

Paper 454-2013

The Box-Jenkins Methodology for Time Series Models

Theresa Hoang Diem Ngo, Warner Bros. Entertainment Group, Burbank, CA

ABSTRACT

A time series is a set of values of a particular variable that occur over a period of time in a certain pattern. The most common patterns are increasing or decreasing trend, cycle, seasonality, and irregular fluctuations (Bowerman, O'Connell, and Koehler 2005). To model a time series event as a function of its past values, analysts identify the pattern with the assumption that the pattern will persist in the future. Applying the Box-Jenkins methodology, this paper emphasizes how to identify an appropriate time series model by matching behaviors of the sample autocorrelation function (ACF) and partial autocorrelation function (PACF) to the theoretical autocorrelation functions. In addition to model identification, the paper examines the significance of the parameter estimates, checks the diagnostics, and validates the forecasts.

INTRODUCTION

This paper is an introduction to applied time series modeling for analysts who have minimum experience in model building, but are not very familiar with time series models. It would help to have a basic understanding of regression analysis such as simple linear regression or multiple regressions. The challenge of modeling is to diagnose the problem and decide on an appropriate model to help answer the real-world questions. It takes experience to develop an ability to formulate appropriate statistical models and to interpret the results, but this paper gives a head start on practicing these techniques.

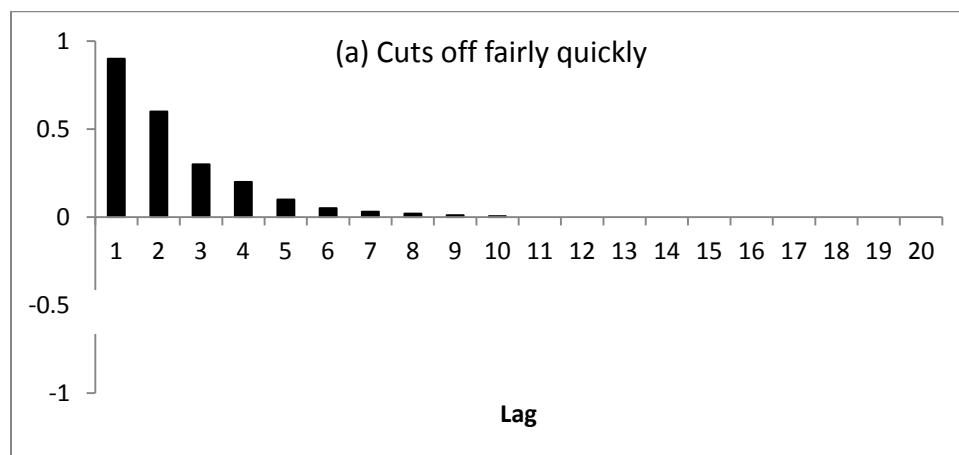
NON-SEASONAL BOX-JENKINS MODEL IDENTIFICATION

Before identifying the pattern, the time series values y_1, y_2, \dots, y_n must be stationary where its mean and variance are constant through time. The constant mean and variance can be achieved by removing the pattern caused by the time dependent autocorrelation. Besides looking at the plot of the time series values over time to determine stationary or non-stationary, the sample autocorrelation function (ACF) also gives visibility to the data. If the ACF of the time series values either cuts off or dies down fairly quickly (Figure 1(a)), then the time series values should be considered stationary. On the other hand, if the ACF of the time series values either cuts off or dies down extremely slowly (Figure 1(b)), then it should be considered non-stationary. In general, if the original time series values are *non-stationary* and *non-seasonal*, perform the first or second differencing transformation on the original data will usually produce stationary time series values.

First Difference: $z_t = y_t - y_{t-1}$ where $t = 2, 3, \dots, n$

Second Difference: $z_t = (y_t - y_{t-1}) - (y_{t-1} - y_{t-2})$ where $t = 3, 4, \dots, n$

Figure 1. The ACF (PACF) cuts off fairly quickly versus dies down extremely slowly (Bowerman, O'Connell, and Koehler p. 413)



The Box-Jenkins Methodology for Time Series Models, continued



Given that either the original time series y_1, y_2, \dots, y_n or the transformed time series z_b, z_{b+1}, \dots, z_n are stationary, we can now look at the sample autocorrelation function (ACF) and partial autocorrelation function (PACF) for particular behaviors that indicate a theoretical non-seasonal Box-Jenkins model. **Figure 2** shows different behaviors of the ACF and PACF. There are three types of non-seasonal theoretical Box-Jenkins models summarized in **Table 1**: moving average model of order q , autoregressive model of order p , and mixed autoregressive – moving average model of (p, q) . Please note that there is no theoretical Box-Jenkins model when the ACF cuts off quickly after lag q and the PACF cuts off quickly after lag p . However, we could look at which of the ACF or PACF is cutting off more abruptly to tentatively identify the model and look at the estimations, diagnostics, and forecasts to select the best model.

