

## CHAPTER 5 - Modeling and Forecasting with ARMA Processes .

### 5.1 Preliminary Estimation

The Yule-Walker and Burg procedures apply to the fitting of pure AR models (Although the former can be adapted to models with  $q > 0$ , it will be less efficient.)

The innovation and Hannan-Rissanen algorithms are used to provide preliminary estimates of the ARMA when  $q > 0$ .

#### Yule-Walker

The Yule-Walker coefficients  $\hat{\phi}_1, \dots, \hat{\phi}_P$  are precisely the coefficients of the best linear predictor of  $x_{t+1}$  in terms of  $x_t, \dots, x_1$  under the assumption the ACF of  $\{x_t\}$  coincides with the sample ACF at lags  $1, \dots, n$ .

For a pure AR model,  $\Phi(B)x_t = z_t \quad (*)$

$$\text{for } h \geq 0, E(x_t x_{t+h}) = E[(z_t + \sum_{r=1}^P \phi_r x_{t-r}) x_{t+h}] \\ = \sigma^2 + \sum_{r=1}^P \phi_r E(x_{t-r} x_{t+h})$$

$$\Rightarrow r(h) = \sum_{r=1}^P \phi_r r(1|t-r|)$$

$$\text{for } h=0, r(0) = \sigma^2 + \sum_{r=1}^P \phi_r r(r)$$

Then, we obtain the Yule-Walker equations

$$P_P \phi = \delta_P \text{ and } \sigma^2 = r(0) - \phi^T \delta_P$$

where  $P_P$  is the covariance matrix  $[r(i-j)]_{i,j=1}^P$  and  $\delta_P = (r(1), \dots, r(P))^T$ .

If we replace covariances  $r(ij)$  by the corresponding sample covariances  $\hat{r}(ij)$ , we obtain a set of equations for the so-called Yule-Walker estimators  $\hat{\phi}$  and  $\hat{\sigma}^2$  of  $\phi$  and  $\sigma^2$ .

$$\hat{P}_P \hat{\phi} = \hat{\delta}_P \text{ and } \hat{\sigma}^2 = \hat{r}(0) - \hat{\phi}^T \hat{\delta}_P$$

where  $\hat{P}_P = [\hat{r}(i-j)]_{i,j=1}^P$  and  $\hat{\delta}_P = (\hat{r}(1), \dots, \hat{r}(P))^T$ .

If  $\hat{\gamma}(0) > 0$ , then  $\hat{P}_m$  is nonsingular for every  $m=1, 2, \dots$

Sample Yule-Walker Equations:

$$\hat{\phi} = (\hat{\phi}_1, \dots, \hat{\phi}_p)' = \hat{R}_p^{-1} \hat{\rho}_p \quad (5.1.7)$$

and

$$\hat{\sigma}^2 = \hat{\gamma}(0) [1 - \hat{\rho}_p' \hat{R}_p^{-1} \hat{\rho}_p], \quad (5.1.8)$$

where  $\hat{\rho}_p = (\hat{\rho}(1), \dots, \hat{\rho}(p))'$  =  $\hat{\gamma}_p / \hat{\gamma}(0)$ .

It is often the case that moment estimators, i.e. estimators that like  $\hat{\phi}$  are obtained by equating theoretical and sample moments, have much higher variances than estimators obtained by alternative methods such as maximum likelihood. However, the Yule-Walker estimators of the coefficients  $\phi_1, \dots, \phi_p$  of an AR(p) process have approximately the same distribution for large samples as the corresponding maximum likelihood estimators.

### Large-Sample Distribution of Yule-Walker Estimators

For large  $n$  and true AR(p) model,  $\hat{\phi}_p \xrightarrow{d} N(\phi, \frac{1}{n} \sigma^2 \hat{P}_p^{-1})$

- If  $m < p$ , missing  $\phi$ 's!

*fitting truth*

$\Rightarrow \hat{\phi}_m$  will not be consistent for the truth.

