the same in (3.53) and (3.54); we will explain this result in the next section, but it will be evident from the examples.

**Definition 3.9** The **partial autocorrelation function (PACF)** of a stationary process,  $x_t$ , denoted  $\phi_{hh}$ , for  $h = 1, 2, \dots$ , is

$$\phi_{11} = \text{corr}(x_{t+1}, x_t) = \rho(1) \quad (3.55)$$

and

$$\phi_{hh} = \text{corr}(x_{t+h} - \hat{x}_{t+h}, x_t - \hat{x}_t), \quad h \geq 2. \quad (3.56)$$

The reason for using a double subscript will become evident in the next section. The PACF,  $\phi_{hh}$ , is the correlation between  $x_{t+h}$  and  $x_t$  with the linear dependence of  $\{x_{t+1}, \dots, x_{t+h-1}\}$  on each, removed. If the process  $x_t$  is Gaussian, then  $\phi_{hh} = \text{corr}(x_{t+h}, x_t \mid x_{t+1}, \dots, x_{t+h-1})$ ; that is,  $\phi_{hh}$  is the correlation coefficient between  $x_{t+h}$  and  $x_t$  in the bivariate distribution of  $(x_{t+h}, x_t)$  conditional on  $\{x_{t+1}, \dots, x_{t+h-1}\}$ .

<sup>3.3</sup> The term regression here refers to regression in the population sense. That is,  $\hat{x}_{t+h}$  is the linear combination of  $\{x_{t+h-1}, x_{t+h-2}, \dots, x_{t+1}\}$  that minimizes the mean squared error  $E(x_{t+h} - \sum_{j=1}^{h-1} \alpha_j x_{t+j})^2$ .

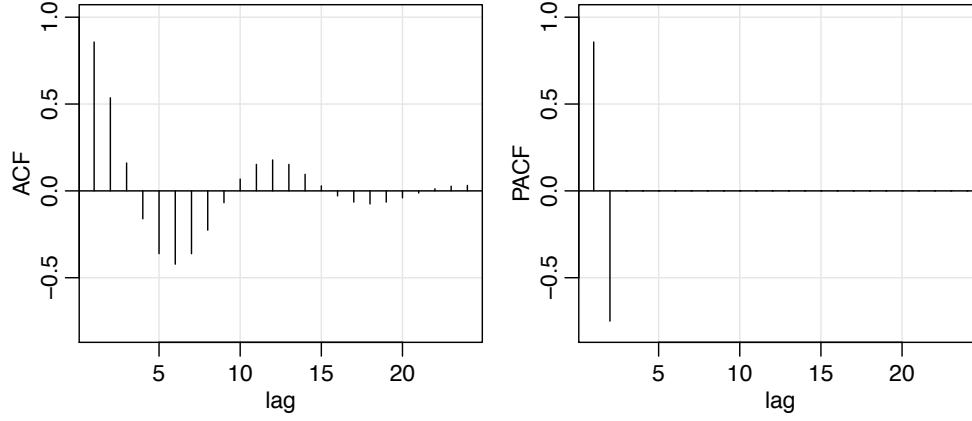


Fig. 3.5. The ACF and PACF of an AR(2) model with  $\phi_1 = 1.5$  and  $\phi_2 = -.75$ .

### Example 3.15 The PACF of an AR(1)

Consider the PACF of the AR(1) process given by  $x_t = \phi x_{t-1} + w_t$ , with  $|\phi| < 1$ . By definition,  $\phi_{11} = \rho(1) = \phi$ . To calculate  $\phi_{22}$ , consider the regression of  $x_{t+2}$  on  $x_{t+1}$ , say,  $\hat{x}_{t+2} = \beta x_{t+1}$ . We choose  $\beta$  to minimize

$$E(x_{t+2} - \hat{x}_{t+2})^2 = E(x_{t+2} - \beta x_{t+1})^2 = \gamma(0) - 2\beta\gamma(1) + \beta^2\gamma(0).$$

Taking derivatives with respect to  $\beta$  and setting the result equal to zero, we have  $\beta = \gamma(1)/\gamma(0) = \rho(1) = \phi$ . Next, consider the regression of  $x_t$  on  $x_{t+1}$ , say  $\hat{x}_t = \beta x_{t+1}$ . We choose  $\beta$  to minimize

$$E(x_t - \hat{x}_t)^2 = E(x_t - \beta x_{t+1})^2 = \gamma(0) - 2\beta\gamma(1) + \beta^2\gamma(0).$$

This is the same equation as before, so  $\beta = \phi$ . Hence,

$$\begin{aligned} \phi_{22} &= \text{corr}(x_{t+2} - \hat{x}_{t+2}, x_t - \hat{x}_t) = \text{corr}(x_{t+2} - \phi x_{t+1}, x_t - \phi x_{t+1}) \\ &= \text{corr}(w_{t+2}, x_t - \phi x_{t+1}) = 0 \end{aligned}$$

by causality. Thus,  $\phi_{22} = 0$ . In the next example, we will see that in this case,  $\phi_{hh} = 0$  for all  $h > 1$ .

### Example 3.16 The PACF of an AR(p)

The model implies  $x_{t+h} = \sum_{j=1}^p \phi_j x_{t+h-j} + w_{t+h}$ , where the roots of  $\phi(z)$  are outside the unit circle. When  $h > p$ , the regression of  $x_{t+h}$  on  $\{x_{t+1}, \dots, x_{t+h-1}\}$ , is

$$\hat{x}_{t+h} = \sum_{j=1}^p \phi_j x_{t+h-j}.$$

We have not proved this obvious result yet, but we will prove it in the next section. Thus, when  $h > p$ ,



**Table 3.1.** Behavior of the ACF and PACF for ARMA Models

	AR( $p$ )	MA( $q$ )	ARMA( $p, q$ )
ACF	Tails off	Cuts off after lag $q$	Tails off
PACF	Cuts off after lag $p$	Tails off	Tails off

$$\phi_{hh} = \text{corr}(x_{t+h} - \hat{x}_{t+h}, x_t - \hat{x}_t) = \text{corr}(w_{t+h}, x_t - \hat{x}_t) = 0,$$

because, by causality,  $x_t - \hat{x}_t$  depends only on  $\{w_{t+h-1}, w_{t+h-2}, \dots\}$ ; recall equation (3.54). When  $h \leq p$ ,  $\phi_{pp}$  is not zero, and  $\phi_{11}, \dots, \phi_{p-1,p-1}$  are not necessarily zero. We will see later that, in fact,  $\phi_{pp} = \phi_p$ . Figure 3.5 shows the ACF and the PACF of the AR(2) model presented in Example 3.11. To reproduce Figure 3.5 in R, use the following commands:

```
ACF = ARMAacf(ar=c(1.5,-.75), ma=0, 24)[-1]
PACF = ARMAacf(ar=c(1.5,-.75), ma=0, 24, pacf=TRUE)
par(mfrow=c(1,2))
plot(ACF, type="h", xlab="lag", ylim=c(-.8,1)); abline(h=0)
plot(PACF, type="h", xlab="lag", ylim=c(-.8,1)); abline(h=0)
```

### Example 3.17 The PACF of an Invertible MA( $q$ )

For an invertible MA( $q$ ), we can write  $x_t = -\sum_{j=1}^{\infty} \pi_j x_{t-j} + w_t$ . Moreover, no finite representation exists. From this result, it should be apparent that the PACF will never cut off, as in the case of an AR( $p$ ).

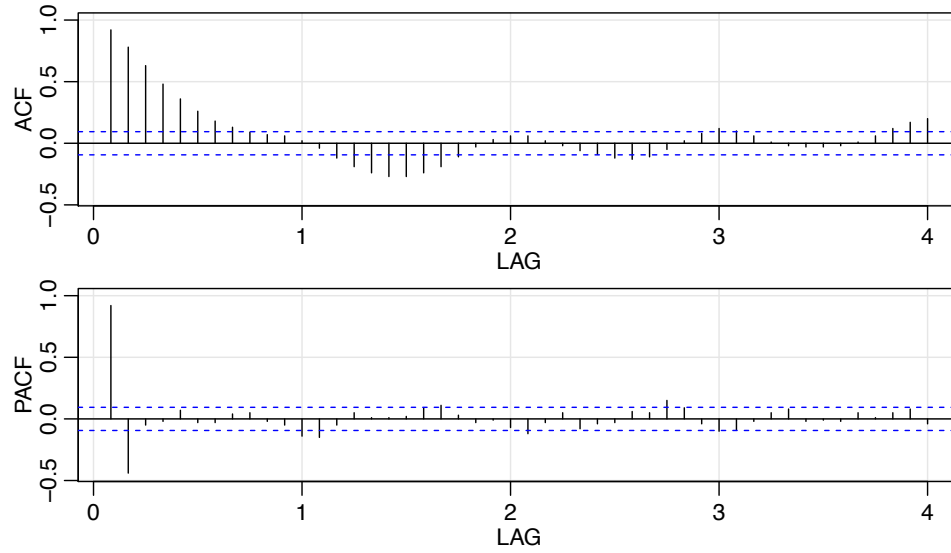
For an MA(1),  $x_t = w_t + \theta w_{t-1}$ , with  $|\theta| < 1$ , calculations similar to Example 3.15 will yield  $\phi_{22} = -\theta^2/(1 + \theta^2 + \theta^4)$ . For the MA(1) in general, we can show that

$$\phi_{hh} = -\frac{(-\theta)^h(1 - \theta^2)}{1 - \theta^{2(h+1)}}, \quad h \geq 1.$$

In the next section, we will discuss methods of calculating the PACF. The PACF for MA models behaves much like the ACF for AR models. Also, the PACF for AR models behaves much like the ACF for MA models. Because an invertible ARMA model has an infinite AR representation, the PACF will not cut off. We may summarize these results in Table 3.1.

### Example 3.18 Preliminary Analysis of the Recruitment Series

We consider the problem of modeling the Recruitment series shown in Figure 1.5. There are 453 months of observed recruitment ranging over the years 1950-1987. The ACF and the PACF given in Figure 3.6 are consistent with the behavior of an AR(2). The ACF has cycles corresponding roughly to a 12-month period, and the PACF has large values for  $h = 1, 2$  and then is essentially zero for higher order lags. Based on Table 3.1, these results suggest that a second-order ( $p = 2$ ) autoregressive model might provide a good fit. Although we will discuss estimation



**Fig. 3.6.** ACF and PACF of the Recruitment series. Note that the lag axes are in terms of season (12 months in this case).

in detail in [Section 3.5](#), we ran a regression (see [Section 2.1](#)) using the data triplets  $\{(x; z_1, z_2) : (x_3; x_2, x_1), (x_4; x_3, x_2), \dots, (x_{453}; x_{452}, x_{451})\}$  to fit a model of the form

$$x_t = \phi_0 + \phi_1 x_{t-1} + \phi_2 x_{t-2} + w_t$$

for  $t = 3, 4, \dots, 453$ . The estimates and standard errors (in parentheses) are  $\hat{\phi}_0 = 6.74_{(1.11)}$ ,  $\hat{\phi}_1 = 1.35_{(.04)}$ ,  $\hat{\phi}_2 = -.46_{(.04)}$ , and  $\hat{\sigma}_w^2 = 89.72$ .

The following R code can be used for this analysis. We use `acf2` from `astsa` to print and plot the ACF and PACF.

```
acf2(rec, 48)      # will produce values and a graphic
(regr = ar.ols(rec, order=2, demean=FALSE, intercept=TRUE))
regr$asy.se.coef  # standard errors of the estimates
```

### 3.4 Forecasting

In forecasting, the goal is to predict future values of a time series,  $x_{n+m}$ ,  $m = 1, 2, \dots$ , based on the data collected to the present,  $x_{1:n} = \{x_1, x_2, \dots, x_n\}$ . Throughout this section, we will assume  $x_t$  is stationary and the model parameters are known. The problem of forecasting when the model parameters are unknown will be discussed in the next section; also, see [Problem 3.26](#). The minimum mean square error predictor of  $x_{n+m}$  is

$$x_{n+m}^n = E(x_{n+m} \mid x_{1:n}) \quad (3.57)$$

because the conditional expectation minimizes the mean square error

$$E[x_{n+m} - g(x_{1:n})]^2, \quad (3.58)$$

where  $g(x_{1:n})$  is a function of the observations  $x_{1:n}$ ; see [Problem 3.14](#).

First, we will restrict attention to predictors that are linear functions of the data, that is, predictors of the form

$$x_{n+m}^n = \alpha_0 + \sum_{k=1}^n \alpha_k x_k, \quad (3.59)$$

where  $\alpha_0, \alpha_1, \dots, \alpha_n$  are real numbers. We note that the  $\alpha$ s depend on  $n$  and  $m$ , but for now we drop the dependence from the notation. For example, if  $n = m = 1$ , then  $x_2^1$  is the one-step-ahead linear forecast of  $x_2$  given  $x_1$ . In terms of (3.59),  $x_2^1 = \alpha_0 + \alpha_1 x_1$ . But if  $n = 2$ ,  $x_3^2$  is the one-step-ahead linear forecast of  $x_3$  given  $x_1$  and  $x_2$ . In terms of (3.59),  $x_3^2 = \alpha_0 + \alpha_1 x_1 + \alpha_2 x_2$ , and in general, the  $\alpha$ s in  $x_2^1$  and  $x_3^2$  will be different.

Linear predictors of the form (3.59) that minimize the mean square prediction error (3.58) are called *best linear predictors* (BLPs). As we shall see, linear prediction depends only on the second-order moments of the process, which are easy to estimate from the data. Much of the material in this section is enhanced by the theoretical material presented in **Appendix B**. For example, **Theorem B.3** states that if the process is Gaussian, minimum mean square error predictors and best linear predictors are the same. The following property, which is based on the Projection Theorem, **Theorem B.1**, is a key result.

### Property 3.3 Best Linear Prediction for Stationary Processes

Given data  $x_1, \dots, x_n$ , the best linear predictor,  $x_{n+m}^n = \alpha_0 + \sum_{k=1}^n \alpha_k x_k$ , of  $x_{n+m}$ , for  $m \geq 1$ , is found by solving

$$E[(x_{n+m} - x_{n+m}^n) x_k] = 0, \quad k = 0, 1, \dots, n, \quad (3.60)$$

where  $x_0 = 1$ , for  $\alpha_0, \alpha_1, \dots, \alpha_n$ .

The equations specified in (3.60) are called the *prediction equations*, and they are used to solve for the coefficients  $\{\alpha_0, \alpha_1, \dots, \alpha_n\}$ . The results of **Property 3.3** can also be obtained via least squares; i.e., to minimize  $Q = E(x_{n+m} - \sum_{k=0}^n \alpha_k x_k)^2$  with respect to the  $\alpha$ s, solve  $\partial Q / \partial \alpha_j = 0$  for the  $\alpha_j$ ,  $j = 0, 1, \dots, n$ . This leads to (3.60).

If  $E(x_t) = \mu$ , the first equation ( $k = 0$ ) of (3.60) implies

$$E(x_{n+m}^n) = E(x_{n+m}) = \mu.$$

Thus, taking expectation in (3.59), we have

$$\mu = \alpha_0 + \sum_{k=1}^n \alpha_k \mu \quad \text{or} \quad \alpha_0 = \mu \left( 1 - \sum_{k=1}^n \alpha_k \right).$$

Hence, the form of the BLP is

$$x_{n+m}^n = \mu + \sum_{k=1}^n \alpha_k (x_k - \mu).$$

Thus, until we discuss estimation, there is no loss of generality in considering the case that  $\mu = 0$ , in which case,  $\alpha_0 = 0$ .

First, consider *one-step-ahead prediction*. That is, given  $\{x_1, \dots, x_n\}$ , we wish to forecast the value of the time series at the next time point,  $x_{n+1}$ . The BLP of  $x_{n+1}$  is of the form

$$x_{n+1}^n = \phi_{n1}x_n + \phi_{n2}x_{n-1} + \dots + \phi_{nn}x_1, \quad (3.61)$$

where we now display the dependence of the coefficients on  $n$ ; in this case,  $\alpha_k$  in (3.59) is  $\phi_{n,n+1-k}$  in (3.61), for  $k = 1, \dots, n$ . Using **Property 3.3**, the coefficients  $\{\phi_{n1}, \phi_{n2}, \dots, \phi_{nn}\}$  satisfy

$$E\left[\left(x_{n+1} - \sum_{j=1}^n \phi_{nj}x_{n+1-j}\right)x_{n+1-k}\right] = 0, \quad k = 1, \dots, n,$$

or

$$\sum_{j=1}^n \phi_{nj}\gamma(k-j) = \gamma(k), \quad k = 1, \dots, n. \quad (3.62)$$

The prediction equations (3.62) can be written in matrix notation as

$$\Gamma_n \phi_n = \gamma_n, \quad (3.63)$$

where  $\Gamma_n = \{\gamma(k-j)\}_{j,k=1}^n$  is an  $n \times n$  matrix,  $\phi_n = (\phi_{n1}, \dots, \phi_{nn})'$  is an  $n \times 1$  vector, and  $\gamma_n = (\gamma(1), \dots, \gamma(n))'$  is an  $n \times 1$  vector.

The matrix  $\Gamma_n$  is nonnegative definite. If  $\Gamma_n$  is singular, there are many solutions to (3.63), but, by the Projection Theorem (**Theorem B.1**),  $x_{n+1}^n$  is unique. If  $\Gamma_n$  is nonsingular, the elements of  $\phi_n$  are unique, and are given by

$$\phi_n = \Gamma_n^{-1} \gamma_n. \quad (3.64)$$

For ARMA models, the fact that  $\sigma_w^2 > 0$  and  $\gamma(h) \rightarrow 0$  as  $h \rightarrow \infty$  is enough to ensure that  $\Gamma_n$  is positive definite (**Problem 3.12**). It is sometimes convenient to write the one-step-ahead forecast in vector notation

$$x_{n+1}^n = \phi_n' x, \quad (3.65)$$

where  $x = (x_n, x_{n-1}, \dots, x_1)'$ .

The *mean square one-step-ahead prediction error* is

$$P_{n+1}^n = E(x_{n+1} - x_{n+1}^n)^2 = \gamma(0) - \gamma_n' \Gamma_n^{-1} \gamma_n. \quad (3.66)$$

To verify (3.66) using (3.64) and (3.65),

$$\begin{aligned} E(x_{n+1} - x_{n+1}^n)^2 &= E(x_{n+1} - \phi_n' x)^2 = E(x_{n+1} - \gamma_n' \Gamma_n^{-1} x)^2 \\ &= E(x_{n+1}^2 - 2\gamma_n' \Gamma_n^{-1} x x_{n+1} + \gamma_n' \Gamma_n^{-1} x x' \Gamma_n^{-1} \gamma_n) \\ &= \gamma(0) - 2\gamma_n' \Gamma_n^{-1} \gamma_n + \gamma_n' \Gamma_n^{-1} \Gamma_n \Gamma_n^{-1} \gamma_n \\ &= \gamma(0) - \gamma_n' \Gamma_n^{-1} \gamma_n. \end{aligned}$$

**Example 3.19 Prediction for an AR(2)**

Suppose we have a causal AR(2) process  $x_t = \phi_1 x_{t-1} + \phi_2 x_{t-2} + w_t$ , and one observation  $x_1$ . Then, using equation (3.64), the one-step-ahead prediction of  $x_2$  based on  $x_1$  is

$$x_2^1 = \phi_{11} x_1 = \frac{\gamma(1)}{\gamma(0)} x_1 = \rho(1) x_1.$$

Now, suppose we want the one-step-ahead prediction of  $x_3$  based on two observations  $x_1$  and  $x_2$ ; i.e.,  $x_3^2 = \phi_{21} x_2 + \phi_{22} x_1$ . We could use (3.62)

$$\begin{aligned}\phi_{21} \gamma(0) + \phi_{22} \gamma(1) &= \gamma(1) \\ \phi_{21} \gamma(1) + \phi_{22} \gamma(0) &= \gamma(2)\end{aligned}$$

to solve for  $\phi_{21}$  and  $\phi_{22}$ , or use the matrix form in (3.64) and solve

$$\begin{pmatrix} \phi_{21} \\ \phi_{22} \end{pmatrix} = \begin{pmatrix} \gamma(0) & \gamma(1) \\ \gamma(1) & \gamma(0) \end{pmatrix}^{-1} \begin{pmatrix} \gamma(1) \\ \gamma(2) \end{pmatrix},$$

but, it should be apparent from the model that  $x_3^2 = \phi_1 x_2 + \phi_2 x_1$ . Because  $\phi_1 x_2 + \phi_2 x_1$  satisfies the prediction equations (3.60),

$$E\{[x_3 - (\phi_1 x_2 + \phi_2 x_1)]x_1\} = E(w_3 x_1) = 0,$$

$$E\{[x_3 - (\phi_1 x_2 + \phi_2 x_1)]x_2\} = E(w_3 x_2) = 0,$$

it follows that, indeed,  $x_3^2 = \phi_1 x_2 + \phi_2 x_1$ , and by the uniqueness of the coefficients in this case, that  $\phi_{21} = \phi_1$  and  $\phi_{22} = \phi_2$ . Continuing in this way, it is easy to verify that, for  $n \geq 2$ ,

$$x_{n+1}^n = \phi_1 x_n + \phi_2 x_{n-1}.$$

That is,  $\phi_{n1} = \phi_1$ ,  $\phi_{n2} = \phi_2$ , and  $\phi_{nj} = 0$ , for  $j = 3, 4, \dots, n$ .

From Example 3.19, it should be clear (Problem 3.45) that, if the time series is a causal AR( $p$ ) process, then, for  $n \geq p$ ,

$$x_{n+1}^n = \phi_1 x_n + \phi_2 x_{n-1} + \dots + \phi_p x_{n-p+1}. \quad (3.67)$$

For ARMA models in general, the prediction equations will not be as simple as the pure AR case. In addition, for  $n$  large, the use of (3.64) is prohibitive because it requires the inversion of a large matrix. There are, however, iterative solutions that do not require any matrix inversion. In particular, we mention the recursive solution due to Levinson (1947) and Durbin (1960).

**Property 3.4 The Durbin–Levinson Algorithm**

Equations (3.64) and (3.66) can be solved iteratively as follows:

$$\phi_{00} = 0, \quad P_1^0 = \gamma(0). \quad (3.68)$$

For  $n \geq 1$ ,

$$\phi_{nn} = \frac{\rho(n) - \sum_{k=1}^{n-1} \phi_{n-1,k} \rho(n-k)}{1 - \sum_{k=1}^{n-1} \phi_{n-1,k} \rho(k)}, \quad P_{n+1}^n = P_n^{n-1} (1 - \phi_{nn}^2), \quad (3.69)$$

where, for  $n \geq 2$ ,

$$\phi_{nk} = \phi_{n-1,k} - \phi_{nn} \phi_{n-1,n-k}, \quad k = 1, 2, \dots, n-1. \quad (3.70)$$

The proof of **Property 3.4** is left as an exercise; see **Problem 3.13**.

### Example 3.20 Using the Durbin–Levinson Algorithm

To use the algorithm, start with  $\phi_{00} = 0$ ,  $P_1^0 = \gamma(0)$ . Then, for  $n = 1$ ,

$$\phi_{11} = \rho(1), \quad P_2^1 = \gamma(0)[1 - \phi_{11}^2].$$

For  $n = 2$ ,

$$\begin{aligned} \phi_{22} &= \frac{\rho(2) - \phi_{11} \rho(1)}{1 - \phi_{11} \rho(1)}, \quad \phi_{21} = \phi_{11} - \phi_{22} \phi_{11}, \\ P_3^2 &= P_2^1 [1 - \phi_{22}^2] = \gamma(0)[1 - \phi_{11}^2][1 - \phi_{22}^2]. \end{aligned}$$

For  $n = 3$ ,

$$\begin{aligned} \phi_{33} &= \frac{\rho(3) - \phi_{21} \rho(2) - \phi_{22} \rho(1)}{1 - \phi_{21} \rho(1) - \phi_{22} \rho(2)}, \\ \phi_{32} &= \phi_{22} - \phi_{33} \phi_{21}, \quad \phi_{31} = \phi_{21} - \phi_{33} \phi_{22}, \\ P_4^3 &= P_3^2 [1 - \phi_{33}^2] = \gamma(0)[1 - \phi_{11}^2][1 - \phi_{22}^2][1 - \phi_{33}^2], \end{aligned}$$

and so on. Note that, in general, the standard error of the one-step-ahead forecast is the square root of

$$P_{n+1}^n = \gamma(0) \prod_{j=1}^n [1 - \phi_{jj}^2]. \quad (3.71)$$

An important consequence of the Durbin–Levinson algorithm is (see **Problem 3.13**) as follows.

### Property 3.5 Iterative Solution for the PACF

The PACF of a stationary process  $x_t$ , can be obtained iteratively via (3.69) as  $\phi_{nn}$ , for  $n = 1, 2, \dots$ .

Using **Property 3.5** and putting  $n = p$  in (3.61) and (3.67), it follows that for an  $\text{AR}(p)$  model,

$$\begin{aligned} x_{p+1}^p &= \phi_{p1} x_p + \phi_{p2} x_{p-1} + \dots + \phi_{pp} x_1 \\ &= \phi_1 x_p + \phi_2 x_{p-1} + \dots + \phi_p x_1. \end{aligned} \quad (3.72)$$

Result (3.72) shows that for an  $\text{AR}(p)$  model, the partial autocorrelation coefficient at lag  $p$ ,  $\phi_{pp}$ , is also the last coefficient in the model,  $\phi_p$ , as was claimed in **Example 3.16**.

**Example 3.21 The PACF of an AR(2)**

We will use the results of [Example 3.20](#) and [Property 3.5](#) to calculate the first three values,  $\phi_{11}$ ,  $\phi_{22}$ ,  $\phi_{33}$ , of the PACF. Recall from [Example 3.10](#) that  $\rho(h) - \phi_1\rho(h-1) - \phi_2\rho(h-2) = 0$  for  $h \geq 1$ . When  $h = 1, 2, 3$ , we have  $\rho(1) = \phi_1/(1 - \phi_2)$ ,  $\rho(2) = \phi_1\rho(1) + \phi_2$ ,  $\rho(3) - \phi_1\rho(2) - \phi_2\rho(1) = 0$ . Thus,

$$\begin{aligned}\phi_{11} &= \rho(1) = \frac{\phi_1}{1 - \phi_2} \\ \phi_{22} &= \frac{\rho(2) - \rho(1)^2}{1 - \rho(1)^2} = \frac{\left[\phi_1\left(\frac{\phi_1}{1 - \phi_2}\right) + \phi_2\right] - \left(\frac{\phi_1}{1 - \phi_2}\right)^2}{1 - \left(\frac{\phi_1}{1 - \phi_2}\right)^2} = \phi_2 \\ \phi_{21} &= \rho(1)[1 - \phi_2] = \phi_1 \\ \phi_{33} &= \frac{\rho(3) - \phi_1\rho(2) - \phi_2\rho(1)}{1 - \phi_1\rho(1) - \phi_2\rho(2)} = 0.\end{aligned}$$

Notice that, as shown in [\(3.72\)](#),  $\phi_{22} = \phi_2$  for an AR(2) model.