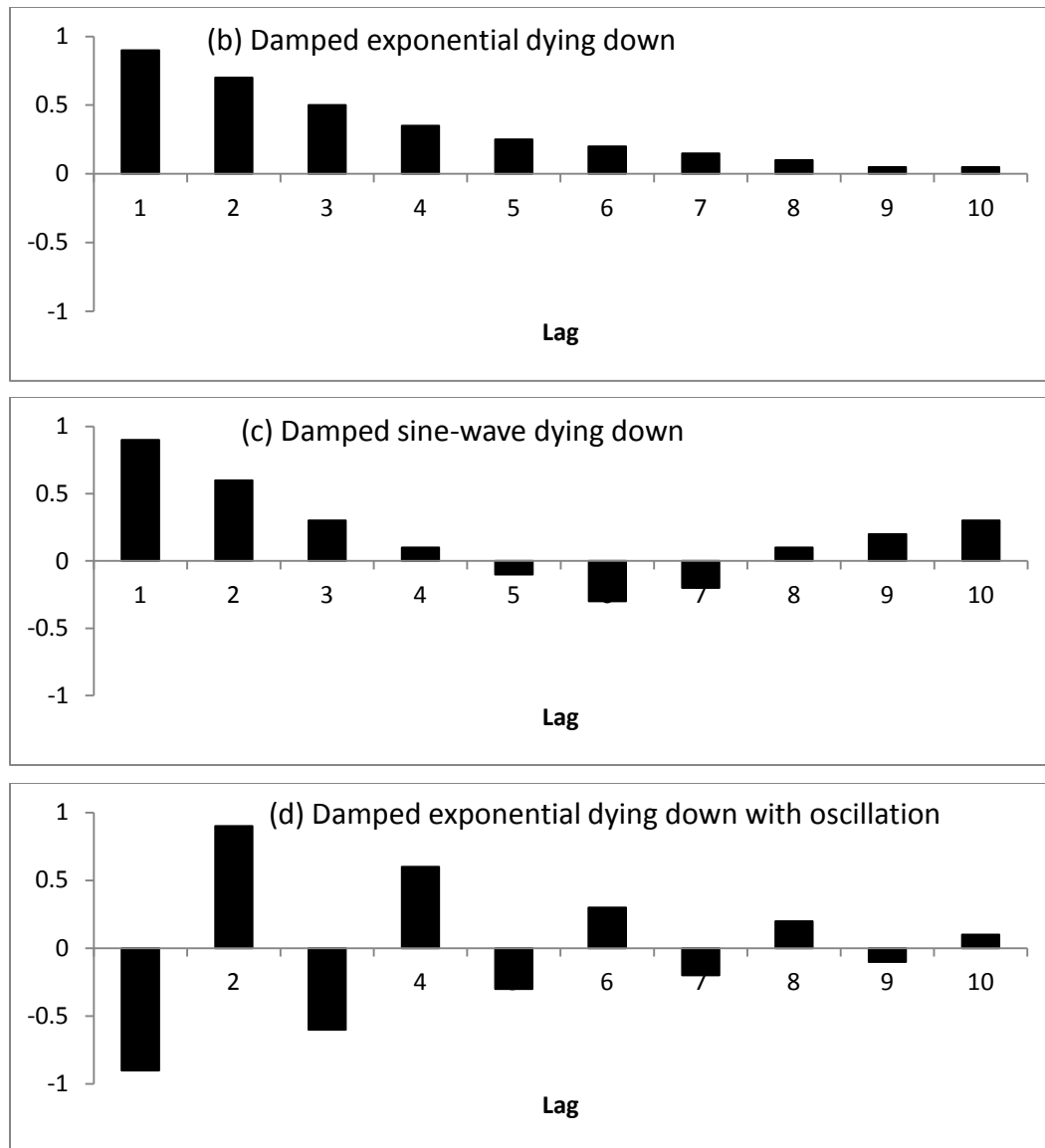
Table 1. Non-Seasonal Theoretical Box-Jenkins Models (Bowerman, O'Connell, and Koehler p. 436)

| Model | ACF | PACF |
|--|------------------------|------------------------|
| Moving average of order q $z_t = \delta + a_t - \theta_1 a_{t-1} - \theta_2 a_{t-2} - \dots - \theta_q a_{t-q}$ | Cuts off after lag q | Dies down |
| Autoregressive of order p $z_t = \delta + \varphi_1 z_{t-1} + \varphi_2 z_{t-2} + \dots + \varphi_p z_{t-p} + a_t$ | Dies down | Cuts off after lag p |
| Mixed autoregressive-moving average of order (p, q) $z_t = \delta + \varphi_1 z_{t-1} + \varphi_2 z_{t-2} + \dots + \varphi_p z_{t-p} + a_t - \theta_1 a_{t-1} - \theta_2 a_{t-2} - \dots - \theta_q a_{t-q}$ | Dies down | Dies down |

Figure 2. The ACF and PACF Behaviors (Bowerman, O'Connell, and Koehler p. 412)



The Box-Jenkins Methodology for Time Series Models, continued



SEASONAL BOX-JENKINS MODEL IDENTIFICATION

If the original time series values are *non-stationary* and *seasonal*, more complex differencing transformations are required. Before using differencing to transform seasonal non-stationary time series values into stationary time series values, we need to check if the data over time shows constant seasonal variance. To stabilize the variance, apply an appropriate predifferencing transformation – log, square root, and etc – on the original time series values. Next we examine the ACF of the transformed data both at non-seasonal and seasonal levels for any indication of stationary. The seasonal behaviors appear at the **exact seasonal lags** L , $2L$, $3L$, and $4L$. For daily data ($L = 7$), the exact seasonal lags are 7, 14, 21, and 28. For monthly data ($L = 12$), the exact seasonal lags would be 12, 24, 36, and 48. For quarterly data ($L = 4$), the exact seasonal lags would be 4, 8, 12, and 16. Beside looking for spikes greater than two standard errors at the exact seasonal lags, we should also look for spikes at the **near seasonal lags**, which are $L - 2$, $L - 1$, $L + 1$, $L + 2$, $2L - 2$, $2L - 1$, $2L + 1$, $2L + 2$, $3L - 2$, $3L - 1$, $3L + 1$, $3L + 2$, $4L - 2$, $4L - 1$, $4L + 1$, and $4L + 2$. In general, the transformed time series values z_b, z_{b+1}, \dots, z_n are considered stationary if the ACF shows **both** of the following behaviors:

1. Cuts off or dies down fairly quickly at the non-seasonal level.
2. Cuts off or dies down fairly quickly at the seasonal level (exact seasonal lags or near seasonal lags).

Otherwise, these values are considered non-stationary. The four stationary transformations on the time series values y_1, y_2, \dots, y_n are shown in **Table 2**. Sometimes the transformation (1) $z_t = y_t^*$ does not need differencing to produce

stationary time series values. The transformation (2) $z_t = y_t^* - y_{t-1}^*$ is the first non-seasonal differencing that sometimes transforms seasonal time series values into stationary time series values. The transformations (3) $z_t = y_t^* - y_{t-L}^*$ and (4) $z_t = y_t^* - y_{t-1}^* - y_{t-L}^* + y_{t-L-1}^*$ are the first seasonal differencing and the mixed of first non-seasonal differencing and first seasonal differencing. These transformations frequently produce stationary time series values.