- If  $m > p$ , NOT missing  $\phi$ 's!

$$X_t = \phi_1 X_{t-1} + \dots + \phi_p X_{t-p} + \epsilon_t \quad (t=p+1, \dots, n)$$

"True" model with  $m$   $\phi$ 's just needs  $m-p$  0-valued  $\phi$ 's for the remaining terms.

## Burg's algorithm

Burg's algorithm estimates the PACF  $\{\phi_{11}, \phi_{22}, \dots\}$  by successively minimizing sum of squares of forward and backward one-step prediction errors with respect to the coefficients  $\phi_{ii}$ .

Let  $u_{it}(t) = x_{n+i-t} - \hat{x}_{n+i-t}^{(i)}$  for  $t=1, \dots, n$  (error) where  
 $\hat{x}_{n+i-t}^{(i)}$  is the BLF of observation  $x_{n+i-t}$  given the preceding  $i$  observations.

e.g.  $n=15, i=3, t=5$

$$u_5(5) = x_{14} - \hat{x}_{14}^{(3)}$$

• BLF for  $x_{14}$  using  $x_3, x_4, x_5$

Let  $v_{it}(t) = x_{n+i-t} - \hat{x}_{n+i-t}^{(i)}$  for  $t=1, \dots, n$  where  
 $\hat{x}_{n+i-t}^{(i)}$  is BLF of  $x_{n+i-t}$  using  $i$  subsequent (immediately following) observations.

e.g.  $v_3(5) = x_1 - \hat{x}_1^{(3)}$

• BLF for  $x_1$  using  $x_2, x_3, x_4$

$u_{it}(t)$  and  $v_{it}(t)$  satisfy

$$u_{it}(t) = v_{it}(t) = x_{n+i-t}$$

$$u_{it}(t) = u_{i-1}(t-1) - \phi_{ii} v_{i-1}(t)$$

$$v_{it}(t) = v_{i-1}(t) - \phi_{ii} u_{i-1}(t-1)$$

equal (PACF we want)

Note  $u_{it}(t)$  and  $v_{it}(t)$  are functions of the data!

Burg's estimate is computed by innovation

$$\phi_{11}^{(B)} = \min_{\phi_{11}} \frac{1}{2(n-1)} \sum_{t=2}^n (u_{it}(t)^2 + v_{it}(t)^2)$$

$$\phi_{22}^{(B)} = \min_{\phi_{22}} \frac{1}{2(n-2)} \sum_{t=3}^n (u_{it}(t)^2 + v_{it}(t)^2) \rightarrow \text{use } \phi_{11}^{(B)} \text{ in } u_{it}(t) \text{ and } v_{it}(t)$$

until we get  $\phi_{pp}^{(B)}$ .

The above recursion is equivalent to solving the recursion

$$\begin{aligned}
 d(i) &= \sum_{t=2}^n (u_t^{(t-1)} + v_t^{(t)}) \\
 \phi_{ii}^{(B)} &= \frac{\sum_{t=i+1}^n v_t^{(t)} u_t^{(t-1)}}{d(i)} \\
 d(i+1) &= (1 - \phi_{ii}^{(B)}) d(i) - v_i^{(i+1)} - u_i^{(i)} \\
 \alpha_i^{(B)} &= \frac{(1 - \phi_{ii}^{(B)}) d(i)}{d(i+1)}
 \end{aligned}$$

PACF coefficients:  $\hat{\phi}_{11}^{(B)}, \hat{\phi}_{22}^{(B)}, \dots, \hat{\phi}_{pp}^{(B)}$

The coefficients for AR(p) are  $\phi_1, \dots, \phi_p$

⇒ Durbin-Levinson algorithm

① Compute  $\phi_{nn}$ .

② Compute  $\phi_{nn}, \dots, \phi_{11}$  using  $\phi_{nn}$  and  $\phi_{n+1,n}, \dots, \phi_{1,n}$

Replace the step where we compute  $\phi_{nn}$  with  $\hat{\phi}_{nn}$  from Burg's algorithm.

