



# SAS/ETS® 13.2 User's Guide

## The ARIMA Procedure

This document is an individual chapter from *SAS/ETS® 13.2 User's Guide*.

The correct bibliographic citation for the complete manual is as follows: SAS Institute Inc. 2014. *SAS/ETS® 13.2 User's Guide*. Cary, NC: SAS Institute Inc.

Copyright © 2014, SAS Institute Inc., Cary, NC, USA

All rights reserved. Produced in the United States of America.

**For a hard-copy book:** No part of this publication may be reproduced, stored in a retrieval system, or transmitted, in any form or by any means, electronic, mechanical, photocopying, or otherwise, without the prior written permission of the publisher, SAS Institute Inc.

**For a Web download or e-book:** Your use of this publication shall be governed by the terms established by the vendor at the time you acquire this publication.

The scanning, uploading, and distribution of this book via the Internet or any other means without the permission of the publisher is illegal and punishable by law. Please purchase only authorized electronic editions and do not participate in or encourage electronic piracy of copyrighted materials. Your support of others' rights is appreciated.

**U.S. Government License Rights; Restricted Rights:** The Software and its documentation is commercial computer software developed at private expense and is provided with RESTRICTED RIGHTS to the United States Government. Use, duplication or disclosure of the Software by the United States Government is subject to the license terms of this Agreement pursuant to, as applicable, FAR 12.212, DFAR 227.7202-1(a), DFAR 227.7202-3(a) and DFAR 227.7202-4 and, to the extent required under U.S. federal law, the minimum restricted rights as set out in FAR 52.227-19 (DEC 2007). If FAR 52.227-19 is applicable, this provision serves as notice under clause (c) thereof and no other notice is required to be affixed to the Software or documentation. The Government's rights in Software and documentation shall be only those set forth in this Agreement.

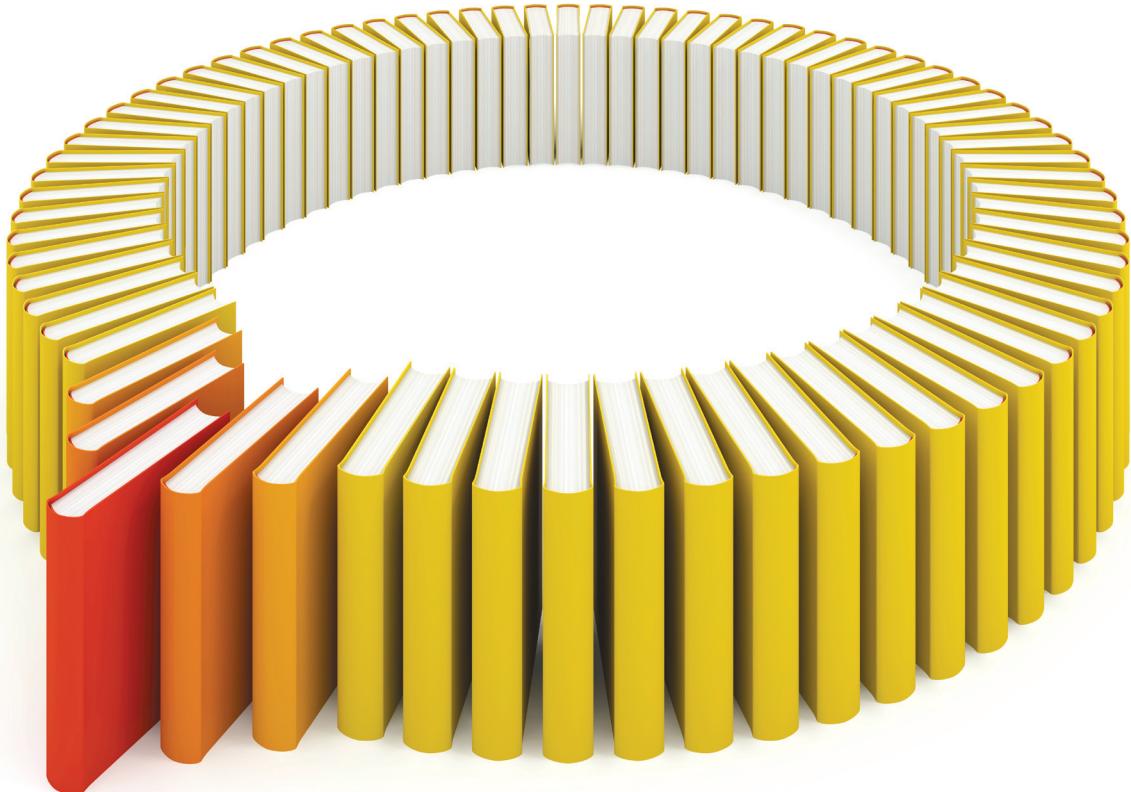
SAS Institute Inc., SAS Campus Drive, Cary, North Carolina 27513.

August 2014

SAS provides a complete selection of books and electronic products to help customers use SAS® software to its fullest potential. For more information about our offerings, visit [support.sas.com/bookstore](http://support.sas.com/bookstore) or call 1-800-727-3228.

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.



# Gain Greater Insight into Your SAS® Software with SAS Books.

Discover all that you need on your journey to knowledge and empowerment.

 support.sas.com/bookstore  
for additional books and resources.

  
sas®  
THE POWER TO KNOW®



# Chapter 7

## The ARIMA Procedure

### Contents

---

Overview: ARIMA Procedure . . . . .	186
Getting Started: ARIMA Procedure . . . . .	187
The Three Stages of ARIMA Modeling . . . . .	187
Identification Stage . . . . .	188
Estimation and Diagnostic Checking Stage . . . . .	193
Forecasting Stage . . . . .	199
Using ARIMA Procedure Statements . . . . .	201
General Notation for ARIMA Models . . . . .	201
Stationarity . . . . .	204
Differencing . . . . .	204
Subset, Seasonal, and Factored ARMA Models . . . . .	206
Input Variables and Regression with ARMA Errors . . . . .	207
Intervention Models and Interrupted Time Series . . . . .	210
Rational Transfer Functions and Distributed Lag Models . . . . .	211
Forecasting with Input Variables . . . . .	213
Data Requirements . . . . .	214
Syntax: ARIMA Procedure . . . . .	215
Functional Summary . . . . .	215
PROC ARIMA Statement . . . . .	218
BY Statement . . . . .	220
IDENTIFY Statement . . . . .	221
ESTIMATE Statement . . . . .	225
OUTLIER Statement . . . . .	229
FORECAST Statement . . . . .	230
Details: ARIMA Procedure . . . . .	231
The Inverse Autocorrelation Function . . . . .	231
The Partial Autocorrelation Function . . . . .	232
The Cross-Correlation Function . . . . .	232
The ESACF Method . . . . .	233
The MINIC Method . . . . .	235
The SCAN Method . . . . .	236
Stationarity Tests . . . . .	238
Prewhitening . . . . .	238
Identifying Transfer Function Models . . . . .	239
Missing Values and Autocorrelations . . . . .	239
Estimation Details . . . . .	240

Specifying Inputs and Transfer Functions . . . . .	244
Initial Values . . . . .	245
Stationarity and Invertibility . . . . .	246
Naming of Model Parameters . . . . .	247
Missing Values and Estimation and Forecasting . . . . .	247
Forecasting Details . . . . .	248
Forecasting Log Transformed Data . . . . .	249
Specifying Series Periodicity . . . . .	250
Detecting Outliers . . . . .	250
OUT= Data Set . . . . .	252
OUTCOV= Data Set . . . . .	253
OUTTEST= Data Set . . . . .	254
OUTMODEL= SAS Data Set . . . . .	257
OUTSTAT= Data Set . . . . .	258
Printed Output . . . . .	259
ODS Table Names . . . . .	262
Statistical Graphics . . . . .	263
Examples: ARIMA Procedure . . . . .	<b>267</b>
Example 7.1: Simulated IMA Model . . . . .	267
Example 7.2: Seasonal Model for the Airline Series . . . . .	271
Example 7.3: Model for Series J Data from Box and Jenkins . . . . .	278
Example 7.4: An Intervention Model for Ozone Data . . . . .	286
Example 7.5: Using Diagnostics to Identify ARIMA Models . . . . .	288
Example 7.6: Detection of Level Changes in the Nile River Data . . . . .	293
Example 7.7: Iterative Outlier Detection . . . . .	294
References . . . . .	<b>296</b>

---

## Overview: ARIMA Procedure

The ARIMA procedure analyzes and forecasts equally spaced univariate time series data, transfer function data, and intervention data by using the autoregressive integrated moving-average (ARIMA) or autoregressive moving-average (ARMA) model. An ARIMA model predicts a value in a response time series as a linear combination of its own past values, past errors (also called shocks or innovations), and current and past values of other time series.

The ARIMA approach was first popularized by Box and Jenkins, and ARIMA models are often referred to as Box-Jenkins models. The general transfer function model employed by the ARIMA procedure was discussed by Box and Tiao (1975). When an ARIMA model includes other time series as input variables, the model is sometimes referred to as an ARIMAX model. Pankratz (1991) refers to the ARIMAX model as *dynamic regression*.

The ARIMA procedure provides a comprehensive set of tools for univariate time series model identification, parameter estimation, and forecasting, and it offers great flexibility in the kinds of ARIMA or ARIMAX

models that can be analyzed. The ARIMA procedure supports seasonal, subset, and factored ARIMA models; intervention or interrupted time series models; multiple regression analysis with ARMA errors; and rational transfer function models of any complexity.

The design of PROC ARIMA closely follows the Box-Jenkins strategy for time series modeling with features for the identification, estimation and diagnostic checking, and forecasting steps of the Box-Jenkins method.

Before you use PROC ARIMA, you should be familiar with Box-Jenkins methods, and you should exercise care and judgment when you use the ARIMA procedure. The ARIMA class of time series models is complex and powerful, and some degree of expertise is needed to use them correctly.

## Getting Started: ARIMA Procedure

This section outlines the use of the ARIMA procedure and gives a cursory description of the ARIMA modeling process for readers who are less familiar with these methods.