Table 2. Four Stationarity Transformations (Bowerman, O'Connell, and Koehler p. 492)

| (1) $z_t = y_t^*$ | (2) $z_t = y_t^* - y_{t-1}^*$ | (3) $z_t = y_t^* - y_{t-L}^*$ | (4) $z_t = y_t^* - y_{t-1}^* - y_{t-L}^* + y_{t-L-1}^*$ |
|---|---|---|---|
| $z_1 = y_1^*$ $z_2 = y_2^*$ \vdots \vdots \vdots \vdots \vdots $z_n = y_n^*$ | $z_2 = y_2^* - y_1^*$ $z_3 = y_3^* - y_2^*$ \vdots \vdots \vdots \vdots $z_n = y_n^* - y_{n-1}^*$ | $z_{L+1} = y_{L+1}^* - y_1^*$ $z_{L+2} = y_{L+2}^* - y_2^*$ \vdots \vdots $z_n = y_n^* - y_{n-L}^*$ | $z_{L+2} = y_{L+2}^* - y_{L+1}^* - y_2^* + y_1^*$ $z_{L+3} = y_{L+3}^* - y_{L+2}^* - y_3^* + y_2^*$ \vdots $z_n = y_n^* - y_{n-1}^* - y_{n-L}^* + y_{n-L-1}^*$ |

Given that the time series values are considered stationary and exhibit behaviors (spikes) described in **Table 1 Non-Seasonal Theoretical Box-Jenkins Models** at the non-seasonal and seasonal levels, here is the three-step procedure for tentatively identifying a model:

STEP 1: Look at the behaviors (spikes) of the ACF and PACF at the non-seasonal level to identify a non-seasonal model.

STEP 2: Look at the behaviors (spikes) of the ACF and PACF at the seasonal level to identify a seasonal model.

STEP 3: Combine models from STEP 1 and STEP 2 to identify an overall tentatively model.

Once we obtain the overall tentatively model, we need to determine whether or not to include a constant term δ in a Box-Jenkins model. A constant term δ is the mean (μ) of the stationary time series values z_b, z_{b+1}, \dots, z_n , which μ is equal (or nearly equal) or not equal to zero. In general, the rule of thumb is to include a constant mean δ in the model if the absolute value of

$$\frac{\bar{z}}{s_z/\sqrt{n-b+1}} \text{ where } \bar{z} = \frac{\sum_{t=b}^n z_t}{n-b+1} \text{ and } s_z = \left[\frac{\sum_{t=b}^n (z_t - \bar{z})^2}{(n-b+1)-1} \right]^{1/2}$$

is greater than 2 (Bowerman, O'Connell, and Koehler p. 427). Equivalently, if the p-value associated to a constant mean δ is less than a significant level α , we should include a constant mean δ in the model. Otherwise, we do not include it in the model.

EXAMPLE 1

Using SASHELP.USECON data set, we will apply the Box-Jenkins methodology to examine the trend and develop a time series model that could be used to forecast the Airline Revenue Passenger Miles Domestic (AIRRPMD) on a monthly basis. The historical data started on January 1971 to December 1991. We will build the model based on the first 19 years of historical data and forecast the monthly AIRRPMD for the 20th year (1991) to validate the model.

Output 1.1 shows that the monthly AIRRPMDs follow an increasing trend and have a seasonal pattern with one major peak and several minor peaks during the year. The major peaks appear to be bigger over the years indicating non-constant seasonal variation. We need to perform a transformation on the data to stabilize the seasonal variation. The square roots of the AIRRPMD (Output 1.2) show that the transformation is not strong enough to equalize the seasonal variation. The quartic roots of the AIRRPMD (Output 1.3) and the natural logarithms of the AIRRPMD (Output 1.4) both seem to produce constant seasonal variation. However, the quartic root transformation best equalizes the seasonal variation because the major peaks' scales are consistent over the range of time (note the variations of the major peaks in 1971 and 1991 appear more equalize).

The Box-Jenkins Methodology for Time Series Models, continued

```

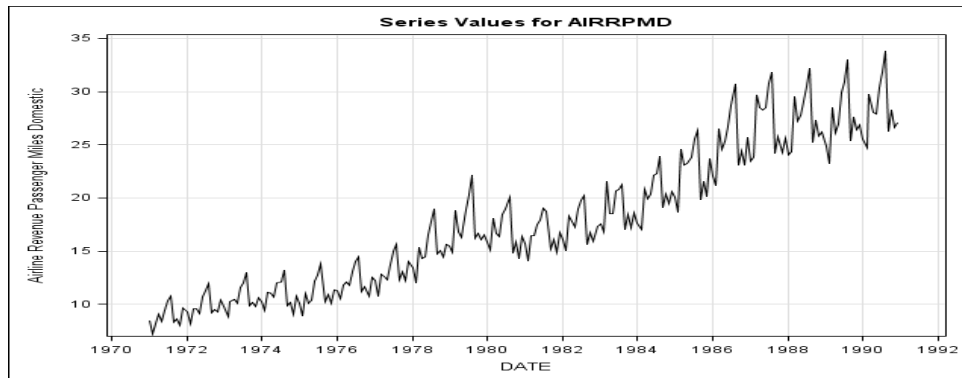
DATA airrpmd;
SET sashelp.usecon (WHERE=(Date < '01Jan1991'd));
  ln_AIRRPMD=log(AIRRPMD);
  sqrt_AIRRPMD=AIRRPMD**.5;
  Qtroot_AIRRPMD=AIRRPMD**.25;

RUN;

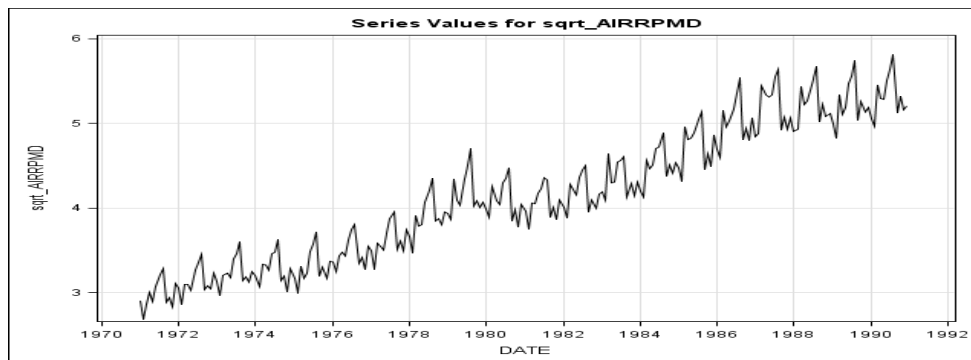
ODS GRAPHICS ON;
PROC TIMESERIES data=airrpmd plot=series;
ID date interval=month;
VAR AIRRPMD ln_AIRRPMD sqrt_AIRRPMD Qtroot_AIRRPMD;
RUN;