Both algorithms are consistent for true  $\phi_1, \dots, \phi_p$  and they have the same asymptotic distribution

$$\hat{\phi}_p^{(B)} \sim N(\phi_p, \frac{1}{n} \sigma^2 P_p^{-1})$$

### Innovation algorithm

Substitute  $\hat{\epsilon}_{t,j}$  for  $\epsilon_{t,j}$  in the innovation algorithm (KCF,j)

For MA(q),  $X_t = Z_t + \theta_1 Z_{t-1} + \dots + \theta_q Z_{t-q} \quad Z_t \sim WN(0, \sigma^2)$

If we apply innovation algorithm to obtain  $\hat{x}_n$ , we have that  
 $(\hat{\theta}_{m1}, \dots, \hat{\theta}_{mq})$  at n-step

$\hat{\theta}_{m1}, \dots, \hat{\theta}_{mq}$  for sufficiently large m will be approximately equal to  $\theta_1, \dots, \theta_q$

Now we can plug in  $\hat{R}_{(i,j)} = \hat{\epsilon}_{(i,j)}$  and obtain  $\hat{\theta}_{m1}, \dots, \hat{\theta}_{mq}$   
and then for large enough m

$$(\hat{\theta}_{m1}, \dots, \hat{\theta}_{mq}) \sim N(\varnothing, \frac{1}{n} \text{Var}(Z))$$

### Hannan - Rissanen algorithm

Basic idea:

- ① Use Yule-Walker as fit AR(p), then use AR coefficients to estimate S<sub>t,p</sub> process. (if invertible)

$$\hat{Z}_t = X_t - \hat{\phi}_1 X_{t-1} - \hat{\phi}_2 X_{t-2} - \dots - \hat{\phi}_p X_{t-p}$$

- ②  $X_t = \phi_1 X_{t-1} + \phi_2 X_{t-2} + \dots + \phi_p X_{t-p} + \theta_1 \hat{Z}_{t-1} + \dots + \theta_q \hat{Z}_{t-q} + \hat{Z}_t$

Solve for new values of  $\phi$  and  $\theta$  via least squares

→ use multiple linear regression to estimate  $\phi$  and  $\theta$  by minimizing

$$\hat{\Sigma}(\phi, \theta) = \sum_{t=1}^n (X_t - \phi_1 X_{t-1} - \dots - \phi_p X_{t-p} - \theta_1 \hat{Z}_{t-1} - \dots - \theta_q \hat{Z}_{t-q})^2$$

Problem Sometimes  $\hat{Z}_t$  is not good estimator

Moment estimator (don't need assumption for distribution)

- ① Yule-Walker for AR(p) (ACF)
- ② Burg's algorithm for AR(p) (PACF)
- ③ Innovation algorithm for MA(q) (ACF)
- ④ Hannan-Rissanen for ARMA(p,q) (ACF/PACF + MLR)

Problem Might not perform well if the model is not correctly specified.

## 5.2 Maximum Likelihood Estimation

Assumption Assume  $\{X_t\}$  is a Gaussian time series with mean zero and  $K_{(i,j)} = E(X_i X_j)$

Let  $\Omega_n = (\theta_1, \dots, \theta_n)$ . Then the likelihood for a mean-zero Gaussian time series is the multivariate normal density.

$$L(\Omega_n) = (2\pi)^{-n/2} (\det \Gamma_n)^{-1/2} \exp \left( -\frac{1}{2} \mathbf{X}'_n \Gamma_n^{-1} \mathbf{X}_n \right)$$

where  $\Gamma_n = E(\mathbf{X}_n \mathbf{X}_n^T) = \begin{pmatrix} K_{(1,1)} & \cdots & K_{(1,n)} \\ \vdots & & \vdots \\ K_{(n,1)} & \cdots & K_{(n,n)} \end{pmatrix}$

Recall that  $K_{(i,j)}$  is a function of  $\alpha_1, \dots, \alpha_p$  and  $\sigma^2$  of ARMA(p,q).

In general, maximizing  $L(\Omega_n)$  will be impossible due to the presence of  $\Gamma_n^{-1}$  and  $\det(\Gamma_n)$ . However, the innovation algorithm allows us to simplify things.