### The Three Stages of ARIMA Modeling

The analysis performed by PROC ARIMA is divided into three stages, corresponding to the stages described by Box and Jenkins (1976).

1. In the *identification* stage, you use the IDENTIFY statement to specify the response series and identify candidate ARIMA models for it. The IDENTIFY statement reads time series that are to be used in later statements, possibly differencing them, and computes autocorrelations, inverse autocorrelations, partial autocorrelations, and cross-correlations. Stationarity tests can be performed to determine if differencing is necessary. The analysis of the IDENTIFY statement output usually suggests one or more ARIMA models that could be fit. Options enable you to test for stationarity and tentative ARMA order identification.
2. In the *estimation and diagnostic checking* stage, you use the ESTIMATE statement to specify the ARIMA model to fit to the variable specified in the previous IDENTIFY statement and to estimate the parameters of that model. The ESTIMATE statement also produces diagnostic statistics to help you judge the adequacy of the model.

Significance tests for parameter estimates indicate whether some terms in the model might be unnecessary. Goodness-of-fit statistics aid in comparing this model to others. Tests for white noise residuals indicate whether the residual series contains additional information that might be used by a more complex model. The OUTLIER statement provides another useful tool to check whether the currently estimated model accounts for all the variation in the series. If the diagnostic tests indicate problems with the model, you try another model and then repeat the estimation and diagnostic checking stage.

3. In the *forecasting* stage, you use the FORECAST statement to forecast future values of the time series and to generate confidence intervals for these forecasts from the ARIMA model produced by the preceding ESTIMATE statement.

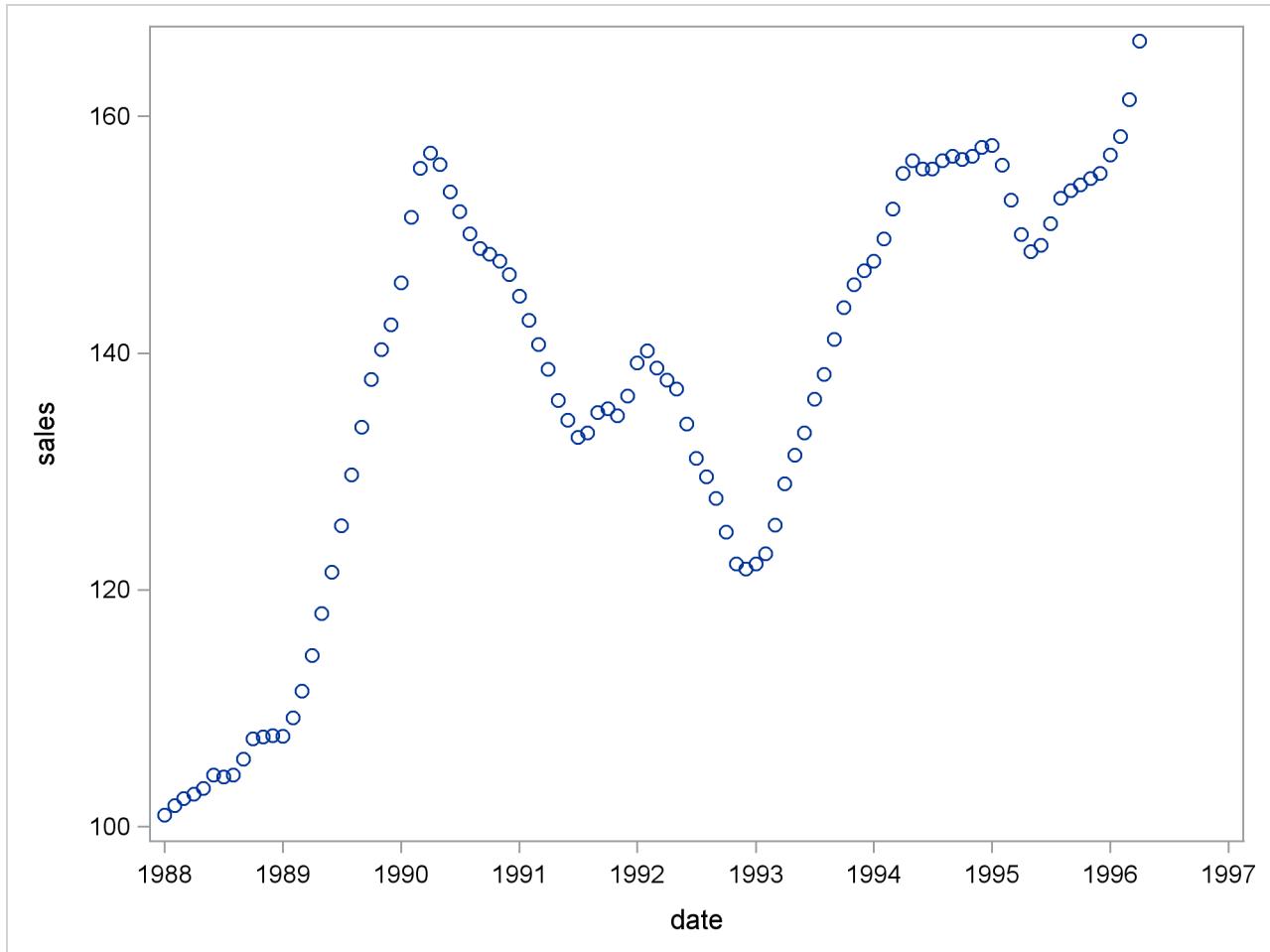
These three steps are explained further and illustrated through an extended example in the following sections.

## Identification Stage

Suppose you have a variable called SALES that you want to forecast. The following example illustrates ARIMA modeling and forecasting by using a simulated data set TEST that contains a time series SALES generated by an ARIMA(1,1,1) model. The output produced by this example is explained in the following sections. The simulated SALES series is shown in Figure 7.1.

```
proc sgplot data=test;
  scatter y=sales x=date;
run;
```

**Figure 7.1** Simulated ARIMA(1,1,1) Series SALES



## Using the IDENTIFY Statement

You first specify the input data set in the PROC ARIMA statement. Then, you use an IDENTIFY statement to read in the SALES series and analyze its correlation properties. You do this by using the following statements:

```
proc arima data=test ;
    identify var=sales nlag=24;
run;
```

### Descriptive Statistics

The IDENTIFY statement first prints descriptive statistics for the SALES series. This part of the IDENTIFY statement output is shown in Figure 7.2.

**Figure 7.2** IDENTIFY Statement Descriptive Statistics Output

#### The ARIMA Procedure

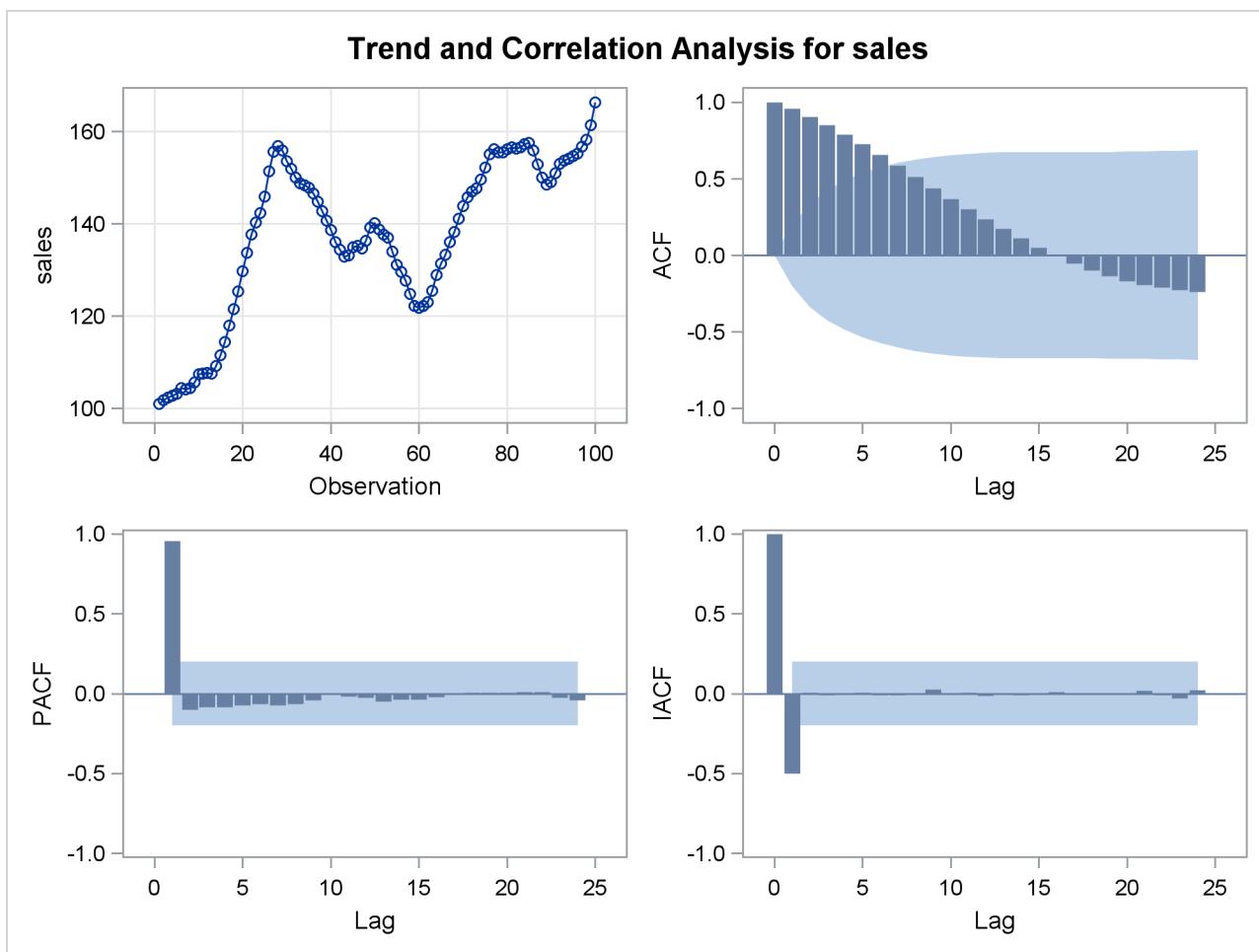
Name of Variable = sales	
Mean of Working Series	137.3662
Standard Deviation	17.36385
Number of Observations	100

### Autocorrelation Function Plots

The IDENTIFY statement next produces a panel of plots used for its autocorrelation and trend analysis. The panel contains the following plots:

- the time series plot of the series
- the sample autocorrelation function plot (ACF)
- the sample inverse autocorrelation function plot (IACF)
- the sample partial autocorrelation function plot (PACF)

This correlation analysis panel is shown in Figure 7.3.