```

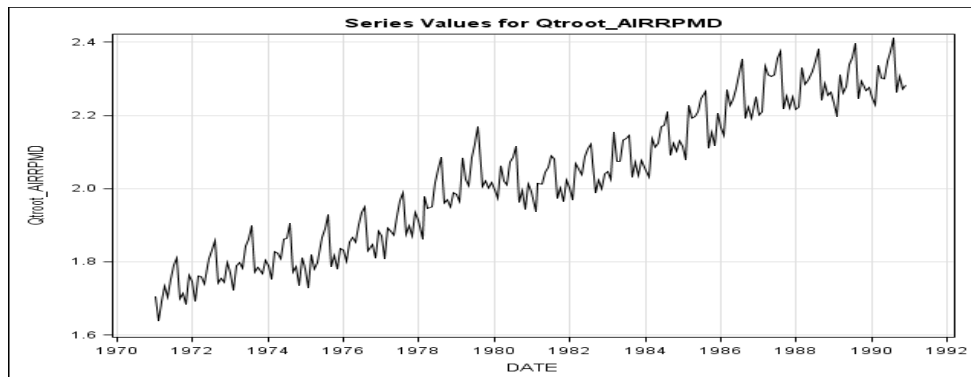
Output 1.1 AIRRPMDs



Output 1.2 SQUARE ROOT AIRRPMDs

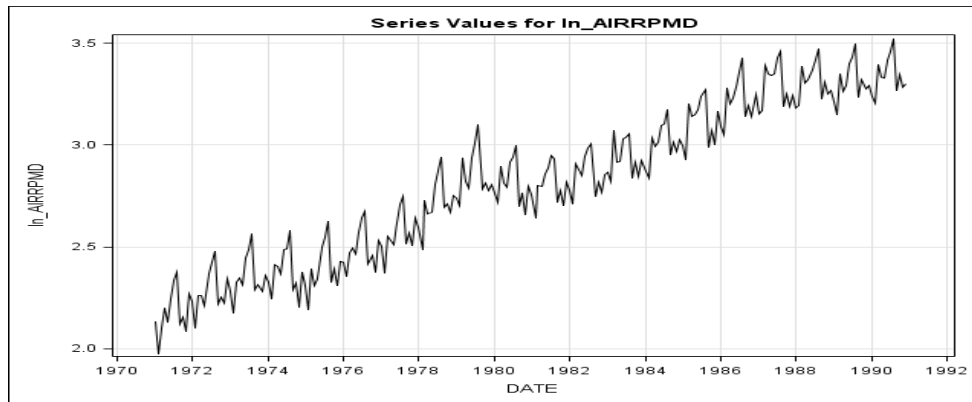


Output 1.3 QUARTIC ROOT AIRRPMDs



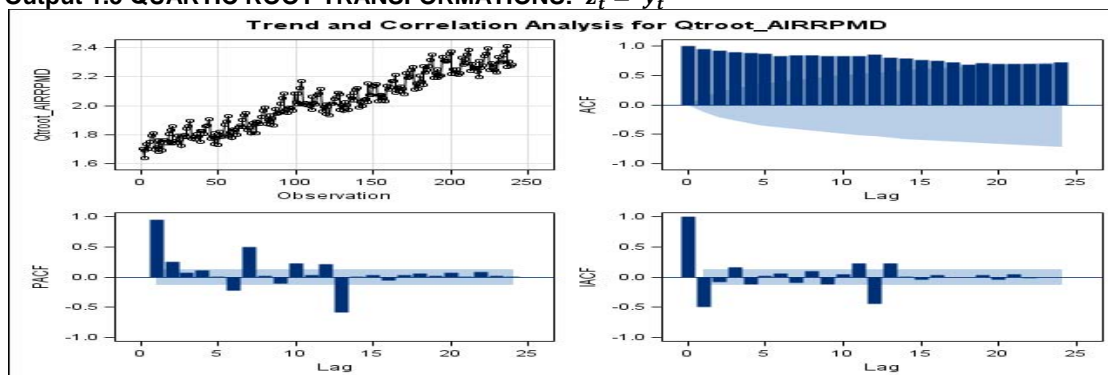
The Box-Jenkins Methodology for Time Series Models, continued