Let  $\hat{X}_j = E(X_j | X_{j-1}, \dots, X_1) = P_{j-1} X_j$

Let  $\{\theta_{i,j}\}$  be the coefficients from the innovations algorithm and  $X_j - \hat{X}_j$  is the prediction error of the  $j$ th observation.

We know that  $X_n = C_n (X_n - \hat{X}_n)$  where

$$C_n = \begin{pmatrix} 1 & 0 & 0 & \cdots \\ \theta_{n1} & 1 & 0 & \cdots \\ \theta_{n2} & \theta_{n1} & 1 & \cdots \\ \vdots & \vdots & \vdots & \ddots \end{pmatrix}$$

Now we also have that

$$E[(x_i - \hat{x}_i)(x_j - \hat{x}_j)] = 0 \quad (\text{induction})$$

i.e. innovation for  $i \neq j$  are uncorrelated  $\forall i, j, i \neq j$

Therefore,  $\underline{x}_n - \hat{\underline{x}}_n$  has a diagonal covariance matrix

$$D_n = \text{diag}(v_0, \dots, v_{n-1})$$

where  $v_i$  comes from the innovation algorithm.

because the variance of the innovations is

$$\begin{aligned} \text{Var}(\underline{x}_n - \hat{\underline{x}}_n) &= E[(\underline{x}_n - \hat{\underline{x}}_n)^2] \quad \text{b.c. } E(\underline{x}_n - \hat{\underline{x}}_n) = 0 \\ &= \text{MSE} \end{aligned}$$

$\Gamma_n$

$$\begin{aligned} \text{Var}(\underline{x}_n) &= E(\underline{x}_n \underline{x}_n^T) = E([C_n(\underline{x}_n - \hat{\underline{x}}_n)][C_n(\underline{x}_n - \hat{\underline{x}}_n)]^T) \\ &= E[C_n(\underline{x}_n - \hat{\underline{x}}_n)(\underline{x}_n - \hat{\underline{x}}_n)^T C_n^T] \\ &= C_n E[(\underline{x}_n - \hat{\underline{x}}_n)(\underline{x}_n - \hat{\underline{x}}_n)^T] C_n^T \\ &= C_n \text{Var}(\underline{x}_n - \hat{\underline{x}}_n) C_n^T \\ &= C_n D_n C_n^T \end{aligned}$$

$$\begin{aligned} \text{So in } L(\Gamma_n), \quad \underline{x}_n^T \Gamma_n^{-1} \underline{x}_n &= (\underline{x}_n - \hat{\underline{x}}_n)^T C_n^T \Gamma_n^{-1} C_n (\underline{x}_n - \hat{\underline{x}}_n) \\ &= (\underline{x}_n - \hat{\underline{x}}_n)^T C_n^T (C_n^T)^{-1} D_n^{-1} C_n^T C_n (\underline{x}_n - \hat{\underline{x}}_n) \\ &= (\underline{x}_n - \hat{\underline{x}}_n)^T D_n^{-1} (\underline{x}_n - \hat{\underline{x}}_n) \\ &= \sum_{j=1}^n \frac{(x_j - \hat{x}_j)^2}{v_{j-1}} \quad \text{since } D_n \text{ diagonal}. \end{aligned}$$

$$\text{and } \det(\Gamma_n) = \det(C_n D_n C_n^T) = \det(C_n) \det(D_n) \det(C_n^T)$$

$$\text{lower diagonal} = \det(C_n) \det(D_n)$$

$$\text{matrix} = v_0 \cdot v_1 \cdots v_{n-1}$$

$$\Rightarrow \det(d_1 \cdots d_n) = v_0 \cdot v_1 \cdots v_{n-1}$$

$$= \prod_{j=0}^{n-1} v_j$$

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

If  $\Gamma_n$  depends on a finite number of parameters, the MLE maximizes  $L$  for any given datasets.