**Figure 7.3** Correlation Analysis of SALES

These autocorrelation function plots show the degree of correlation with past values of the series as a function of the number of periods in the past (that is, the lag) at which the correlation is computed.

The NLAG= option controls the number of lags for which the autocorrelations are shown. By default, the autocorrelation functions are plotted to lag 24.

Most books on time series analysis explain how to interpret the autocorrelation and the partial autocorrelation plots. See the section “[The Inverse Autocorrelation Function](#)” on page 231 for a discussion of the inverse autocorrelation plots.

By examining these plots, you can judge whether the series is *stationary* or *nonstationary*. In this case, a visual inspection of the autocorrelation function plot indicates that the SALES series is nonstationary, since the ACF decays very slowly. For more formal stationarity tests, use the STATIONARITY= option. (See the section “[Stationarity](#)” on page 204.)

### **White Noise Test**

The last part of the default IDENTIFY statement output is the check for white noise. This is an approximate statistical test of the hypothesis that none of the autocorrelations of the series up to a given lag are significantly different from 0. If this is true for all lags, then there is no information in the series to model, and no ARIMA model is needed for the series.

The autocorrelations are checked in groups of six, and the number of lags checked depends on the NLAG= option. The check for white noise output is shown in Figure 7.4.

**Figure 7.4** IDENTIFY Statement Check for White Noise

Autocorrelation Check for White Noise									
To	Lag	Chi-Square	DF	Pr > ChiSq	Autocorrelations				
	<b>6</b>	426.44	6	<.0001	0.957	0.907	0.852	0.791	0.726
	<b>12</b>	547.82	12	<.0001	0.588	0.514	0.440	0.370	0.303
	<b>18</b>	554.70	18	<.0001	0.174	0.112	0.052	-0.004	-0.054
	<b>24</b>	585.73	24	<.0001	-0.135	-0.167	-0.192	-0.211	-0.227

In this case, the white noise hypothesis is rejected very strongly, which is expected since the series is nonstationary. The *p*-value for the test of the first six autocorrelations is printed as <0.0001, which means the *p*-value is less than 0.0001.

### Identification of the Differenced Series

Since the series is nonstationary, the next step is to transform it to a stationary series by differencing. That is, instead of modeling the SALES series itself, you model the change in SALES from one period to the next. To difference the SALES series, use another IDENTIFY statement and specify that the first difference of SALES be analyzed, as shown in the following statements:

```
proc arima data=test;
  identify var=sales(1);
run;
```

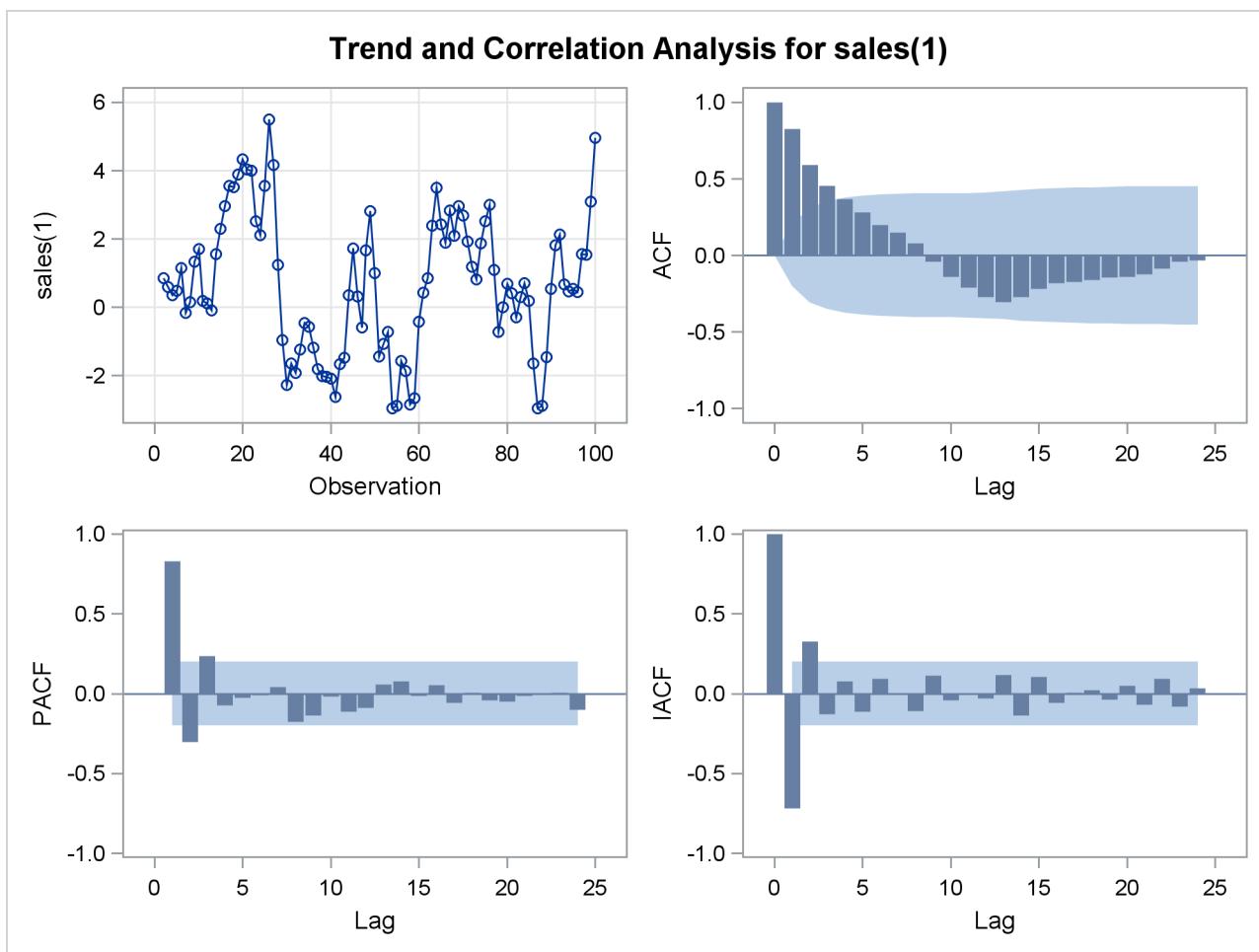
The second IDENTIFY statement produces the same information as the first, but for the change in SALES from one period to the next rather than for the total SALES in each period. The summary statistics output from this IDENTIFY statement is shown in Figure 7.5. Note that the period of differencing is given as 1, and one observation was lost through the differencing operation.

**Figure 7.5** IDENTIFY Statement Output for Differenced Series

### The ARIMA Procedure

Name of Variable = sales	
Period(s) of Differencing	1
Mean of Working Series	0.660589
Standard Deviation	2.011543
Number of Observations	99
Observation(s) eliminated by differencing	1

The autocorrelation plots for the differenced series are shown in Figure 7.6.

**Figure 7.6** Correlation Analysis of the Change in SALES

The autocorrelations decrease rapidly in this plot, indicating that the change in SALES is a stationary time series.

The next step in the Box-Jenkins methodology is to examine the patterns in the autocorrelation plot to choose candidate ARMA models to the series. The partial and inverse autocorrelation function plots are also useful aids in identifying appropriate ARMA models for the series.

In the usual Box-Jenkins approach to ARIMA modeling, the sample autocorrelation function, inverse autocorrelation function, and partial autocorrelation function are compared with the theoretical correlation functions expected from different kinds of ARMA models. This matching of theoretical autocorrelation functions of different ARMA models to the sample autocorrelation functions computed from the response series is the heart of the identification stage of Box-Jenkins modeling. Most textbooks on time series analysis, such as Pankratz (1983), discuss the theoretical autocorrelation functions for different kinds of ARMA models.

Since the input data are only a limited sample of the series, the sample autocorrelation functions computed from the input series only approximate the true autocorrelation function of the process that generates the series. This means that the sample autocorrelation functions do not exactly match the theoretical autocorrelation functions for any ARMA model and can have a pattern similar to that of several different ARMA models. If the series is white noise (a purely random process), then there is no need to fit a model. The check for

white noise, shown in Figure 7.7, indicates that the change in SALES is highly autocorrelated. Thus, an autocorrelation model, for example an AR(1) model, might be a good candidate model to fit to this process.

**Figure 7.7** IDENTIFY Statement Check for White Noise

Autocorrelation Check for White Noise									
To Lag	Chi-Square	DF	Pr > ChiSq	Autocorrelations					
6	154.44	6	<.0001	0.828	0.591	0.454	0.369	0.281	0.198
12	173.66	12	<.0001	0.151	0.081	-0.039	-0.141	-0.210	-0.274
18	209.64	18	<.0001	-0.305	-0.271	-0.218	-0.183	-0.174	-0.161
24	218.04	24	<.0001	-0.144	-0.141	-0.125	-0.085	-0.040	-0.032

## Estimation and Diagnostic Checking Stage

The autocorrelation plots for this series, as shown in the previous section, suggest an AR(1) model for the change in SALES. You should check the diagnostic statistics to see if the AR(1) model is adequate. Other candidate models include an MA(1) model and low-order mixed ARMA models. In this example, the AR(1) model is tried first.