Output 1.4 NATURAL LOGARITHMS AIRRPMDS

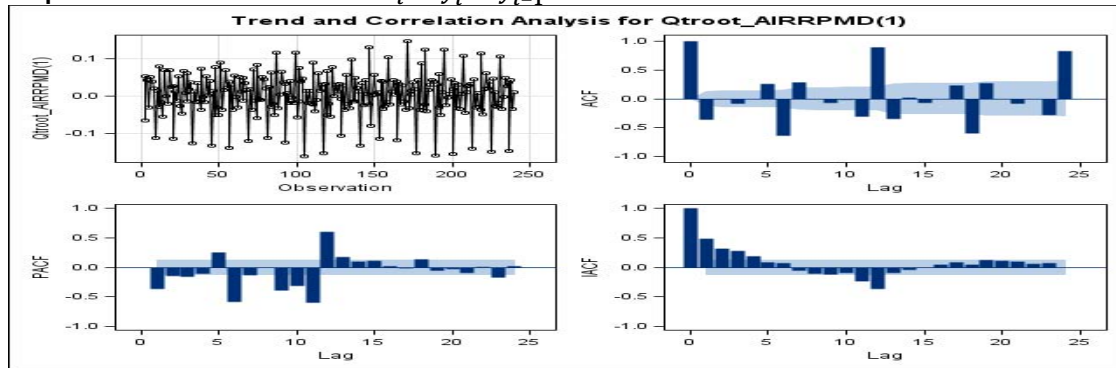
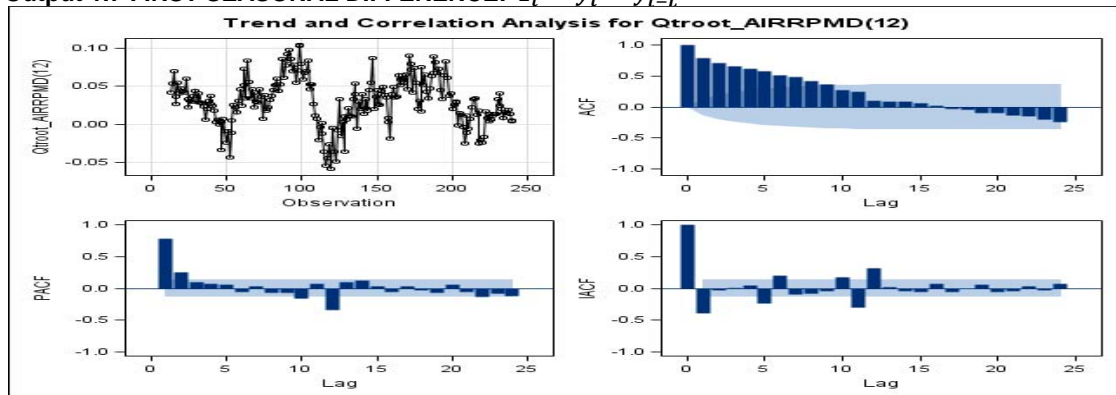
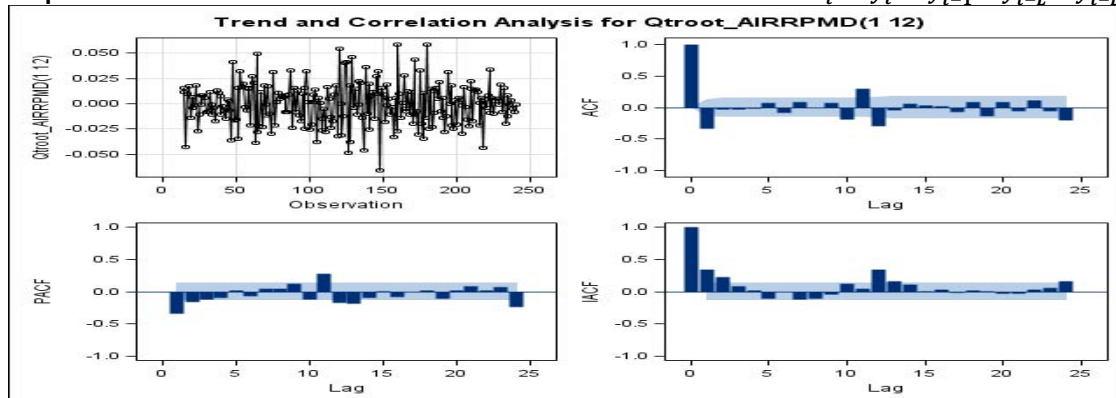


Let's examine the following sample autocorrelation functions (ACF) of quartic roots time series values and determine which one of the four stationary transformations in **Table 2** produces stationary time series values. The time series values z_b, z_{b+1}, \dots, z_n are considered stationary if the ACF cuts off fairly quickly both at the non-seasonal and seasonal levels. When we examine the time series values at the seasonal level, make sure to look at the exact seasonal lags ($L=12, 24, 36$, and etc.) and near seasonal lags ($L=10, 11, 13, 14, 22, 23, 25, 26$, and etc.) that are greater than two standard errors. Since the ACF in Output 1.5 dies down extremely slowly both at the non-seasonal and seasonal levels, the quartic roots time series values are considered non-stationary. In Output 1.6, the ACF dies down fairly quicker at the non-seasonal level than at the seasonal level. The transformed values are still non-stationary. The ACF in Output 1.7 shows similar behaviors as Output 1.5 indicating non-stationary time series values. In Output 1.8, the ACF cuts off quickly after lag 1 at the non-seasonal level and dies down fairly quickly after lag 12 at the seasonal level. We can conclude that the quartic roots time series values transformed by the first difference and first seasonal difference are stationary. Please note that the values can be over-differenced. If that is the case, the sample inverse autocorrelation functions (IACF) will die down extremely slowly.

```
ODS GRAPHICS ON;
PROC ARIMA DATA=airrpm;
/* QUARTIC ROOT TRANSFORMATIONS */
IDENTIFY VAR=Qtroot_AIRRPMD NLAG=24 ;
/* FIRST DIFFERENCE */
IDENTIFY VAR=Qtroot_AIRRPMD(1) NLAG=24 ;
/* FIRST SEASONAL DIFFERENCE */
IDENTIFY VAR=Qtroot_AIRRPMD(12) NLAG=24 ;
/* FIRST DIFFERENCE AND FIRST SEASONAL DIFFERENCE */
IDENTIFY VAR=Qtroot_AIRRPMD(1,12) NLAG=24 ;
RUN;
```