So far, we have concentrated on one-step-ahead prediction, but [Property 3.3](#) allows us to calculate the BLP of  $x_{n+m}$  for any  $m \geq 1$ . Given data,  $\{x_1, \dots, x_n\}$ , the  $m$ -step-ahead predictor is

$$x_{n+m}^n = \phi_{n1}^{(m)} x_n + \phi_{n2}^{(m)} x_{n-1} + \dots + \phi_{nn}^{(m)} x_1, \quad (3.73)$$

where  $\{\phi_{n1}^{(m)}, \phi_{n2}^{(m)}, \dots, \phi_{nn}^{(m)}\}$  satisfy the prediction equations,

$$\sum_{j=1}^n \phi_{nj}^{(m)} E(x_{n+1-j} x_{n+1-k}) = E(x_{n+m} x_{n+1-k}), \quad k = 1, \dots, n,$$

or

$$\sum_{j=1}^n \phi_{nj}^{(m)} \gamma(k-j) = \gamma(m+k-1), \quad k = 1, \dots, n. \quad (3.74)$$

The prediction equations can again be written in matrix notation as

$$\Gamma_n \phi_n^{(m)} = \gamma_n^{(m)}, \quad (3.75)$$

where  $\gamma_n^{(m)} = (\gamma(m), \dots, \gamma(m+n-1))'$ , and  $\phi_n^{(m)} = (\phi_{n1}^{(m)}, \dots, \phi_{nn}^{(m)})'$  are  $n \times 1$  vectors. The *mean square  $m$ -step-ahead prediction error* is

$$P_{n+m}^n = E(x_{n+m} - x_{n+m}^n)^2 = \gamma(0) - \gamma_n^{(m)'} \Gamma_n^{-1} \gamma_n^{(m)}. \quad (3.76)$$

Another useful algorithm for calculating forecasts was given by Brockwell and Davis (1991, Chapter 5). This algorithm follows directly from applying the projection theorem ([Theorem B.1](#)) to the *innovations*,  $x_t - x_t^{t-1}$ , for  $t = 1, \dots, n$ , using the fact that the innovations  $x_t - x_t^{t-1}$  and  $x_s - x_s^{s-1}$  are uncorrelated for  $s \neq t$  (see [Problem 3.46](#)). We present the case in which  $x_t$  is a mean-zero stationary time series.

**Property 3.6 The Innovations Algorithm**

The one-step-ahead predictors,  $x_{t+1}^t$ , and their mean-squared errors,  $P_{t+1}^t$ , can be calculated iteratively as

$$x_1^0 = 0, \quad P_1^0 = \gamma(0)$$

$$x_{t+1}^t = \sum_{j=1}^t \theta_{tj} (x_{t+1-j} - x_{t+1-j}^{t-j}), \quad t = 1, 2, \dots \quad (3.77)$$

$$P_{t+1}^t = \gamma(0) - \sum_{j=0}^{t-1} \theta_{t,t-j}^2 P_{j+1}^j \quad t = 1, 2, \dots, \quad (3.78)$$

where, for  $j = 0, 1, \dots, t-1$ ,

$$\theta_{t,t-j} = \left( \gamma(t-j) - \sum_{k=0}^{j-1} \theta_{j,j-k} \theta_{t,t-k} P_{k+1}^k \right) / P_{j+1}^j. \quad (3.79)$$

Given data  $x_1, \dots, x_n$ , the innovations algorithm can be calculated successively for  $t = 1$ , then  $t = 2$  and so on, in which case the calculation of  $x_{n+1}^n$  and  $P_{n+1}^n$  is made at the final step  $t = n$ . The  $m$ -step-ahead predictor and its mean-square error based on the innovations algorithm (Problem 3.46) are given by

$$x_{n+m}^n = \sum_{j=m}^{n+m-1} \theta_{n+m-1,j} (x_{n+m-j} - x_{n+m-j}^{n+m-j-1}), \quad (3.80)$$

$$P_{n+m}^n = \gamma(0) - \sum_{j=m}^{n+m-1} \theta_{n+m-1,j}^2 P_{n+m-j}^{n+m-j-1}, \quad (3.81)$$

where the  $\theta_{n+m-1,j}$  are obtained by continued iteration of (3.79).

**Example 3.22 Prediction for an MA(1)**

The innovations algorithm lends itself well to prediction for moving average processes. Consider an MA(1) model,  $x_t = w_t + \theta w_{t-1}$ . Recall that  $\gamma(0) = (1 + \theta^2)\sigma_w^2$ ,  $\gamma(1) = \theta\sigma_w^2$ , and  $\gamma(h) = 0$  for  $h > 1$ . Then, using Property 3.6, we have

$$\theta_{n1} = \theta\sigma_w^2 / P_n^{n-1}$$

$$\theta_{nj} = 0, \quad j = 2, \dots, n$$

$$P_1^0 = (1 + \theta^2)\sigma_w^2$$

$$P_{n+1}^n = (1 + \theta^2 - \theta\theta_{n1})\sigma_w^2.$$

Finally, from (3.77), the one-step-ahead predictor is

$$x_{n+1}^n = \theta (x_n - x_n^{n-1}) \sigma_w^2 / P_n^{n-1}.$$



## FORECASTING ARMA PROCESSES

The general prediction equations (3.60) provide little insight into forecasting for ARMA models in general. There are a number of different ways to express these forecasts, and each aids in understanding the special structure of ARMA prediction. Throughout, we assume  $x_t$  is a causal and invertible ARMA( $p, q$ ) process,  $\phi(B)x_t = \theta(B)w_t$ , where  $w_t \sim \text{iid } N(0, \sigma_w^2)$ . In the non-zero mean case,  $E(x_t) = \mu_x$ , simply replace  $x_t$  with  $x_t - \mu_x$  in the model. First, we consider two types of forecasts. We write  $x_{n+m}^n$  to mean the minimum mean square error predictor of  $x_{n+m}$  based on the data  $\{x_n, \dots, x_1\}$ , that is,

$$x_{n+m}^n = E(x_{n+m} \mid x_n, \dots, x_1).$$

For ARMA models, it is easier to calculate the predictor of  $x_{n+m}$ , assuming we have the complete history of the process  $\{x_n, x_{n-1}, \dots, x_1, x_0, x_{-1}, \dots\}$ . We will denote the predictor of  $x_{n+m}$  based on the infinite past as

$$\tilde{x}_{n+m} = E(x_{n+m} \mid x_n, x_{n-1}, \dots, x_1, x_0, x_{-1}, \dots).$$

In general,  $x_{n+m}^n$  and  $\tilde{x}_{n+m}$  are not the same, but the idea here is that, for large samples,  $\tilde{x}_{n+m}$  will provide a good approximation to  $x_{n+m}^n$ .

Now, write  $x_{n+m}$  in its causal and invertible forms:

$$x_{n+m} = \sum_{j=0}^{\infty} \psi_j w_{n+m-j}, \quad \psi_0 = 1 \quad (3.82)$$

$$w_{n+m} = \sum_{j=0}^{\infty} \pi_j x_{n+m-j}, \quad \pi_0 = 1. \quad (3.83)$$

Then, taking conditional expectations in (3.82), we have

$$\tilde{x}_{n+m} = \sum_{j=0}^{\infty} \psi_j \tilde{w}_{n+m-j} = \sum_{j=m}^{\infty} \psi_j w_{n+m-j}, \quad (3.84)$$

because, by causality and invertibility,

$$\tilde{w}_t = E(w_t \mid x_n, x_{n-1}, \dots, x_0, x_{-1}, \dots) = \begin{cases} 0 & t > n \\ w_t & t \leq n. \end{cases}$$

Similarly, taking conditional expectations in (3.83), we have

$$0 = \tilde{x}_{n+m} + \sum_{j=1}^{\infty} \pi_j \tilde{x}_{n+m-j},$$

or

$$\tilde{x}_{n+m} = - \sum_{j=1}^{m-1} \pi_j \tilde{x}_{n+m-j} - \sum_{j=m}^{\infty} \pi_j x_{n+m-j}, \quad (3.85)$$

using the fact  $E(x_t \mid x_n, x_{n-1}, \dots, x_0, x_{-1}, \dots) = x_t$ , for  $t \leq n$ . Prediction is accomplished recursively using (3.85), starting with the one-step-ahead predictor,  $m = 1$ , and then continuing for  $m = 2, 3, \dots$ . Using (3.84), we can write

$$x_{n+m} - \tilde{x}_{n+m} = \sum_{j=0}^{m-1} \psi_j w_{n+m-j},$$

so the *mean-square prediction error* can be written as

$$P_{n+m}^n = E(x_{n+m} - \tilde{x}_{n+m})^2 = \sigma_w^2 \sum_{j=0}^{m-1} \psi_j^2. \quad (3.86)$$

Also, we note, for a fixed sample size,  $n$ , the prediction errors are correlated. That is, for  $k \geq 1$ ,

$$E\{(x_{n+m} - \tilde{x}_{n+m})(x_{n+m+k} - \tilde{x}_{n+m+k})\} = \sigma_w^2 \sum_{j=0}^{m-1} \psi_j \psi_{j+k}. \quad (3.87)$$

### Example 3.23 Long-Range Forecasts

Consider forecasting an ARMA process with mean  $\mu_x$ . Replacing  $x_{n+m}$  with  $x_{n+m} - \mu_x$  in (3.82), and taking conditional expectation as in (3.84), we deduce that the  $m$ -step-ahead forecast can be written as

$$\tilde{x}_{n+m} = \mu_x + \sum_{j=m}^{\infty} \psi_j w_{n+m-j}. \quad (3.88)$$

Noting that the  $\psi$ -weights dampen to zero exponentially fast, it is clear that

$$\tilde{x}_{n+m} \rightarrow \mu_x \quad (3.89)$$

exponentially fast (in the mean square sense) as  $m \rightarrow \infty$ . Moreover, by (3.86), the mean square prediction error

$$P_{n+m}^n \rightarrow \sigma_w^2 \sum_{j=0}^{\infty} \psi_j^2 = \gamma_x(0) = \sigma_x^2, \quad (3.90)$$

exponentially fast as  $m \rightarrow \infty$ .

It should be clear from (3.89) and (3.90) that ARMA forecasts quickly settle to the mean with a constant prediction error as the forecast horizon,  $m$ , grows. This effect can be seen in Figure 3.7 where the Recruitment series is forecast for 24 months; see Example 3.25.

When  $n$  is small, the general prediction equations (3.60) can be used easily. When  $n$  is large, we would use (3.85) by truncating, because we do not observe

$x_0, x_{-1}, x_{-2}, \dots$ , and only the data  $x_1, x_2, \dots, x_n$  are available. In this case, we can truncate (3.85) by setting  $\sum_{j=n+m}^{\infty} \pi_j x_{n+m-j} = 0$ . The *truncated predictor* is then written as

$$\tilde{x}_{n+m}^n = - \sum_{j=1}^{m-1} \pi_j \tilde{x}_{n+m-j}^n - \sum_{j=m}^{n+m-1} \pi_j x_{n+m-j}, \quad (3.91)$$

which is also calculated recursively,  $m = 1, 2, \dots$ . The mean square prediction error, in this case, is approximated using (3.86).

For AR( $p$ ) models, and when  $n > p$ , equation (3.67) yields the exact predictor,  $x_{n+m}^n$ , of  $x_{n+m}$ , and there is no need for approximations. That is, for  $n > p$ ,  $\tilde{x}_{n+m}^n = x_{n+m}^n = x_{n+m}^n$ . Also, in this case, the one-step-ahead prediction error is  $E(x_{n+1} - x_{n+1}^n)^2 = \sigma_w^2$ . For pure MA( $q$ ) or ARMA( $p, q$ ) models, truncated prediction has a fairly nice form.

### Property 3.7 Truncated Prediction for ARMA

For ARMA( $p, q$ ) models, the truncated predictors for  $m = 1, 2, \dots$ , are

$$\tilde{x}_{n+m}^n = \phi_1 \tilde{x}_{n+m-1}^n + \dots + \phi_p \tilde{x}_{n+m-p}^n + \theta_1 \tilde{w}_{n+m-1}^n + \dots + \theta_q \tilde{w}_{n+m-q}^n, \quad (3.92)$$

where  $\tilde{x}_t^n = x_t$  for  $1 \leq t \leq n$  and  $\tilde{x}_t^n = 0$  for  $t \leq 0$ . The truncated prediction errors are given by:  $\tilde{w}_t^n = 0$  for  $t \leq 0$  or  $t > n$ , and

$$\tilde{w}_t^n = \phi(B)\tilde{x}_t^n - \theta_1 \tilde{w}_{t-1}^n - \dots - \theta_q \tilde{w}_{t-q}^n$$

for  $1 \leq t \leq n$ .

### Example 3.24 Forecasting an ARMA(1, 1) Series

Given data  $x_1, \dots, x_n$ , for forecasting purposes, write the model as

$$x_{n+1} = \phi x_n + w_{n+1} + \theta w_n.$$

Then, based on (3.92), the one-step-ahead truncated forecast is

$$\tilde{x}_{n+1}^n = \phi x_n + 0 + \theta \tilde{w}_n^n.$$

For  $m \geq 2$ , we have

$$\tilde{x}_{n+m}^n = \phi \tilde{x}_{n+m-1}^n,$$

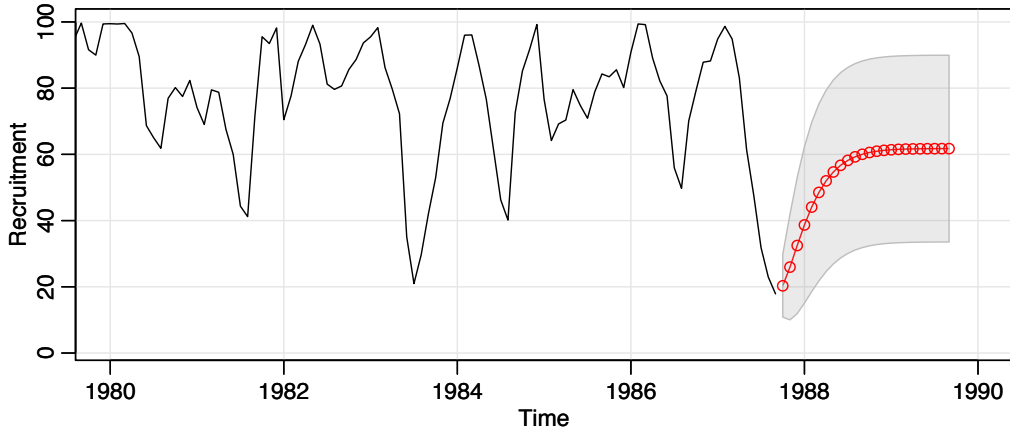
which can be calculated recursively,  $m = 2, 3, \dots$ .

To calculate  $\tilde{w}_n^n$ , which is needed to initialize the successive forecasts, the model can be written as  $w_t = x_t - \phi x_{t-1} - \theta w_{t-1}$  for  $t = 1, \dots, n$ . For truncated forecasting using (3.92), put  $\tilde{w}_0^n = 0$ ,  $x_0 = 0$ , and then iterate the errors forward in time

$$\tilde{w}_t^n = x_t - \phi x_{t-1} - \theta \tilde{w}_{t-1}^n, \quad t = 1, \dots, n.$$

The approximate forecast variance is computed from (3.86) using the  $\psi$ -weights determined as in Example 3.12. In particular, the  $\psi$ -weights satisfy  $\psi_j = (\phi + \theta)\phi^{j-1}$ , for  $j \geq 1$ . This result gives

$$P_{n+m}^n = \sigma_w^2 \left[ 1 + (\phi + \theta)^2 \sum_{j=1}^{m-1} \phi^{2(j-1)} \right] = \sigma_w^2 \left[ 1 + \frac{(\phi + \theta)^2 (1 - \phi^{2(m-1)})}{(1 - \phi^2)} \right].$$



**Fig. 3.7.** Twenty-four month forecasts for the Recruitment series. The actual data shown are from about January 1980 to September 1987, and then the forecasts plus and minus one standard error are displayed.

To assess the precision of the forecasts, *prediction intervals* are typically calculated along with the forecasts. In general,  $(1 - \alpha)$  prediction intervals are of the form

$$x_{n+m}^n \pm c_{\alpha/2} \sqrt{P_{n+m}^n}, \quad (3.93)$$

where  $c_{\alpha/2}$  is chosen to get the desired degree of confidence. For example, if the process is Gaussian, then choosing  $c_{\alpha/2} = 2$  will yield an approximate 95% prediction interval for  $x_{n+m}$ . If we are interested in establishing prediction intervals over more than one time period, then  $c_{\alpha/2}$  should be adjusted appropriately, for example, by using Bonferroni's inequality [see (4.63) in Chapter 4 or Johnson and Wichern, 1992, Chapter 5].

### Example 3.25 Forecasting the Recruitment Series

Using the parameter estimates as the actual parameter values, Figure 3.7 shows the result of forecasting the Recruitment series given in Example 3.18 over a 24-month horizon,  $m = 1, 2, \dots, 24$ . The actual forecasts are calculated as

$$x_{n+m}^n = 6.74 + 1.35x_{n+m-1}^n - .46x_{n+m-2}^n$$

for  $n = 453$  and  $m = 1, 2, \dots, 12$ . Recall that  $x_t^s = x_t$  when  $t \leq s$ . The forecasts errors  $P_{n+m}^n$  are calculated using (3.86). Recall that  $\hat{\sigma}_w^2 = 89.72$ , and using (3.40) from Example 3.12, we have  $\psi_j = 1.35\psi_{j-1} - .46\psi_{j-2}$  for  $j \geq 2$ , where  $\psi_0 = 1$  and  $\psi_1 = 1.35$ . Thus, for  $n = 453$ ,

$$\begin{aligned} P_{n+1}^n &= 89.72, \\ P_{n+2}^n &= 89.72(1 + 1.35^2), \\ P_{n+3}^n &= 89.72(1 + 1.35^2 + [1.35^2 - .46]^2), \end{aligned}$$

and so on.

Note how the forecast levels off quickly and the prediction intervals are wide, even though in this case the forecast limits are only based on one standard error; that is,  $x_{n+m}^n \pm \sqrt{P_{n+m}^n}$ .

To reproduce the analysis and Figure 3.7, use the following commands:

```
regr = ar.ols(rec, order=2, demean=FALSE, intercept=TRUE)
fore = predict(regr, n.ahead=24)
ts.plot(rec, fore$pred, col=1:2, xlim=c(1980,1990), ylab="Recruitment")
U = fore$pred+fore$se; L = fore$pred-fore$se
xx = c(time(U), rev(time(U))); yy = c(L, rev(U))
polygon(xx, yy, border = 8, col = gray(.6, alpha = .2))
lines(fore$pred, type="p", col=2)
```

We complete this section with a brief discussion of *backcasting*. In backcasting, we want to predict  $x_{1-m}$ , for  $m = 1, 2, \dots$ , based on the data  $\{x_1, \dots, x_n\}$ . Write the backcast as

$$x_{1-m}^n = \sum_{j=1}^n \alpha_j x_j. \quad (3.94)$$

Analogous to (3.74), the prediction equations (assuming  $\mu_x = 0$ ) are

$$\sum_{j=1}^n \alpha_j E(x_j x_k) = E(x_{1-m} x_k), \quad k = 1, \dots, n, \quad (3.95)$$

or

$$\sum_{j=1}^n \alpha_j \gamma(k-j) = \gamma(m+k-1), \quad k = 1, \dots, n. \quad (3.96)$$

These equations are precisely the prediction equations for forward prediction. That is,  $\alpha_j \equiv \phi_{nj}^{(m)}$ , for  $j = 1, \dots, n$ , where the  $\phi_{nj}^{(m)}$  are given by (3.75). Finally, the backcasts are given by

$$x_{1-m}^n = \phi_{n1}^{(m)} x_1 + \dots + \phi_{nn}^{(m)} x_n, \quad m = 1, 2, \dots \quad (3.97)$$

### Example 3.26 Backcasting an ARMA(1, 1)

Consider an ARMA(1, 1) process,  $x_t = \phi x_{t-1} + \theta w_{t-1} + w_t$ ; we will call this the *forward model*. We have just seen that best linear prediction backward in time is the same as best linear prediction forward in time for stationary models. Assuming the models are Gaussian, we also have that minimum mean square error prediction backward in time is the same as forward in time for ARMA models.<sup>3.4</sup> Thus, the process can equivalently be generated by the *backward model*,

$$x_t = \phi x_{t+1} + \theta v_{t+1} + v_t,$$

<sup>3.4</sup> In the stationary Gaussian case, (a) the distribution of  $\{x_{n+1}, x_n, \dots, x_1\}$  is the same as (b) the distribution of  $\{x_0, x_1, \dots, x_n\}$ . In forecasting we use (a) to obtain  $E(x_{n+1} | x_n, \dots, x_1)$ ; in backcasting we use (b) to obtain  $E(x_0 | x_1, \dots, x_n)$ . Because (a) and (b) are the same, the two problems are equivalent.

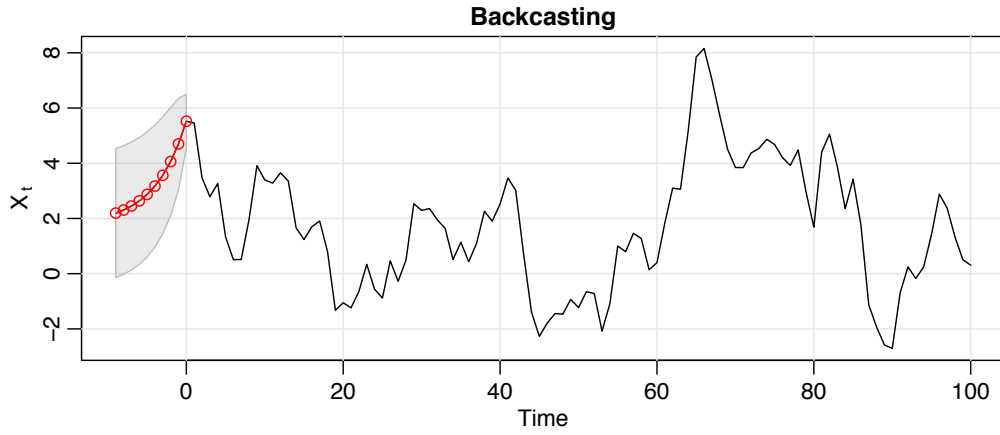


Fig. 3.8. Display for *Example 3.26*; backcasts from a simulated ARMA(1, 1).

where  $\{v_t\}$  is a Gaussian white noise process with variance  $\sigma_w^2$ . We may write  $x_t = \sum_{j=0}^{\infty} \psi_j v_{t+j}$ , where  $\psi_0 = 1$ ; this means that  $x_t$  is uncorrelated with  $\{v_{t-1}, v_{t-2}, \dots\}$ , in analogy to the forward model.

Given data  $\{x_1, \dots, x_n\}$ , truncate  $v_n^n = E(v_n | x_1, \dots, x_n)$  to zero and then iterate backward. That is, put  $\tilde{v}_n^n = 0$ , as an initial approximation, and then generate the errors backward

$$\tilde{v}_t^n = x_t - \phi x_{t+1} - \theta \tilde{v}_{t+1}^n, \quad t = (n-1), (n-2), \dots, 1.$$

Then,

$$\tilde{x}_0^n = \phi x_1 + \theta \tilde{v}_1^n + \tilde{v}_0^n = \phi x_1 + \theta \tilde{v}_1^n,$$

because  $\tilde{v}_t^n = 0$  for  $t \leq 0$ . Continuing, the general truncated backcasts are given by

$$\tilde{x}_{1-m}^n = \phi \tilde{x}_{2-m}^n, \quad m = 2, 3, \dots$$

To backcast data in R, simply reverse the data, fit the model and predict. In the following, we backcasted a simulated ARMA(1,1) process; see *Figure 3.8*.

```
set.seed(90210)
x = arima.sim(list(order = c(1,0,1), ar = .9, ma=.5), n = 100)
xr = rev(x) # xr is the reversed data
pxr = predict(arima(xr, order=c(1,0,1)), 10) # predict the reversed data
pxrp = rev(pxr$pred) # reorder the predictors (for plotting)
pxrse = rev(pxr$se) # reorder the SEs
nx = ts(c(pxrp, x), start=-9) # attach the backcasts to the data
plot(nx, ylab=expression(X[~t]), main='Backcasting')
U = nx[1:10] + pxrse; L = nx[1:10] - pxrse
xx = c(-9:0, 0:-9); yy = c(L, rev(U))
polygon(xx, yy, border = 8, col = gray(0.6, alpha = 0.2))
lines(-9:0, nx[1:10], col=2, type='o')
```

### 3.5 Estimation

Throughout this section, we assume we have  $n$  observations,  $x_1, \dots, x_n$ , from a causal and invertible Gaussian ARMA( $p, q$ ) process in which, initially, the order parameters,  $p$  and  $q$ , are known. Our goal is to estimate the parameters,  $\phi_1, \dots, \phi_p, \theta_1, \dots, \theta_q$ , and  $\sigma_w^2$ . We will discuss the problem of determining  $p$  and  $q$  later in this section.

We begin with *method of moments* estimators. The idea behind these estimators is that of equating population moments to sample moments and then solving for the parameters in terms of the sample moments. We immediately see that, if  $E(x_t) = \mu$ , then the method of moments estimator of  $\mu$  is the sample average,  $\bar{x}$ . Thus, while discussing method of moments, we will assume  $\mu = 0$ . Although the method of moments can produce good estimators, they can sometimes lead to suboptimal estimators. We first consider the case in which the method leads to optimal (efficient) estimators, that is, **AR( $p$ ) models**,

$$x_t = \phi_1 x_{t-1} + \dots + \phi_p x_{t-p} + w_t,$$

where the first  $p + 1$  equations of (3.47) and (3.48) lead to the following:

**Definition 3.10** The **Yule–Walker equations** are given by

$$\gamma(h) = \phi_1 \gamma(h-1) + \dots + \phi_p \gamma(h-p), \quad h = 1, 2, \dots, p, \quad (3.98)$$

$$\sigma_w^2 = \gamma(0) - \phi_1 \gamma(1) - \dots - \phi_p \gamma(p). \quad (3.99)$$

In matrix notation, the Yule–Walker equations are

$$\Gamma_p \phi = \gamma_p, \quad \sigma_w^2 = \gamma(0) - \phi' \gamma_p, \quad (3.100)$$

where  $\Gamma_p = \{\gamma(k-j)\}_{j,k=1}^p$  is a  $p \times p$  matrix,  $\phi = (\phi_1, \dots, \phi_p)'$  is a  $p \times 1$  vector, and  $\gamma_p = (\gamma(1), \dots, \gamma(p))'$  is a  $p \times 1$  vector. Using the method of moments, we replace  $\gamma(h)$  in (3.100) by  $\hat{\gamma}(h)$  [see equation (1.36)] and solve

$$\hat{\phi} = \hat{\Gamma}_p^{-1} \hat{\gamma}_p, \quad \hat{\sigma}_w^2 = \hat{\gamma}(0) - \hat{\gamma}_p' \hat{\Gamma}_p^{-1} \hat{\gamma}_p. \quad (3.101)$$

These estimators are typically called the *Yule–Walker estimators*. For calculation purposes, it is sometimes more convenient to work with the sample ACF. By factoring  $\hat{\gamma}(0)$  in (3.101), we can write the Yule–Walker estimates as

$$\hat{\phi} = \hat{R}_p^{-1} \hat{\rho}_p, \quad \hat{\sigma}_w^2 = \hat{\gamma}(0) [1 - \hat{\rho}_p' \hat{R}_p^{-1} \hat{\rho}_p], \quad (3.102)$$

where  $\hat{R}_p = \{\hat{\rho}(k-j)\}_{j,k=1}^p$  is a  $p \times p$  matrix and  $\hat{\rho}_p = (\hat{\rho}(1), \dots, \hat{\rho}(p))'$  is a  $p \times 1$  vector.