### Estimating an AR(1) Model

The following statements fit an AR(1) model (an autoregressive model of order 1), which predicts the change in SALES as an average change, plus some fraction of the previous change, plus a random error. To estimate an AR model, you specify the order of the autoregressive model with the P= option in an ESTIMATE statement:

```
estimate p=1;
run;
```

The ESTIMATE statement fits the model to the data and prints parameter estimates and various diagnostic statistics that indicate how well the model fits the data. The first part of the ESTIMATE statement output, the table of parameter estimates, is shown in Figure 7.8.

**Figure 7.8** Parameter Estimates for AR(1) Model

### The ARIMA Procedure

Conditional Least Squares Estimation					
Parameter	Standard Estimate	Standard Error	Approx t Value	Pr >  t	Lag
MU	0.90280	0.65984	1.37	0.1744	0
AR1,1	0.86847	0.05485	15.83	<.0001	1

The table of parameter estimates is titled “Conditional Least Squares Estimation,” which indicates the estimation method used. You can request different estimation methods with the METHOD= option.

The table of parameter estimates lists the parameters in the model; for each parameter, the table shows the estimated value and the standard error and *t* value for the estimate. The table also indicates the lag at which the parameter appears in the model.

In this case, there are two parameters in the model. The mean term is labeled MU; its estimated value is 0.90280. The autoregressive parameter is labeled AR1,1; this is the coefficient of the lagged value of the change in SALES, and its estimate is 0.86847.

The *t* values provide significance tests for the parameter estimates and indicate whether some terms in the model might be unnecessary. In this case, the *t* value for the autoregressive parameter is 15.83, so this term is highly significant. The *t* value for MU indicates that the mean term adds little to the model.

The standard error estimates are based on large sample theory. Thus, the standard errors are labeled as approximate, and the standard errors and *t* values might not be reliable in small samples.

The next part of the ESTIMATE statement output is a table of goodness-of-fit statistics, which aid in comparing this model to other models. This output is shown in Figure 7.9.

**Figure 7.9** Goodness-of-Fit Statistics for AR(1) Model

<b>Constant Estimate</b>	0.118749
<b>Variance Estimate</b>	1.15794
<b>Std Error Estimate</b>	1.076076
<b>AIC</b>	297.4469
<b>SBC</b>	302.6372
<b>Number of Residuals</b>	99

The “Constant Estimate” is a function of the mean term MU and the autoregressive parameters. This estimate is computed only for AR or ARMA models, but not for strictly MA models. See the section “General Notation for ARIMA Models” on page 201 for an explanation of the constant estimate.

The “Variance Estimate” is the variance of the residual series, which estimates the innovation variance. The item labeled “Std Error Estimate” is the square root of the variance estimate. In general, when you are comparing candidate models, smaller AIC and SBC statistics indicate the better fitting model. The section “Estimation Details” on page 240 explains the AIC and SBC statistics.

The ESTIMATE statement next prints a table of correlations of the parameter estimates, as shown in Figure 7.10. This table can help you assess the extent to which collinearity might have influenced the results. If two parameter estimates are very highly correlated, you might consider dropping one of them from the model.

**Figure 7.10** Correlations of the Estimates for AR(1) Model

<b>Parameter</b>	<b>Correlations of Parameter Estimates</b>	
	MU	AR1,1
MU	1.000	0.114
AR1,1	0.114	1.000

The next part of the ESTIMATE statement output is a check of the autocorrelations of the residuals. This output has the same form as the autocorrelation check for white noise that the IDENTIFY statement prints for the response series. The autocorrelation check of residuals is shown in Figure 7.11.

**Figure 7.11** Check for White Noise Residuals for AR(1) Model

Autocorrelation Check of Residuals										
To	Lag	Chi-Square	DF	Pr > ChiSq	Autocorrelations					
6	6	19.09	5	0.0019	0.327	-0.220	-0.128	0.068	-0.002	-0.096
12	12	22.90	11	0.0183	0.072	0.116	-0.042	-0.066	0.031	-0.091
18	18	31.63	17	0.0167	-0.233	-0.129	-0.024	0.056	-0.014	-0.008
24	24	32.83	23	0.0841	0.009	-0.057	-0.057	-0.001	0.049	-0.015

The  $\chi^2$  test statistics for the residuals series indicate whether the residuals are uncorrelated (white noise) or contain additional information that might be used by a more complex model. In this case, the test statistics reject the no-autocorrelation hypothesis at a high level of significance ( $p = 0.0019$  for the first six lags.) This means that the residuals are not white noise, and so the AR(1) model is not a fully adequate model for this series. The ESTIMATE statement output also includes graphical analysis of the residuals. It is not shown here. The graphical analysis also reveals the inadequacy of the AR(1) model.

The final part of the ESTIMATE statement output is a listing of the estimated model, using the backshift notation. This output is shown in [Figure 7.12](#).

**Figure 7.12** Estimated ARIMA(1, 1, 0) Model for SALES

Model for variable sales	
Estimated Mean	0.902799
Period(s) of Differencing	1
Autoregressive Factors	
Factor 1: 1 - 0.86847 B**(1)	

This listing combines the differencing specification given in the IDENTIFY statement with the parameter estimates of the model for the change in SALES. Since the AR(1) model is for the change in SALES, the final model for SALES is an ARIMA(1,1,0) model. Using  $B$ , the backshift operator, the mathematical form of the estimated model shown in this output is as follows:

$$(1 - B)sales_t = 0.902799 + \frac{1}{(1 - 0.86847B)}a_t$$

See the section “[General Notation for ARIMA Models](#)” on page 201 for further explanation of this notation.

### Estimating an ARMA(1,1) Model

The IDENTIFY statement plots suggest a mixed autoregressive and moving-average model, and the previous ESTIMATE statement check of residuals indicates that an AR(1) model is not sufficient. You now try estimating an ARMA(1,1) model for the change in SALES.

An ARMA(1,1) model predicts the change in SALES as an average change, plus some fraction of the previous change, plus a random error, plus some fraction of the random error in the preceding period. An ARMA(1,1) model for the change in SALES is the same as an ARIMA(1,1,1) model for the level of SALES.

To estimate a mixed autoregressive moving-average model, you specify the order of the moving-average part of the model with the Q= option in an ESTIMATE statement in addition to specifying the order of the autoregressive part with the P= option. The following statements fit an ARMA(1,1) model to the differenced SALES series:

```
estimate p=1 q=1;
run;
```

The parameter estimates table and goodness-of-fit statistics for this model are shown in Figure 7.13.

**Figure 7.13** Estimated ARMA(1, 1) Model for Change in SALES

### The ARIMA Procedure

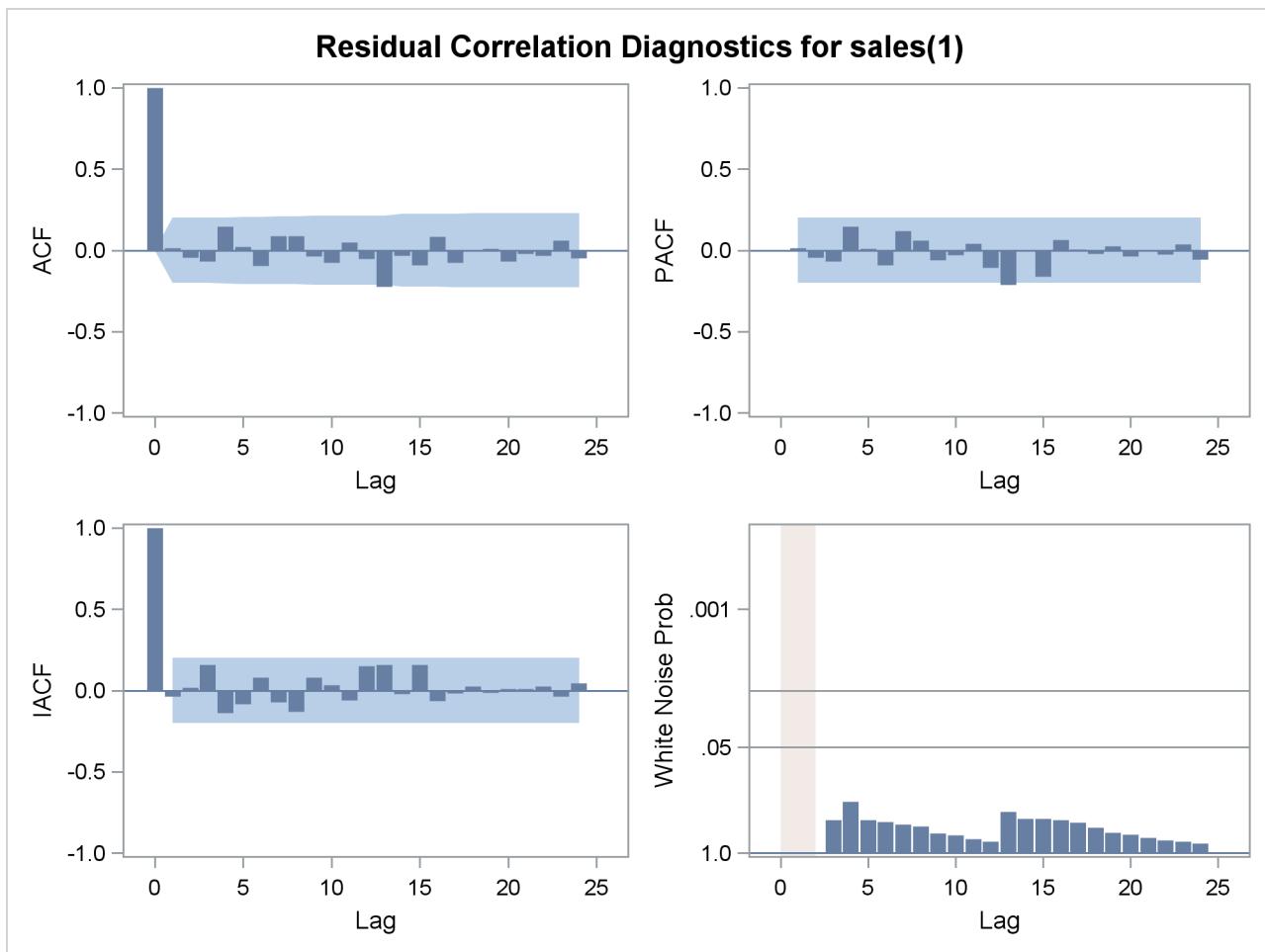
Conditional Least Squares Estimation					
Parameter	Estimate	Standard Error	t Value	Approx Pr >  t	Lag
MU	0.89288	0.49391	1.81	0.0738	0
MA1,1	-0.58935	0.08988	-6.56	<.0001	1
AR1,1	0.74755	0.07785	9.60	<.0001	1