Output 1.5 QUARTIC ROOT TRANSFORMATIONS: $z_t = y_t^*$ 

The Box-Jenkins Methodology for Time Series Models, continued

Output 1.6 FIRST DIFFERENCE: $z_t = y_t^* - y_{t-1}^*$ **Output 1.7 FIRST SEASONAL DIFFERENCE:** $z_t = y_t^* - y_{t-L}^*$ **Output 1.8 FIRST DIFFERENCE AND FIRST SEASONAL DIFFERENCE:** $z_t = y_t^* - y_{t-1}^* - y_{t-L}^* + y_{t-L-1}^*$ 

❶ Autocorrelation Check for White Noise

| To Lag | Chi- Square | DF | Pr > ChiSq | -----Autocorrelations----- | | | | | |
|-----------|----------------|----|---------------|----------------------------|--------|--------|--------|--------|--------|
| 6 | 29.66 | 6 | <.0001 | -0.337 | -0.027 | -0.037 | -0.023 | 0.071 | -0.087 |
| 12 | 85.07 | 12 | <.0001 | 0.090 | -0.004 | 0.077 | -0.184 | 0.305 | -0.299 |
| 18 | 90.41 | 18 | <.0001 | -0.050 | 0.061 | 0.032 | 0.023 | -0.078 | 0.088 |
| 24 | 113.81 | 24 | <.0001 | -0.136 | 0.094 | -0.063 | 0.118 | -0.057 | -0.209 |

❶ In addition to the autocorrelation plots, the white noise output is provided to test the null hypothesis: the autocorrelations of the time series values are equal to zero. If H_0 is rejected, the autocorrelations of the time series values are nonzero. In this case, the white noise output shows that the null hypotheses for lags up to 6, 12, 18, and 24 are strongly rejected indicating that a time series model is needed.

The Box-Jenkins Methodology for Time Series Models, continued

Looking at the ACF and PACF of the stationary time series in Output 1.8, we can identify an appropriate time series model by following the three-step procedure discussed in **Seasonal Box-Jenkins Model Identification**.

STEP 1: The ACF cuts off after lag 1 and the PACF dies down at the non-seasonal level indicate a first-order moving average model. Therefore the tentatively non-seasonal model is

$$z_t = \delta + a_t - \theta_1 a_{t-1}$$

STEP 2: The ACF cuts off after lag 12 and the PACF dies down at the seasonal level indicate a seasonal moving average model with lag 12. Therefore the tentatively seasonal model is

$$z_t = \delta + a_t - \theta_{1,12} a_{t-12}$$

STEP 3: Combining models from STEP 1 and 2, the tentatively overall model is

$$z_t = \delta - \theta_1 a_{t-1} - \theta_{1,12} a_{t-12} + a_t$$

ESTIMATIONS

Not only does the Box-Jenkins model have to be stationary, it also has to be invertible. Invertible means recent observations are more heavily weighted than more remote observations; the parameters $(\varphi_1, \varphi_2, \dots, \varphi_p, \theta_1, \theta_2, \dots, \theta_q)$ used in the model decline from the most recent observations down to the further past observations. By default, PROC ARIMA in SAS® requires that the preliminary and final parameter estimates for the autoregressive and moving-average models satisfy the stationarity and invertibility conditions. If analysts use other software packages, make sure that the parameter estimates meet the stationarity and invertibility conditions (Bowerman, O'Connell, and Koehler p. 450).

The t -values and approximate p -values test the following hypothesis (Bowerman, O'Connell, and Koehler pp. 455 – 456). Let θ be any particular parameter in a Box-Jenkins model.

$$H_0: \theta = 0 \text{ versus } H_a: \theta \neq 0$$

We can reject the null hypothesis H_0 if and only if either of the following conditions holds:

1. $|t| > t_{\alpha/2}^{(n-n_p)}$ where n_p is the number of parameters in the model.
2. $p\text{-value} < \alpha$

Please note that if the null hypothesis is rejected at the smaller significant level α , the stronger the evidence indicates that the parameter is important in the model. We also look at the estimated correlation between parameters in the correlation matrix. The parameters are always correlated, but very highly correlations greater than 0.9 suggests poor parameter estimates. If that is the case, drop one of the parameters. ❸ Lastly, when we are comparing models, the best model has smaller standard error, Akaike's Information Criterion (AIC), and Schwarz's Bayesian Criterion (SBC) values.

EXAMPLE 1 CONTINUED

To determine where or not δ should be included in the combined model, here is the calculation:

❶

$$\frac{\bar{z}}{s_z/\sqrt{n-b+1}} = \frac{-0.00017}{0.020485/\sqrt{240-13+1}} = -0.1253 \text{ is less than } 2.$$

We conclude that δ should not be included in the combined model. Equivalently, ❷ the p -value (0.6980) associated with δ is greater than the significant level $\alpha = 0.05$ indicating that δ is insignificant. The p -values associated with θ_1 and $\theta_{1,12}$ are greater than 0.05 indicating that these parameters are significant in the model. The correlation matrix ❸ shows that the parameters, θ_1 and $\theta_{1,12}$, are not highly correlated.

```
PROC ARIMA DATA=airrpm;
/* NONSEASONAL FIRST DIFFERENCE AND SEASONAL DIFFERENCE*/
IDENTIFY VAR=Qtrout_AIRRPMD(1,12) NLAG=24 ;
/* FIRST-ORDER MOVING AVERAGE AND SEASONAL MOVING AVERAGE AT LAG 12*/
ESTIMATE Q=(1) (12) ;
RUN;
```

| | | |
|---|-----------------------------------|----------|
| ❶ | Name of Variable = Qtrout_AIRRPMD | |
| Period(s) of Differencing | | 1,12 |
| Mean of Working Series | | -0.00017 |
| Standard Deviation | | 0.020485 |
| Number of Observations | | 227 |
| Observation(s) eliminated by differencing | | 13 |