For AR( $p$ ) models, if the sample size is large, the Yule–Walker estimators are approximately normally distributed, and  $\hat{\sigma}_w^2$  is close to the true value of  $\sigma_w^2$ . We state these results in **Property 3.8**; for details, see **Section B.3**.

**Property 3.8 Large Sample Results for Yule–Walker Estimators**

The asymptotic ( $n \rightarrow \infty$ ) behavior of the Yule–Walker estimators in the case of causal AR( $p$ ) processes is as follows:

$$\sqrt{n}(\hat{\phi} - \phi) \xrightarrow{d} N(0, \sigma_w^2 \Gamma_p^{-1}), \quad \hat{\sigma}_w^2 \xrightarrow{p} \sigma_w^2. \quad (3.103)$$

The Durbin–Levinson algorithm, (3.68)–(3.70), can be used to calculate  $\hat{\phi}$  without inverting  $\hat{\Gamma}_p$  or  $\hat{R}_p$ , by replacing  $\gamma(h)$  by  $\hat{\gamma}(h)$  in the algorithm. In running the algorithm, we will iteratively calculate the  $h \times 1$  vector,  $\hat{\phi}_h = (\hat{\phi}_{h1}, \dots, \hat{\phi}_{hh})'$ , for  $h = 1, 2, \dots$ . Thus, in addition to obtaining the desired forecasts, the Durbin–Levinson algorithm yields  $\hat{\phi}_{hh}$ , the sample PACF. Using (3.103), we can show the following property.

**Property 3.9 Large Sample Distribution of the PACF**

For a causal AR( $p$ ) process, asymptotically ( $n \rightarrow \infty$ ),

$$\sqrt{n} \hat{\phi}_{hh} \xrightarrow{d} N(0, 1), \quad \text{for } h > p. \quad (3.104)$$

**Example 3.27 Yule–Walker Estimation for an AR(2) Process**

The data shown in Figure 3.4 were  $n = 144$  simulated observations from the AR(2) model

$$x_t = 1.5x_{t-1} - .75x_{t-2} + w_t,$$

where  $w_t \sim \text{iid } N(0, 1)$ . For these data,  $\hat{\gamma}(0) = 8.903$ ,  $\hat{\rho}(1) = .849$ , and  $\hat{\rho}(2) = .519$ . Thus,

$$\hat{\phi} = \begin{pmatrix} \hat{\phi}_1 \\ \hat{\phi}_2 \end{pmatrix} = \begin{bmatrix} 1 & .849 \\ .849 & 1 \end{bmatrix}^{-1} \begin{pmatrix} .849 \\ .519 \end{pmatrix} = \begin{pmatrix} 1.463 \\ -.723 \end{pmatrix}$$

and

$$\hat{\sigma}_w^2 = 8.903 \left[ 1 - (.849, .519) \begin{pmatrix} 1.463 \\ -.723 \end{pmatrix} \right] = 1.187.$$

By Property 3.8, the asymptotic variance–covariance matrix of  $\hat{\phi}$  is

$$\frac{1}{144} \frac{1.187}{8.903} \begin{bmatrix} 1 & .849 \\ .849 & 1 \end{bmatrix}^{-1} = \begin{bmatrix} .058^2 & -.003 \\ -.003 & .058^2 \end{bmatrix},$$

and it can be used to get confidence regions for, or make inferences about  $\hat{\phi}$  and its components. For example, an approximate 95% confidence interval for  $\phi_2$  is  $-.723 \pm 2(.058)$ , or  $(-.838, -.608)$ , which contains the true value of  $\phi_2 = -.75$ .

For these data, the first three sample partial autocorrelations are  $\hat{\phi}_{11} = \hat{\rho}(1) = .849$ ,  $\hat{\phi}_{22} = \hat{\phi}_2 = -.721$ , and  $\hat{\phi}_{33} = -.085$ . According to Property 3.9, the asymptotic standard error of  $\hat{\phi}_{33}$  is  $1/\sqrt{144} = .083$ , and the observed value,  $-.085$ , is about only one standard deviation from  $\phi_{33} = 0$ .



**Example 3.28 Yule–Walker Estimation of the Recruitment Series**

In [Example 3.18](#) we fit an AR(2) model to the recruitment series using ordinary least squares (OLS). For AR models, the estimators obtained via OLS and Yule–Walker are nearly identical; we will see this when we discuss conditional sum of squares estimation in [\(3.111\)–\(3.116\)](#).

Below are the results of fitting the same model using Yule–Walker estimation in R, which are nearly identical to the values in [Example 3.18](#).

```
rec.yw = ar.yw(rec, order=2)
rec.yw$x.mean      # = 62.26 (mean estimate)
rec.yw$ar          # = 1.33, -.44 (coefficient estimates)
sqrt(diag(rec.yw$asy.var.coef)) # = .04, .04 (standard errors)
rec.yw$var.pred    # = 94.80 (error variance estimate)
```

To obtain the 24 month ahead predictions and their standard errors, and then plot the results (not shown) as in [Example 3.25](#), use the R commands:

```
rec.pr = predict(rec.yw, n.ahead=24)
ts.plot(rec, rec.pr$pred, col=1:2)
lines(rec.pr$pred + rec.pr$se, col=4, lty=2)
lines(rec.pr$pred - rec.pr$se, col=4, lty=2)
```

In the case of AR( $p$ ) models, the Yule–Walker estimators given in [\(3.102\)](#) are optimal in the sense that the asymptotic distribution, [\(3.103\)](#), is the best asymptotic normal distribution. This is because, given initial conditions, AR( $p$ ) models are linear models, and the Yule–Walker estimators are essentially least squares estimators. If we use method of moments for MA or ARMA models, we will not get optimal estimators because such processes are nonlinear in the parameters.

**Example 3.29 Method of Moments Estimation for an MA(1)**

Consider the time series

$$x_t = w_t + \theta w_{t-1},$$

where  $|\theta| < 1$ . The model can then be written as

$$x_t = \sum_{j=1}^{\infty} (-\theta)^j x_{t-j} + w_t,$$

which is nonlinear in  $\theta$ . The first two population autocovariances are  $\gamma(0) = \sigma_w^2(1 + \theta^2)$  and  $\gamma(1) = \sigma_w^2\theta$ , so the estimate of  $\theta$  is found by solving:

$$\hat{\rho}(1) = \frac{\hat{\gamma}(1)}{\hat{\gamma}(0)} = \frac{\hat{\theta}}{1 + \hat{\theta}^2}.$$

Two solutions exist, so we would pick the invertible one. If  $|\hat{\rho}(1)| \leq \frac{1}{2}$ , the solutions are real, otherwise, a real solution does not exist. Even though  $|\rho(1)| < \frac{1}{2}$  for an invertible MA(1), it may happen that  $|\hat{\rho}(1)| \geq \frac{1}{2}$  because it is an estimator. For example, the following simulation in R produces a value of  $\hat{\rho}(1) = .507$  when the true value is  $\rho(1) = .9/(1 + .9^2) = .497$ .

```
set.seed(2)
mal = arima.sim(list(order = c(0,0,1), ma = 0.9), n = 50)
acf(mal, plot=FALSE)[1] # = .507 (lag 1 sample ACF)
```

When  $|\hat{\rho}(1)| < \frac{1}{2}$ , the invertible estimate is

$$\hat{\theta} = \frac{1 - \sqrt{1 - 4\hat{\rho}(1)^2}}{2\hat{\rho}(1)}. \quad (3.105)$$

It can be shown that<sup>3.5</sup>

$$\hat{\theta} \sim \text{AN} \left( \theta, \frac{1 + \theta^2 + 4\theta^4 + \theta^6 + \theta^8}{n(1 - \theta^2)^2} \right);$$

AN is read *asymptotically normal* and is defined in [Definition A.5](#). The maximum likelihood estimator (which we discuss next) of  $\theta$ , in this case, has an asymptotic variance of  $(1 - \theta^2)/n$ . When  $\theta = .5$ , for example, the ratio of the asymptotic variance of the method of moments estimator to the maximum likelihood estimator of  $\theta$  is about 3.5. That is, for large samples, the variance of the method of moments estimator is about 3.5 times larger than the variance of the MLE of  $\theta$  when  $\theta = .5$ .

### MAXIMUM LIKELIHOOD AND LEAST SQUARES ESTIMATION

To fix ideas, we first focus on the causal AR(1) case. Let

$$x_t = \mu + \phi(x_{t-1} - \mu) + w_t \quad (3.106)$$

where  $|\phi| < 1$  and  $w_t \sim \text{iid } N(0, \sigma_w^2)$ . Given data  $x_1, x_2, \dots, x_n$ , we seek the likelihood

$$L(\mu, \phi, \sigma_w^2) = f(x_1, x_2, \dots, x_n \mid \mu, \phi, \sigma_w^2).$$

In the case of an AR(1), we may write the likelihood as

$$L(\mu, \phi, \sigma_w^2) = f(x_1)f(x_2 \mid x_1) \cdots f(x_n \mid x_{n-1}),$$

where we have dropped the parameters in the densities,  $f(\cdot)$ , to ease the notation. Because  $x_t \mid x_{t-1} \sim N(\mu + \phi(x_{t-1} - \mu), \sigma_w^2)$ , we have

$$f(x_t \mid x_{t-1}) = f_w[(x_t - \mu) - \phi(x_{t-1} - \mu)],$$

where  $f_w(\cdot)$  is the density of  $w_t$ , that is, the normal density with mean zero and variance  $\sigma_w^2$ . We may then write the likelihood as

$$L(\mu, \phi, \sigma_w) = f(x_1) \prod_{t=2}^n f_w[(x_t - \mu) - \phi(x_{t-1} - \mu)].$$

<sup>3.5</sup> The result follows from [Theorem A.7](#) and the delta method. See the proof of [Theorem A.7](#) for details on the delta method.

To find  $f(x_1)$ , we can use the causal representation

$$x_1 = \mu + \sum_{j=0}^{\infty} \phi^j w_{1-j}$$

to see that  $x_1$  is normal, with mean  $\mu$  and variance  $\sigma_w^2/(1-\phi^2)$ . Finally, for an AR(1), the likelihood is

$$L(\mu, \phi, \sigma_w^2) = (2\pi\sigma_w^2)^{-n/2} (1-\phi^2)^{1/2} \exp \left[ -\frac{S(\mu, \phi)}{2\sigma_w^2} \right], \quad (3.107)$$

where

$$S(\mu, \phi) = (1-\phi^2)(x_1 - \mu)^2 + \sum_{t=2}^n [(x_t - \mu) - \phi(x_{t-1} - \mu)]^2. \quad (3.108)$$

Typically,  $S(\mu, \phi)$  is called the *unconditional sum of squares*. We could have also considered the estimation of  $\mu$  and  $\phi$  using *unconditional least squares*, that is, estimation by minimizing  $S(\mu, \phi)$ .

Taking the partial derivative of the log of (3.107) with respect to  $\sigma_w^2$  and setting the result equal to zero, we get the typical normal result that for any given values of  $\mu$  and  $\phi$  in the parameter space,  $\sigma_w^2 = n^{-1}S(\mu, \phi)$  maximizes the likelihood. Thus, the maximum likelihood estimate of  $\sigma_w^2$  is

$$\hat{\sigma}_w^2 = n^{-1}S(\hat{\mu}, \hat{\phi}), \quad (3.109)$$

where  $\hat{\mu}$  and  $\hat{\phi}$  are the MLEs of  $\mu$  and  $\phi$ , respectively. If we replace  $n$  in (3.109) by  $n-2$ , we would obtain the unconditional least squares estimate of  $\sigma_w^2$ .

If, in (3.107), we take logs, replace  $\sigma_w^2$  by  $\hat{\sigma}_w^2$ , and ignore constants,  $\hat{\mu}$  and  $\hat{\phi}$  are the values that minimize the criterion function

$$l(\mu, \phi) = \log [n^{-1}S(\mu, \phi)] - n^{-1} \log(1-\phi^2); \quad (3.110)$$

that is,  $l(\mu, \phi) \propto -2 \log L(\mu, \phi, \hat{\sigma}_w^2)$ .<sup>3.6</sup> Because (3.108) and (3.110) are complicated functions of the parameters, the minimization of  $l(\mu, \phi)$  or  $S(\mu, \phi)$  is accomplished numerically. In the case of AR models, we have the advantage that, conditional on initial values, they are linear models. That is, we can drop the term in the likelihood that causes the nonlinearity. Conditioning on  $x_1$ , the *conditional likelihood* becomes

$$\begin{aligned} L(\mu, \phi, \sigma_w^2 \mid x_1) &= \prod_{t=2}^n f_w [(x_t - \mu) - \phi(x_{t-1} - \mu)] \\ &= (2\pi\sigma_w^2)^{-(n-1)/2} \exp \left[ -\frac{S_c(\mu, \phi)}{2\sigma_w^2} \right], \end{aligned} \quad (3.111)$$

where the *conditional sum of squares* is

<sup>3.6</sup> The criterion function is sometimes called the profile or concentrated likelihood.

$$S_c(\mu, \phi) = \sum_{t=2}^n [(x_t - \mu) - \phi(x_{t-1} - \mu)]^2. \quad (3.112)$$

The conditional MLE of  $\sigma_w^2$  is

$$\hat{\sigma}_w^2 = S_c(\hat{\mu}, \hat{\phi}) / (n - 1), \quad (3.113)$$

and  $\hat{\mu}$  and  $\hat{\phi}$  are the values that minimize the conditional sum of squares,  $S_c(\mu, \phi)$ . Letting  $\alpha = \mu(1 - \phi)$ , the conditional sum of squares can be written as

$$S_c(\mu, \phi) = \sum_{t=2}^n [x_t - (\alpha + \phi x_{t-1})]^2. \quad (3.114)$$

The problem is now the linear regression problem stated in [Section 2.1](#). Following the results from least squares estimation, we have  $\hat{\alpha} = \bar{x}_{(2)} - \hat{\phi}\bar{x}_{(1)}$ , where  $\bar{x}_{(1)} = (n - 1)^{-1} \sum_{t=1}^{n-1} x_t$ , and  $\bar{x}_{(2)} = (n - 1)^{-1} \sum_{t=2}^n x_t$ , and the conditional estimates are then

$$\hat{\mu} = \frac{\bar{x}_{(2)} - \hat{\phi}\bar{x}_{(1)}}{1 - \hat{\phi}} \quad (3.115)$$

$$\hat{\phi} = \frac{\sum_{t=2}^n (x_t - \bar{x}_{(2)})(x_{t-1} - \bar{x}_{(1)})}{\sum_{t=2}^n (x_{t-1} - \bar{x}_{(1)})^2}. \quad (3.116)$$

From (3.115) and (3.116), we see that  $\hat{\mu} \approx \bar{x}$  and  $\hat{\phi} \approx \hat{\rho}(1)$ . That is, the Yule–Walker estimators and the conditional least squares estimators are approximately the same. The only difference is the inclusion or exclusion of terms involving the endpoints,  $x_1$  and  $x_n$ . We can also adjust the estimate of  $\sigma_w^2$  in (3.113) to be equivalent to the least squares estimator, that is, divide  $S_c(\hat{\mu}, \hat{\phi})$  by  $(n - 3)$  instead of  $(n - 1)$  in (3.113).

For general AR( $p$ ) models, maximum likelihood estimation, unconditional least squares, and conditional least squares follow analogously to the AR(1) example. For general ARMA models, it is difficult to write the likelihood as an explicit function of the parameters. Instead, it is advantageous to write the likelihood in terms of the *innovations*, or one-step-ahead prediction errors,  $x_t - x_t^{t-1}$ . This will also be useful in [Chapter 6](#) when we study state-space models.

For a normal ARMA( $p, q$ ) model, let  $\beta = (\mu, \phi_1, \dots, \phi_p, \theta_1, \dots, \theta_q)'$  be the  $(p + q + 1)$ -dimensional vector of the model parameters. The likelihood can be written as

$$L(\beta, \sigma_w^2) = \prod_{t=1}^n f(x_t \mid x_{t-1}, \dots, x_1).$$

The conditional distribution of  $x_t$  given  $x_{t-1}, \dots, x_1$  is Gaussian with mean  $x_t^{t-1}$  and variance  $P_t^{t-1}$ . Recall from (3.71) that  $P_t^{t-1} = \gamma(0) \prod_{j=1}^{t-1} (1 - \phi_{jj}^2)$ . For ARMA models,  $\gamma(0) = \sigma_w^2 \sum_{j=0}^{\infty} \psi_j^2$ , in which case we may write

$$P_t^{t-1} = \sigma_w^2 \left\{ \left[ \sum_{j=0}^{\infty} \psi_j^2 \right] \left[ \prod_{j=1}^{t-1} (1 - \phi_{jj}^2) \right] \right\} \stackrel{\text{def}}{=} \sigma_w^2 r_t,$$

where  $r_t$  is the term in the braces. Note that the  $r_t$  terms are functions only of the regression parameters and that they may be computed recursively as  $r_{t+1} = (1 - \phi_{tt}^2)r_t$  with initial condition  $r_1 = \sum_{j=0}^{\infty} \psi_j^2$ . The likelihood of the data can now be written as

$$L(\beta, \sigma_w^2) = (2\pi\sigma_w^2)^{-n/2} [r_1(\beta)r_2(\beta) \cdots r_n(\beta)]^{-1/2} \exp \left[ -\frac{S(\beta)}{2\sigma_w^2} \right], \quad (3.117)$$

where

$$S(\beta) = \sum_{t=1}^n \left[ \frac{(x_t - x_t^{t-1}(\beta))^2}{r_t(\beta)} \right]. \quad (3.118)$$

Both  $x_t^{t-1}$  and  $r_t$  are functions of  $\beta$  alone, and we make that fact explicit in (3.117)-(3.118). Given values for  $\beta$  and  $\sigma_w^2$ , the likelihood may be evaluated using the techniques of Section 3.4. Maximum likelihood estimation would now proceed by maximizing (3.117) with respect to  $\beta$  and  $\sigma_w^2$ . As in the AR(1) example, we have

$$\hat{\sigma}_w^2 = n^{-1} S(\hat{\beta}), \quad (3.119)$$

where  $\hat{\beta}$  is the value of  $\beta$  that minimizes the concentrated likelihood

$$l(\beta) = \log [n^{-1} S(\beta)] + n^{-1} \sum_{t=1}^n \log r_t(\beta). \quad (3.120)$$

For the AR(1) model (3.106) discussed previously, recall that  $x_1^0 = \mu$  and  $x_t^{t-1} = \mu + \phi(x_{t-1} - \mu)$ , for  $t = 2, \dots, n$ . Also, using the fact that  $\phi_{11} = \phi$  and  $\phi_{hh} = 0$  for  $h > 1$ , we have  $r_1 = \sum_{j=0}^{\infty} \phi^{2j} = (1 - \phi^2)^{-1}$ ,  $r_2 = (1 - \phi^2)^{-1}(1 - \phi^2) = 1$ , and in general,  $r_t = 1$  for  $t = 2, \dots, n$ . Hence, the likelihood presented in (3.107) is identical to the innovations form of the likelihood given by (3.117). Moreover, the generic  $S(\beta)$  in (3.118) is  $S(\mu, \phi)$  given in (3.108) and the generic  $l(\beta)$  in (3.120) is  $l(\mu, \phi)$  in (3.110).

Unconditional least squares would be performed by minimizing (3.118) with respect to  $\beta$ . Conditional least squares estimation would involve minimizing (3.118) with respect to  $\beta$  but where, to ease the computational burden, the predictions and their errors are obtained by conditioning on initial values of the data. In general, numerical optimization routines are used to obtain the actual estimates and their standard errors.

### Example 3.30 The Newton–Raphson and Scoring Algorithms

Two common numerical optimization routines for accomplishing maximum likelihood estimation are Newton–Raphson and scoring. We will give a brief account of the mathematical ideas here. The actual implementation of these algorithms is much more complicated than our discussion might imply. For details, the reader is referred to any of the *Numerical Recipes* books, for example, Press et al. (1993).

Let  $l(\beta)$  be a criterion function of  $k$  parameters  $\beta = (\beta_1, \dots, \beta_k)$  that we wish to minimize with respect to  $\beta$ . For example, consider the likelihood function given by (3.110) or by (3.120). Suppose  $l(\hat{\beta})$  is the extremum that we are interested in

finding, and  $\hat{\beta}$  is found by solving  $\partial l(\beta)/\partial \beta_j = 0$ , for  $j = 1, \dots, k$ . Let  $l^{(1)}(\beta)$  denote the  $k \times 1$  vector of partials

$$l^{(1)}(\beta) = \left( \frac{\partial l(\beta)}{\partial \beta_1}, \dots, \frac{\partial l(\beta)}{\partial \beta_k} \right)'.$$

Note,  $l^{(1)}(\hat{\beta}) = 0$ , the  $k \times 1$  zero vector. Let  $l^{(2)}(\beta)$  denote the  $k \times k$  matrix of second-order partials

$$l^{(2)}(\beta) = \left\{ -\frac{\partial^2 l(\beta)}{\partial \beta_i \partial \beta_j} \right\}_{i,j=1}^k,$$

and assume  $l^{(2)}(\beta)$  is nonsingular. Let  $\beta_{(0)}$  be a “sufficiently good” initial estimator of  $\beta$ . Then, using a Taylor expansion, we have the following approximation:

$$0 = l^{(1)}(\hat{\beta}) \approx l^{(1)}(\beta_{(0)}) - l^{(2)}(\beta_{(0)}) [\hat{\beta} - \beta_{(0)}].$$

Setting the right-hand side equal to zero and solving for  $\hat{\beta}$  [call the solution  $\beta_{(1)}$ ], we get

$$\beta_{(1)} = \beta_{(0)} + \left[ l^{(2)}(\beta_{(0)}) \right]^{-1} l^{(1)}(\beta_{(0)}).$$

The Newton–Raphson algorithm proceeds by iterating this result, replacing  $\beta_{(0)}$  by  $\beta_{(1)}$  to get  $\beta_{(2)}$ , and so on, until convergence. Under a set of appropriate conditions, the sequence of estimators,  $\beta_{(1)}, \beta_{(2)}, \dots$ , will converge to  $\hat{\beta}$ , the MLE of  $\beta$ .

For maximum likelihood estimation, the criterion function used is  $l(\beta)$  given by (3.120);  $l^{(1)}(\beta)$  is called the score vector, and  $l^{(2)}(\beta)$  is called the *Hessian*. In the method of scoring, we replace  $l^{(2)}(\beta)$  by  $E[l^{(2)}(\beta)]$ , the *information* matrix. Under appropriate conditions, the inverse of the information matrix is the asymptotic variance–covariance matrix of the estimator  $\hat{\beta}$ . This is sometimes approximated by the inverse of the Hessian at  $\hat{\beta}$ . If the derivatives are difficult to obtain, it is possible to use quasi-maximum likelihood estimation where numerical techniques are used to approximate the derivatives.

### Example 3.31 MLE for the Recruitment Series

So far, we have fit an AR(2) model to the Recruitment series using ordinary least squares (Example 3.18) and using Yule–Walker (Example 3.28). The following is an R session used to fit an AR(2) model via maximum likelihood estimation to the Recruitment series; these results can be compared to the results in Example 3.18 and Example 3.28.

```
rec.mle = ar.mle(rec, order=2)
rec.mle$x.mean # 62.26
rec.mle$ar      # 1.35, -.46
sqrt(diag(rec.mle$asy.var.coef)) # .04, .04
rec.mle$var.pred # 89.34
```

We now discuss least squares for ARMA( $p, q$ ) models via *Gauss–Newton*. For general and complete details of the Gauss–Newton procedure, the reader is referred

to Fuller (1996). As before, write  $\beta = (\phi_1, \dots, \phi_p, \theta_1, \dots, \theta_q)'$ , and for the ease of discussion, we will put  $\mu = 0$ . We write the model in terms of the errors

$$w_t(\beta) = x_t - \sum_{j=1}^p \phi_j x_{t-j} - \sum_{k=1}^q \theta_k w_{t-k}(\beta), \quad (3.121)$$

emphasizing the dependence of the errors on the parameters.

For conditional least squares, we approximate the residual sum of squares by conditioning on  $x_1, \dots, x_p$  (if  $p > 0$ ) and  $w_p = w_{p-1} = w_{p-2} = \dots = w_{1-q} = 0$  (if  $q > 0$ ), in which case, given  $\beta$ , we may evaluate (3.121) for  $t = p+1, p+2, \dots, n$ . Using this conditioning argument, the conditional error sum of squares is

$$S_c(\beta) = \sum_{t=p+1}^n w_t^2(\beta). \quad (3.122)$$

Minimizing  $S_c(\beta)$  with respect to  $\beta$  yields the conditional least squares estimates. If  $q = 0$ , the problem is linear regression and no iterative technique is needed to minimize  $S_c(\phi_1, \dots, \phi_p)$ . If  $q > 0$ , the problem becomes nonlinear regression and we will have to rely on numerical optimization.

When  $n$  is large, conditioning on a few initial values will have little influence on the final parameter estimates. In the case of small to moderate sample sizes, one may wish to rely on unconditional least squares. The unconditional least squares problem is to choose  $\beta$  to minimize the unconditional sum of squares, which we have generically denoted by  $S(\beta)$  in this section. The unconditional sum of squares can be written in various ways, and one useful form in the case of ARMA( $p, q$ ) models is derived in Box et al. (1994, Appendix A7.3). They showed (see Problem 3.19) the unconditional sum of squares can be written as

$$S(\beta) = \sum_{t=-\infty}^n \tilde{w}_t^2(\beta), \quad (3.123)$$

where  $\tilde{w}_t(\beta) = E(w_t \mid x_1, \dots, x_n)$ . When  $t \leq 0$ , the  $\tilde{w}_t(\beta)$  are obtained by backcasting. As a practical matter, we approximate  $S(\beta)$  by starting the sum at  $t = -M+1$ , where  $M$  is chosen large enough to guarantee  $\sum_{t=-\infty}^{-M} \tilde{w}_t^2(\beta) \approx 0$ . In the case of unconditional least squares estimation, a numerical optimization technique is needed even when  $q = 0$ .