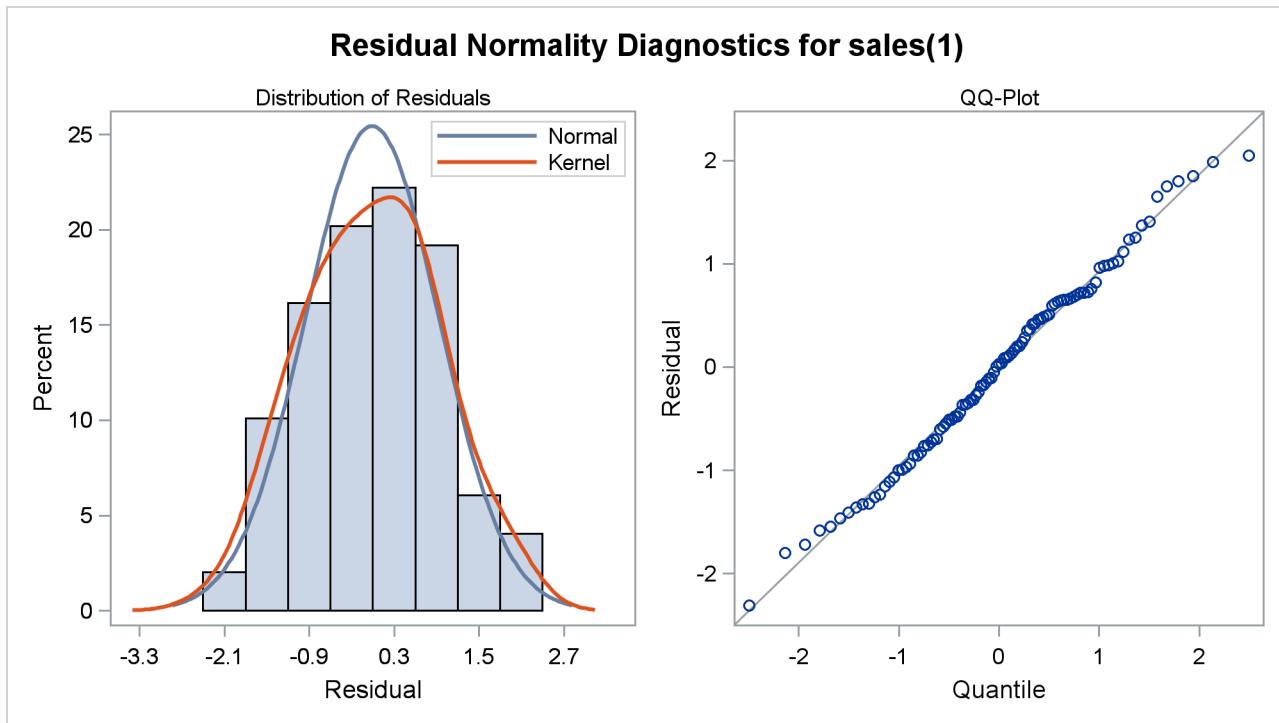
  

Constant Estimate	0.225409
Variance Estimate	0.904034
Std Error Estimate	0.950807
AIC	273.9155
SBC	281.7009
Number of Residuals	99

The moving-average parameter estimate, labeled “MA1,1”, is -0.58935. Both the moving-average and the autoregressive parameters have significant *t* values. Note that the variance estimate, AIC, and SBC are all smaller than they were for the AR(1) model, indicating that the ARMA(1,1) model fits the data better without over-parameterizing.

The graphical check of the residuals from this model is shown in Figure 7.14 and Figure 7.15. The residual correlation and white noise test plots show that you cannot reject the hypothesis that the residuals are uncorrelated. The normality plots also show no departure from normality. Thus, you conclude that the ARMA(1,1) model is adequate for the change in SALES series, and there is no point in trying more complex models.

**Figure 7.14** White Noise Check of Residuals for the ARMA(1,1) Model

**Figure 7.15** Normality Check of Residuals for the ARMA(1,1) Model

The form of the estimated ARIMA(1,1,1) model for SALES is shown in [Figure 7.16](#).

**Figure 7.16** Estimated ARIMA(1,1,1) Model for SALES

Model for variable sales	
Estimated Mean	0.892875
Period(s) of Differencing	1
Autoregressive Factors	
Factor 1: $1 - 0.74755 B^{**}(1)$	
Moving Average Factors	
Factor 1: $1 + 0.58935 B^{**}(1)$	

The estimated model shown in this output is

$$(1 - B)sales_t = 0.892875 + \frac{(1 + 0.58935B)}{(1 - 0.74755B)}a_t$$

In addition to the residual analysis of a model, it is often useful to check whether there are any changes in the time series that are not accounted for by the currently estimated model. The OUTLIER statement enables you to detect such changes. For a long series, this task can be computationally burdensome. Therefore, in general, it is better done after a model that fits the data reasonably well has been found. [Figure 7.17](#) shows the output of the simplest form of the OUTLIER statement:

```
outlier;
run;
```

Two possible outliers have been found for the model in question. See the section “[Detecting Outliers](#)” on page 250, and the examples [Example 7.6](#) and [Example 7.7](#), for more details about modeling in the presence of outliers. In this illustration these outliers are not discussed any further.

**Figure 7.17** Outliers for the ARIMA(1,1,1) Model for SALES

The ARIMA Procedure			
Outlier Detection Summary			
Maximum number searched	2		
Number found	2		
Significance used	0.05		

Outlier Details				
Obs	Type	Estimate	Chi-Square	Approx Prob>ChiSq
10	Additive	0.56879	4.20	0.0403
67	Additive	0.55698	4.42	0.0355

Since the model diagnostic tests show that all the parameter estimates are significant and the residual series is white noise, the estimation and diagnostic checking stage is complete. You can now proceed to forecasting the SALES series with this ARIMA(1,1,1) model.

## Forecasting Stage

To produce the forecast, use a FORECAST statement after the ESTIMATE statement for the model you decide is best. If the last model fit is not the best, then repeat the ESTIMATE statement for the best model before you use the FORECAST statement.

Suppose that the SALES series is monthly, that you want to forecast one year ahead from the most recently available SALES figure, and that the dates for the observations are given by a variable DATE in the input data set TEST. You use the following FORECAST statement:

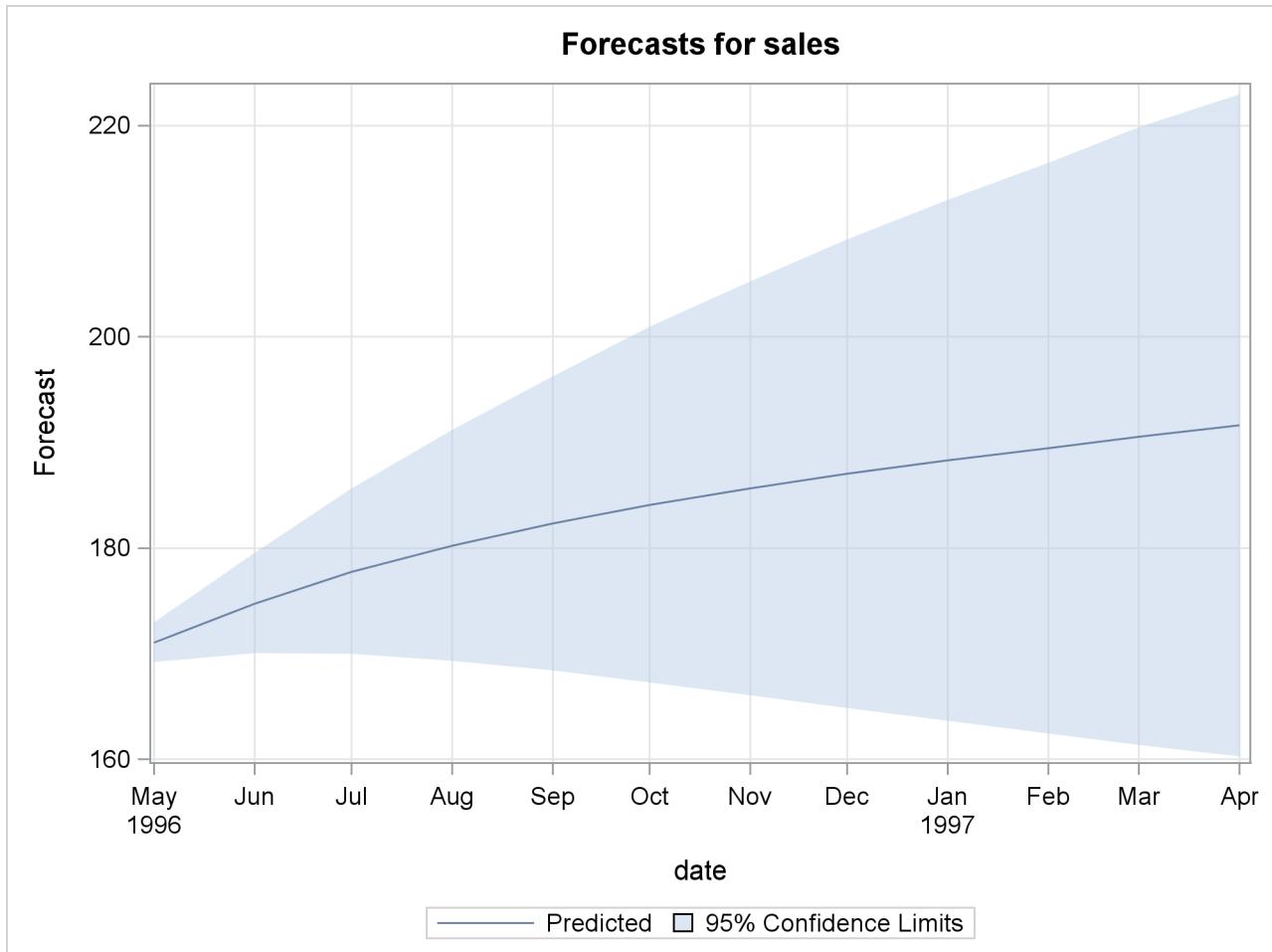
```
forecast lead=12 interval=month id=date out=results;
run;
```

The LEAD= option specifies how many periods ahead to forecast (12 months, in this case). The ID= option specifies the ID variable, which is typically a SAS *date*, *time*, or *datetime* variable, used to date the observations of the SALES time series. The INTERVAL= option indicates that data are monthly and enables PROC ARIMA to extrapolate DATE values for forecast periods. The OUT= option writes the forecasts to the output data set RESULTS. See the section “[OUT= Data Set](#)” on page 252 for information about the contents of the output data set.

