

## Estimation

### method of moments estimators:

AR(p):

When the process is AR(p),

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

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

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

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

In matrix notation, the Yule–Walker equations are

$$\Gamma_p \boldsymbol{\phi} = \boldsymbol{\gamma}_p, \quad \sigma_w^2 = \gamma(0) - \boldsymbol{\phi}' \boldsymbol{\gamma}_p, \quad (3.100)$$

where  $\Gamma_p = \{\gamma(k-j)\}_{j,k=1}^p$  is a  $p \times p$  matrix,  $\boldsymbol{\phi} = (\phi_1, \dots, \phi_p)'$  is a  $p \times 1$  vector, and  $\boldsymbol{\gamma}_p = (\gamma(1), \dots, \gamma(p))'$  is a  $p \times 1$  vector. Using the method of moments, we replace  $\gamma(h)$  in (3.100) by  $\hat{\gamma}(h)$  [see equation (1.34)] and solve

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

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

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

where  $\hat{\mathbf{R}}_p = \{\hat{\rho}(k-j)\}_{j,k=1}^p$  is a  $p \times p$  matrix and  $\hat{\boldsymbol{\rho}}_p = (\hat{\rho}(1), \dots, \hat{\rho}(p))'$  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  $\sigma_w^2$ .

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 3.28 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  $\rho(1) = .9/(1 + .9^2) = .497$ .

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

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

It can be shown that<sup>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, page 515, of Appendix A. 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

$$x_t = \mu + \phi(x_{t-1} - \mu) + w_t \quad (3.105)$$

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

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

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

$$L(\mu, \phi, \sigma_w^2) = (2\pi\sigma_w^2)^{-n/2} (1 - \phi^2)^{1/2} \exp \left[ -\frac{S(\mu, \phi)}{2\sigma_w^2} \right], \quad (3.106)$$

where

$$S(\mu, \phi) = (1 - \phi^2)(x_1 - \mu)^2 + \sum_{t=2}^n [(x_t - \mu) - \phi(x_{t-1} - \mu)]^2. \quad (3.107)$$

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 (3.106) with respect to  $\sigma_w^2$  and setting the result equal to zero, we see 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

$$\hat{\sigma}_w^2 = n^{-1}S(\hat{\mu}, \hat{\phi}), \quad (3.108)$$

where  $\hat{\mu}$  and  $\hat{\phi}$  are the MLEs of  $\mu$  and  $\phi$ , respectively. If we replace  $n$  in (3.108) by  $n - 2$ , we would obtain the unconditional least squares estimate of  $\sigma_w^2$ .

If, in (3.106), 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

$$l(\mu, \phi) = \log [n^{-1} S(\mu, \phi)] - n^{-1} \log(1 - \phi^2); \quad (3.109)$$

that is,  $l(\mu, \phi) \propto -2 \log L(\mu, \phi, \hat{\sigma}_w^2)$ .<sup>6</sup> Because (3.107) and (3.109) 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

$$\begin{aligned} L(\mu, \phi, \sigma_w^2 \mid x_1) &= \prod_{t=2}^n f_w [(x_t - \mu) - \phi(x_{t-1} - \mu)] \\ &= (2\pi\sigma_w^2)^{-(n-1)/2} \exp \left[ -\frac{S_c(\mu, \phi)}{2\sigma_w^2} \right], \end{aligned} \quad (3.110)$$

where the conditional sum of squares is

$$S_c(\mu, \phi) = \sum_{t=2}^n [(x_t - \mu) - \phi(x_{t-1} - \mu)]^2. \quad (3.111)$$

The conditional MLE of  $\sigma_w^2$  is

$$\hat{\sigma}_w^2 = S_c(\hat{\mu}, \hat{\phi}) / (n - 1), \quad (3.112)$$

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

$$S_c(\mu, \phi) = \sum_{t=2}^n [x_t - (\alpha + \phi x_{t-1})]^2. \quad (3.113)$$

The problem is now the linear regression problem stated in §2.2. 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

$$\hat{\mu} = \frac{\bar{x}_{(2)} - \hat{\phi}\bar{x}_{(1)}}{1 - \hat{\phi}} \quad (3.114)$$

$$\hat{\phi} = \frac{\sum_{t=2}^n (x_t - \bar{x}_{(2)})(x_{t-1} - \bar{x}_{(1)})}{\sum_{t=2}^n (x_{t-1} - \bar{x}_{(1)})^2}. \quad (3.115)$$

From (3.114) and (3.115), 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 (3.112) 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 (3.112).

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

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

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

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

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

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

$$L(\boldsymbol{\beta}, \sigma_w^2) = (2\pi\sigma_w^2)^{-n/2} [r_1(\boldsymbol{\beta})r_2(\boldsymbol{\beta}) \cdots r_n(\boldsymbol{\beta})]^{-1/2} \exp \left[ -\frac{S(\boldsymbol{\beta})}{2\sigma_w^2} \right], \quad (3.116)$$

where

$$S(\boldsymbol{\beta}) = \sum_{t=1}^n \left[ \frac{(x_t - x_t^{t-1}(\boldsymbol{\beta}))^2}{r_t(\boldsymbol{\beta})} \right]. \quad (3.117)$$

Both  $x_t^{t-1}$  and  $r_t$  are functions of  $\boldsymbol{\beta}$  alone, and we make that fact explicit in (3.116)-(3.117). Given values for  $\boldsymbol{\beta}$  and  $\sigma_w^2$ , the likelihood may be evaluated using the techniques of §3.5. Maximum likelihood estimation would now proceed by maximizing (3.116) with respect to  $\boldsymbol{\beta}$  and  $\sigma_w^2$ . As in the AR(1) example, we have

$$\hat{\sigma}_w^2 = n^{-1} S(\hat{\boldsymbol{\beta}}), \quad (3.118)$$

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

$$l(\boldsymbol{\beta}) = \log [n^{-1} S(\boldsymbol{\beta})] + n^{-1} \sum_{t=1}^n \log r_t(\boldsymbol{\beta}). \quad (3.119)$$