To employ Gauss–Newton, let  $\beta_{(0)} = (\phi_1^{(0)}, \dots, \phi_p^{(0)}, \theta_1^{(0)}, \dots, \theta_q^{(0)})'$  be an initial estimate of  $\beta$ . For example, we could obtain  $\beta_{(0)}$  by method of moments. The first-order Taylor expansion of  $w_t(\beta)$  is

$$w_t(\beta) \approx w_t(\beta_{(0)}) - (\beta - \beta_{(0)})' z_t(\beta_{(0)}), \quad (3.124)$$

where

$$z_t'(\beta_{(0)}) = \left( -\frac{\partial w_t(\beta)}{\partial \beta_1}, \dots, -\frac{\partial w_t(\beta)}{\partial \beta_{p+q}} \right) \bigg|_{\beta=\beta_{(0)}}, \quad t = 1, \dots, n.$$

The linear approximation of  $S_c(\beta)$  is

$$Q(\beta) = \sum_{t=p+1}^n [w_t(\beta_{(0)}) - (\beta - \beta_{(0)})' z_t(\beta_{(0)})]^2 \quad (3.125)$$

and this is the quantity that we will minimize. For approximate unconditional least squares, we would start the sum in (3.125) at  $t = -M + 1$ , for a large value of  $M$ , and work with the backcasted values.

Using the results of ordinary least squares (Section 2.1), we know

$$\widehat{(\beta - \beta_{(0)})} = \left( n^{-1} \sum_{t=p+1}^n z_t(\beta_{(0)}) z_t'(\beta_{(0)}) \right)^{-1} \left( n^{-1} \sum_{t=p+1}^n z_t(\beta_{(0)}) w_t(\beta_{(0)}) \right) \quad (3.126)$$

minimizes  $Q(\beta)$ . From (3.126), we write the *one-step Gauss–Newton estimate* as

$$\beta_{(1)} = \beta_{(0)} + \Delta(\beta_{(0)}), \quad (3.127)$$

where  $\Delta(\beta_{(0)})$  denotes the right-hand side of (3.126). Gauss–Newton estimation is accomplished by replacing  $\beta_{(0)}$  by  $\beta_{(1)}$  in (3.127). This process is repeated by calculating, at iteration  $j = 2, 3, \dots$ ,

$$\beta_{(j)} = \beta_{(j-1)} + \Delta(\beta_{(j-1)})$$

until convergence.

### Example 3.32 Gauss–Newton for an MA(1)

Consider an invertible MA(1) process,  $x_t = w_t + \theta w_{t-1}$ . Write the truncated errors as

$$w_t(\theta) = x_t - \theta w_{t-1}(\theta), \quad t = 1, \dots, n, \quad (3.128)$$

where we condition on  $w_0(\theta) = 0$ . Taking derivatives and negating,

$$-\frac{\partial w_t(\theta)}{\partial \theta} = w_{t-1}(\theta) + \theta \frac{\partial w_{t-1}(\theta)}{\partial \theta}, \quad t = 1, \dots, n, \quad (3.129)$$

where  $\partial w_0(\theta)/\partial \theta = 0$ . We can also write (3.129) as

$$z_t(\theta) = w_{t-1}(\theta) - \theta z_{t-1}(\theta), \quad t = 1, \dots, n, \quad (3.130)$$

where  $z_t(\theta) = -\partial w_t(\theta)/\partial \theta$  and  $z_0(\theta) = 0$ .

Let  $\theta_{(0)}$  be an initial estimate of  $\theta$ , for example, the estimate given in Example 3.29. Then, the Gauss–Newton procedure for conditional least squares is given by

$$\theta_{(j+1)} = \theta_{(j)} + \frac{\sum_{t=1}^n z_t(\theta_{(j)}) w_t(\theta_{(j)})}{\sum_{t=1}^n z_t^2(\theta_{(j)})}, \quad j = 0, 1, 2, \dots, \quad (3.131)$$

where the values in (3.131) are calculated recursively using (3.128) and (3.130). The calculations are stopped when  $|\theta_{(j+1)} - \theta_{(j)}|$ , or  $|Q(\theta_{(j+1)}) - Q(\theta_{(j)})|$ , are smaller than some preset amount.



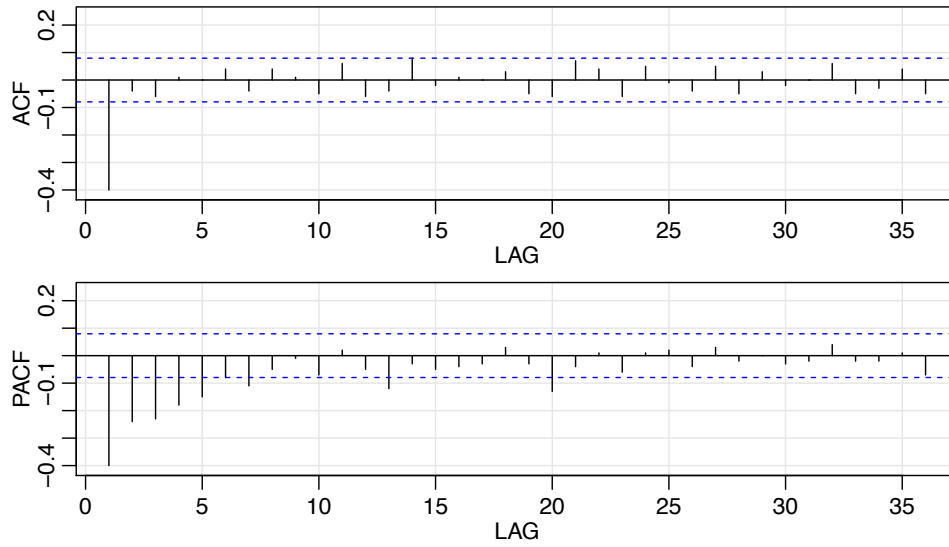


Fig. 3.9. ACF and PACF of transformed glacial varves.

### Example 3.33 Fitting the Glacial Varve Series

Consider the series of glacial varve thicknesses from Massachusetts for  $n = 634$  years, as analyzed in Example 2.7 and in Problem 2.8, where it was argued that a first-order moving average model might fit the logarithmically transformed and differenced varve series, say,

$$\nabla \log(x_t) = \log(x_t) - \log(x_{t-1}) = \log\left(\frac{x_t}{x_{t-1}}\right),$$

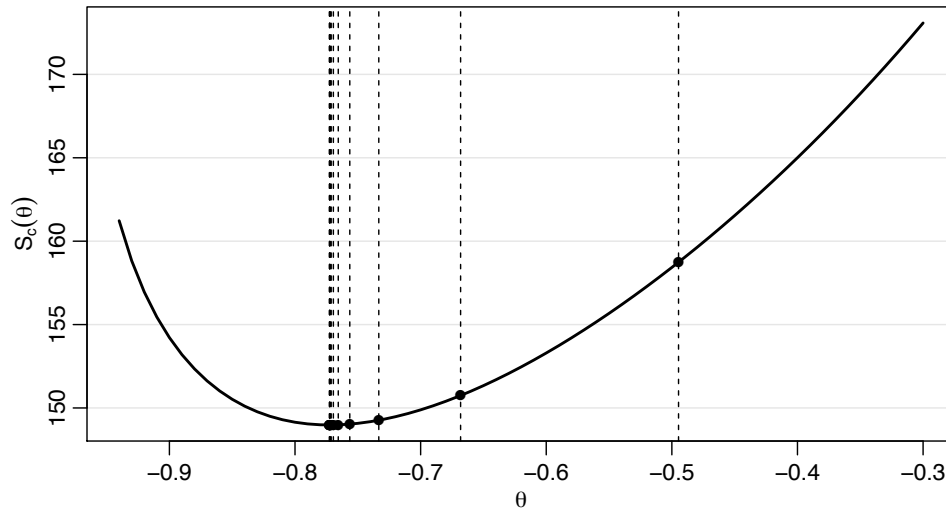
which can be interpreted as being approximately the percentage change in the thickness.

The sample ACF and PACF, shown in Figure 3.9, confirm the tendency of  $\nabla \log(x_t)$  to behave as a first-order moving average process as the ACF has only a significant peak at lag one and the PACF decreases exponentially. Using Table 3.1, this sample behavior fits that of the MA(1) very well.

Since  $\hat{\rho}(1) = -.397$ , our initial estimate is  $\theta_{(0)} = -.495$  using (3.105). The results of eleven iterations of the Gauss–Newton procedure, (3.131), starting with  $\theta_{(0)}$  are given in Table 3.2. The final estimate is  $\hat{\theta} = \theta_{(11)} = -.773$ ; interim values and the corresponding value of the conditional sum of squares,  $S_c(\theta)$  given in (3.122), are also displayed in the table. The final estimate of the error variance is  $\hat{\sigma}_w^2 = 148.98/632 = .236$  with 632 degrees of freedom (one is lost in differencing). The value of the sum of the squared derivatives at convergence is  $\sum_{t=1}^n z_t^2(\theta_{(11)}) = 368.741$ , and consequently, the estimated standard error of  $\hat{\theta}$  is  $\sqrt{.236/368.741} = .025$ ,<sup>3.7</sup> this leads to a  $t$ -value of  $-.773/.025 = -30.92$  with 632 degrees of freedom.

Figure 3.10 displays the conditional sum of squares,  $S_c(\theta)$  as a function of  $\theta$ , as well as indicating the values of each step of the Gauss–Newton algorithm. Note

<sup>3.7</sup> To estimate the standard error, we are using the standard regression results from (2.6) as an approximation



**Fig. 3.10.** Conditional sum of squares versus values of the moving average parameter for the glacial varve example, [Example 3.33](#). Vertical lines indicate the values of the parameter obtained via Gauss–Newton; see [Table 3.2](#) for the actual values.

that the Gauss–Newton procedure takes large steps toward the minimum initially, and then takes very small steps as it gets close to the minimizing value. When there is only one parameter, as in this case, it would be easy to evaluate  $S_c(\theta)$  on a grid of points, and then choose the appropriate value of  $\theta$  from the grid search. It would be difficult, however, to perform grid searches when there are many parameters.

The following code was used in this example.

```
x = diff(log(varve))
# Evaluate Sc on a Grid
c(0) -> w -> z
c() -> Sc -> Sz -> Szw
num = length(x)
th = seq(-.3, -.94, -.01)
for (p in 1:length(th)){
  for (i in 2:num){ w[i] = x[i]-th[p]*w[i-1] }
  Sc[p] = sum(w^2) }
plot(th, Sc, type="l", ylab=expression(S[c](theta)), xlab=expression(theta),
      lwd=2)
# Gauss-Newton Estimation
r = acf(x, lag=1, plot=FALSE)$acf[-1]
rstart = (1-sqrt(1-4*(r^2)))/(2*r) # from (3.105)
c(0) -> w -> z
c() -> Sc -> Sz -> Szw -> para
niter = 12
para[1] = rstart
for (p in 1:niter){
  for (i in 2:num){ w[i] = x[i]-para[p]*w[i-1]
                    z[i] = w[i-1]-para[p]*z[i-1] }
  Sc[p] = sum(w^2)
  Sz[p] = sum(z^2)
  Szw[p] = sum(z*w)
  para[p+1] = para[p] + Szw[p]/Sz[p] }
```

**Table 3.2.** Gauss–Newton Results for *Example 3.33*

$j$	$\theta_{(j)}$	$S_c(\theta_{(j)})$	$\sum_{t=1}^n z_t^2(\theta_{(j)})$
0	-0.495	158.739	171.240
1	-0.668	150.747	235.266
2	-0.733	149.264	300.562
3	-0.756	149.031	336.823
4	-0.766	148.990	354.173
5	-0.769	148.982	362.167
6	-0.771	148.980	365.801
7	-0.772	148.980	367.446
8	-0.772	148.980	368.188
9	-0.772	148.980	368.522
10	-0.773	148.980	368.673
11	-0.773	148.980	368.741

```
round(cbind(iteration=0:(niter-1), thetahat=para[1:niter] , Sc , Sz ), 3)
abline(v = para[1:12], lty=2)
points(para[1:12], Sc[1:12], pch=16)
```

In the general case of causal and invertible ARMA( $p, q$ ) models, maximum likelihood estimation and conditional and unconditional least squares estimation (and Yule–Walker estimation in the case of AR models) all lead to optimal estimators. The proof of this general result can be found in a number of texts on theoretical time series analysis (for example, Brockwell and Davis, 1991, or Hannan, 1970, to mention a few). We will denote the ARMA coefficient parameters by  $\beta = (\phi_1, \dots, \phi_p, \theta_1, \dots, \theta_q)'$ .

### Property 3.10 Large Sample Distribution of the Estimators

Under appropriate conditions, for causal and invertible ARMA processes, the maximum likelihood, the unconditional least squares, and the conditional least squares estimators, each initialized by the method of moments estimator, all provide optimal estimators of  $\sigma_w^2$  and  $\beta$ , in the sense that  $\hat{\sigma}_w^2$  is consistent, and the asymptotic distribution of  $\hat{\beta}$  is the best asymptotic normal distribution. In particular, as  $n \rightarrow \infty$ ,

$$\sqrt{n}(\hat{\beta} - \beta) \xrightarrow{d} N(0, \sigma_w^2 \Gamma_{p,q}^{-1}). \quad (3.132)$$

The asymptotic variance–covariance matrix of the estimator  $\hat{\beta}$  is the inverse of the information matrix. In particular, the  $(p+q) \times (p+q)$  matrix  $\Gamma_{p,q}$ , has the form

$$\Gamma_{p,q} = \begin{pmatrix} \Gamma_{\phi\phi} & \Gamma_{\phi\theta} \\ \Gamma_{\theta\phi} & \Gamma_{\theta\theta} \end{pmatrix}. \quad (3.133)$$

The  $p \times p$  matrix  $\Gamma_{\phi\phi}$  is given by (3.100), that is, the  $ij$ -th element of  $\Gamma_{\phi\phi}$ , for  $i, j = 1, \dots, p$ , is  $\gamma_x(i-j)$  from an AR( $p$ ) process,  $\phi(B)x_t = w_t$ . Similarly,  $\Gamma_{\theta\theta}$  is a  $q \times q$  matrix with the  $ij$ -th element, for  $i, j = 1, \dots, q$ , equal to  $\gamma_y(i-j)$  from an AR( $q$ ) process,  $\theta(B)y_t = w_t$ . The  $p \times q$  matrix  $\Gamma_{\phi\theta} = \{\gamma_{xy}(i-j)\}$ , for  $i = 1, \dots, p$ ;  $j = 1, \dots, q$ ; that is, the  $ij$ -th element is the cross-covariance between the two AR processes given by  $\phi(B)x_t = w_t$  and  $\theta(B)y_t = w_t$ . Finally,  $\Gamma_{\theta\phi} = \Gamma'_{\phi\theta}$  is  $q \times p$ .

Further discussion of **Property 3.10**, including a proof for the case of least squares estimators for AR( $p$ ) processes, can be found in **Section B.3**.

### Example 3.34 Some Specific Asymptotic Distributions

The following are some specific cases of **Property 3.10**.

**AR(1):**  $\gamma_x(0) = \sigma_w^2 / (1 - \phi^2)$ , so  $\sigma_w^2 \Gamma_{1,0}^{-1} = (1 - \phi^2)$ . Thus,

$$\hat{\phi} \sim \text{AN} \left[ \phi, n^{-1}(1 - \phi^2) \right]. \quad (3.134)$$

**AR(2):** The reader can verify that

$$\gamma_x(0) = \left( \frac{1 - \phi_2}{1 + \phi_2} \right) \frac{\sigma_w^2}{(1 - \phi_2)^2 - \phi_1^2}$$

and  $\gamma_x(1) = \phi_1 \gamma_x(0) + \phi_2 \gamma_x(1)$ . From these facts, we can compute  $\Gamma_{2,0}^{-1}$ . In particular, we have

$$\begin{pmatrix} \hat{\phi}_1 \\ \hat{\phi}_2 \end{pmatrix} \sim \text{AN} \left[ \begin{pmatrix} \phi_1 \\ \phi_2 \end{pmatrix}, n^{-1} \begin{pmatrix} 1 - \phi_2^2 & -\phi_1(1 + \phi_2) \\ \text{sym} & 1 - \phi_2^2 \end{pmatrix} \right]. \quad (3.135)$$

**MA(1):** In this case, write  $\theta(B)y_t = w_t$ , or  $y_t + \theta y_{t-1} = w_t$ . Then, analogous to the AR(1) case,  $\gamma_y(0) = \sigma_w^2 / (1 - \theta^2)$ , so  $\sigma_w^2 \Gamma_{0,1}^{-1} = (1 - \theta^2)$ . Thus,

$$\hat{\theta} \sim \text{AN} \left[ \theta, n^{-1}(1 - \theta^2) \right]. \quad (3.136)$$

**MA(2):** Write  $y_t + \theta_1 y_{t-1} + \theta_2 y_{t-2} = w_t$ , so, analogous to the AR(2) case, we have

$$\begin{pmatrix} \hat{\theta}_1 \\ \hat{\theta}_2 \end{pmatrix} \sim \text{AN} \left[ \begin{pmatrix} \theta_1 \\ \theta_2 \end{pmatrix}, n^{-1} \begin{pmatrix} 1 - \theta_2^2 & \theta_1(1 + \theta_2) \\ \text{sym} & 1 - \theta_2^2 \end{pmatrix} \right]. \quad (3.137)$$

**ARMA(1,1):** To calculate  $\Gamma_{\phi\theta}$ , we must find  $\gamma_{xy}(0)$ , where  $x_t - \phi x_{t-1} = w_t$  and  $y_t + \theta y_{t-1} = w_t$ . We have

$$\begin{aligned} \gamma_{xy}(0) &= \text{cov}(x_t, y_t) = \text{cov}(\phi x_{t-1} + w_t, -\theta y_{t-1} + w_t) \\ &= -\phi\theta \gamma_{xy}(0) + \sigma_w^2. \end{aligned}$$

Solving, we find,  $\gamma_{xy}(0) = \sigma_w^2 / (1 + \phi\theta)$ . Thus,

$$\begin{pmatrix} \hat{\phi} \\ \hat{\theta} \end{pmatrix} \sim \text{AN} \left[ \begin{pmatrix} \phi \\ \theta \end{pmatrix}, n^{-1} \begin{bmatrix} (1 - \phi^2)^{-1} & (1 + \phi\theta)^{-1} \\ \text{sym} & (1 - \theta^2)^{-1} \end{bmatrix}^{-1} \right]. \quad (3.138)$$

### Example 3.35 Overfitting Caveat

The asymptotic behavior of the parameter estimators gives us an additional insight into the problem of fitting ARMA models to data. For example, suppose a time series follows an AR(1) process and we decide to fit an AR(2) to the data. Do any problems occur in doing this? More generally, why not simply fit large-order AR models to make sure that we capture the dynamics of the process? After all,

if the process is truly an AR(1), the other autoregressive parameters will not be significant. The answer is that if we *overfit*, we obtain less efficient, or less precise parameter estimates. For example, if we fit an AR(1) to an AR(1) process, for large  $n$ ,  $\text{var}(\hat{\phi}_1) \approx n^{-1}(1 - \phi_1^2)$ . But, if we fit an AR(2) to the AR(1) process, for large  $n$ ,  $\text{var}(\hat{\phi}_1) \approx n^{-1}(1 - \phi_2^2) = n^{-1}$  because  $\phi_2 = 0$ . Thus, the variance of  $\phi_1$  has been inflated, making the estimator less precise.

We do want to mention, however, that overfitting can be used as a diagnostic tool. For example, if we fit an AR(2) model to the data and are satisfied with that model, then adding one more parameter and fitting an AR(3) should lead to approximately the same model as in the AR(2) fit. We will discuss model diagnostics in more detail in [Section 3.7](#).

The reader might wonder, for example, why the asymptotic distributions of  $\hat{\phi}$  from an AR(1) and  $\hat{\theta}$  from an MA(1) are of the same form; compare (3.134) to (3.136). It is possible to explain this unexpected result heuristically using the intuition of linear regression. That is, for the normal regression model presented in [Section 2.1](#) with no intercept term,  $x_t = \beta z_t + w_t$ , we know  $\hat{\beta}$  is normally distributed with mean  $\beta$ , and from (2.6),

$$\text{var} \left\{ \sqrt{n} (\hat{\beta} - \beta) \right\} = n \sigma_w^2 \left( \sum_{t=1}^n z_t^2 \right)^{-1} = \sigma_w^2 \left( n^{-1} \sum_{t=1}^n z_t^2 \right)^{-1}.$$

For the causal AR(1) model given by  $x_t = \phi x_{t-1} + w_t$ , the intuition of regression tells us to expect that, for  $n$  large,

$$\sqrt{n} (\hat{\phi} - \phi)$$

is approximately normal with mean zero and with variance given by

$$\sigma_w^2 \left( n^{-1} \sum_{t=2}^n x_{t-1}^2 \right)^{-1}.$$

Now,  $n^{-1} \sum_{t=2}^n x_{t-1}^2$  is the sample variance (recall that the mean of  $x_t$  is zero) of the  $x_t$ , so as  $n$  becomes large we would expect it to approach  $\text{var}(x_t) = \gamma(0) = \sigma_w^2 / (1 - \phi^2)$ . Thus, the large sample variance of  $\sqrt{n} (\hat{\phi} - \phi)$  is

$$\sigma_w^2 \gamma_x(0)^{-1} = \sigma_w^2 \left( \frac{\sigma_w^2}{1 - \phi^2} \right)^{-1} = (1 - \phi^2);$$

that is, (3.134) holds.

In the case of an MA(1), we may use the discussion of [Example 3.32](#) to write an approximate regression model for the MA(1). That is, consider the approximation (3.130) as the regression model

$$z_t(\hat{\theta}) = -\theta z_{t-1}(\hat{\theta}) + w_{t-1},$$

where now,  $z_{t-1}(\hat{\theta})$  as defined in [Example 3.32](#), plays the role of the regressor. Continuing with the analogy, we would expect the asymptotic distribution of  $\sqrt{n}(\hat{\theta} - \theta)$  to be normal, with mean zero, and approximate variance

$$\sigma_w^2 \left( n^{-1} \sum_{t=2}^n z_{t-1}^2(\hat{\theta}) \right)^{-1}.$$

As in the AR(1) case,  $n^{-1} \sum_{t=2}^n z_{t-1}^2(\hat{\theta})$  is the sample variance of the  $z_t(\hat{\theta})$  so, for large  $n$ , this should be  $\text{var}\{z_t(\theta)\} = \gamma_z(0)$ , say. But note, as seen from [\(3.130\)](#),  $z_t(\theta)$  is approximately an AR(1) process with parameter  $-\theta$ . Thus,

$$\sigma_w^2 \gamma_z(0)^{-1} = \sigma_w^2 \left( \frac{\sigma_w^2}{1 - (-\theta)^2} \right)^{-1} = (1 - \theta^2),$$

which agrees with [\(3.136\)](#). Finally, the asymptotic distributions of the AR parameter estimates and the MA parameter estimates are of the same form because in the MA case, the “regressors” are the differential processes  $z_t(\theta)$  that have AR structure, and it is this structure that determines the asymptotic variance of the estimators. For a rigorous account of this approach for the general case, see Fuller (1996, Theorem 5.5.4).

In [Example 3.33](#), the estimated standard error of  $\hat{\theta}$  was .025. In that example, we used regression results to estimate the standard error as the square root of

$$n^{-1} \hat{\sigma}_w^2 \left( n^{-1} \sum_{t=1}^n z_t^2(\hat{\theta}) \right)^{-1} = \frac{\hat{\sigma}_w^2}{\sum_{t=1}^n z_t^2(\hat{\theta})},$$

where  $n = 632$ ,  $\hat{\sigma}_w^2 = .236$ ,  $\sum_{t=1}^n z_t^2(\hat{\theta}) = 368.74$  and  $\hat{\theta} = -.773$ . Using [\(3.136\)](#), we could have also calculated this value using the asymptotic approximation, the square root of  $(1 - (-.773)^2)/632$ , which is also .025.

If  $n$  is small, or if the parameters are close to the boundaries, the asymptotic approximations can be quite poor. The *bootstrap* can be helpful in this case; for a broad treatment of the bootstrap, see Efron and Tibshirani (1994). We discuss the case of an AR(1) here and leave the general discussion for [Chapter 6](#). For now, we give a simple example of the bootstrap for an AR(1) process.

### Example 3.36 Bootstrapping an AR(1)

We consider an AR(1) model with a regression coefficient near the boundary of causality and an error process that is symmetric but not normal. Specifically, consider the causal model

$$x_t = \mu + \phi(x_{t-1} - \mu) + w_t, \quad (3.139)$$

where  $\mu = 50$ ,  $\phi = .95$ , and  $w_t$  are iid double exponential (Laplace) with location zero, and scale parameter  $\beta = 2$ . The density of  $w_t$  is given by

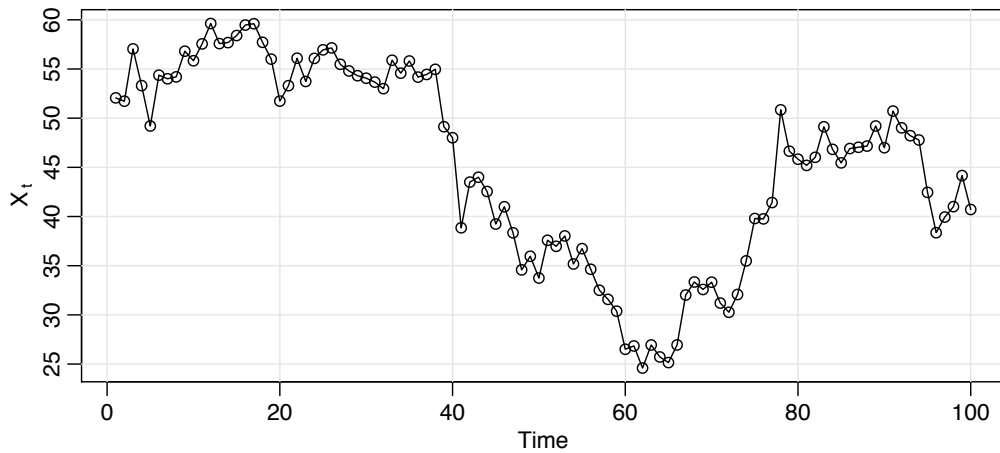


Fig. 3.11. One hundred observations generated from the model in Example 3.36.

$$f(w) = \frac{1}{2\beta} \exp\{-|w|/\beta\} \quad -\infty < w < \infty.$$

In this example,  $E(w_t) = 0$  and  $\text{var}(w_t) = 2\beta^2 = 8$ . Figure 3.11 shows  $n = 100$  simulated observations from this process. This particular realization is interesting: the data look like they were generated from a nonstationary process with three different mean levels. In fact, the data were generated from a well-behaved, albeit non-normal, stationary and causal model. To show the advantages of the bootstrap, we will act as if we do not know the actual error distribution. The data in Figure 3.11 were generated as follows.

```
set.seed(101010)
e = rexp(150, rate=.5); u = runif(150,-1,1); de = e*sign(u)
dex = 50 + arima.sim(n=100, list(ar=.95), innov=de, n.start=50)
plot.ts(dex, type='o', ylab=expression(X[~t]))
```

Using these data, we obtained the Yule–Walker estimates  $\hat{\mu} = 45.25$ ,  $\hat{\phi} = .96$ , and  $\hat{\sigma}_w^2 = 7.88$ , as follows.

```
fit = ar.yw(dex, order=1)
round(cbind(fit$x.mean, fit$ar, fit$var.pred), 2)
[1,] 45.25 0.96 7.88
```

To assess the finite sample distribution of  $\hat{\phi}$  when  $n = 100$ , we simulated 1000 realizations of this AR(1) process and estimated the parameters via Yule–Walker. The finite sampling density of the Yule–Walker estimate of  $\phi$ , based on the 1000 repeated simulations, is shown in Figure 3.12. Based on Property 3.10, we would say that  $\hat{\phi}$  is approximately normal with mean  $\phi$  (which we supposedly do not know) and variance  $(1 - \phi^2)/100$ , which we would approximate by  $(1 - .96^2)/100 = .032^2$ ; this distribution is superimposed on Figure 3.12. Clearly the sampling distribution is not close to normality for this sample size. The R code to perform the simulation is as follows. We use the results at the end of the example

```
set.seed(111)
phi.yw = rep(NA, 1000)
for (i in 1:1000){
```

```
e = rexp(150, rate=.5); u = runif(150,-1,1); de = e*sign(u)
x = 50 + arima.sim(n=100,list(ar=.95), innov=de, n.start=50)
phi.yw[i] = ar.yw(x, order=1)$ar }
```

The preceding simulation required full knowledge of the model, the parameter values and the noise distribution. Of course, in a sampling situation, we would not have the information necessary to do the preceding simulation and consequently would not be able to generate a figure like Figure 3.12. The bootstrap, however, gives us a way to attack the problem.