The Box-Jenkins Methodology for Time Series Models, continued

| Autocorrelation Check for White Noise | | | | | | | | | |
|---------------------------------------|------------|----|------------|----------------------------|--------|--------|--------|--------|--------|
| To Lag | Chi-Square | DF | Pr > ChiSq | -----Autocorrelations----- | | | | | |
| 6 | 29.66 | 6 | <.0001 | -0.337 | -0.027 | -0.037 | -0.023 | 0.071 | -0.087 |
| 12 | 85.07 | 12 | <.0001 | 0.090 | -0.004 | 0.077 | -0.184 | 0.305 | -0.299 |
| 18 | 90.41 | 18 | <.0001 | -0.050 | 0.061 | 0.032 | 0.023 | -0.078 | 0.088 |
| 24 | 113.81 | 24 | <.0001 | -0.136 | 0.094 | -0.063 | 0.118 | -0.057 | -0.209 |

| ② Conditional Least Squares Estimation | | | | | |
|--|------------|----------------|---------|----------------|-----|
| Parameter | Estimate | Standard Error | t Value | Approx Pr > t | Lag |
| MU | -0.0001057 | 0.0002720 | -0.39 | 0.6980 | 0 |
| MA1,1 | 0.43886 | 0.06022 | 7.29 | <.0001 | 1 |
| MA2,1 | 0.60252 | 0.05499 | 10.96 | <.0001 | 12 |

| | |
|---------------------|----------|
| Constant Estimate | -0.00011 |
| Variance Estimate | 0.000291 |
| Std Error Estimate | 0.017072 |
| ⑤ AIC | -1200.74 |
| SBC | -1190.46 |
| Number of Residuals | 227 |

* AIC and SBC do not include log determinant.

| ⑥ Correlations of Parameter Estimates | | | | |
|---------------------------------------|-------|--------|--------|--|
| Parameter | MU | MA1,1 | MA2,1 | |
| MU | 1.000 | 0.005 | 0.020 | |
| MA1,1 | 0.005 | 1.000 | -0.051 | |
| MA2,1 | 0.020 | -0.051 | 1.000 | |

DIAGNOSTICS CHECKING

To test the adequacy of an overall model, the null and alternative hypotheses are H_0 : Model is adequate versus H_a : Model is inadequate. We perform the test by using the Ljung-Box statistic (Q^*) given below:

$$\text{The Ljung - Box Statistic : } Q^* = n'(n' + 2) \sum_{l=1}^K (n' - l)^{-1} r_l^2(\hat{a})$$

Please note that n_c is the number of parameters in the model excluded the constant mean δ and $n' = n - d$ where n is the number of observations and d is the degree of non-seasonal differencing used to transform the original time series values into stationary. Also $r_l^2(\hat{a})$ is the square of the autocorrelation of the residuals at lag l (Bowerman, O'Connell, and Koehler p. 459).

If the p-value is greater than significant level α or equivalently Q^* is less than chi-square distribution with $K - n_c$ degree of freedom, the null hypothesis cannot be rejected concluding that the model is adequate. The greater the p-value is, the stronger the evidence indicates that the model is adequate. Furthermore, we can improve the model by examining the autocorrelations and partial autocorrelation of the residuals. If there are spikes exceeding two standard errors, it is possibly an indication of a more adequate model.

EXAMPLE 1 CONTINUED

④ The p-values associated with Q^* are greater than 0.05 indicating that the model is adequate. Looking at the autocorrelation function and partial autocorrelation function of the residuals, there is only one spike at lag 11 that is slightly over two standard errors. Since the standard error, Akaike's Information Criteria (AIC), and Schwarz Bayesian

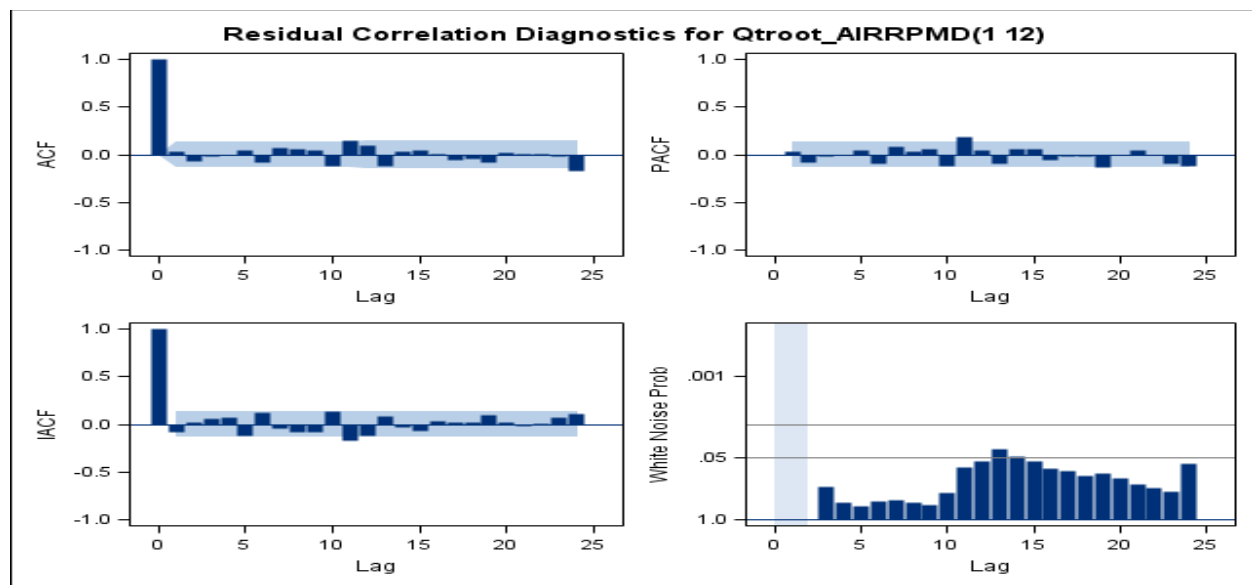
The Box-Jenkins Methodology for Time Series Models, continued

Criteria (SBC) values do not significantly improved from a revised model that accounts for lag 11 (output omitted), the original overall model is considered adequate.

