

Time Series (MATH5845)

Dr. Atefeh Zamani

Based on the notes by Prof. William T.M. Dunsmuir

T2 2025

Chapter 6

Maximum Likelihood Estimation for ARMA models

Contents

6.1	Innovations form of the Gaussian Likelihood.	90
6.2	Optimizing the likelihood	91
6.3	Large Sample Distribution of MLE's	93
6.4	Exercises	95

As mentioned before, Yule-Walker estimator are efficient estimator for Autoregressive processes but not for general ARMA processes. For more on this topic, you can refer to Shumway et al. [2000], P. 113-116. In this chapter, we are going to consider maximum likelihood estimators for ARMA models in time series.

6.1 Innovations form of the Gaussian Likelihood.

Let $\{X_t\}$ be a Gaussian time series with zero mean function and covariance function $\gamma(i, j) = E(X_i, X_j)$. Let \mathbf{X}_n and $\hat{\mathbf{X}}_n$ be as before. Let

$$\Gamma_n = E(\mathbf{X}_n \mathbf{X}_n^T) = [\gamma(i, j)]_{i,j=1}^n$$

and assume Γ_n^{-1} exists. The likelihood of \mathbf{X}_n is

$$L(\Gamma_n) = (2\pi)^{-n/2} (\det \Gamma_n)^{-1/2} \exp(-\mathbf{x}_n^T \Gamma_n^{-1} \mathbf{x}_n / 2).$$

In general this is a difficult expression to work with. Using the ideas used to develop the innovations algorithm we note that

$$\mathbf{X}_n = C_n(\mathbf{X}_n - \hat{\mathbf{X}}_n)$$

so that

$$\begin{aligned} \Gamma_n &= E(\mathbf{X}_n \mathbf{X}_n^T) \\ &= E[C_n(\mathbf{X}_n - \hat{\mathbf{X}}_n)(\mathbf{X}_n - \hat{\mathbf{X}}_n)^T C_n^T] \\ &= C_n D_n C_n^T \end{aligned}$$

where $D_n = \text{diag}(v_0, \dots, v_{n-1})$. The quadratic form in the likelihood exponent is therefore

$$\begin{aligned} \mathbf{x}_n^T \Gamma_n^{-1} \mathbf{x}_n &= (\mathbf{x}_n - \hat{\mathbf{x}}_n)^T D_n^{-1} (\mathbf{x}_n - \hat{\mathbf{x}}_n) \\ &= \sum_{j=1}^n (x_j - \hat{x}_j)^2 / v_{j-1} \end{aligned}$$

and the determinant term is $\det \Gamma_n = (\det C_n)^2 \det(D_n) = v_0 v_1 \dots v_{n-1}$. Using these expressions the likelihood becomes

$$L(\Gamma_n) = \frac{1}{\sqrt{(2\pi)^n v_0 v_1 \dots v_{n-1}}} \exp\left\{-\frac{1}{2} \sum_{j=1}^n \frac{(x_j - \hat{x}_j)^2}{v_{j-1}}\right\}.$$

This is the innovations form of the likelihood. It is a product of the distributions of $X_j - \hat{X}_j$ for $j = 1, \dots, n$ and can be re-written as

$$L(\Gamma_n) = \prod_{j=1}^n \frac{1}{\sqrt{2\pi v_{j-1}}} \exp\left\{-\frac{1}{2} \frac{(x_j - \hat{x}_j)^2}{v_{j-1}}\right\}.$$

Up to this point we have not expressed the $n(n+1)/2$ covariances $\gamma(i, j)$ as functions of a smaller set of parameters. Without this there are too many “unknown parameters” to estimate using the n successive values x_1, \dots, x_n observed on the time series. Further progress with

6.2. OPTIMIZING THE LIKELIHOOD

estimating a model requires us to model the covariances in some fashion. Let $\psi = (\psi_1, \dots, \psi_r)^T$ where r is fixed and finite and assume that all covariances, $\gamma(i, j)$, are defined in terms of ψ . Then the value $\hat{\psi}$ that maximises $L(\psi) = L(\Gamma_n(\psi))$ is called the maximum likelihood estimate of ψ .