To simplify the discussion and the notation, we condition on  $x_1$  throughout the example. In this case, the one-step-ahead predictors have a simple form,

$$x_t^{t-1} = \mu + \phi(x_{t-1} - \mu), \quad t = 2, \dots, 100.$$

Consequently, the innovations,  $\epsilon_t = x_t - x_t^{t-1}$ , are given by

$$\epsilon_t = (x_t - \mu) - \phi(x_{t-1} - \mu), \quad t = 2, \dots, 100, \quad (3.140)$$

each with MSPE  $P_t^{t-1} = E(\epsilon_t^2) = E(w_t^2) = \sigma_w^2$  for  $t = 2, \dots, 100$ . We can use (3.140) to write the model in terms of the innovations,

$$x_t = x_t^{t-1} + \epsilon_t = \mu + \phi(x_{t-1} - \mu) + \epsilon_t \quad t = 2, \dots, 100. \quad (3.141)$$

To perform the bootstrap simulation, we replace the parameters with their estimates in (3.141), that is,  $\hat{\mu} = 45.25$  and  $\hat{\phi} = .96$ , and denote the resulting sample innovations as  $\{\hat{\epsilon}_2, \dots, \hat{\epsilon}_{100}\}$ . To obtain one bootstrap sample, first randomly sample, with replacement,  $n = 99$  values from the set of sample innovations; call the sampled values  $\{\epsilon_2^*, \dots, \epsilon_{100}^*\}$ . Now, generate a bootstrapped data set sequentially by setting

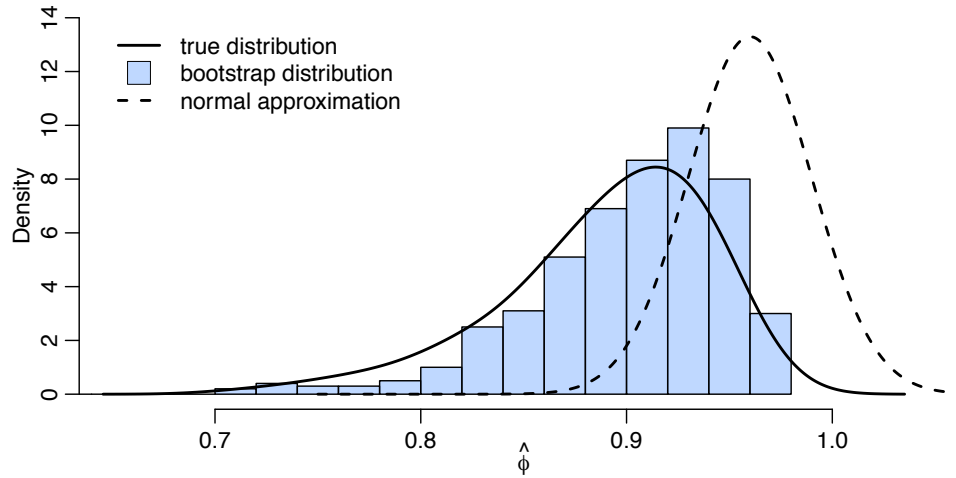
$$x_t^* = 45.25 + .96(x_{t-1}^* - 45.25) + \epsilon_t^*, \quad t = 2, \dots, 100. \quad (3.142)$$

with  $x_1^*$  held fixed at  $x_1$ . Next, estimate the parameters as if the data were  $x_t^*$ . Call these estimates  $\hat{\mu}(1)$ ,  $\hat{\phi}(1)$ , and  $\sigma_w^2(1)$ . Repeat this process a large number,  $B$ , of times, generating a collection of bootstrapped parameter estimates,  $\{\hat{\mu}(b), \hat{\phi}(b), \sigma_w^2(b); b = 1, \dots, B\}$ . We can then approximate the finite sample distribution of an estimator from the bootstrapped parameter values. For example, we can approximate the distribution of  $\hat{\phi} - \phi$  by the empirical distribution of  $\hat{\phi}(b) - \hat{\phi}$ , for  $b = 1, \dots, B$ .

Figure 3.12 shows the bootstrap histogram of 500 bootstrapped estimates of  $\phi$  using the data shown in Figure 3.11. Note that the bootstrap distribution of  $\hat{\phi}$  is close to the distribution of  $\hat{\phi}$  shown in Figure 3.12. The following code was used to perform the bootstrap.

```
set.seed(666) # not that 666
fit = ar.yw(dex, order=1) # assumes the data were retained
m = fit$x.mean # estimate of mean
phi = fit$ar # estimate of phi
nboot = 500 # number of bootstrap replicates
resids = fit$resid[-1] # the 99 innovations
```





**Fig. 3.12.** Finite sample density of the Yule-Walker estimate of  $\phi$  (solid line) in [Example 3.36](#) and the corresponding asymptotic normal density (dashed line). Bootstrap histogram of  $\hat{\phi}$  based on 500 bootstrapped samples.

```

x.star = dex                                # initialize x*
phi.star.yw = rep(NA, nboot)
# Bootstrap
for (i in 1:nboot) {
  resid.star = sample(resids, replace=TRUE)
  for (t in 1:99){ x.star[t+1] = m + phi*(x.star[t]-m) + resid.star[t] }
  phi.star.yw[i] = ar.yw(x.star, order=1)$ar
}
# Picture
culer = rgb(.5,.7,1,.5)
hist(phi.star.yw, 15, main="", prob=TRUE, xlim=c(.65,1.05), ylim=c(0,14),
      col=culer, xlab=expression(hat(phi)))
lines(density(phi.yw, bw=.02), lwd=2)      # from previous simulation
u = seq(.75, 1.1, by=.001)                # normal approximation
lines(u, dnorm(u, mean=.96, sd=.03), lty=2, lwd=2)
legend(.65, 14, legend=c('true distribution', 'bootstrap distribution',
  'normal approximation'), bty='n', lty=c(1,0,2), lwd=c(2,0,2),
  col=1, pch=c(NA,22,NA), pt.bg=c(NA,culer,NA), pt.cex=2.5)

```

## 3.6 Integrated Models for Nonstationary Data

In [Chapter 1](#) and [Chapter 2](#), we saw that if  $x_t$  is a random walk,  $x_t = x_{t-1} + w_t$ , then by differencing  $x_t$ , we find that  $\nabla x_t = w_t$  is stationary. In many situations, time series can be thought of as being composed of two components, a nonstationary trend component and a zero-mean stationary component. For example, in [Section 2.1](#) we considered the model

$$x_t = \mu_t + y_t, \quad (3.143)$$

where  $\mu_t = \beta_0 + \beta_1 t$  and  $y_t$  is stationary. Differencing such a process will lead to a stationary process:

$$\nabla x_t = x_t - x_{t-1} = \beta_1 + y_t - y_{t-1} = \beta_1 + \nabla y_t.$$

Another model that leads to first differencing is the case in which  $\mu_t$  in (3.143) is stochastic and slowly varying according to a random walk. That is,

$$\mu_t = \mu_{t-1} + v_t$$

where  $v_t$  is stationary. In this case,

$$\nabla x_t = v_t + \nabla y_t,$$

is stationary. If  $\mu_t$  in (3.143) is a  $k$ -th order polynomial,  $\mu_t = \sum_{j=0}^k \beta_j t^j$ , then (Problem 3.27) the differenced series  $\nabla^k x_t$  is stationary. Stochastic trend models can also lead to higher order differencing. For example, suppose

$$\mu_t = \mu_{t-1} + v_t \quad \text{and} \quad v_t = v_{t-1} + e_t,$$

where  $e_t$  is stationary. Then,  $\nabla x_t = v_t + \nabla y_t$  is not stationary, but

$$\nabla^2 x_t = e_t + \nabla^2 y_t$$

is stationary.

The *integrated* ARMA, or ARIMA, model is a broadening of the class of ARMA models to include differencing.

**Definition 3.11** A process  $x_t$  is said to be **ARIMA**( $p, d, q$ ) if

$$\nabla^d x_t = (1 - B)^d x_t$$

is ARMA( $p, q$ ). In general, we will write the model as

$$\phi(B)(1 - B)^d x_t = \theta(B)w_t. \quad (3.144)$$

If  $E(\nabla^d x_t) = \mu$ , we write the model as

$$\phi(B)(1 - B)^d x_t = \delta + \theta(B)w_t,$$

where  $\delta = \mu(1 - \phi_1 - \cdots - \phi_p)$ .

Because of the nonstationarity, care must be taken when deriving forecasts. For the sake of completeness, we discuss this issue briefly here, but we stress the fact that both the theoretical and computational aspects of the problem are best handled via state-space models. We discuss the theoretical details in [Chapter 6](#). For information on the state-space based computational aspects in R, see the ARIMA help files ([?arima](#) and [?predict.Arima](#)); our scripts [sarima](#) and [sarima.for](#) are basically wrappers for these R scripts.

It should be clear that, since  $y_t = \nabla^d x_t$  is ARMA, we can use [Section 3.4](#) methods to obtain forecasts of  $y_t$ , which in turn lead to forecasts for  $x_t$ . For example, if  $d = 1$ , given forecasts  $y_{n+m}^n$  for  $m = 1, 2, \dots$ , we have  $y_{n+m}^n = x_{n+m}^n - x_{n+m-1}^n$ , so that

$$x_{n+m}^n = y_{n+m}^n + x_{n+m-1}^n$$

with initial condition  $x_{n+1}^n = y_{n+1}^n + x_n$  (noting  $x_n^n = x_n$ ).

It is a little more difficult to obtain the prediction errors  $P_{n+m}^n$ , but for large  $n$ , the approximation used in Section 3.4, equation (3.86), works well. That is, the mean-squared prediction error can be approximated by

$$P_{n+m}^n = \sigma_w^2 \sum_{j=0}^{m-1} \psi_j^{*2}, \quad (3.145)$$

where  $\psi_j^*$  is the coefficient of  $z^j$  in  $\psi^*(z) = \theta(z)/\phi(z)(1-z)^d$ .

To better understand integrated models, we examine the properties of some simple cases; Problem 3.29 covers the ARIMA(1, 1, 0) case.

### Example 3.37 Random Walk with Drift

To fix ideas, we begin by considering the random walk with drift model first presented in Example 1.11, that is,

$$x_t = \delta + x_{t-1} + w_t,$$

for  $t = 1, 2, \dots$ , and  $x_0 = 0$ . Technically, the model is not ARIMA, but we could include it trivially as an ARIMA(0, 1, 0) model. Given data  $x_1, \dots, x_n$ , the one-step-ahead forecast is given by

$$x_{n+1}^n = E(x_{n+1} \mid x_n, \dots, x_1) = E(\delta + x_n + w_{n+1} \mid x_n, \dots, x_1) = \delta + x_n.$$

The two-step-ahead forecast is given by  $x_{n+2}^n = \delta + x_{n+1}^n = 2\delta + x_n$ , and consequently, the  $m$ -step-ahead forecast, for  $m = 1, 2, \dots$ , is

$$x_{n+m}^n = m\delta + x_n, \quad (3.146)$$

To obtain the forecast errors, it is convenient to recall equation (1.4); i.e.,  $x_n = n\delta + \sum_{j=1}^n w_j$ , in which case we may write

$$x_{n+m} = (n+m)\delta + \sum_{j=1}^{n+m} w_j = m\delta + x_n + \sum_{j=n+1}^{n+m} w_j.$$

From this it follows that the  $m$ -step-ahead prediction error is given by

$$P_{n+m}^n = E(x_{n+m} - x_{n+m}^n)^2 = E\left(\sum_{j=n+1}^{n+m} w_j\right)^2 = m\sigma_w^2. \quad (3.147)$$

Hence, unlike the stationary case (see Example 3.23), as the forecast horizon grows, the prediction errors, (3.147), increase without bound and the forecasts follow a straight line with slope  $\delta$  emanating from  $x_n$ . We note that (3.145) is exact in this case because  $\psi^*(z) = 1/(1-z) = \sum_{j=0}^{\infty} z^j$  for  $|z| < 1$ , so that  $\psi_j^* = 1$  for all  $j$ .

The  $w_t$  are Gaussian, so estimation is straightforward because the differenced data, say  $y_t = \nabla x_t$ , are independent and identically distributed normal variates with mean  $\delta$  and variance  $\sigma_w^2$ . Consequently, optimal estimates of  $\delta$  and  $\sigma_w^2$  are the sample mean and variance of the  $y_t$ , respectively.

**Example 3.38 IMA(1, 1) and EWMA**

The ARIMA(0,1,1), or IMA(1,1) model is of interest because many economic time series can be successfully modeled this way. In addition, the model leads to a frequently used, and abused, forecasting method called exponentially weighted moving averages (EWMA). We will write the model as

$$x_t = x_{t-1} + w_t - \lambda w_{t-1}, \quad (3.148)$$

with  $|\lambda| < 1$ , for  $t = 1, 2, \dots$ , and  $x_0 = 0$ , because this model formulation is easier to work with here, and it leads to the standard representation for EWMA. We could have included a drift term in (3.148), as was done in the previous example, but for the sake of simplicity, we leave it out of the discussion. If we write

$$y_t = w_t - \lambda w_{t-1},$$

we may write (3.148) as  $x_t = x_{t-1} + y_t$ . Because  $|\lambda| < 1$ ,  $y_t$  has an invertible representation,  $y_t = \sum_{j=1}^{\infty} \lambda^j y_{t-j} + w_t$ , and substituting  $y_t = x_t - x_{t-1}$ , we may write

$$x_t = \sum_{j=1}^{\infty} (1 - \lambda) \lambda^{j-1} x_{t-j} + w_t. \quad (3.149)$$

as an approximation for large  $t$  (put  $x_t = 0$  for  $t \leq 0$ ). Verification of (3.149) is left to the reader (Problem 3.28). Using the approximation (3.149), we have that the approximate one-step-ahead predictor, using the notation of Section 3.4, is

$$\begin{aligned} \tilde{x}_{n+1} &= \sum_{j=1}^{\infty} (1 - \lambda) \lambda^{j-1} x_{n+1-j} \\ &= (1 - \lambda) x_n + \lambda \sum_{j=1}^{\infty} (1 - \lambda) \lambda^{j-1} x_{n-j} \\ &= (1 - \lambda) x_n + \lambda \tilde{x}_n. \end{aligned} \quad (3.150)$$

From (3.150), we see that the new forecast is a linear combination of the old forecast and the new observation. Based on (3.150) and the fact that we only observe  $x_1, \dots, x_n$ , and consequently  $y_1, \dots, y_n$  (because  $y_t = x_t - x_{t-1}$ ;  $x_0 = 0$ ), the truncated forecasts are

$$\tilde{x}_{n+1}^n = (1 - \lambda) x_n + \lambda \tilde{x}_n^{n-1}, \quad n \geq 1, \quad (3.151)$$

with  $\tilde{x}_1^0 = x_1$  as an initial value. The mean-square prediction error can be approximated using (3.145) by noting that  $\psi^*(z) = (1 - \lambda z)/(1 - z) = 1 + (1 - \lambda) \sum_{j=1}^{\infty} z^j$  for  $|z| < 1$ ; consequently, for large  $n$ , (3.145) leads to

$$P_{n+m}^n \approx \sigma_w^2 [1 + (m - 1)(1 - \lambda)^2].$$

In EWMA, the parameter  $1 - \lambda$  is often called the smoothing parameter and is restricted to be between zero and one. Larger values of  $\lambda$  lead to smoother forecasts.

This method of forecasting is popular because it is easy to use; we need only retain the previous forecast value and the current observation to forecast the next time period. Unfortunately, as previously suggested, the method is often abused because some forecasters do not verify that the observations follow an IMA(1, 1) process, and often arbitrarily pick values of  $\lambda$ . In the following, we show how to generate 100 observations from an IMA(1,1) model with  $\lambda = -\theta = .8$  and then calculate and display the fitted EWMA superimposed on the data. This is accomplished using the Holt-Winters command in R (see the help file `?HoltWinters` for details; no output is shown):

```
set.seed(666)
x = arima.sim(list(order = c(0,1,1), ma = -0.8), n = 100)
(x.ima = HoltWinters(x, beta=FALSE, gamma=FALSE)) #  $\alpha$  below is  $1 - \lambda$ 
  Smoothing parameter:  alpha:    0.1663072
plot(x.ima)
```

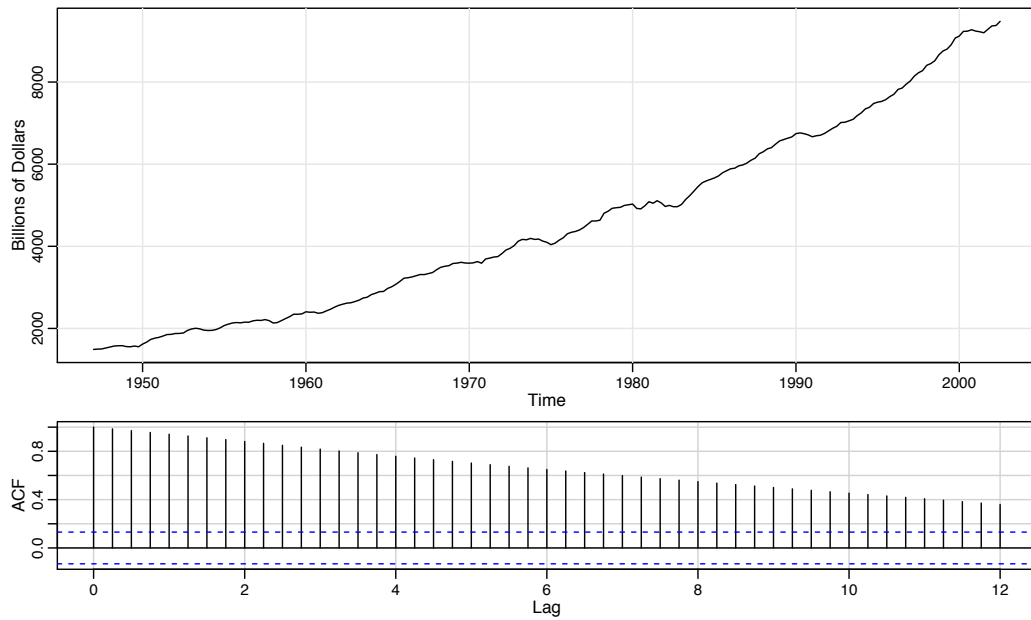
### 3.7 Building ARIMA Models

There are a few basic steps to fitting ARIMA models to time series data. These steps involve

- plotting the data,
- possibly transforming the data,
- identifying the dependence orders of the model,
- parameter estimation,
- diagnostics, and
- model choice.

First, as with any data analysis, we should construct a time plot of the data, and inspect the graph for any anomalies. If, for example, the variability in the data grows with time, it will be necessary to transform the data to stabilize the variance. In such cases, the Box–Cox class of power transformations, equation (2.34), could be employed. Also, the particular application might suggest an appropriate transformation. For example, we have seen numerous examples where the data behave as  $x_t = (1 + p_t)x_{t-1}$ , where  $p_t$  is a small percentage change from period  $t - 1$  to  $t$ , which may be negative. If  $p_t$  is a relatively stable process, then  $\nabla \log(x_t) \approx p_t$  will be relatively stable. Frequently,  $\nabla \log(x_t)$  is called the *return* or *growth rate*. This general idea was used in Example 3.33, and we will use it again in Example 3.39.

After suitably transforming the data, the next step is to identify preliminary values of the autoregressive order,  $p$ , the order of differencing,  $d$ , and the moving average order,  $q$ . A time plot of the data will typically suggest whether any differencing is needed. If differencing is called for, then difference the data once,  $d = 1$ , and inspect the time plot of  $\nabla x_t$ . If additional differencing is necessary, then try differencing again and inspect a time plot of  $\nabla^2 x_t$ . Be careful not to overdifference because this may introduce dependence where none exists. For example,  $x_t = w_t$  is serially uncorrelated, but  $\nabla x_t = w_t - w_{t-1}$  is MA(1). In addition to time plots, the sample



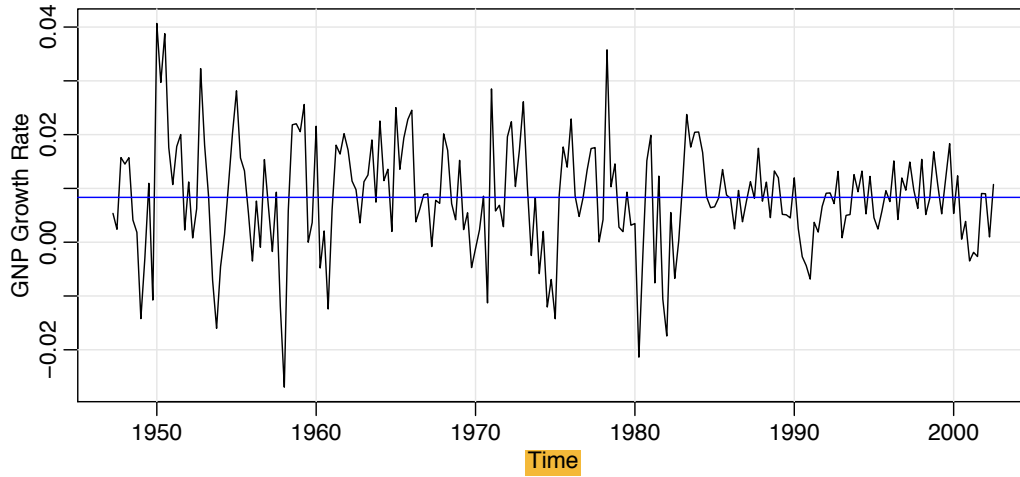
**Fig. 3.13.** Top: Quarterly U.S. GNP from 1947(1) to 2002(3). Bottom: Sample ACF of the GNP data. Lag is in terms of years.

ACF can help in indicating whether differencing is needed. Because the polynomial  $\phi(z)(1 - z)^d$  has a unit root, the sample ACF,  $\hat{\rho}(h)$ , will not decay to zero fast as  $h$  increases. Thus, a slow decay in  $\hat{\rho}(h)$  is an indication that differencing may be needed.

When preliminary values of  $d$  have been settled, the next step is to look at the sample ACF and PACF of  $\nabla^d x_t$  for whatever values of  $d$  have been chosen. Using Table 3.1 as a guide, preliminary values of  $p$  and  $q$  are chosen. Note that it cannot be the case that both the ACF and PACF cut off. Because we are dealing with estimates, it will not always be clear whether the sample ACF or PACF is tailing off or cutting off. Also, two models that are seemingly different can actually be very similar. With this in mind, we should not worry about being so precise at this stage of the model fitting. At this point, a few preliminary values of  $p$ ,  $d$ , and  $q$  should be at hand, and we can start estimating the parameters.

### Example 3.39 Analysis of GNP Data

In this example, we consider the analysis of quarterly U.S. GNP from 1947(1) to 2002(3),  $n = 223$  observations. The data are real U.S. gross national product in billions of chained 1996 dollars and have been seasonally adjusted. The data were obtained from the Federal Reserve Bank of St. Louis (<http://research.stlouisfed.org/>). Figure 3.13 shows a plot of the data, say,  $y_t$ . Because strong trend tends to obscure other effects, it is difficult to see any other variability in data except for periodic large dips in the economy. When reports of GNP and similar economic indicators are given, it is often in growth rate (percent change) rather than in actual (or adjusted) values that is of interest. The growth rate, say,  $x_t = \nabla \log(y_t)$ , is plotted in Figure 3.14, and it appears to be a stable process.



**Fig. 3.14.** U.S. GNP quarterly growth rate. The horizontal line displays the average growth of the process, which is close to 1%.

The sample ACF and PACF of the quarterly growth rate are plotted in Figure 3.15. Inspecting the sample ACF and PACF, we might feel that the ACF is cutting off at lag 2 and the PACF is tailing off. This would suggest the GNP growth rate follows an MA(2) process, or log GNP follows an ARIMA(0, 1, 2) model. Rather than focus on one model, we will also suggest that it appears that the ACF is tailing off and the PACF is cutting off at lag 1. This suggests an AR(1) model for the growth rate, or ARIMA(1, 1, 0) for log GNP. As a preliminary analysis, we will fit both models.

Using MLE to fit the MA(2) model for the growth rate,  $x_t$ , the estimated model is

$$\hat{x}_t = .008_{(.001)} + .303_{(.065)}\hat{w}_{t-1} + .204_{(.064)}\hat{w}_{t-2} + \hat{w}_t, \quad (3.152)$$

where  $\hat{\sigma}_w = .0094$  is based on 219 degrees of freedom. The values in parentheses are the corresponding estimated standard errors. All of the regression coefficients are significant, including the constant. *We make a special note of this because, as a default, some computer packages do not fit a constant in a differenced model.* That is, these packages assume, by default, that there is no drift. In this example, not including a constant leads to the wrong conclusions about the nature of the U.S. economy. Not including a constant assumes the average quarterly growth rate is zero, whereas the U.S. GNP average quarterly growth rate is about 1% (which can be seen easily in Figure 3.14). We leave it to the reader to investigate what happens when the constant is not included.