By default, the FORECAST statement also prints and plots the forecast values, as shown in [Figure 7.18](#) and [Figure 7.19](#). The forecast table shows for each forecast period the observation number, forecast value, standard error estimate for the forecast value, and lower and upper limits for a 95% confidence interval for the forecast.

**Figure 7.18** Forecasts for ARIMA(1,1,1) Model for SALES**The ARIMA Procedure**

Obs	Forecast	Std Error	95% Confidence Limits	
			Lower	Upper
101	171.0320	0.9508	169.1684	172.8955
102	174.7534	2.4168	170.0165	179.4903
103	177.7608	3.9879	169.9445	185.5770
104	180.2343	5.5658	169.3256	191.1430
105	182.3088	7.1033	168.3866	196.2310
106	184.0850	8.5789	167.2707	200.8993
107	185.6382	9.9841	166.0698	205.2066
108	187.0247	11.3173	164.8433	209.2061
109	188.2866	12.5807	163.6289	212.9443
110	189.4553	13.7784	162.4501	216.4605
111	190.5544	14.9153	161.3209	219.7879
112	191.6014	15.9964	160.2491	222.9538

**Figure 7.19** Forecasts for the ARMA(1,1,1) Model

Normally, you want the forecast values stored in an output data set, and you are not interested in seeing this printed list of the forecast. You can use the NOPRINT option in the FORECAST statement to suppress this output.

## Using ARIMA Procedure Statements

The IDENTIFY, ESTIMATE, and FORECAST statements are related in a hierarchy. An IDENTIFY statement brings in a time series to be modeled; several ESTIMATE statements can follow to estimate different ARIMA models for the series; for each model estimated, several FORECAST statements can be used. Thus, a FORECAST statement must be preceded at some point by an ESTIMATE statement, and an ESTIMATE statement must be preceded at some point by an IDENTIFY statement. Additional IDENTIFY statements can be used to switch to modeling a different response series or to change the degree of differencing used.

The ARIMA procedure can be used interactively in the sense that all ARIMA procedure statements can be executed any number of times without reinvoking PROC ARIMA. You can execute ARIMA procedure statements singly or in groups by following the single statement or group of statements with a RUN statement. The output for each statement or group of statements is produced when the RUN statement is entered.

A RUN statement does not terminate the PROC ARIMA step but tells the procedure to execute the statements given so far. You can end PROC ARIMA by submitting a QUIT statement, a DATA step, another PROC step, or an ENDSAS statement.

The example in the preceding section illustrates the interactive use of ARIMA procedure statements. The complete PROC ARIMA program for that example is as follows:

```
proc arima data=test;
  identify var=sales nlag=24;
  run;
  identify var=sales(1);
  run;
  estimate p=1;
  run;
  estimate p=1 q=1;
  run;
  outlier;
  run;
  forecast lead=12 interval=month id=date out=results;
  run;
quit;
```

## General Notation for ARIMA Models

The order of an ARIMA (autoregressive integrated moving-average) model is usually denoted by the notation ARIMA( $p,d,q$ ), where

- $p$  is the order of the autoregressive part
- $d$  is the order of the differencing

$q$  is the order of the moving-average process

If no differencing is done ( $d = 0$ ), the models are usually referred to as ARMA( $p, q$ ) models. The final model in the preceding example is an ARIMA(1,1,1) model since the IDENTIFY statement specified  $d = 1$ , and the final ESTIMATE statement specified  $p = 1$  and  $q = 1$ .

## Notation for Pure ARIMA Models

Mathematically the pure ARIMA model is written as

$$W_t = \mu + \frac{\theta(B)}{\phi(B)} a_t$$

where

$t$	indexes time
$W_t$	is the response series $Y_t$ or a difference of the response series
$\mu$	is the mean term
$B$	is the backshift operator; that is, $BX_t = X_{t-1}$
$\phi(B)$	is the autoregressive operator, represented as a polynomial in the backshift operator: $\phi(B) = 1 - \phi_1 B - \dots - \phi_p B^p$
$\theta(B)$	is the moving-average operator, represented as a polynomial in the backshift operator: $\theta(B) = 1 - \theta_1 B - \dots - \theta_q B^q$
$a_t$	is the independent disturbance, also called the random error

The series  $W_t$  is computed by the IDENTIFY statement and is the series processed by the ESTIMATE statement. Thus,  $W_t$  is either the response series  $Y_t$  or a difference of  $Y_t$  specified by the differencing operators in the IDENTIFY statement.

For simple (nonseasonal) differencing,  $W_t = (1 - B)^d Y_t$ . For seasonal differencing  $W_t = (1 - B)^d (1 - B^s)^D Y_t$ , where  $d$  is the degree of nonseasonal differencing,  $D$  is the degree of seasonal differencing, and  $s$  is the length of the seasonal cycle.

For example, the mathematical form of the ARIMA(1,1,1) model estimated in the preceding example is

$$(1 - B)Y_t = \mu + \frac{(1 - \theta_1 B)}{(1 - \phi_1 B)} a_t$$

## Model Constant Term

The ARIMA model can also be written as

$$\phi(B)(W_t - \mu) = \theta(B)a_t$$

or

$$\phi(B)W_t = const + \theta(B)a_t$$

where

$$\text{const} = \phi(B)\mu = \mu - \phi_1\mu - \phi_2\mu - \dots - \phi_p\mu$$

Thus, when an autoregressive operator and a mean term are both included in the model, the constant term for the model can be represented as  $\phi(B)\mu$ . This value is printed with the label “Constant Estimate” in the ESTIMATE statement output.

## Notation for Transfer Function Models

The general ARIMA model with input series, also called the ARIMAX model, is written as

$$W_t = \mu + \sum_i \frac{\omega_i(B)}{\delta_i(B)} B^{k_i} X_{i,t} + \frac{\theta(B)}{\phi(B)} a_t$$

where

$X_{i,t}$	is the $i$ th input time series or a difference of the $i$ th input series at time $t$
$k_i$	is the pure time delay for the effect of the $i$ th input series
$\omega_i(B)$	is the numerator polynomial of the transfer function for the $i$ th input series
$\delta_i(B)$	is the denominator polynomial of the transfer function for the $i$ th input series.

The model can also be written more compactly as

$$W_t = \mu + \sum_i \Psi_i(B) X_{i,t} + n_t$$

where

$\Psi_i(B)$	is the transfer function for the $i$ th input series modeled as a ratio of the $\omega$ and $\delta$ polynomials: $\Psi_i(B) = (\omega_i(B)/\delta_i(B))B^{k_i}$
$n_t$	is the noise series: $n_t = (\theta(B)/\phi(B))a_t$

This model expresses the response series as a combination of past values of the random shocks and past values of other input series. The response series is also called the *dependent series* or *output series*. An input time series is also referred to as an *independent series* or a *predictor series*. Response variable, dependent variable, independent variable, or predictor variable are other terms often used.

## Notation for Factored Models

ARIMA models are sometimes expressed in a factored form. This means that the  $\phi$ ,  $\theta$ ,  $\omega$ , or  $\delta$  polynomials are expressed as products of simpler polynomials. For example, you could express the pure ARIMA model as

$$W_t = \mu + \frac{\theta_1(B)\theta_2(B)}{\phi_1(B)\phi_2(B)} a_t$$

where  $\phi_1(B)\phi_2(B) = \phi(B)$  and  $\theta_1(B)\theta_2(B) = \theta(B)$ .

When an ARIMA model is expressed in factored form, the order of the model is usually expressed by using a factored notation also. The order of an ARIMA model expressed as the product of two factors is denoted as ARIMA( $p,d,q$ ) $\times(P,D,Q)$ .

## Notation for Seasonal Models

ARIMA models for time series with regular seasonal fluctuations often use differencing operators and autoregressive and moving-average parameters at lags that are multiples of the length of the seasonal cycle. When all the terms in an ARIMA model factor refer to lags that are a multiple of a constant  $s$ , the constant is factored out and suffixed to the ARIMA( $p,d,q$ ) notation.

Thus, the general notation for the order of a seasonal ARIMA model with both seasonal and nonseasonal factors is ARIMA( $p,d,q$ ) $\times(P,D,Q)_s$ . The term ( $p,d,q$ ) gives the order of the nonseasonal part of the ARIMA model; the term ( $P,D,Q$ ) $_s$  gives the order of the seasonal part. The value of  $s$  is the number of observations in a seasonal cycle: 12 for monthly series, 4 for quarterly series, 7 for daily series with day-of-week effects, and so forth.

For example, the notation ARIMA(0,1,2) $\times(0,1,1)_{12}$  describes a seasonal ARIMA model for monthly data with the following mathematical form:

$$(1 - B)(1 - B^{12})Y_t = \mu + (1 - \theta_{1,1}B - \theta_{1,2}B^2)(1 - \theta_{2,1}B^{12})a_t$$

## Stationarity

The noise (or residual) series for an ARMA model must be *stationary*, which means that both the expected values of the series and its autocovariance function are independent of time.

The standard way to check for nonstationarity is to plot the series and its autocorrelation function. You can visually examine a graph of the series over time to see if it has a visible trend or if its variability changes noticeably over time. If the series is nonstationary, its autocorrelation function will usually decay slowly.