We now apply these general results to estimation for ARMA(p, q) time series. Let

$$r_j = v_j / \sigma^2$$

and note that r_j and \hat{x}_j depend on $\phi = (\phi_1, \dots, \phi_p)^T$ and $\theta = (\theta_1, \dots, \theta_q)^T$ but are independent from σ^2 . Using this we can re-write $-2/n$ times the log-likelihood as follows.

$$l(\phi, \theta, \sigma^2) = \log \sigma^2 + \frac{1}{\sigma^2} \frac{1}{n} S(\phi, \theta) + \frac{1}{n} \sum_{j=1}^n \log r_{j-1} + \log(2\pi)$$

where

$$S(\phi, \theta) = \sum_{j=1}^n \frac{(x_j - \hat{x}_j)^2}{r_{j-1}}.$$

For fixed (ϕ, θ) , $l(\phi, \theta, \sigma^2)$ is minimised by

$$\hat{\sigma}^2(\phi, \theta) = \frac{1}{n} S(\phi, \theta)$$

and substituting this back we get the reduced expression

$$\begin{aligned} l(\phi, \theta) &= l(\phi, \theta, \hat{\sigma}^2(\phi, \theta)) \\ &= \log\left(\frac{1}{n} S(\phi, \theta)\right) + \frac{1}{n} \sum_{j=1}^n \log r_{j-1} + 1 + \log(2\pi). \end{aligned}$$

The values of (ϕ, θ) that minimise this are the maximum likelihood estimates $(\hat{\phi}, \hat{\theta})$ and the maximum likelihood estimate of σ^2 is $\hat{\sigma}^2 = \hat{\sigma}^2(\hat{\phi}, \hat{\theta})$.

For many applications the “determinant” term $\frac{1}{n} \sum_{j=1}^n \log r_{j-1}$ is negligible for large n and $l(\phi, \theta)$ can be minimised by minimising the sum of squares of scaled innovations $S(\phi, \theta)$. For example, for the AR(1) time series, $r_0 = 1/(1 - \phi^2)$ and $r_n = 1$ for $n \geq 1$ giving

$$\frac{1}{n} \sum_{j=1}^n \log r_{j-1} = -\frac{1}{n} \log(1 - \phi^2) \rightarrow 0$$

provided $|\phi| < 1$ and uniformly on any closed subset of this interval.

Least squares estimates are those obtained by choosing $(\tilde{\phi}, \tilde{\theta})$ to minimise $S(\phi, \theta)$ and forming

$$\tilde{\sigma}^2 = \frac{1}{n - p - q} S(\tilde{\phi}, \tilde{\theta}).$$

6.2 Optimizing the likelihood

As we learnt from Exercise 6.2, even for the simplest ARMA models, the likelihood is not a quadratic function of the parameters and equating first derivatives with respect to parameters

6.2. OPTIMIZING THE LIKELIHOOD

to zero will not result in easily solved linear equations. Compare this with linear regression for example.

There are many alternative ways to optimize non-linear functions and a proper discussion of this is beyond the scope of this course. However, we will review some key underlying concepts because these provide some insights also about how to obtain standard errors for the maximum likelihood estimates.

For this section, we refer to Shumway et al. [2000] (3rd edition, Section 3.6) for more details and, as they do, let $\beta = (\phi, \theta, \sigma^2)$ denote all the parameters of the ARMA equation. Let $l(\beta)$ denote $-1/n$ times the log-likelihood which is to be *minimized* over β .

Assume that there is unique global extremum of $l(\beta)$ in the interior of the parameter space allowed for β to range over – this is typically a subset of $\mathbb{R}^{r=p+q+1}$ in the ARMA(p, q) case. The negative of the score function is the vector of first derivatives

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

To find the minimum we can solve the first order condition that $l^{(1)} = 0$ to obtain the maximum likelihood estimator $\hat{\beta}$. In a neighbourhood of $\hat{\beta}$ the matrix of second derivatives (called the Hessian) given by

$$l^{(2)}(\beta) = \frac{\partial^2 l(\beta)}{\partial \beta \partial \beta'} = \left\{ \frac{\partial^2 l(\beta)}{\partial \beta_j \partial \beta_k} \right\}_{j,k=1}^r.$$

is strictly negative definite and therefore invertible.