The estimated AR(1) model is

$$\hat{x}_t = .008_{(.001)}(1 - .347) + .347_{(.063)}\hat{x}_{t-1} + \hat{w}_t, \quad (3.153)$$

where  $\hat{\sigma}_w = .0095$  on 220 degrees of freedom; note that the constant in (3.153) is  $.008(1 - .347) = .005$ .

We will discuss diagnostics next, but assuming both of these models fit well, how are we to reconcile the apparent differences of the estimated models (3.152)

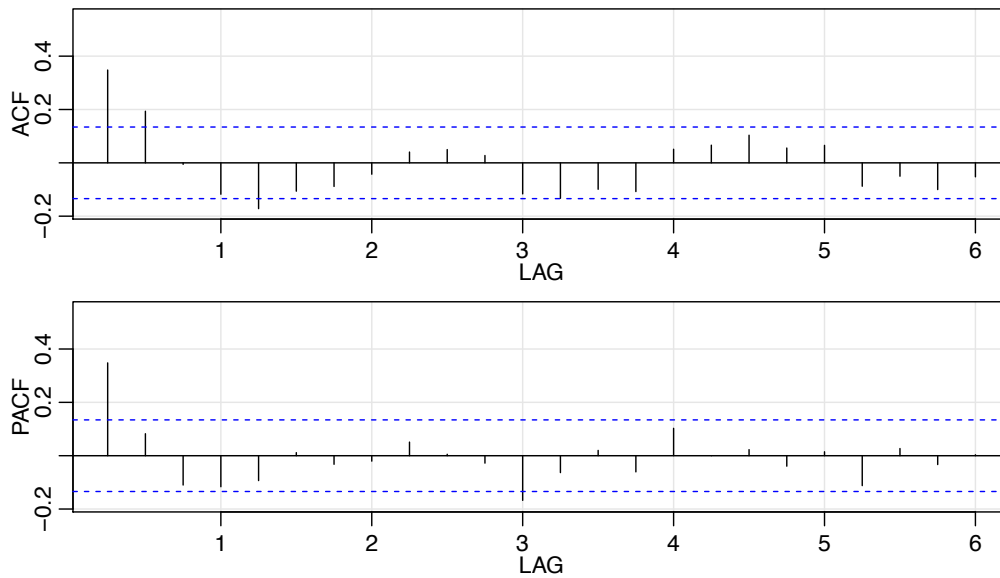


Fig. 3.15. Sample ACF and PACF of the GNP quarterly growth rate. Lag is in terms of years.

and (3.153)? In fact, the fitted models are nearly the same. To show this, consider an AR(1) model of the form in (3.153) without a constant term; that is,

$$x_t = .35x_{t-1} + w_t,$$

and write it in its causal form,  $x_t = \sum_{j=0}^{\infty} \psi_j w_{t-j}$ , where we recall  $\psi_j = .35^j$ . Thus,  $\psi_0 = 1, \psi_1 = .350, \psi_2 = .123, \psi_3 = .043, \psi_4 = .015, \psi_5 = .005, \psi_6 = .002, \psi_7 = .001, \psi_8 = 0, \psi_9 = 0, \psi_{10} = 0$ , and so forth. Thus,

$$x_t \approx .35w_{t-1} + .12w_{t-2} + w_t,$$

which is similar to the fitted MA(2) model in (3.153).

The analysis can be performed in R as follows.

```
plot(gnp)
acf2(gnp, 50)
gnpgr = diff(log(gnp)) # growth rate
plot(gnpgr)
acf2(gnpgr, 24)
sarima(gnpgr, 1, 0, 0) # AR(1)
sarima(gnpgr, 0, 0, 2) # MA(2)
ARMAtoMA(ar=.35, ma=0, 10) # prints psi-weights
```

The next step in model fitting is diagnostics. This investigation includes the analysis of the residuals as well as model comparisons. Again, the first step involves a time plot of the *innovations* (or residuals),  $x_t - \hat{x}_t^{t-1}$ , or of the *standardized innovations*

$$e_t = (x_t - \hat{x}_t^{t-1}) / \sqrt{\hat{P}_t^{t-1}}, \quad (3.154)$$

where  $\hat{x}_t^{t-1}$  is the one-step-ahead prediction of  $x_t$  based on the fitted model and  $\hat{P}_t^{t-1}$  is the estimated one-step-ahead error variance. If the model fits well, the standardized



residuals should behave as an iid sequence with mean zero and variance one. The time plot should be inspected for any obvious departures from this assumption. Unless the time series is Gaussian, it is not enough that the residuals are uncorrelated. For example, it is possible in the non-Gaussian case to have an uncorrelated process for which values contiguous in time are highly dependent. As an example, we mention the family of GARCH models that are discussed in Chapter 5.

Investigation of marginal normality can be accomplished visually by looking at a histogram of the residuals. In addition to this, a normal probability plot or a Q-Q plot can help in identifying departures from normality. See Johnson and Wichern (1992, Chapter 4) for details of this test as well as additional tests for multivariate normality.

There are several tests of randomness, for example the runs test, that could be applied to the residuals. We could also inspect the sample autocorrelations of the residuals, say,  $\hat{\rho}_e(h)$ , for any patterns or large values. Recall that, for a white noise sequence, the sample autocorrelations are approximately independently and normally distributed with zero means and variances  $1/n$ . Hence, a good check on the correlation structure of the residuals is to plot  $\hat{\rho}_e(h)$  versus  $h$  along with the error bounds of  $\pm 2/\sqrt{n}$ . The residuals from a model fit, however, will not quite have the properties of a white noise sequence and the variance of  $\hat{\rho}_e(h)$  can be much less than  $1/n$ . Details can be found in Box and Pierce (1970) and McLeod (1978). This part of the diagnostics can be viewed as a visual inspection of  $\hat{\rho}_e(h)$  with the main concern being the detection of obvious departures from the independence assumption.

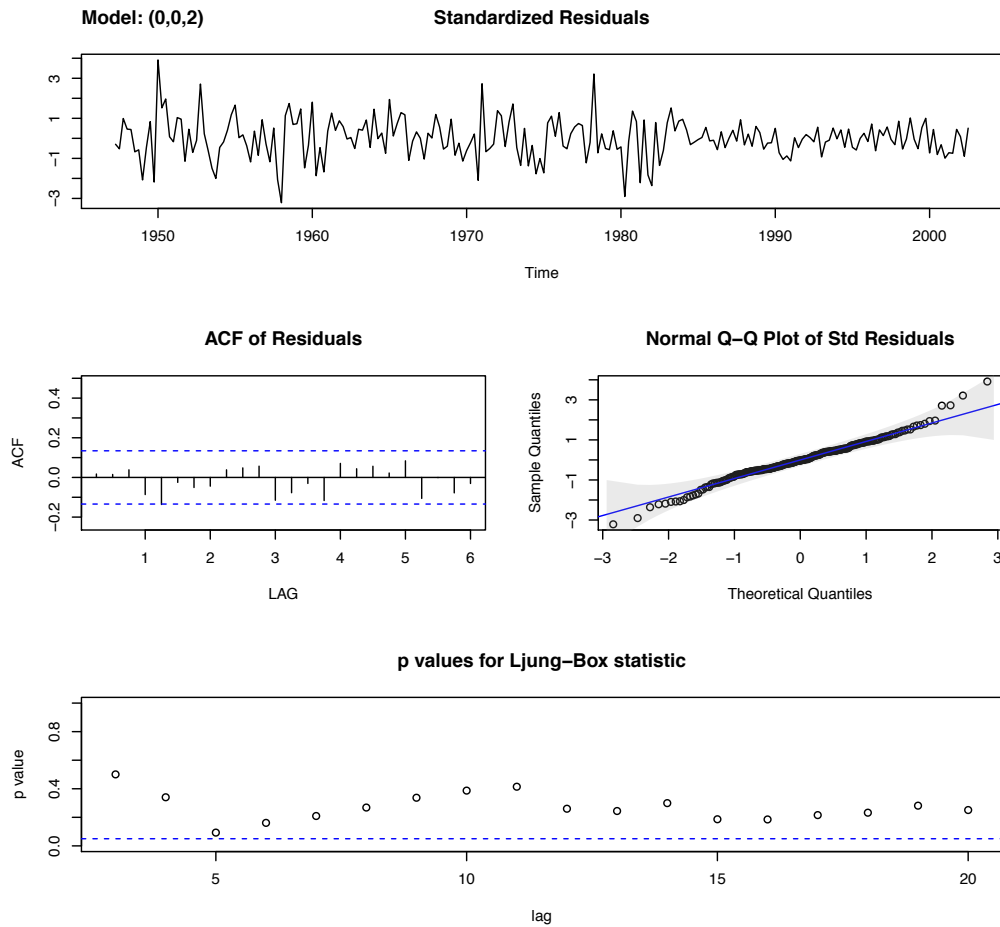
In addition to plotting  $\hat{\rho}_e(h)$ , we can perform a general test that takes into consideration the magnitudes of  $\hat{\rho}_e(h)$  as a group. For example, it may be the case that, individually, each  $\hat{\rho}_e(h)$  is small in magnitude, say, each one is just slightly less than  $2/\sqrt{n}$  in magnitude, but, collectively, the values are large. The *Ljung–Box–Pierce Q-statistic* given by

$$Q = n(n+2) \sum_{h=1}^H \frac{\hat{\rho}_e^2(h)}{n-h} \quad (3.155)$$

can be used to perform such a test. The value  $H$  in (3.155) is chosen somewhat arbitrarily, typically,  $H = 20$ . Under the null hypothesis of model adequacy, asymptotically ( $n \rightarrow \infty$ ),  $Q \sim \chi_{H-p-q}^2$ . Thus, we would reject the null hypothesis at level  $\alpha$  if the value of  $Q$  exceeds the  $(1 - \alpha)$ -quantile of the  $\chi_{H-p-q}^2$  distribution. Details can be found in Box and Pierce (1970), Ljung and Box (1978), and Davies et al. (1977). The basic idea is that if  $w_t$  is white noise, then by Property 1.2,  $n\hat{\rho}_w^2(h)$ , for  $h = 1, \dots, H$ , are asymptotically independent  $\chi_1^2$  random variables. This means that  $n \sum_{h=1}^H \hat{\rho}_w^2(h)$  is approximately a  $\chi_H^2$  random variable. Because the test involves the ACF of residuals from a model fit, there is a loss of  $p + q$  degrees of freedom; the other values in (3.155) are used to adjust the statistic to better match the asymptotic chi-squared distribution.

#### Example 3.40 Diagnostics for GNP Growth Rate Example

We will focus on the MA(2) fit from Example 3.39; the analysis of the AR(1) residuals is similar. Figure 3.16 displays a plot of the standardized residuals, the ACF of the residuals, a boxplot of the standardized residuals, and the p-values



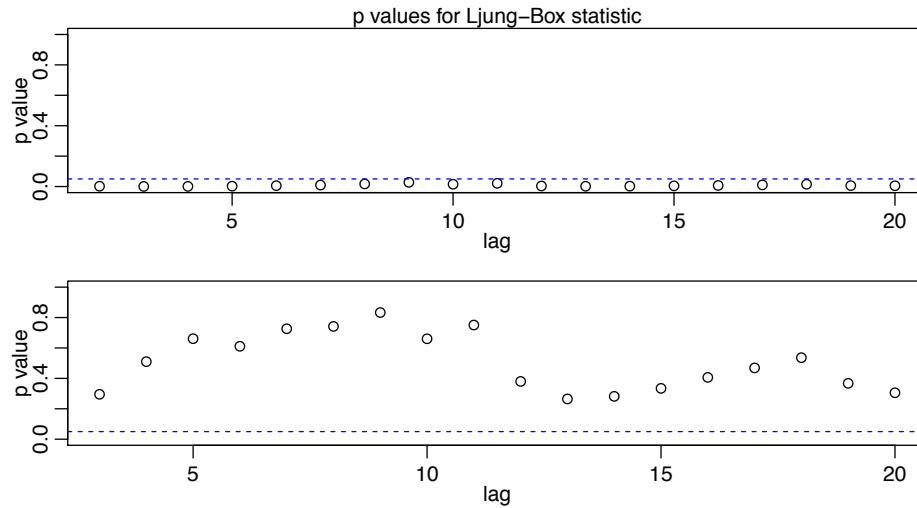
**Fig. 3.16.** Diagnostics of the residuals from MA(2) fit on GNP growth rate.

associated with the Q-statistic, (3.155), at lags  $H = 3$  through  $H = 20$  (with corresponding degrees of freedom  $H - 2$ ).

Inspection of the time plot of the standardized residuals in Figure 3.16 shows no obvious patterns. Notice that there may be outliers, with a few values exceeding 3 standard deviations in magnitude. The ACF of the standardized residuals shows no apparent departure from the model assumptions, and the Q-statistic is never significant at the lags shown. The normal Q-Q plot of the residuals shows that the assumption of normality is reasonable, with the exception of the possible outliers.

The model appears to fit well. The diagnostics shown in Figure 3.16 are a by-product of the `sarima` command from the previous example.<sup>3.8</sup>

<sup>3.8</sup> The script `tsdiag` is available in R to run diagnostics for an ARIMA object, however, the script has errors and we do not recommend using it.



**Fig. 3.17.** *Q-statistic p-values for the ARIMA(0, 1, 1) fit (top) and the ARIMA(1, 1, 1) fit (bottom) to the logged varve data.*

### Example 3.41 Diagnostics for the Glacial Varve Series

In [Example 3.33](#), we fit an ARIMA(0, 1, 1) model to the logarithms of the glacial varve data and there appears to be a small amount of autocorrelation left in the residuals and the Q-tests are all significant; see [Figure 3.17](#).

To adjust for this problem, we fit an ARIMA(1, 1, 1) to the logged varve data and obtained the estimates

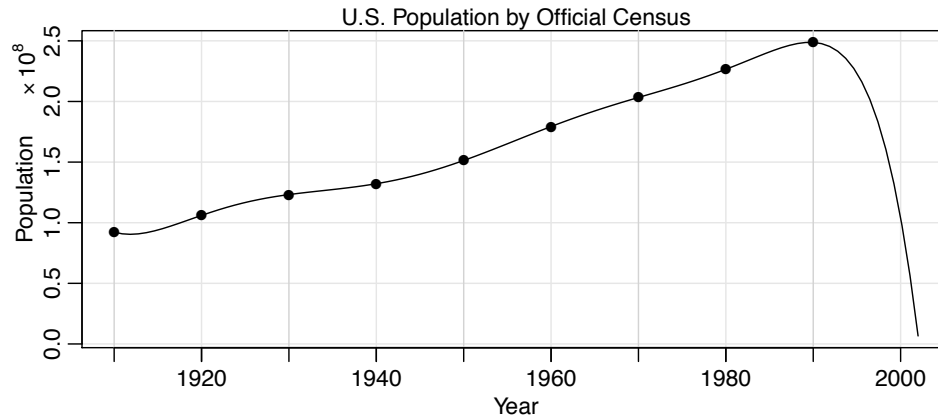
$$\hat{\phi} = .23_{(.05)}, \hat{\theta} = -.89_{(.03)}, \hat{\sigma}_w^2 = .23.$$

Hence the AR term is significant. The Q-statistic p-values for this model are also displayed in [Figure 3.17](#), and it appears this model fits the data well.

As previously stated, the diagnostics are byproducts of the individual `sarima` runs. We note that we did not fit a constant in either model because there is no apparent drift in the differenced, logged varve series. This fact can be verified by noting the constant is not significant when the command `no.constant=TRUE` is removed in the code:

```
sarima(log(varve), 0, 1, 1, no.constant=TRUE) # ARIMA(0,1,1)
sarima(log(varve), 1, 1, 1, no.constant=TRUE) # ARIMA(1,1,1)
```

In [Example 3.39](#), we have two competing models, an AR(1) and an MA(2) on the GNP growth rate, that each appear to fit the data well. In addition, we might also consider that an AR(2) or an MA(3) might do better for forecasting. Perhaps combining both models, that is, fitting an ARMA(1, 2) to the GNP growth rate, would be the best. As previously mentioned, we have to be concerned with *overfitting* the model; it is not always the case that more is better. Overfitting leads to less-precise estimators, and adding more parameters may fit the data better but may also lead to bad forecasts. This result is illustrated in the following example.



*Fig. 3.18. A perfect fit and a terrible forecast.*

### Example 3.42 A Problem with Overfitting

Figure 3.18 shows the U.S. population by official census, every ten years from 1910 to 1990, as points. If we use these nine observations to predict the future population, we can use an eight-degree polynomial so the fit to the nine observations is perfect. The model in this case is

$$x_t = \beta_0 + \beta_1 t + \beta_2 t^2 + \cdots + \beta_8 t^8 + w_t.$$

The fitted line, which is plotted in the figure, passes through the nine observations. The model predicts that the population of the United States will be close to zero in the year 2000, and will cross zero sometime in the year 2002!

The final step of model fitting is model choice or model selection. That is, we must decide which model we will retain for forecasting. The most popular techniques, AIC, AICc, and BIC, were described in Section 2.1 in the context of regression models.

### Example 3.43 Model Choice for the U.S. GNP Series

Returning to the analysis of the U.S. GNP data presented in Example 3.39 and Example 3.40, recall that two models, an AR(1) and an MA(2), fit the GNP growth rate well. To choose the final model, we compare the AIC, the AICc, and the BIC for both models. These values are a byproduct of the `sarima` runs displayed at the end of Example 3.39, but for convenience, we display them again here (recall the growth rate data are in `gnpgr`):

```
sarima(gnpgr, 1, 0, 0) # AR(1)
$AIC: -8.294403 $AICc: -8.284898 $BIC: -9.263748
sarima(gnpgr, 0, 0, 2) # MA(2)
$AIC: -8.297693 $AICc: -8.287854 $BIC: -9.251711
```

The AIC and AICc both prefer the MA(2) fit, whereas the BIC prefers the simpler AR(1) model. It is often the case that the BIC will select a model of smaller order than the AIC or AICc. In either case, it is not unreasonable to retain the AR(1) because pure autoregressive models are easier to work with.

### 3.8 Regression with Autocorrelated Errors

In [Section 2.1](#), we covered the classical regression model with uncorrelated errors  $w_t$ . In this section, we discuss the modifications that might be considered when the errors are correlated. That is, consider the regression model

$$y_t = \sum_{j=1}^r \beta_j z_{tj} + x_t \quad (3.156)$$

where  $x_t$  is a process with some covariance function  $\gamma_x(s, t)$ . In ordinary least squares, the assumption is that  $x_t$  is white Gaussian noise, in which case  $\gamma_x(s, t) = 0$  for  $s \neq t$  and  $\gamma_x(t, t) = \sigma^2$ , independent of  $t$ . If this is not the case, then weighted least squares should be used.

Write the model in vector notation,  $y = Z\beta + x$ , where  $y = (y_1, \dots, y_n)'$  and  $x = (x_1, \dots, x_n)'$  are  $n \times 1$  vectors,  $\beta = (\beta_1, \dots, \beta_r)'$  is  $r \times 1$ , and  $Z = [z_1 | z_2 | \dots | z_n]'$  is the  $n \times r$  matrix composed of the input variables. Let  $\Gamma = \{\gamma_x(s, t)\}$ , then  $\Gamma^{-1/2}y = \Gamma^{-1/2}Z\beta + \Gamma^{-1/2}x$ , so that we can write the model as

$$y^* = Z^*\beta + \delta,$$

where  $y^* = \Gamma^{-1/2}y$ ,  $Z^* = \Gamma^{-1/2}Z$ , and  $\delta = \Gamma^{-1/2}x$ . Consequently, the covariance matrix of  $\delta$  is the identity and the model is in the classical linear model form. It follows that the weighted estimate of  $\beta$  is  $\hat{\beta}_w = (Z^{*'}Z^*)^{-1}Z^{*'}y^* = (Z'\Gamma^{-1}Z)^{-1}Z'\Gamma^{-1}y$ , and the variance-covariance matrix of the estimator is  $\text{var}(\hat{\beta}_w) = (Z'\Gamma^{-1}Z)^{-1}$ . If  $x_t$  is white noise, then  $\Gamma = \sigma^2 I$  and these results reduce to the usual least squares results.

In the time series case, it is often possible to assume a stationary covariance structure for the error process  $x_t$  that corresponds to a linear process and try to find an ARMA representation for  $x_t$ . For example, if we have a pure AR( $p$ ) error, then

$$\phi(B)x_t = w_t,$$

and  $\phi(B) = 1 - \phi_1 B - \dots - \phi_p B^p$  is the linear transformation that, when applied to the error process, produces the white noise  $w_t$ . Multiplying the regression equation through by the transformation  $\phi(B)$  yields,

$$\underbrace{\phi(B)y_t}_{y_t^*} = \sum_{j=1}^r \beta_j \underbrace{\phi(B)z_{tj}}_{z_{tj}^*} + \underbrace{\phi(B)x_t}_{w_t},$$

and we are back to the linear regression model where the observations have been transformed so that  $y_t^* = \phi(B)y_t$  is the dependent variable,  $z_{tj}^* = \phi(B)z_{tj}$  for  $j = 1, \dots, r$ , are the independent variables, but the  $\beta$ s are the same as in the original model. For example, if  $p = 1$ , then  $y_t^* = y_t - \phi y_{t-1}$  and  $z_{tj}^* = z_{tj} - \phi z_{t-1,j}$ .

In the AR case, we may set up the least squares problem as minimizing the error sum of squares

$$S(\phi, \beta) = \sum_{t=1}^n w_t^2 = \sum_{t=1}^n \left[ \phi(B)y_t - \sum_{j=1}^r \beta_j \phi(B)z_{tj} \right]^2$$

with respect to all the parameters,  $\phi = \{\phi_1, \dots, \phi_p\}$  and  $\beta = \{\beta_1, \dots, \beta_r\}$ . Of course, the optimization is performed using numerical methods.

If the error process is ARMA( $p, q$ ), i.e.,  $\phi(B)x_t = \theta(B)w_t$ , then in the above discussion, we transform by  $\pi(B)x_t = w_t$ , where  $\pi(B) = \theta(B)^{-1}\phi(B)$ . In this case the error sum of squares also depends on  $\theta = \{\theta_1, \dots, \theta_q\}$ :

$$S(\phi, \theta, \beta) = \sum_{t=1}^n w_t^2 = \sum_{t=1}^n \left[ \pi(B)y_t - \sum_{j=1}^r \beta_j \pi(B)z_{tj} \right]^2$$

At this point, the main problem is that we do not typically know the behavior of the noise  $x_t$  prior to the analysis. An easy way to tackle this problem was first presented in Cochran and Orcutt (1949), and with the advent of cheap computing is modernized below:

- (i) First, run an ordinary regression of  $y_t$  on  $z_{t1}, \dots, z_{tr}$  (acting as if the errors are uncorrelated). Retain the residuals,  $\hat{x}_t = y_t - \sum_{j=1}^r \hat{\beta}_j z_{tj}$ .
- (ii) Identify ARMA model(s) for the residuals  $\hat{x}_t$ .
- (iii) Run weighted least squares (or MLE) on the regression model with autocorrelated errors using the model specified in step (ii).
- (iv) Inspect the residuals  $\hat{w}_t$  for whiteness, and adjust the model if necessary.

### Example 3.44 Mortality, Temperature and Pollution

We consider the analyses presented in [Example 2.2](#), relating mean adjusted temperature  $T_t$ , and particulate levels  $P_t$  to cardiovascular mortality  $M_t$ . We consider the regression model

$$M_t = \beta_1 + \beta_2 t + \beta_3 T_t + \beta_4 T_t^2 + \beta_5 P_t + x_t, \quad (3.157)$$

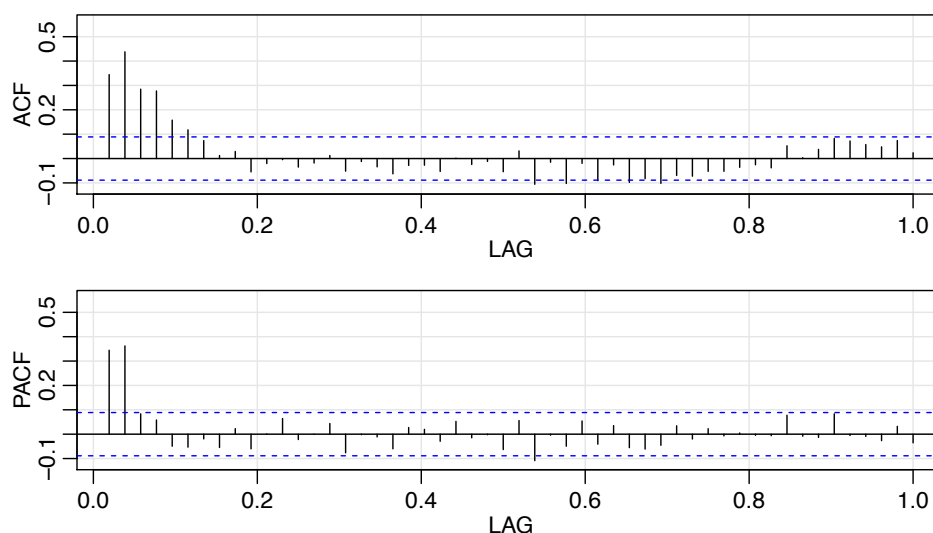
where, for now, we assume that  $x_t$  is white noise. The sample ACF and PACF of the residuals from the ordinary least squares fit of (3.157) are shown in [Figure 3.19](#), and the results suggest an AR(2) model for the residuals.

Our next step is to fit the correlated error model (3.157), but where  $x_t$  is AR(2),

$$x_t = \phi_1 x_{t-1} + \phi_2 x_{t-2} + w_t$$

and  $w_t$  is white noise. The model can be fit using the `sarima` function as follows (partial output shown).

```
trend = time(cmort); temp = tempr - mean(tempr); temp2 = temp^2
summary(fit <- lm(cmort~trend + temp + temp2 + part, na.action=NULL))
acf2(resid(fit), 52) # implies AR2
sarima(cmort, 2,0,0, xreg=cbind(trend,temp,temp2,part))
Coefficients:
      ar1      ar2  intercept      trend      temp      temp2      part
    0.3848  0.4326   80.2116  -1.5165  -0.0190   0.0154   0.1545
s.e.  0.0436  0.0400    1.8072   0.4226   0.0495   0.0020   0.0272
sigma^2 estimated as 26.01: loglikelihood = -1549.04, aic = 3114.07
```



**Fig. 3.19.** Sample ACF and PACF of the mortality residuals indicating an AR(2) process.

The residual analysis output from `sarima` (not shown) shows no obvious departure of the residuals from whiteness.