| ④ Autocorrelation Check of Residuals | | | | | | | | | |
|--------------------------------------|------------|----|------------|----------------------------|--------|--------|--------|--------|--------|
| To Lag | Chi-Square | DF | Pr > ChiSq | -----Autocorrelations----- | | | | | |
| 6 | 4.00 | 4 | 0.4061 | 0.035 | -0.075 | -0.019 | 0.000 | 0.050 | -0.086 |
| 12 | 17.81 | 10 | 0.0582 | 0.073 | 0.055 | 0.047 | -0.120 | 0.152 | 0.097 |
| 18 | 22.87 | 16 | 0.1172 | -0.114 | 0.038 | 0.041 | 0.006 | -0.054 | -0.039 |
| 24 | 32.65 | 22 | 0.0670 | -0.080 | 0.016 | 0.007 | 0.013 | -0.016 | -0.177 |
| 30 | 37.37 | 28 | 0.1109 | 0.008 | 0.115 | 0.012 | -0.021 | -0.046 | -0.047 |
| 36 | 42.25 | 34 | 0.1566 | -0.033 | -0.114 | 0.005 | 0.010 | 0.042 | -0.047 |
| 42 | 51.49 | 40 | 0.1053 | -0.031 | -0.089 | -0.088 | 0.031 | 0.065 | 0.107 |

| | |
|-----------------------------------|----------|
| Model for variable Qtroot_AIRRPMD | |
| Estimated Mean | -0.00011 |
| Period(s) of Differencing | 1,12 |

| | |
|------------------------|---------------------|
| Moving Average Factors | |
| Factor 1: | 1 - 0.43886 B**(1) |
| Factor 2: | 1 - 0.60252 B**(12) |



FORECASTING

We validate the forecast by splitting the data in two parts: one part of the data is used for modeling and the other part of the data is used for forecasting. We look at the residuals to determine how accurate the model predicts. The desired accuracy of the forecasts depends on the analyst's goal.

EXAMPLE 1 CONTINUED

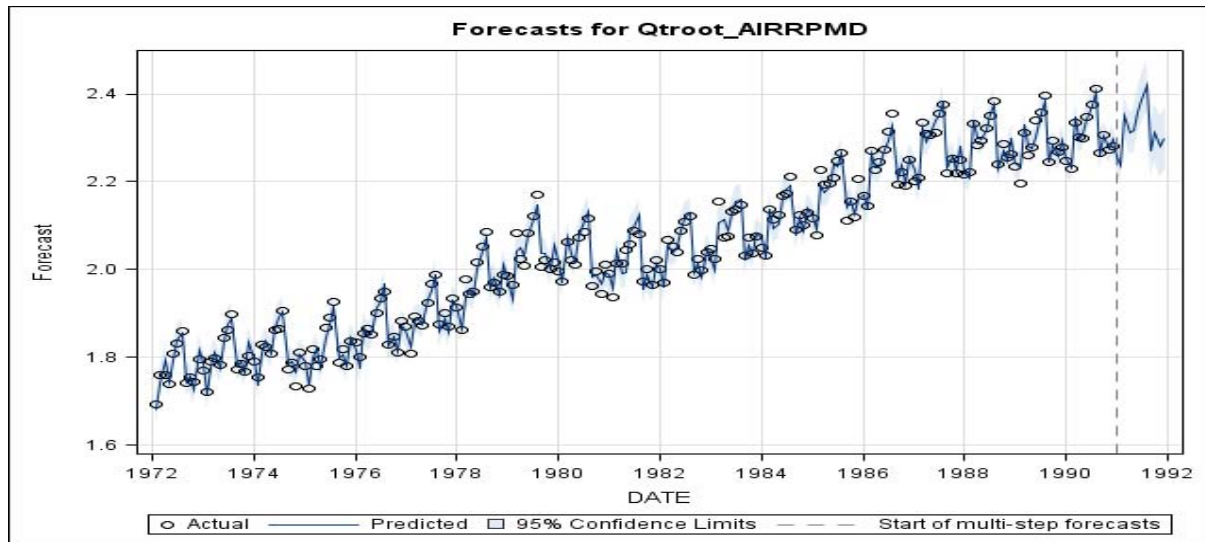
The first 19 years of the monthly AIRRPMD values are used to build the model forecasting the 20th year. Note that the forecast values for the 20th year are still in quartic roots transformation. To get the forecast values in the original scales, do the reverse transformation; simply calculate each value to the fourth power.

The Box-Jenkins Methodology for Time Series Models, continued

```

ODS GRAPHICS ON;
PROC ARIMA DATA=airrpmdd PLOTS = ALL;
/* NONSEASONAL 1ST DIFFERENCE AND SEASONAL DIFFERENCE*/
IDENTIFY VAR=Qtrroot_AIRRPMD(1,12) NLAG=24 ;
/* FIRST-ORDER MOVING AVERAGE AND SEASONAL MOVING AVERAGE AT LAG 12*/
ESTIMATE Q=(1) (12) ;
FORECAST ID=date LEAD=12 INTERVAL=month PRINTALL OUT=predictions;
RUN;

```



CONCLUSION

The Box-Jenkins methodology has four steps: model identification, estimation, diagnostics checking, and forecasting. Of all the four steps, model identification is the most complex step that analysts usually struggle with. The best way to learn is to understand the methodology and apply it to real-world problems. Practices help gain the experience you need to become efficient in model building.

REFERENCE

Bowerman, Bruce L., Richard T. O'Connell, and Anne B. Koehler. Forecasting, Time Series, and Regression, 4th ed. Belmont, CA: Thomson Brooks/Cole, 2005.

ACKNOWLEDGMENTS

A special thank you goes to my mentor, Art Carpenter, for enthusiastically reviewing my paper and providing comments and questions that make me think twice. I appreciate Professor Subir Ghosh for educating me on Time Series Models with so much patience.

CONTACT INFORMATION

Your comments and questions are valued and encouraged. Contact the author at:

Name: Theresa Hoang Diem Ngo
E-mail: theresa.ngo1120@gmail.com

SAS and all other SAS Institute Inc. product or service names are registered trademarks or trademarks of SAS Institute Inc. in the USA and other countries. ® indicates USA registration.

Other brand and product names are trademarks of their respective companies.