Since the score equation is non-linear, methods such as Newton-Raphson and Gauss-Newton can be used. Consider Newton-Raphson. Let $\hat{\beta}^{(k)}$ be the current attempt to solve the score equation. Let $\hat{\beta}^{(k+1)}$ be an improvement on that attempt. If this improved solution was perfect we would have

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

where the approximation is based on Taylor expansion. Assuming equality throughout and solving for $\hat{\beta}^{(k+1)}$ gives

$$\hat{\beta}^{(k+1)} = \hat{\beta}^{(k)} + \left[-l^{(2)}(\hat{\beta}^{(k)}) \right]^{-1} l^{(1)}(\hat{\beta}^{(k)}) \quad (6.1)$$

Of course, for non-linear score functions, we will not get perfection in a single step and (6.1) needs to be iterated by replacing the old guess by the new one until convergence to the MLE $\hat{\beta}$ to some required accuracy is achieved. Note that, at convergence, we must have $l^{(1)}(\hat{\beta}) = 0$ as required.

For implementation of the Newton-Raphson method first and second derivatives of the likelihood with respect to parameters β need to be derived. For ARMA models these can be readily computed using recursions based on the ARMA model. Variations in methods typically involve how starting values for the recursions are calculated and whether or not the log determinant term is included. All these variations provide the same estimates in the limit as $n \rightarrow \infty$ and have the same asymptotic properties – see next section. Further details on the implementation of the Gauss-Newton method for ARMA models can be found in standard books on time series and in particular in Shumway and Stoffer (3rd edition, pages 129-131).

6.3 Large Sample Distribution of MLE's

Under specific regularity conditions it can be shown that, as $n \rightarrow \infty$,

$$\hat{\beta} \rightarrow \beta_0$$

$$\sqrt{n}(\hat{\beta} - \beta_0) \rightarrow^d \text{MVN}(0, \Omega(\beta_0))$$

where β_0 is the true parameter vector.

Also, it can be shown that

$$\hat{\Omega}(\beta_0) = \left[-l^{(2)}(\hat{\beta}) \right]^{-1} \rightarrow \Omega(\beta_0)$$

so that standard errors of the estimates can be calculated as

$$\text{s.e.}(\hat{\beta}_j) = \left[\frac{1}{n} \hat{\Omega}(\beta_0)_{j,j} \right]^{1/2}$$

The large sample multivariate normal distribution applies to the maximum likelihood estimates $(\hat{\phi}, \hat{\theta})$ as well as to the least squares estimates $(\tilde{\phi}, \tilde{\theta})$ and other variations of approximation to the likelihood.

Example 6.1 $AR(p)$.

$$\sqrt{n}(\hat{\phi} - \phi_0) \rightarrow^d \text{MVN}(0, \Omega(\phi_0))$$

where

$$\Omega(\phi_0) = \sigma_0^2 \Gamma_p^{-1}(\phi_0)$$

and in particular when $p = 1$

$$\Omega(\phi_0) = (1 - \phi_{0,1}^2)$$

and when $p = 2$

$$\Omega(\phi_0) = \begin{bmatrix} (1 - \phi_{0,2}^2) & -\phi_{0,1}(1 + \phi_{0,2}) \\ -\phi_{0,1}(1 + \phi_{0,2}) & (1 - \phi_{0,2}^2) \end{bmatrix}$$

Example 6.2 $MA(q)$.

$$\sqrt{n}(\hat{\theta} - \theta_0) \rightarrow^d \text{MVN}(0, \Omega(\theta_0))$$

where

$$\Omega(\theta_0) = \sigma_0^2 \Gamma_q^{-1}(\theta_0)$$

and $\Gamma_q(\theta_0)$ is the covariance matrix of an $AR(q)$ process

$$Y_t + \theta_{0,1}Y_{t-1} + \cdots + \theta_{0,q}Y_{t-q} = Z_t$$

in particular when $q = 1$

$$\Omega(\theta_0) = (1 - \theta_{0,1}^2)$$

and when $q = 2$

$$\Omega(\theta_0) = \begin{bmatrix} (1 - \theta_{0,2}^2) & \theta_{0,1}(1 + \theta_{0,2}) \\ \theta_{0,1}(1 + \theta_{0,2}) & (1 - \theta_{0,2}^2) \end{bmatrix}$$

6.3. LARGE SAMPLE DISTRIBUTION OF MLE'S

An intuitive explanation of why the form of the asymptotic variance for estimating the AR(1) and the MA(1) are of the same form is given on page 135 of Shumway and Stoffer 3rd Ed.

Note also that the precision of estimation improves as the parameters get closer to the limits of stationarity or invertibility.