### Example 3.45 Regression with Lagged Variables (cont)

In [Example 2.9](#) we fit the model

$$R_t = \beta_0 + \beta_1 S_{t-6} + \beta_2 D_{t-6} + \beta_3 D_{t-6} S_{t-6} + w_t,$$

where  $R_t$  is Recruitment,  $S_t$  is SOI, and  $D_t$  is a dummy variable that is 0 if  $S_t < 0$  and 1 otherwise. However, residual analysis indicates that the residuals are not white noise. The sample (P)ACF of the residuals indicates that an AR(2) model might be appropriate, which is similar to the results of [Example 3.44](#). We display partial results of the final model below.

```
dummy = ifelse(soi < 0, 0, 1)
fish = ts.intersect(rec, soil6=lag(soi,-6), dL6=lag(dummy,-6), dframe=TRUE)
summary(fit <- lm(rec ~soil6*dL6, data=fish, na.action=NULL))
attach(fish)
plot(resid(fit))
acf2(resid(fit)) # indicates AR(2)
intract = soil6*dL6 # interaction term
sarima(rec,2,0,0, xreg = cbind(soil6, dL6, intract))
$ttable
```

	Estimate	SE	t.value	p.value
ar1	1.3624	0.0440	30.9303	0.0000
ar2	-0.4703	0.0444	-10.5902	0.0000
intercept	64.8028	4.1121	15.7590	0.0000
soil6	8.6671	2.2205	3.9033	0.0001
dL6	-2.5945	0.9535	-2.7209	0.0068
intract	-10.3092	2.8311	-3.6415	0.0003

### 3.9 Multiplicative Seasonal ARIMA Models

In this section, we introduce several modifications made to the ARIMA model to account for seasonal and nonstationary behavior. Often, the dependence on the past tends to occur most strongly at multiples of some underlying seasonal lag  $s$ . For example, with monthly economic data, there is a strong yearly component occurring at lags that are multiples of  $s = 12$ , because of the strong connections of all activity to the calendar year. Data taken quarterly will exhibit the yearly repetitive period at  $s = 4$  quarters. Natural phenomena such as temperature also have strong components corresponding to seasons. Hence, the natural variability of many physical, biological, and economic processes tends to match with seasonal fluctuations. Because of this, it is appropriate to introduce autoregressive and moving average polynomials that identify with the seasonal lags. The resulting *pure seasonal autoregressive moving average model*, say,  $\text{ARMA}(P, Q)_s$ , then takes the form

$$\Phi_P(B^s)x_t = \Theta_Q(B^s)w_t, \quad (3.158)$$

where the operators

$$\Phi_P(B^s) = 1 - \Phi_1 B^s - \Phi_2 B^{2s} - \dots - \Phi_P B^{Ps} \quad (3.159)$$

and

$$\Theta_Q(B^s) = 1 + \Theta_1 B^s + \Theta_2 B^{2s} + \dots + \Theta_Q B^{Qs} \quad (3.160)$$

are the **seasonal autoregressive operator** and the **seasonal moving average operator** of orders  $P$  and  $Q$ , respectively, with seasonal period  $s$ .

Analogous to the properties of nonseasonal ARMA models, the pure seasonal  $\text{ARMA}(P, Q)_s$  is *causal* only when the roots of  $\Phi_P(z^s)$  lie outside the unit circle, and it is *invertible* only when the roots of  $\Theta_Q(z^s)$  lie outside the unit circle.

#### Example 3.46 A Seasonal AR Series

A first-order seasonal autoregressive series that might run over months could be written as

$$(1 - \Phi B^{12})x_t = w_t$$

or

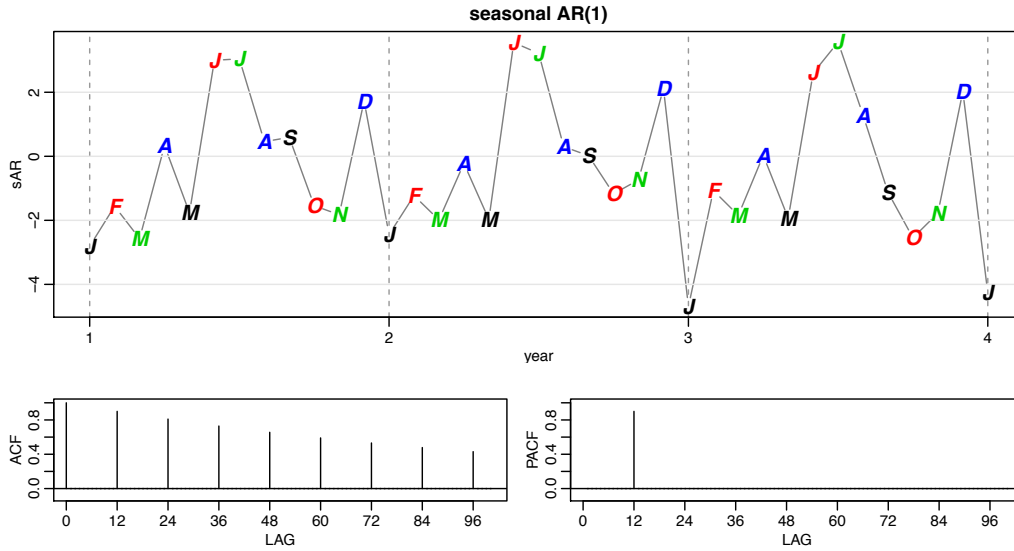
$$x_t = \Phi x_{t-12} + w_t.$$

This model exhibits the series  $x_t$  in terms of past lags at the multiple of the yearly seasonal period  $s = 12$  months. It is clear from the above form that estimation and forecasting for such a process involves only straightforward modifications of the unit lag case already treated. In particular, the causal condition requires  $|\Phi| < 1$ .

We simulated 3 years of data from the model with  $\Phi = .9$ , and exhibit the *theoretical* ACF and PACF of the model. See [Figure 3.20](#).

```
set.seed(666)
phi = c(rep(0,11),.9)
sAR = arima.sim(list(order=c(12,0,0), ar=phi), n=37)
sAR = ts(sAR, freq=12)
layout(matrix(c(1,1,2, 1,1,3), nc=2))
```





**Fig. 3.20.** Data generated from a seasonal ( $s = 12$ ) AR(1), and the true ACF and PACF of the model  $x_t = .9x_{t-12} + w_t$ .

```
par(mar=c(3,3,2,1), mgp=c(1.6,.6,0))
plot(sAR, axes=FALSE, main='seasonal AR(1)', xlab="year", type='c')
Months = c("J","F","M","A","M","J","J","A","S","O","N","D")
points(sAR, pch=Months, cex=1.25, font=4, col=1:4)
axis(1, 1:4); abline(v=1:4, lty=2, col=gray(.7))
axis(2); box()
ACF = ARMAacf(ar=phi, ma=0, 100)
PACF = ARMAacf(ar=phi, ma=0, 100, pacf=TRUE)
plot(ACF,type="h", xlab="LAG", ylim=c(-.1,1)); abline(h=0)
plot(PACF, type="h", xlab="LAG", ylim=c(-.1,1)); abline(h=0)
```

For the first-order seasonal ( $s = 12$ ) MA model,  $x_t = w_t + \Theta w_{t-12}$ , it is easy to verify that

$$\begin{aligned}\gamma(0) &= (1 + \Theta^2)\sigma^2 \\ \gamma(\pm 12) &= \Theta\sigma^2 \\ \gamma(h) &= 0, \quad \text{otherwise.}\end{aligned}$$

Thus, the only nonzero correlation, aside from lag zero, is

$$\rho(\pm 12) = \Theta/(1 + \Theta^2).$$

For the first-order seasonal ( $s = 12$ ) AR model, using the techniques of the nonseasonal AR(1), we have

$$\begin{aligned}\gamma(0) &= \sigma^2/(1 - \Phi^2) \\ \gamma(\pm 12k) &= \sigma^2\Phi^k/(1 - \Phi^2) \quad k = 1, 2, \dots \\ \gamma(h) &= 0, \quad \text{otherwise.}\end{aligned}$$

In this case, the only non-zero correlations are

**Table 3.3.** Behavior of the ACF and PACF for Pure SARMA Models

	$AR(P)_s$	$MA(Q)_s$	$ARMA(P, Q)_s$
ACF*	Tails off at lags $ks$ , $k = 1, 2, \dots$ ,	Cuts off after lag $Qs$	Tails off at lags $ks$
PACF*	Cuts off after lag $Ps$	Tails off at lags $ks$ $k = 1, 2, \dots$ ,	Tails off at lags $ks$

\*The values at nonseasonal lags  $h \neq ks$ , for  $k = 1, 2, \dots$ , are zero.

$$\rho(\pm 12k) = \Phi^k, \quad k = 0, 1, 2, \dots$$

These results can be verified using the general result that  $\gamma(h) = \Phi\gamma(h - 12)$ , for  $h \geq 1$ . For example, when  $h = 1$ ,  $\gamma(1) = \Phi\gamma(11)$ , but when  $h = 11$ , we have  $\gamma(11) = \Phi\gamma(1)$ , which implies that  $\gamma(1) = \gamma(11) = 0$ . In addition to these results, the PACF have the analogous extensions from nonseasonal to seasonal models. These results are demonstrated in Figure 3.20.

As an initial diagnostic criterion, we can use the properties for the pure seasonal autoregressive and moving average series listed in Table 3.3. These properties may be considered as generalizations of the properties for nonseasonal models that were presented in Table 3.1.

In general, we can combine the seasonal and nonseasonal operators into a *multiplicative seasonal autoregressive moving average model*, denoted by  $ARMA(p, q) \times (P, Q)_s$ , and write

$$\Phi_P(B^s)\phi(B)x_t = \Theta_Q(B^s)\theta(B)w_t \quad (3.161)$$

as the overall model. Although the diagnostic properties in Table 3.3 are not strictly true for the overall mixed model, the behavior of the ACF and PACF tends to show rough patterns of the indicated form. In fact, for mixed models, we tend to see a mixture of the facts listed in Table 3.1 and Table 3.3. In fitting such models, focusing on the seasonal autoregressive and moving average components first generally leads to more satisfactory results.

#### Example 3.47 A Mixed Seasonal Model

Consider an  $ARMA(0, 1) \times (1, 0)_{12}$  model

$$x_t = \Phi x_{t-12} + w_t + \theta w_{t-1},$$

where  $|\Phi| < 1$  and  $|\theta| < 1$ . Then, because  $x_{t-12}$ ,  $w_t$ , and  $w_{t-1}$  are uncorrelated, and  $x_t$  is stationary,  $\gamma(0) = \Phi^2\gamma(0) + \sigma_w^2 + \theta^2\sigma_w^2$ , or

$$\gamma(0) = \frac{1 + \theta^2}{1 - \Phi^2} \sigma_w^2.$$

In addition, multiplying the model by  $x_{t-h}$ ,  $h > 0$ , and taking expectations, we have  $\gamma(1) = \Phi\gamma(11) + \theta\sigma_w^2$ , and  $\gamma(h) = \Phi\gamma(h - 12)$ , for  $h \geq 2$ . Thus, the ACF for this model is

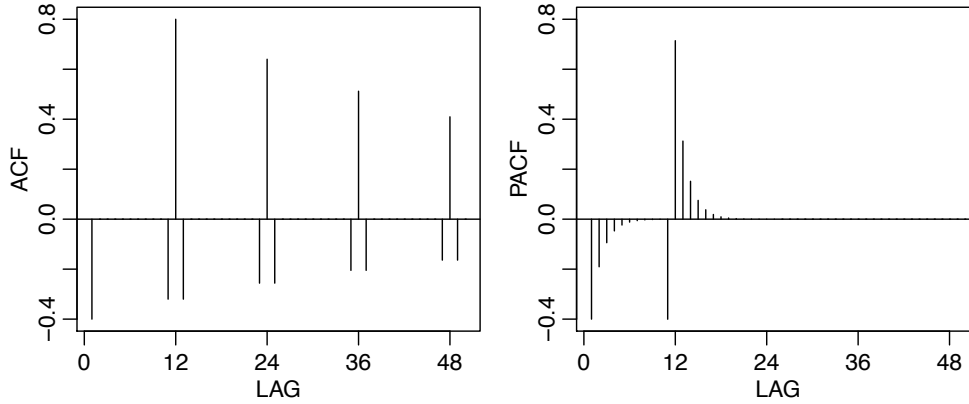


Fig. 3.21. ACF and PACF of the mixed seasonal ARMA model  $x_t = .8x_{t-12} + w_t - .5w_{t-1}$ .

$$\begin{aligned}\rho(12h) &= \Phi^h \quad h = 1, 2, \dots \\ \rho(12h - 1) &= \rho(12h + 1) = \frac{\theta}{1 + \theta^2} \Phi^h \quad h = 0, 1, 2, \dots, \\ \rho(h) &= 0, \quad \text{otherwise.}\end{aligned}$$

The ACF and PACF for this model, with  $\Phi = .8$  and  $\theta = -.5$ , are shown in Figure 3.21. These type of correlation relationships, although idealized here, are typically seen with seasonal data.

To reproduce Figure 3.21 in R, use the following commands:

```
phi = c(rep(0,11),.8)
ACF = ARMAacf(ar=phi, ma=-.5, 50)[-1] # [-1] removes 0 lag
PACF = ARMAacf(ar=phi, ma=-.5, 50, pacf=TRUE)
par(mfrow=c(1,2))
plot(ACF, type="h", xlab="LAG", ylim=c(-.4,.8)); abline(h=0)
plot(PACF, type="h", xlab="LAG", ylim=c(-.4,.8)); abline(h=0)
```

Seasonal persistence occurs when the process is nearly periodic in the season. For example, with average monthly temperatures over the years, each January would be approximately the same, each February would be approximately the same, and so on. In this case, we might think of average monthly temperature  $x_t$  as being modeled as

$$x_t = S_t + w_t,$$

where  $S_t$  is a seasonal component that varies a little from one year to the next, according to a random walk,

$$S_t = S_{t-12} + v_t.$$

In this model,  $w_t$  and  $v_t$  are uncorrelated white noise processes. The tendency of data to follow this type of model will be exhibited in a sample ACF that is large and decays very slowly at lags  $h = 12k$ , for  $k = 1, 2, \dots$ . If we subtract the effect of successive years from each other, we find that

$$(1 - B^{12})x_t = x_t - x_{t-12} = v_t + w_t - w_{t-12}.$$

This model is a stationary  $MA(1)_{12}$ , and its ACF will have a peak only at lag 12. In general, seasonal differencing can be indicated when the ACF decays slowly at multiples of some season  $s$ , but is negligible between the periods. Then, a *seasonal difference of order  $D$*  is defined as

$$\nabla_s^D x_t = (1 - B^s)^D x_t, \quad (3.162)$$

where  $D = 1, 2, \dots$ , takes positive integer values. Typically,  $D = 1$  is sufficient to obtain seasonal stationarity. Incorporating these ideas into a general model leads to the following definition.

**Definition 3.12** *The multiplicative seasonal autoregressive integrated moving average model, or SARIMA model is given by*

$$\Phi_P(B^s)\phi(B)\nabla_s^D\nabla^d x_t = \delta + \Theta_Q(B^s)\theta(B)w_t, \quad (3.163)$$

where  $w_t$  is the usual Gaussian white noise process. The general model is denoted as  $ARIMA(p, d, q) \times (P, D, Q)_s$ . The ordinary autoregressive and moving average components are represented by polynomials  $\phi(B)$  and  $\theta(B)$  of orders  $p$  and  $q$ , respectively, and the seasonal autoregressive and moving average components by  $\Phi_P(B^s)$  and  $\Theta_Q(B^s)$  of orders  $P$  and  $Q$  and ordinary and seasonal difference components by  $\nabla^d = (1 - B)^d$  and  $\nabla_s^D = (1 - B^s)^D$ .

### Example 3.48 An SARIMA Model

Consider the following model, which often provides a reasonable representation for seasonal, nonstationary, economic time series. We exhibit the equations for the model, denoted by  $ARIMA(0, 1, 1) \times (0, 1, 1)_{12}$  in the notation given above, where the seasonal fluctuations occur every 12 months. Then, with  $\delta = 0$ , the model (3.163) becomes

$$\nabla_{12}\nabla x_t = \Theta(B^{12})\theta(B)w_t$$

or

$$(1 - B^{12})(1 - B)x_t = (1 + \Theta B^{12})(1 + \theta B)w_t. \quad (3.164)$$

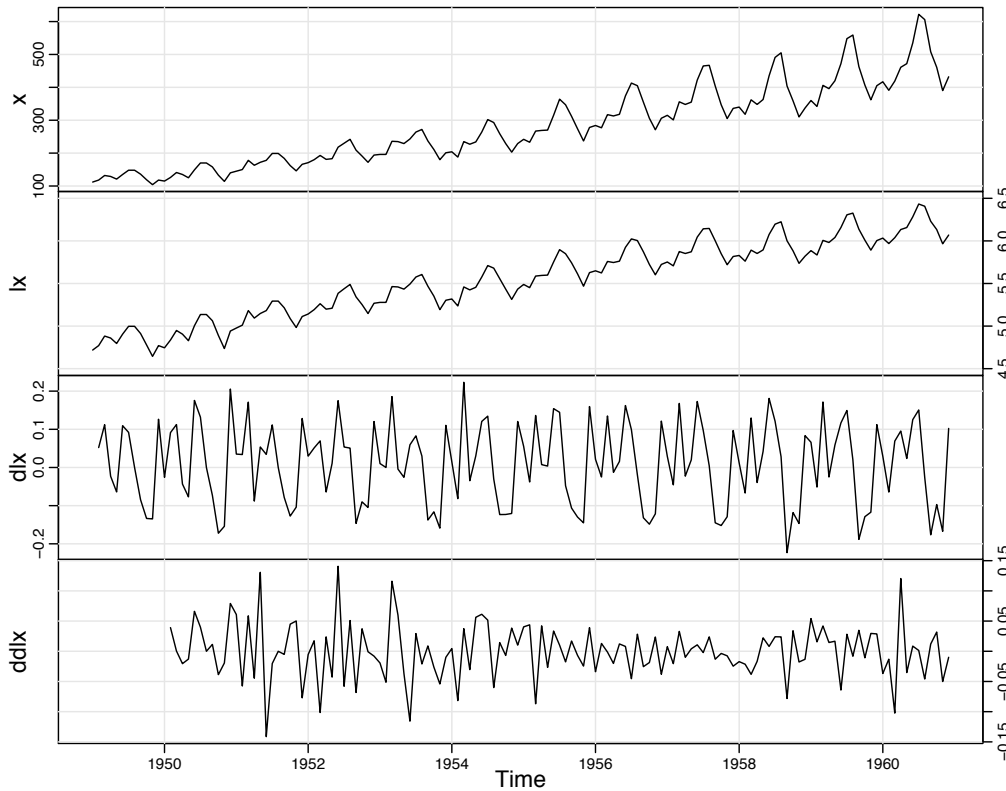
Expanding both sides of (3.164) leads to the representation

$$(1 - B - B^{12} + B^{13})x_t = (1 + \theta B + \Theta B^{12} + \Theta\theta B^{13})w_t,$$

or in difference equation form

$$x_t = x_{t-1} + x_{t-12} - x_{t-13} + w_t + \theta w_{t-1} + \Theta w_{t-12} + \Theta\theta w_{t-13}.$$

Note that the multiplicative nature of the model implies that the coefficient of  $w_{t-13}$  is the product of the coefficients of  $w_{t-1}$  and  $w_{t-12}$  rather than a free parameter. The multiplicative model assumption seems to work well with many seasonal time series data sets while reducing the number of parameters that must be estimated.



**Fig. 3.22.** R data set `AirPassengers`, which are the monthly totals of international airline passengers  $x$ , and the transformed data:  $lx = \log x_t$ ,  $dlx = \nabla \log x_t$ , and  $ddlx = \nabla_{12} \nabla \log x_t$ .

Selecting the appropriate model for a given set of data from all of those represented by the general form (3.163) is a daunting task, and we usually think first in terms of finding difference operators that produce a roughly stationary series and then in terms of finding a set of simple autoregressive moving average or multiplicative seasonal ARMA to fit the resulting residual series. Differencing operations are applied first, and then the residuals are constructed from a series of reduced length. Next, the ACF and the PACF of these residuals are evaluated. Peaks that appear in these functions can often be eliminated by fitting an autoregressive or moving average component in accordance with the general properties of Table 3.1 and Table 3.3. In considering whether the model is satisfactory, the diagnostic techniques discussed in Section 3.7 still apply.

#### Example 3.49 Air Passengers

We consider the R data set `AirPassengers`, which are the monthly totals of international airline passengers, 1949 to 1960, taken from Box & Jenkins (1970). Various plots of the data and transformed data are shown in Figure 3.22 and were obtained as follows:

```
x = AirPassengers
lx = log(x); dlx = diff(lx); ddxl = diff(dlx, 12)
```

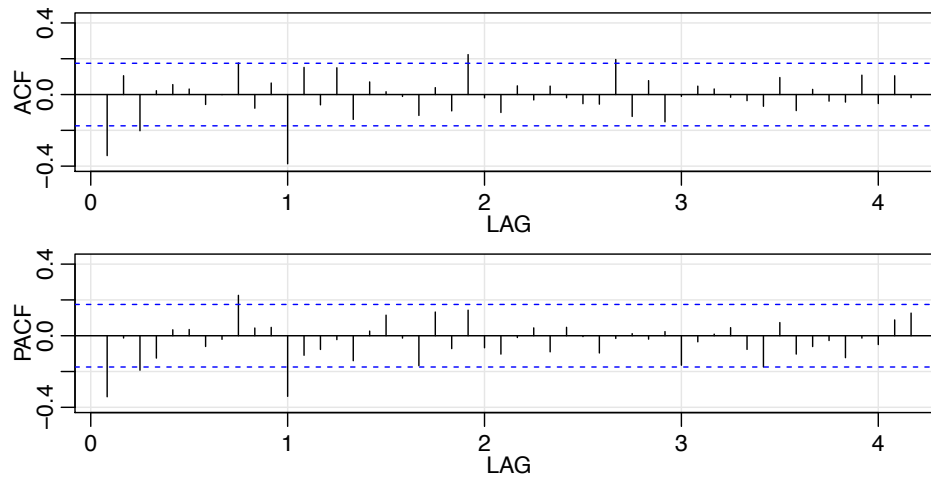


Fig. 3.23. Sample ACF and PACF of  $\text{ddlx}$  ( $\nabla_{12}\nabla \log x_t$ ).

```
plot.ts(cbind(x, lx, dlx, ddx), main="")
# below of interest for showing seasonal RW (not shown here):
par(mfrow=c(2,1))
monthplot(dlx); monthplot(ddlx)
```

Note that  $x$  is the original series, which shows trend plus increasing variance. The logged data are in  $lx$ , and the transformation stabilizes the variance. The logged data are then differenced to remove trend, and are stored in  $dlx$ . It is clear there is still persistence in the seasons (i.e.,  $dlx_t \approx dlx_{t-12}$ ), so that a twelfth-order difference is applied and stored in  $ddlx$ . The transformed data appears to be stationary and we are now ready to fit a model.

The sample ACF and PACF of  $\text{ddlx}$  ( $\nabla_{12}\nabla \log x_t$ ) are shown in Figure 3.23. The R code is:

```
acf2(ddlx, 50)
```

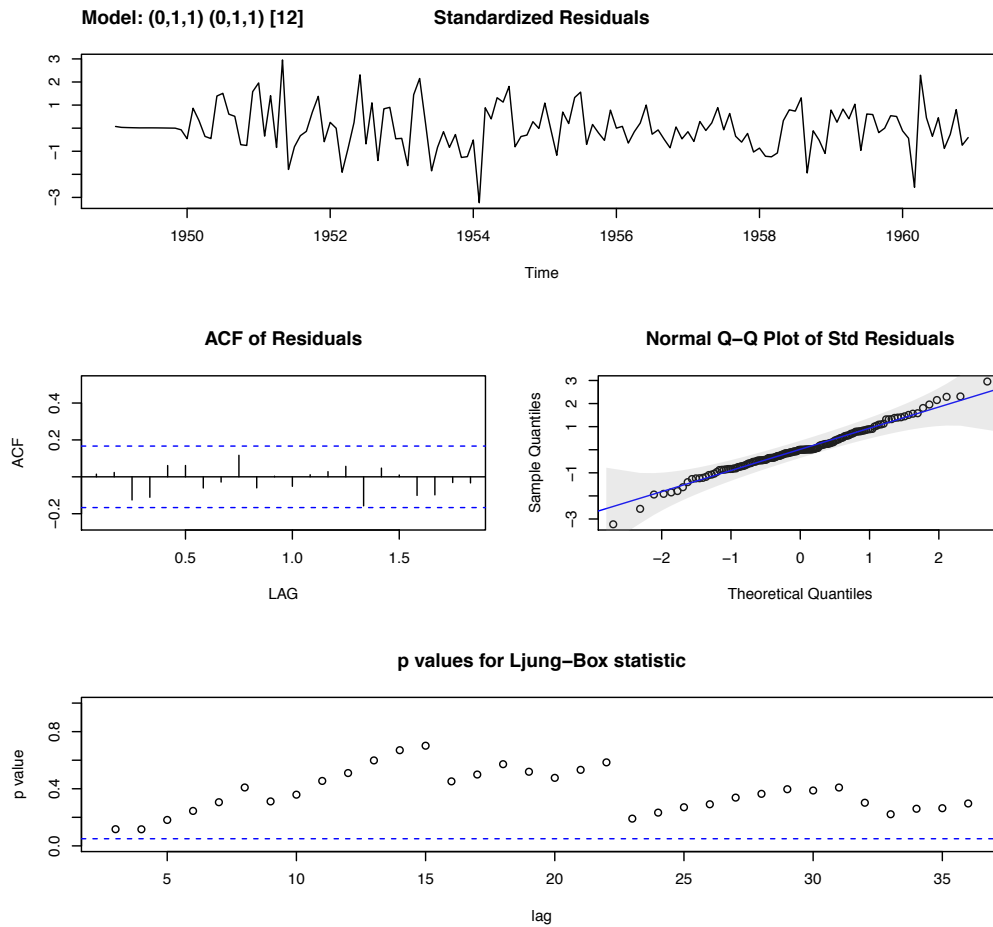
*Seasonal Component:* It appears that at the seasons, the ACF is cutting off a lag  $1s$  ( $s = 12$ ), whereas the PACF is tailing off at lags  $1s, 2s, 3s, 4s, \dots$ . These results implies an  $\text{SMA}(1)$ ,  $P = 0$ ,  $Q = 1$ , in the season ( $s = 12$ ).

*Non-Seasonal Component:* Inspecting the sample ACF and PACF at the lower lags, it appears as though both are tailing off. This suggests an  $\text{ARMA}(1, 1)$  within the seasons,  $p = q = 1$ .