Another way of checking for stationarity is to use the stationarity tests described in the section “[Stationarity Tests](#)” on page 238.

Most time series are nonstationary and must be transformed to a stationary series before the ARIMA modeling process can proceed. If the series has a nonstationary variance, taking the log of the series can help. You can compute the log values in a DATA step and then analyze the log values with PROC ARIMA.

If the series has a trend over time, seasonality, or some other nonstationary pattern, the usual solution is to take the difference of the series from one period to the next and then analyze this differenced series. Sometimes a series might need to be differenced more than once or differenced at lags greater than one period. (If the trend or seasonal effects are very regular, the introduction of explanatory variables can be an appropriate alternative to differencing.)

## Differencing

Differencing of the response series is specified with the VAR= option of the IDENTIFY statement by placing a list of differencing periods in parentheses after the variable name. For example, to take a simple first difference of the series SALES, use the statement

```
identify var=sales(1);
```

In this example, the change in SALES from one period to the next is analyzed.

A deterministic seasonal pattern also causes the series to be nonstationary, since the expected value of the series is not the same for all time periods but is higher or lower depending on the season. When the series has a seasonal pattern, you might want to difference the series at a lag that corresponds to the length of the seasonal cycle. For example, if SALES is a monthly series, the statement

```
identify var=sales(12);
```

takes a seasonal difference of SALES, so that the series analyzed is the change in SALES from its value in the same month one year ago.

To take a second difference, add another differencing period to the list. For example, the following statement takes the second difference of SALES:

```
identify var=sales(1,1);
```

That is, SALES is differenced once at lag 1 and then differenced again, also at lag 1. The statement

```
identify var=sales(2);
```

creates a 2-span difference—that is, current period SALES minus SALES from two periods ago. The statement

```
identify var=sales(1,12);
```

takes a second-order difference of SALES, so that the series analyzed is the difference between the current period-to-period change in SALES and the change 12 periods ago. You might want to do this if the series had both a trend over time and a seasonal pattern.

There is no limit to the order of differencing and the degree of lagging for each difference.

Differencing not only affects the series used for the IDENTIFY statement output but also applies to any following ESTIMATE and FORECAST statements. ESTIMATE statements fit ARMA models to the differenced series. FORECAST statements forecast the differences and automatically sum these differences back to undo the differencing operation specified by the IDENTIFY statement, thus producing the final forecast result.

Differencing of input series is specified by the CROSSCORR= option and works just like differencing of the response series. For example, the statement

```
identify var=y(1) crosscorr=(x1(1) x2(1));
```

takes the first difference of Y, the first difference of X1, and the first difference of X2. Whenever X1 and X2 are used in INPUT= options in following ESTIMATE statements, these names refer to the differenced series.

---

## Subset, Seasonal, and Factored ARMA Models

The simplest way to specify an ARMA model is to give the order of the AR and MA parts with the P= and Q= options. When you do this, the model has parameters for the AR and MA parts for all lags through the order specified. However, you can control the form of the ARIMA model exactly as shown in the following section.

### Subset Models

You can control which lags have parameters by specifying the P= or Q= option as a list of lags in parentheses. A model that includes parameters for only some lags is sometimes called a *subset* or *additive model*. For example, consider the following two ESTIMATE statements:

```
identify var=sales;
estimate p=4;
estimate p=(1 4);
```

Both specify AR(4) models, but the first has parameters for lags 1, 2, 3, and 4, while the second has parameters for lags 1 and 4, with the coefficients for lags 2 and 3 constrained to 0. The mathematical form of the autoregressive models produced by these two specifications is shown in [Table 7.1](#).

**Table 7.1** Saturated versus Subset Models

Option	Autoregressive Operator
P=4	$(1 - \phi_1 B - \phi_2 B^2 - \phi_3 B^3 - \phi_4 B^4)$
P=(1 4)	$(1 - \phi_1 B - \phi_4 B^4)$

### Seasonal Models

One particularly useful kind of subset model is a *seasonal model*. When the response series has a seasonal pattern, the values of the series at the same time of year in previous years can be important for modeling the series. For example, if the series SALES is observed monthly, the statements

```
identify var=sales;
estimate p=(12);
```

model SALES as an average value plus some fraction of its deviation from this average value a year ago, plus a random error. Although this is an AR(12) model, it has only one autoregressive parameter.

### Factored Models

A factored model (also referred to as a multiplicative model) represents the ARIMA model as a product of simpler ARIMA models. For example, you might model SALES as a combination of an AR(1) process that reflects short term dependencies and an AR(12) model that reflects the seasonal pattern.

It might seem that the way to do this is with the option P=(1 12), but the AR(1) process also operates in past years; you really need autoregressive parameters at lags 1, 12, and 13. You can specify a subset model

with separate parameters at these lags, or you can specify a factored model that represents the model as the product of an AR(1) model and an AR(12) model. Consider the following two ESTIMATE statements:

```
identify var=sales;
estimate p=(1 12 13);
estimate p=(1) (12);
```

The mathematical form of the autoregressive models produced by these two specifications are shown in Table 7.2.

**Table 7.2** Subset versus Factored Models

Option	Autoregressive Operator
P=(1 12 13)	$(1 - \phi_1 B - \phi_{12} B^{12} - \phi_{13} B^{13})$
P=(1)(12)	$(1 - \phi_1 B)(1 - \phi_{12} B^{12})$

Both models fit by these two ESTIMATE statements predict SALES from its values 1, 12, and 13 periods ago, but they use different parameterizations. The first model has three parameters, whose meanings may be hard to interpret.

The factored specification P=(1)(12) represents the model as the product of two different AR models. It has only two parameters: one that corresponds to recent effects and one that represents seasonal effects. Thus the factored model is more parsimonious, and its parameter estimates are more clearly interpretable.

## Input Variables and Regression with ARMA Errors

In addition to past values of the response series and past errors, you can also model the response series using the current and past values of other series, called *input series*.

Several different names are used to describe ARIMA models with input series. *Transfer function model*, *intervention model*, *interrupted time series model*, *regression model with ARMA errors*, *Box-Tiao model*, and *ARIMAX model* are all different names for ARIMA models with input series. Pankratz (1991) refers to these models as *dynamic regression* models.

### Using Input Series

To use input series, list the input series in a CROSSCORR= option on the IDENTIFY statement and specify how they enter the model with an INPUT= option on the ESTIMATE statement. For example, you might use a series called PRICE to help model SALES, as shown in the following statements:

```
proc arima data=a;
  identify var=sales crosscorr=price;
  estimate input=price;
run;
```

This example performs a simple linear regression of SALES on PRICE; it produces the same results as PROC REG or another SAS regression procedure. The mathematical form of the model estimated by these statements is

$$Y_t = \mu + \omega_0 X_t + a_t$$

The parameter estimates table for this example (using simulated data) is shown in Figure 7.20. The intercept parameter is labeled MU. The regression coefficient for PRICE is labeled NUM1. (See the section “[Naming of Model Parameters](#)” on page 247 for information about how parameters for input series are named.)

**Figure 7.20** Parameter Estimates Table for Regression Model

### The ARIMA Procedure

Conditional Least Squares Estimation							
Parameter	Estimate	Standard Error	t Value	Pr >  t	Approx Lag	Variable	Shift
MU	199.83602	2.99463	66.73	<.0001	0	sales	0
NUM1	-9.99299	0.02885	-346.38	<.0001	0	price	0

Any number of input variables can be used in a model. For example, the following statements fit a multiple regression of SALES on PRICE and INCOME:

```
proc arima data=a;
  identify var=sales crosscorr=(price income);
  estimate input=(price income);
run;
```

The mathematical form of the regression model estimated by these statements is

$$Y_t = \mu + \omega_1 X_{1,t} + \omega_2 X_{2,t} + a_t$$

### Lagging and Differencing Input Series

You can also difference and lag the input series. For example, the following statements regress the change in SALES on the change in PRICE lagged by one period. The difference of PRICE is specified with the CROSSCORR= option and the lag of the change in PRICE is specified by the 1 \$ in the INPUT= option.

```
proc arima data=a;
  identify var=sales(1) crosscorr=price(1);
  estimate input=( 1 $ price );
run;
```

These statements estimate the model

$$(1 - B)Y_t = \mu + \omega_0(1 - B)X_{t-1} + a_t$$

### Regression with ARMA Errors

You can combine input series with ARMA models for the errors. For example, the following statements regress SALES on INCOME and PRICE but with the error term of the regression model (called the *noise series* in ARIMA modeling terminology) assumed to be an ARMA(1,1) process.

```
proc arima data=a;
  identify var=sales crosscorr=(price income);
  estimate p=1 q=1 input=(price income);
run;
```

These statements estimate the model

$$Y_t = \mu + \omega_1 X_{1,t} + \omega_2 X_{2,t} + \frac{(1 - \theta_1 B)}{(1 - \phi_1 B)} a_t$$

## Stationarity and Input Series

Note that the requirement of stationarity applies to the noise series. If there are no input variables, the response series (after differencing and minus the mean term) and the noise series are the same. However, if there are inputs, the noise series is the residual after the effect of the inputs is removed.

There is no requirement that the input series be stationary. If the inputs are nonstationary, the response series will be nonstationary, even though the noise process might be stationary.

When nonstationary input series are used, you can fit the input variables first with no ARMA model for the errors and then consider the stationarity of the residuals before identifying an ARMA model for the noise part.

## Identifying Regression Models with ARMA Errors

Previous sections described the ARIMA modeling identification process that uses the autocorrelation function plots produced by the IDENTIFY statement. This identification process does not apply when the response series depends on input variables. This is because it is the noise process for which you need to identify an ARIMA model, and when input series are involved the response series adjusted for the mean is no longer an estimate of the noise series.

However, if the input series are independent of the noise series, you can use the residuals from the regression model as an estimate of the noise series, then apply the ARIMA modeling identification process to this residual series. This assumes that the noise process is stationary.