Example 6.3 *ARMA(1,1).*

$$\begin{bmatrix} \sqrt{n}(\hat{\phi} - \phi_0) \\ \sqrt{n}(\hat{\theta} - \theta_0) \end{bmatrix} \rightarrow^d MVN(0, \Omega(\phi_0, \theta_0))$$

where

$$\Omega(\phi_0, \theta_0) = \frac{1 + \phi_0\theta_0}{(\phi_0 + \theta_0)^2} \begin{bmatrix} (1 - \phi_0^2)(1 + \phi_0\theta_0) & -(1 - \phi_0^2)(1 - \theta_0^2) \\ -(1 - \phi_0^2)(1 - \theta_0^2) & (1 - \theta_0^2)(1 + \phi_0\theta_0) \end{bmatrix}$$

Recall that in general

$$\Omega(\beta_0)^{-1} = \lim_{n \rightarrow \infty} \left[-l^{(2)}(\hat{\beta}) \right]$$

and applied to this example gives

$$\Omega(\phi_0, \theta_0)^{-1} = \begin{bmatrix} (1 - \phi_0^2)^{-1} & (1 + \phi_0\theta_0)^{-1} \\ (1 + \phi_0\theta_0)^{-1} & (1 - \theta_0^2)^{-1} \end{bmatrix}$$

and when $\phi_0 = -\theta_0$ this matrix has rank 1 (i.e. is non-invertible).

Note that all models with $\phi_0 = -\theta_0$ correspond to white noise (with no autocorrelation) and the ARMA model degenerates under this condition. The implication is that the second derivative matrix is near singular at points on the line $\phi_0 = -\theta_0$ consistent with the likelihood surface being approximately flat along that line. This situation is called ‘non-identifiability’ of parameters and is the simplest such case for ARMA models. For higher order ARMA models, if the two polynomials $\phi(z)$ and $\theta(z)$ have any roots in common the model is overspecified and the likelihood estimation will not converge. Additionally standard errors for the estimates will be very large and untrustworthy.

See also, Shumway et al. [2000], Example 3.35, for further discussion on overfitting ARMA models.

6.4 Exercises

Exercise 6.1 Consider a time series which starts at time $t = 0$ with a fixed value $X_0 = c$ and evolves, for $t > 0$, according to the model

$$X_t = \phi X_{t-1} + Z_t, \quad Z_t \sim \text{i.i.d. } N(0, \sigma^2).$$

In this formulation the quantity $k = \phi c$ is treated as a parameter to be estimated along with the other two parameters (ϕ, σ^2) .

1. Show that $E(X_t) = \phi E(X_{t-1})$ and hence that $\mu_t = E(X_t) = \phi^t c$.

2. Show that

$$\Phi(X - \mu) = Z$$

in which

$$\Phi = \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ -\phi & 1 & 0 & \ddots & \vdots \\ 0 & -\phi & \ddots & \ddots & 0 \\ \vdots & \ddots & \ddots & 1 & 0 \\ 0 & \dots & 0 & -\phi & 1 \end{bmatrix}$$

and $X = (X_1, \dots, X_n)^T$, $\mu = (\mu_1, \dots, \mu_n)^T$, $Z = (Z_1, \dots, Z_n)^T$.

3. Using part b) show that $\Gamma_n = \text{cov}(X)$ satisfies

$$\Phi \Gamma_n \Phi^T = \sigma^2 I_n$$

and hence

$$\Gamma_n^{-1} = \sigma^{-2} \Phi^T \Phi$$

$$\det(\Gamma_n) = \sigma^{2n}$$

and

$$\begin{aligned} -\frac{1}{2}(X - \mu)^T \Gamma_n^{-1} (X - \mu) &= -\frac{1}{2\sigma^2} Z^T Z \\ &= -\frac{1}{2\sigma^2} \left[\sum_{t=2}^n (X_t - \phi X_{t-1})^2 + (X_1 - k)^2 \right] \end{aligned}$$

4. Hence show that the maximum likelihood estimates are $\hat{k} = X_1$ and

$$\hat{\phi} = \frac{\sum_{t=2}^n X_{t-1} X_t}{\sum_{t=2}^n X_{t-1}^2}.$$

5. Assuming the true values of the parameters in the model are known, derive the k -step ahead forecast \hat{X}_{n+k} given observations on X_1, \dots, X_n . Also derive the standard deviation of this forecast and use it to construct a 95% prediction interval for X_{n+k} .

6. State, giving reasons, an appropriate approximate distribution for $\hat{\phi}$ as $n \rightarrow \infty$.

[2 pages total]

Exercise 6.2 For the AR(1) process derive the maximum likelihood estimates, the conditional least squares estimates and the unconditional least squares estimates.