Thus, we first try an  $\text{ARIMA}(1, 1, 1) \times (0, 1, 1)_{12}$  on the logged data:

```
sarima(lx, 1, 1, 1, 0, 1, 1, 12)
Coefficients:
      ar1      ma1      sma1
    0.1960 -0.5784 -0.5643
s.e. 0.2475 0.2132 0.0747
sigma^2 estimated as 0.001341
$AIC -5.5726 $AICc -5.556713 $BIC -6.510729
```

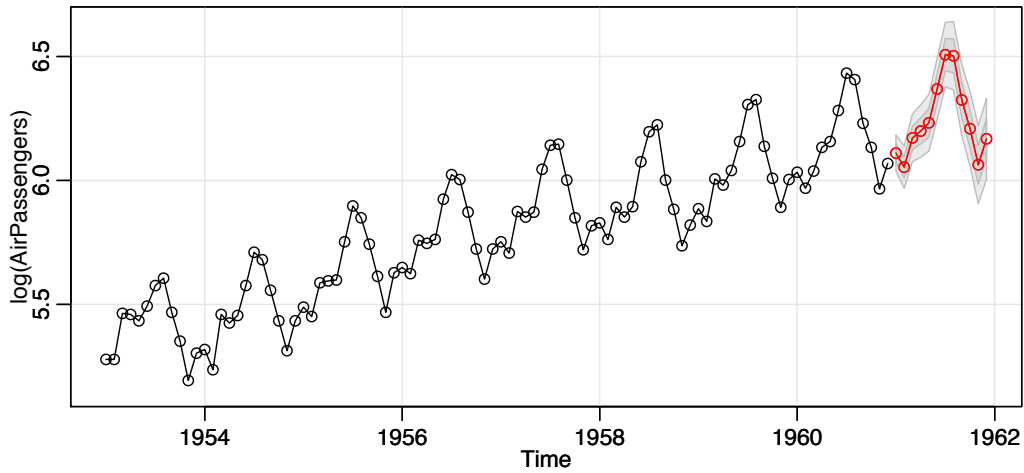
However, the AR parameter is not significant, so we should try dropping one parameter from the within seasons part. In this case, we try both an  $\text{ARIMA}(0, 1, 1) \times (0, 1, 1)_{12}$  and an  $\text{ARIMA}(1, 1, 0) \times (0, 1, 1)_{12}$  model:



**Fig. 3.24.** Residual analysis for the  $ARIMA(0, 1, 1) \times (0, 1, 1)_{12}$  fit to the logged air passengers data set.

```
sarima(lx, 0,1,1, 0,1,1,12)
Coefficients:
      mal      sma1
    -0.4018  -0.5569
s.e.    0.0896   0.0731
sigma^2 estimated as 0.001348
$AIC -5.58133  $AICc -5.56625  $BIC -6.540082
sarima(lx, 1,1,0, 0,1,1,12)
Coefficients:
      ar1      sma1
    -0.3395  -0.5619
s.e.    0.0822   0.0748
sigma^2 estimated as 0.001367
$AIC -5.567081  $AICc -5.552002  $BIC -6.525834
```

All information criteria prefer the  $ARIMA(0, 1, 1) \times (0, 1, 1)_{12}$  model, which is the model displayed in (3.164). The residual diagnostics are shown in Figure 3.24, and except for one or two outliers, the model seems to fit well.



**Fig. 3.25.** Twelve month forecast using the  $ARIMA(0, 1, 1) \times (0, 1, 1)_{12}$  model on the logged air passenger data set.

Finally, we forecast the logged data out twelve months, and the results are shown in [Figure 3.25](#).

```
sarima.for(lx, 12, 0, 1, 1, 0, 1, 1, 12)
```

## Problems

### Section 3.1

**3.1** For an MA(1),  $x_t = w_t + \theta w_{t-1}$ , show that  $|\rho_x(1)| \leq 1/2$  for any number  $\theta$ . For which values of  $\theta$  does  $\rho_x(1)$  attain its maximum and minimum?

**3.2** Let  $\{w_t; t = 0, 1, \dots\}$  be a white noise process with variance  $\sigma_w^2$  and let  $|\phi| < 1$  be a constant. Consider the process  $x_0 = w_0$ , and

$$x_t = \phi x_{t-1} + w_t, \quad t = 1, 2, \dots$$

We might use this method to simulate an AR(1) process from simulated white noise.

- (a) Show that  $x_t = \sum_{j=0}^t \phi^j w_{t-j}$  for any  $t = 0, 1, \dots$
- (b) Find the  $E(x_t)$ .
- (c) Show that, for  $t = 0, 1, \dots$ ,

$$\text{var}(x_t) = \frac{\sigma_w^2}{1 - \phi^2} (1 - \phi^{2(t+1)})$$

- (d) Show that, for  $h \geq 0$ ,

$$\text{cov}(x_{t+h}, x_t) = \phi^h \text{var}(x_t)$$

- (e) Is  $x_t$  stationary?



- (f) Argue that, as  $t \rightarrow \infty$ , the process becomes stationary, so in a sense,  $x_t$  is “asymptotically stationary.”
- (g) Comment on how you could use these results to simulate  $n$  observations of a stationary Gaussian AR(1) model from simulated iid  $N(0, 1)$  values.
- (h) Now suppose  $x_0 = w_0/\sqrt{1 - \phi^2}$ . Is this process stationary? *Hint:* Show  $\text{var}(x_t)$  is constant.

**3.3** Verify the calculations made in **Example 3.4** as follows.

- (a) Let  $x_t = \phi x_{t-1} + w_t$  where  $|\phi| > 1$  and  $w_t \sim \text{iid } N(0, \sigma_w^2)$ . Show  $E(x_t) = 0$  and  $\gamma_x(h) = \sigma_w^2 \phi^{-2} \phi^{-h} / (1 - \phi^{-2})$  for  $h \geq 0$ .
- (b) Let  $y_t = \phi^{-1} y_{t-1} + v_t$  where  $v_t \sim \text{iid } N(0, \sigma_w^2 \phi^{-2})$  and  $\phi$  and  $\sigma_w$  are as in part (a). Argue that  $y_t$  is causal with the same mean function and autocovariance function as  $x_t$ .

**3.4** Identify the following models as ARMA( $p, q$ ) models (watch out for parameter redundancy), and determine whether they are causal and/or invertible:

- (a)  $x_t = .80x_{t-1} - .15x_{t-2} + w_t - .30w_{t-1}$ .
- (b)  $x_t = x_{t-1} - .50x_{t-2} + w_t - w_{t-1}$ .

**3.5** Verify the causal conditions for an AR(2) model given in (3.28). That is, show that an AR(2) is causal if and only if (3.28) holds.

### Section 3.2

**3.6** For the AR(2) model given by  $x_t = -.9x_{t-2} + w_t$ , find the roots of the autoregressive polynomial, and then plot the ACF,  $\rho(h)$ .

**3.7** For the AR(2) series shown below, use the results of **Example 3.10** to determine a set of difference equations that can be used to find the ACF  $\rho(h)$ ,  $h = 0, 1, \dots$ ; solve for the constants in the ACF using the initial conditions. Then plot the ACF values to lag 10 (use `ARMAacf` as a check on your answers).

- (a)  $x_t + 1.6x_{t-1} + .64x_{t-2} = w_t$ .
- (b)  $x_t - .40x_{t-1} - .45x_{t-2} = w_t$ .
- (c)  $x_t - 1.2x_{t-1} + .85x_{t-2} = w_t$ .

### Section 3.3

**3.8** Verify the calculations for the autocorrelation function of an ARMA(1, 1) process given in **Example 3.14**. Compare the form with that of the ACF for the ARMA(1, 0) and the ARMA(0, 1) series. Plot the ACFs of the three series on the same graph for  $\phi = .6$ ,  $\theta = .9$ , and comment on the diagnostic capabilities of the ACF in this case.

**3.9** Generate  $n = 100$  observations from each of the three models discussed in **Problem 3.8**. Compute the sample ACF for each model and compare it to the theoretical values. Compute the sample PACF for each of the generated series and compare the sample ACFs and PACFs with the general results given in **Table 3.1**.

## Section 3.4

**3.10** Let  $x_t$  represent the cardiovascular mortality series (**cmort**) discussed in **Example 2.2**.

- (a) Fit an AR(2) to  $x_t$  using linear regression as in **Example 3.18**.
- (b) Assuming the fitted model in (a) is the true model, find the forecasts over a four-week horizon,  $x_{n+m}^n$ , for  $m = 1, 2, 3, 4$ , and the corresponding 95% prediction intervals.

**3.11** Consider the MA(1) series

$$x_t = w_t + \theta w_{t-1},$$

where  $w_t$  is white noise with variance  $\sigma_w^2$ .

- (a) Derive the minimum mean-square error one-step forecast based on the infinite past, and determine the mean-square error of this forecast.
- (b) Let  $\tilde{x}_{n+1}^n$  be the truncated one-step-ahead forecast as given in (3.92). Show that

$$E[(x_{n+1} - \tilde{x}_{n+1}^n)^2] = \sigma^2(1 + \theta^{2+2n}).$$

Compare the result with (a), and indicate how well the finite approximation works in this case.

**3.12** In the context of equation (3.63), show that, if  $\gamma(0) > 0$  and  $\gamma(h) \rightarrow 0$  as  $h \rightarrow \infty$ , then  $\Gamma_n$  is positive definite.

**3.13** Suppose  $x_t$  is stationary with zero mean and recall the definition of the PACF given by (3.55) and (3.56). That is, let

$$\epsilon_t = x_t - \sum_{i=1}^{h-1} a_i x_{t-i} \quad \text{and} \quad \delta_{t-h} = x_{t-h} - \sum_{j=1}^{h-1} b_j x_{t-j}$$

be the two residuals where  $\{a_1, \dots, a_{h-1}\}$  and  $\{b_1, \dots, b_{h-1}\}$  are chosen so that they minimize the mean-squared errors

$$E[\epsilon_t^2] \quad \text{and} \quad E[\delta_{t-h}^2].$$

The PACF at lag  $h$  was defined as the cross-correlation between  $\epsilon_t$  and  $\delta_{t-h}$ ; that is,

$$\phi_{hh} = \frac{E(\epsilon_t \delta_{t-h})}{\sqrt{E(\epsilon_t^2)E(\delta_{t-h}^2)}}.$$

Let  $R_h$  be the  $h \times h$  matrix with elements  $\rho(i-j)$  for  $i, j = 1, \dots, h$ , and let  $\rho_h = (\rho(1), \rho(2), \dots, \rho(h))'$  be the vector of lagged autocorrelations,  $\rho(h) = \text{corr}(x_{t+h}, x_t)$ . Let  $\tilde{\rho}_h = (\rho(h), \rho(h-1), \dots, \rho(1))'$  be the reversed vector. In addition, let  $x_t^h$  denote the BLP of  $x_t$  given  $\{x_{t-1}, \dots, x_{t-h}\}$ :

$$x_t^h = \alpha_{h1}x_{t-1} + \cdots + \alpha_{hh}x_{t-h},$$

as described in **Property 3.3**. Prove

$$\phi_{hh} = \frac{\rho(h) - \tilde{\rho}'_{h-1} R_{h-1}^{-1} \rho_h}{1 - \tilde{\rho}'_{h-1} R_{h-1}^{-1} \tilde{\rho}_{h-1}} = \alpha_{hh}.$$

In particular, this result proves **Property 3.4**.

*Hint:* Divide the prediction equations [see (3.63)] by  $\gamma(0)$  and write the matrix equation in the partitioned form as

$$\begin{pmatrix} R_{h-1} & \tilde{\rho}_{h-1} \\ \tilde{\rho}'_{h-1} & \rho(0) \end{pmatrix} \begin{pmatrix} \alpha_1 \\ \alpha_{hh} \end{pmatrix} = \begin{pmatrix} \rho_{h-1} \\ \rho(h) \end{pmatrix},$$

where the  $h \times 1$  vector of coefficients  $\alpha = (\alpha_{h1}, \dots, \alpha_{hh})'$  is partitioned as  $\alpha = (\alpha'_1, \alpha_{hh})'$ .

**3.14** Suppose we wish to find a prediction function  $g(x)$  that minimizes

$$MSE = E[(y - g(x))^2],$$

where  $x$  and  $y$  are jointly distributed random variables with density function  $f(x, y)$ .

(a) Show that MSE is minimized by the choice

$$g(x) = E(y \mid x).$$

*Hint:*

$$MSE = EE[(y - g(x))^2 \mid x].$$

(b) Apply the above result to the model

$$y = x^2 + z,$$

where  $x$  and  $z$  are independent zero-mean normal variables with variance one. Show that  $MSE = 1$ .

(c) Suppose we restrict our choices for the function  $g(x)$  to linear functions of the form

$$g(x) = a + bx$$

and determine  $a$  and  $b$  to minimize  $MSE$ . Show that  $a = 1$  and

$$b = \frac{E(xy)}{E(x^2)} = 0$$

and  $MSE = 3$ . What do you interpret this to mean?

**3.15** For an AR(1) model, determine the general form of the  $m$ -step-ahead forecast  $x_{t+m}^t$  and show

$$E[(x_{t+m} - x_{t+m}^t)^2] = \sigma_w^2 \frac{1 - \phi^{2m}}{1 - \phi^2}.$$

**3.16** Consider the ARMA(1,1) model discussed in [Example 3.8](#), equation (3.27); that is,  $x_t = .9x_{t-1} + .5w_{t-1} + w_t$ . Show that truncated prediction as defined in (3.91) is equivalent to truncated prediction using the recursive formula (3.92).

**3.17** Verify statement (3.87), that for a fixed sample size, the ARMA prediction errors are correlated.

### Section 3.5

**3.18** Fit an AR(2) model to the cardiovascular mortality series ([cmort](#)) discussed in [Example 2.2](#). using linear regression and using Yule–Walker.

- (a) Compare the parameter estimates obtained by the two methods.
- (b) Compare the estimated standard errors of the coefficients obtained by linear regression with their corresponding asymptotic approximations, as given in [Property 3.10](#).

**3.19** Suppose  $x_1, \dots, x_n$  are observations from an AR(1) process with  $\mu = 0$ .

- (a) Show the backcasts can be written as  $x_t^n = \phi^{1-t}x_1$ , for  $t \leq 1$ .
- (b) In turn, show, for  $t \leq 1$ , the backcasted errors are

$$\tilde{w}_t(\phi) = x_t^n - \phi x_{t-1}^n = \phi^{1-t}(1 - \phi^2)x_1.$$

- (c) Use the result of (b) to show  $\sum_{t=-\infty}^1 \tilde{w}_t^2(\phi) = (1 - \phi^2)x_1^2$ .
- (d) Use the result of (c) to verify the unconditional sum of squares,  $S(\phi)$ , can be written as  $\sum_{t=-\infty}^n \tilde{w}_t^2(\phi)$ .
- (e) Find  $x_t^{t-1}$  and  $r_t$  for  $1 \leq t \leq n$ , and show that

$$S(\phi) = \sum_{t=1}^n (x_t - x_t^{t-1})^2 / r_t.$$

**3.20** Repeat the following numerical exercise three times. Generate  $n = 500$  observations from the ARMA model given by

$$x_t = .9x_{t-1} + w_t - .9w_{t-1},$$

with  $w_t \sim \text{iid } N(0, 1)$ . Plot the simulated data, compute the sample ACF and PACF of the simulated data, and fit an ARMA(1, 1) model to the data. What happened and how do you explain the results?

**3.21** Generate 10 realizations of length  $n = 200$  each of an ARMA(1,1) process with  $\phi = .9, \theta = .5$  and  $\sigma^2 = 1$ . Find the MLEs of the three parameters in each case and compare the estimators to the true values.

**3.22** Generate  $n = 50$  observations from a Gaussian AR(1) model with  $\phi = .99$  and  $\sigma_w = 1$ . Using an estimation technique of your choice, compare the approximate asymptotic distribution of your estimate (the one you would use for inference) with the results of a bootstrap experiment (use  $B = 200$ ).

**3.23** Using [Example 3.32](#) as your guide, find the Gauss–Newton procedure for estimating the autoregressive parameter,  $\phi$ , from the AR(1) model,  $x_t = \phi x_{t-1} + w_t$ , given data  $x_1, \dots, x_n$ . Does this procedure produce the unconditional or the conditional estimator? *Hint:* Write the model as  $w_t(\phi) = x_t - \phi x_{t-1}$ ; your solution should work out to be a non-recursive procedure.

**3.24** Consider the stationary series generated by

$$x_t = \alpha + \phi x_{t-1} + w_t + \theta w_{t-1},$$

where  $E(x_t) = \mu$ ,  $|\theta| < 1$ ,  $|\phi| < 1$  and the  $w_t$  are iid random variables with zero mean and variance  $\sigma_w^2$ .

- Determine the mean as a function of  $\alpha$  for the above model. Find the autocovariance and ACF of the process  $x_t$ , and show that the process is weakly stationary. Is the process strictly stationary?
- Prove the limiting distribution as  $n \rightarrow \infty$  of the sample mean,

$$\bar{x} = n^{-1} \sum_{t=1}^n x_t,$$

is normal, and find its limiting mean and variance in terms of  $\alpha$ ,  $\phi$ ,  $\theta$ , and  $\sigma_w^2$ . (Note: This part uses results from [Appendix A](#).)

**3.25** A problem of interest in the analysis of geophysical time series involves a simple model for observed data containing a signal and a reflected version of the signal with unknown amplification factor  $a$  and unknown time delay  $\delta$ . For example, the depth of an earthquake is proportional to the time delay  $\delta$  for the P wave and its reflected form pP on a seismic record. Assume the signal, say  $s_t$ , is white and Gaussian with variance  $\sigma_s^2$ , and consider the generating model

$$x_t = s_t + a s_{t-\delta}.$$

- Prove the process  $x_t$  is stationary. If  $|a| < 1$ , show that

$$s_t = \sum_{j=0}^{\infty} (-a)^j x_{t-\delta j}$$

is a mean square convergent representation for the signal  $s_t$ , for  $t = 1, \pm 1, \pm 2, \dots$

- If the time delay  $\delta$  is assumed to be known, suggest an approximate computational method for estimating the parameters  $a$  and  $\sigma_s^2$  using maximum likelihood and the Gauss–Newton method.
- If the time delay  $\delta$  is an unknown integer, specify how we could estimate the parameters including  $\delta$ . Generate a  $n = 500$  point series with  $a = .9$ ,  $\sigma_w^2 = 1$  and  $\delta = 5$ . Estimate the integer time delay  $\delta$  by searching over  $\delta = 3, 4, \dots, 7$ .

**3.26 Forecasting with estimated parameters:** Let  $x_1, x_2, \dots, x_n$  be a sample of size  $n$  from a causal AR(1) process,  $x_t = \phi x_{t-1} + w_t$ . Let  $\hat{\phi}$  be the Yule–Walker estimator of  $\phi$ .

- (a) Show  $\hat{\phi} - \phi = O_p(n^{-1/2})$ . See [Appendix A](#) for the definition of  $O_p(\cdot)$ .  
 (b) Let  $x_{n+1}^n$  be the one-step-ahead forecast of  $x_{n+1}$  given the data  $x_1, \dots, x_n$ , based on the known parameter,  $\phi$ , and let  $\hat{x}_{n+1}^n$  be the one-step-ahead forecast when the parameter is replaced by  $\hat{\phi}$ . Show  $x_{n+1}^n - \hat{x}_{n+1}^n = O_p(n^{-1/2})$ .

### Section 3.6

**3.27** Suppose

$$y_t = \beta_0 + \beta_1 t + \dots + \beta_q t^q + x_t, \quad \beta_q \neq 0,$$

where  $x_t$  is stationary. First, show that  $\nabla^k x_t$  is stationary for any  $k = 1, 2, \dots$ , and then show that  $\nabla^k y_t$  is not stationary for  $k < q$ , but is stationary for  $k \geq q$ .

**3.28** Verify that the IMA(1,1) model given in (3.148) can be inverted and written as (3.149).

**3.29** For the ARIMA(1, 1, 0) model with drift,  $(1 - \phi B)(1 - B)x_t = \delta + w_t$ , let  $y_t = (1 - B)x_t = \nabla x_t$ .

- (a) Noting that  $y_t$  is AR(1), show that, for  $j \geq 1$ ,

$$y_{n+j}^n = \delta [1 + \phi + \dots + \phi^{j-1}] + \phi^j y_n.$$

- (b) Use part (a) to show that, for  $m = 1, 2, \dots$ ,

$$x_{n+m}^n = x_n + \frac{\delta}{1 - \phi} \left[ m - \frac{\phi(1 - \phi^m)}{(1 - \phi)} \right] + (x_n - x_{n-1}) \frac{\phi(1 - \phi^m)}{(1 - \phi)}.$$

*Hint:* From (a),  $x_{n+j}^n - x_{n+j-1}^n = \delta \frac{1 - \phi^j}{1 - \phi} + \phi^j (x_n - x_{n-1})$ . Now sum both sides over  $j$  from 1 to  $m$ .

- (c) Use (3.145) to find  $P_{n+m}^n$  by first showing that  $\psi_0^* = 1$ ,  $\psi_1^* = (1 + \phi)$ , and  $\psi_j^* - (1 + \phi)\psi_{j-1}^* + \phi\psi_{j-2}^* = 0$  for  $j \geq 2$ , in which case  $\psi_j^* = \frac{1 - \phi^{j+1}}{1 - \phi}$ , for  $j \geq 1$ . Note that, as in [Example 3.37](#), equation (3.145) is exact here.

**3.30** For the logarithm of the glacial varve data, say,  $x_t$ , presented in [Example 3.33](#), use the first 100 observations and calculate the EWMA,  $\tilde{x}_{t+1}^t$ , given in (3.151) for  $t = 1, \dots, 100$ , using  $\lambda = .25, .50$ , and  $.75$ , and plot the EWMA's and the data superimposed on each other. Comment on the results.

### Section 3.7

**3.31** In [Example 3.40](#), we presented the diagnostics for the MA(2) fit to the GNP growth rate series. Using that example as a guide, complete the diagnostics for the AR(1) fit.

**3.32** Crude oil prices in dollars per barrel are in [oil](#). Fit an ARIMA( $p, d, q$ ) model to the growth rate performing all necessary diagnostics. Comment.

**3.33** Fit an ARIMA( $p, d, q$ ) model to the global temperature data [globtemp](#) performing all of the necessary diagnostics. After deciding on an appropriate model, forecast (with limits) the next 10 years. Comment.

**3.34** Fit an ARIMA( $p, d, q$ ) model to the sulfur dioxide series, [so2](#), performing all of the necessary diagnostics. After deciding on an appropriate model, forecast the data into the future four time periods ahead (about one month) and calculate 95% prediction intervals for each of the four forecasts. Comment. (Sulfur dioxide is one of the pollutants monitored in the mortality study described in [Example 2.2](#).)

### Section 3.8

**3.35** Let  $S_t$  represent the monthly sales data in [sales](#) ( $n = 150$ ), and let  $L_t$  be the leading indicator in [lead](#).

- Fit an ARIMA model to  $S_t$ , the monthly sales data. Discuss your model fitting in a step-by-step fashion, presenting your (A) initial examination of the data, (B) transformations, if necessary, (C) initial identification of the dependence orders and degree of differencing, (D) parameter estimation, (E) residual diagnostics and model choice.
- Use the CCF and lag plots between  $\nabla S_t$  and  $\nabla L_t$  to argue that a regression of  $\nabla S_t$  on  $\nabla L_{t-3}$  is reasonable. [Note that in `lag2.plot()`, the first named series is the one that gets lagged.]
- Fit the regression model  $\nabla S_t = \beta_0 + \beta_1 \nabla L_{t-3} + x_t$ , where  $x_t$  is an ARMA process (explain how you decided on your model for  $x_t$ ). Discuss your results. [See [Example 3.45](#) for help on coding this problem.]

**3.36** One of the remarkable technological developments in the computer industry has been the ability to store information densely on a hard drive. In addition, the cost of storage has steadily declined causing problems of *too much data* as opposed to *big data*. The data set for this assignment is [cpg](#), which consists of the median annual retail price per GB of hard drives, say  $c_t$ , taken from a sample of manufacturers from 1980 to 2008.

- Plot  $c_t$  and describe what you see.
- Argue that the curve  $c_t$  versus  $t$  behaves like  $c_t \approx \alpha e^{\beta t}$  by fitting a linear regression of  $\log c_t$  on  $t$  and then plotting the fitted line to compare it to the logged data. Comment.

- (c) Inspect the residuals of the linear regression fit and comment.  
 (d) Fit the regression again, but now using the fact that the errors are autocorrelated. Comment.

**3.37** Redo **Problem 2.2** without assuming the error term is white noise.

### Section 3.9

**3.38** Consider the ARIMA model

$$x_t = w_t + \Theta w_{t-2}.$$

- (a) Identify the model using the notation  $\text{ARIMA}(p, d, q) \times (P, D, Q)_s$ .  
 (b) Show that the series is invertible for  $|\Theta| < 1$ , and find the coefficients in the representation

$$w_t = \sum_{k=0}^{\infty} \pi_k x_{t-k}.$$

- (c) Develop equations for the  $m$ -step ahead forecast,  $\tilde{x}_{n+m}$ , and its variance based on the infinite past,  $x_n, x_{n-1}, \dots$ .

**3.39** Plot the ACF of the seasonal  $\text{ARIMA}(0, 1) \times (1, 0)_{12}$  model with  $\Phi = .8$  and  $\theta = .5$ .

**3.40** Fit a seasonal ARIMA model of your choice to the chicken price data in [chicken](#). Use the estimated model to forecast the next 12 months.

**3.41** Fit a seasonal ARIMA model of your choice to the unemployment data in [unemp](#). Use the estimated model to forecast the next 12 months.

**3.42** Fit a seasonal ARIMA model of your choice to the unemployment data in [UnempRate](#). Use the estimated model to forecast the next 12 months.

**3.43** Fit a seasonal ARIMA model of your choice to the U.S. Live Birth Series ([birth](#)). Use the estimated model to forecast the next 12 months.

**3.44** Fit an appropriate seasonal ARIMA model to the log-transformed Johnson and Johnson earnings series ([jj](#)) of **Example 1.1**. Use the estimated model to forecast the next 4 quarters.

*The following problems require supplemental material given in [Appendix B](#).*

**3.45** Suppose  $x_t = \sum_{j=1}^p \phi_j x_{t-j} + w_t$ , where  $\phi_p \neq 0$  and  $w_t$  is white noise such that  $w_t$  is uncorrelated with  $\{x_k; k < t\}$ . Use the Projection Theorem, **Theorem B.1**, to show that, for  $n > p$ , the BLP of  $x_{n+1}$  on  $\overline{\text{sp}}\{x_k, k \leq n\}$  is

$$\hat{x}_{n+1} = \sum_{j=1}^p \phi_j x_{n+1-j}.$$



**3.46** Use the Projection Theorem to derive the Innovations Algorithm, **Property 3.6**, equations (3.77)-(3.79). Then, use **Theorem B.2** to derive the  $m$ -step-ahead forecast results given in (3.80) and (3.81).

**3.47** Consider the series  $x_t = w_t - w_{t-1}$ , where  $w_t$  is a white noise process with mean zero and variance  $\sigma_w^2$ . Suppose we consider the problem of predicting  $x_{n+1}$ , based on only  $x_1, \dots, x_n$ . Use the Projection Theorem to answer the questions below.

(a) Show the best linear predictor is

$$x_{n+1}^n = -\frac{1}{n+1} \sum_{k=1}^n k x_k.$$

(b) Prove the mean square error is

$$E(x_{n+1} - x_{n+1}^n)^2 = \frac{n+2}{n+1} \sigma_w^2.$$

**3.48** Use **Theorem B.2** and **Theorem B.3** to verify (3.117).

**3.49** Prove **Theorem B.2**.

**3.50** Prove **Property 3.2**.