Remember

$$\hat{x}_{n+1} = \begin{cases} \sum_{j=1}^q \theta_j (x_{n+1-j} - \hat{x}_{n+1-j}) & 1 \leq n \leq m \\ \phi_1 x_{n+1} + \dots + \phi_p x_{n+1-p} + \sum_{j=1}^q \theta_j (x_{n+1-j} - \hat{x}_{n+1-j}) & n \geq m \end{cases}$$

$E[(x_{n+1} - \hat{x}_n)^2] = \alpha^2 r_n = v_n$  where  $r_n = E[(w_{n+1} - \hat{w}_n)^2]$  from the recursion in the innovation algorithm.

The Gaussian Likelihood for an ARMA Process:

$$L(\phi, \theta, \sigma^2) = \frac{1}{\sqrt{(2\pi\sigma^2)^n r_0 \cdots r_{n-1}}} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{j=1}^n \frac{(X_j - \hat{X}_j)^2}{r_{j-1}} \right\}. \quad (5.2.9)$$

$$\Rightarrow \ell(\phi, \theta, \sigma^2) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2} \sum_{j=0}^{n-1} \log r_j - \frac{1}{2\sigma^2} \sum_{j=1}^n \frac{(X_j - \hat{X}_j)^2}{r_{j-1}}$$

Note that  $\hat{X}_j$  and  $r_{j-1}$  are both function of model parameters  $\phi, \theta$  but not of  $\sigma^2$ .

Maximum Likelihood Estimators:

$$\hat{\sigma}^2 = n^{-1} S(\hat{\phi}, \hat{\theta}), \quad \text{by setting } \frac{\partial \ell}{\partial \sigma^2} = 0 \quad (5.2.10)$$

where

$$S(\hat{\phi}, \hat{\theta}) = \sum_{j=1}^n (X_j - \hat{X}_j)^2 / r_{j-1}, \quad (5.2.11)$$

and  $\hat{\phi}, \hat{\theta}$  are the values of  $\phi, \theta$  that minimize

$$\ell(\phi, \theta) = \ln(n^{-1} S(\phi, \theta)) + n^{-1} \sum_{j=1}^n \ln r_{j-1}. \quad (5.2.12)$$

No closed form  
use numerical optimization algorithm.

Most optimizer methods require some initial values for  $\phi, \theta$   
(Other estimators can be used as starting point e.g. Burg's)

< Lazy maximization ?

$$\text{try to minimize } S(\phi, \theta) = \sum_{j=1}^n \frac{(X_j - \hat{X}_j)^2}{r_{j-1}} \quad (\text{least square approach})$$

But can lead to noncausal or noninvertible solutions without additional constraints.

And the estimate of  $\sigma^2$  would be

$$\hat{\sigma}^2 = \frac{S(\hat{\phi}^{LS}, \hat{\theta}^{LS})}{n - p - q}$$

The MLE for large  $n$  will be approximately normal.

Let  $\hat{\beta} = (\phi_1, \dots, \phi_p, \theta_1, \dots, \theta_q)$

$\hat{\beta} = (\hat{\phi}_1, \dots, \hat{\phi}_p, \hat{\theta}_1, \dots, \hat{\theta}_q)$  MLE from above

For  $n$  large,  $\hat{\beta} \sim MVN(\hat{\beta}, V(\hat{\beta}))$

$$\text{where } V(\hat{\beta}) = \sigma^2 \begin{pmatrix} E(UU^T) & E(UV^T) \\ E(VU^T) & E(VV^T) \end{pmatrix}^{-1}$$

$P \times P$        $P \times Q$   
 $Q \times P$        $Q \times Q$

where  $\Phi(B)U = Z_U$  and  $\Theta(B)V = Z_V$

$\downarrow$   
AR(CP) with  
 $\phi_1, \dots, \phi_p$

$\downarrow$   
AR(CQ) with  
 $\theta_1, \dots, \theta_q$

In practice, we can approximate  $V(\hat{\beta})$  with  $H^{-1}(B)$  where

$$H = \left( \frac{\partial \ell(B)}{\partial B_i \partial B_j} \right)_{i,j=1}^{P \times Q}$$

NOTE  $H$  is estimated during many optimization algorithms.