The PLOT option in the ESTIMATE statement produces similar plots for the model residuals as the IDENTIFY statement produces for the response series. The PLOT option prints an autocorrelation function plot, an inverse autocorrelation function plot, and a partial autocorrelation function plot for the residual series. Note that these residual correlation plots are produced by default.

The following statements show how the PLOT option is used to identify the ARMA(1,1) model for the noise process used in the preceding example of regression with ARMA errors:

```
proc arima data=a;
  identify var=sales crosscorr=(price income) noprint;
  estimate input=(price income) plot;
  run;
  estimate p=1 q=1 input=(price income);
run;
```

In this example, the IDENTIFY statement includes the NOPRINT option since the autocorrelation plots for the response series are not useful when you know that the response series depends on input series.

The first ESTIMATE statement fits the regression model with no model for the noise process. The PLOT option produces plots of the autocorrelation function, inverse autocorrelation function, and partial autocorrelation function for the residual series of the regression on PRICE and INCOME.

By examining the PLOT option output for the residual series, you verify that the residual series is stationary and identify an ARMA(1,1) model for the noise process. The second ESTIMATE statement fits the final model.

Although this discussion addresses regression models, the same remarks apply to identifying an ARIMA model for the noise process in models that include input series with complex transfer functions.

## Intervention Models and Interrupted Time Series

One special kind of ARIMA model with input series is called an *intervention model* or *interrupted time series* model. In an intervention model, the input series is an indicator variable that contains discrete values that flag the occurrence of an event affecting the response series. This event is an intervention in or an interruption of the normal evolution of the response time series, which, in the absence of the intervention, is usually assumed to be a pure ARIMA process.

Intervention models can be used both to model and forecast the response series and also to analyze the impact of the intervention. When the focus is on estimating the effect of the intervention, the process is often called *intervention analysis* or *interrupted time series analysis*.

### Impulse Interventions

The intervention can be a one-time event. For example, you might want to study the effect of a short-term advertising campaign on the sales of a product. In this case, the input variable has the value of 1 for the period during which the advertising campaign took place and the value 0 for all other periods. Intervention variables of this kind are sometimes called *impulse functions* or *pulse functions*.

Suppose that SALES is a monthly series, and a special advertising effort was made during the month of March 1992. The following statements estimate the effect of this intervention by assuming an ARMA(1,1) model for SALES. The model is specified just like the regression model, but the intervention variable AD is constructed in the DATA step as a zero-one indicator for the month of the advertising effort.

```

data a;
  set a;
  ad = (date = '1mar1992'd);
run;

proc arima data=a;
  identify var=sales crosscorr=ad;
  estimate p=1 q=1 input=ad;
run;

```

## Continuing Interventions

Other interventions can be continuing, in which case the input variable flags periods before and after the intervention. For example, you might want to study the effect of a change in tax rates on some economic measure. Another example is a study of the effect of a change in speed limits on the rate of traffic fatalities. In this case, the input variable has the value 1 after the new speed limit went into effect and the value 0 before. Intervention variables of this kind are called *step functions*.

Another example is the effect of news on product demand. Suppose it was reported in July 1996 that consumption of the product prevents heart disease (or causes cancer), and SALES is consistently higher (or lower) thereafter. The following statements model the effect of this news intervention:

```

data a;
  set a;
  news = (date >= '1jul1996'd);
run;

proc arima data=a;
  identify var=sales crosscorr=news;
  estimate p=1 q=1 input=news;
run;

```

## Interaction Effects

You can include any number of intervention variables in the model. Intervention variables can have any pattern—impulse and continuing interventions are just two possible cases. You can mix discrete valued intervention variables and continuous regressor variables in the same model.

You can also form interaction effects by multiplying input variables and including the product variable as another input. Indeed, as long as the dependent measure is continuous and forms a regular time series, you can use PROC ARIMA to fit any general linear model in conjunction with an ARMA model for the error process by using input variables that correspond to the columns of the design matrix of the linear model.

## Rational Transfer Functions and Distributed Lag Models

How an input series enters the model is called its *transfer function*. Thus, ARIMA models with input series are sometimes referred to as transfer function models.

In the preceding regression and intervention model examples, the transfer function is a single scale parameter. However, you can also specify complex transfer functions composed of numerator and denominator polynomials in the backshift operator. These transfer functions operate on the input series in the same way that the ARMA specification operates on the error term.

## Numerator Factors

For example, suppose you want to model the effect of PRICE on SALES as taking place gradually with the impact distributed over several past lags of PRICE. This is illustrated by the following statements:

```

proc arima data=a;
  identify var=sales crosscorr=price;
  estimate input=( (1 2 3) price );
run;

```

These statements estimate the model

$$Y_t = \mu + (\omega_0 - \omega_1 B - \omega_2 B^2 - \omega_3 B^3) X_t + a_t$$

This example models the effect of PRICE on SALES as a linear function of the current and three most recent values of PRICE. It is equivalent to a multiple linear regression of SALES on PRICE, LAG(PRICE), LAG2(PRICE), and LAG3(PRICE).

This is an example of a transfer function with one *numerator factor*. The numerator factors for a transfer function for an input series are like the MA part of the ARMA model for the noise series.

## Denominator Factors

You can also use transfer functions with *denominator factors*. The denominator factors for a transfer function for an input series are like the AR part of the ARMA model for the noise series. Denominator factors introduce exponentially weighted, infinite distributed lags into the transfer function.

To specify transfer functions with denominator factors, place the denominator factors after a slash (/) in the INPUT= option. For example, the following statements estimate the PRICE effect as an infinite distributed lag model with exponentially declining weights:

```

proc arima data=a;
  identify var=sales crosscorr=price;
  estimate input=( / (1) price );
run;

```

The transfer function specified by these statements is as follows:

$$\frac{\omega_0}{(1 - \delta_1 B)} X_t$$

This transfer function also can be written in the following equivalent form:

$$\omega_0 \left( 1 + \sum_{i=1}^{\infty} \delta_1^i B^i \right) X_t$$

This transfer function can be used with intervention inputs. When it is used with a pulse function input, the result is an intervention effect that dies out gradually over time. When it is used with a step function input, the result is an intervention effect that increases gradually to a limiting value.

## Rational Transfer Functions

By combining various numerator and denominator factors in the INPUT= option, you can specify *rational transfer functions* of any complexity. To specify an input with a general rational transfer function of the form

$$\frac{\omega(B)}{\delta(B)} B^k X_t$$

use an INPUT= option in the ESTIMATE statement of the form

**input=( k \$ ( ω-lags ) / ( δ-lags ) x)**

See the section “Specifying Inputs and Transfer Functions” on page 244 for more information.

## Identifying Transfer Function Models

The CROSSCORR= option of the IDENTIFY statement prints sample cross-correlation functions that show the correlation between the response series and the input series at different lags. The sample cross-correlation function can be used to help identify the form of the transfer function appropriate for an input series. See textbooks on time series analysis for information about using cross-correlation functions to identify transfer function models.

For the cross-correlation function to be meaningful, the input and response series must be filtered with a prewhitening model for the input series. See the section “Prewhitening” on page 238 for more information about this issue.

## Forecasting with Input Variables

To forecast a response series by using an ARIMA model with inputs, you need values of the input series for the forecast periods. You can supply values for the input variables for the forecast periods in the DATA= data set, or you can have PROC ARIMA forecast the input variables.

If you do not have future values of the input variables in the input data set used by the FORECAST statement, the input series must be forecast before the ARIMA procedure can forecast the response series. If you fit an ARIMA model to each of the input series for which you need forecasts before fitting the model for the response series, the FORECAST statement automatically uses the ARIMA models for the input series to generate the needed forecasts of the inputs.

For example, suppose you want to forecast SALES for the next 12 months. In this example, the change in SALES is predicted as a function of the change in PRICE, plus an ARMA(1,1) noise process. To forecast SALES by using PRICE as an input, you also need to fit an ARIMA model for PRICE.

The following statements fit an AR(2) model to the change in PRICE before fitting and forecasting the model for SALES. The FORECAST statement automatically forecasts PRICE using this AR(2) model to get the future inputs needed to produce the forecast of SALES.

```

proc arima data=a;
  identify var=price(1);
  estimate p=2;
  identify var=sales(1) crosscorr=price(1);
  estimate p=1 q=1 input=price;
  forecast lead=12 interval=month id=date out=results;
run;

```

Fitting a model to the input series is also important for identifying transfer functions. (See the section “[Prewhitening](#)” on page 238 for more information.)

Input values from the DATA= data set and input values forecast by PROC ARIMA can be combined. For example, a model for SALES might have three input series: PRICE, INCOME, and TAXRATE. For the forecast, you assume that the tax rate will be unchanged. You have a forecast for INCOME from another source but only for the first few periods of the SALES forecast you want to make. You have no future values for PRICE, which needs to be forecast as in the preceding example.

In this situation, you include observations in the input data set for all forecast periods, with SALES and PRICE set to a missing value, with TAXRATE set to its last actual value, and with INCOME set to forecast values for the periods you have forecasts for and set to missing values for later periods. In the PROC ARIMA step, you estimate ARIMA models for PRICE and INCOME before you estimate the model for SALES, as shown in the following statements:

```

proc arima data=a;
  identify var=price(1);
  estimate p=2;
  identify var=income(1);
  estimate p=2;
  identify var=sales(1) crosscorr=( price(1) income(1) taxrate );
  estimate p=1 q=1 input=( price income taxrate );
  forecast lead=12 interval=month id=date out=results;
run;

```

In forecasting SALES, the ARIMA procedure uses as inputs the value of PRICE forecast by its ARIMA model, the value of TAXRATE found in the DATA= data set, and the value of INCOME found in the DATA= data set, or, when the INCOME variable is missing, the value of INCOME forecast by its ARIMA model. (Because SALES is missing for future time periods, the estimation of model parameters is not affected by the forecast values for PRICE, INCOME, or TAXRATE.)

## Data Requirements

PROC ARIMA can handle time series of moderate size; there should be at least 30 observations. With fewer than 30 observations, the parameter estimates might be poor. With thousands of observations, the method requires considerable computer time and memory.

---

## Syntax: ARIMA Procedure

The ARIMA procedure uses the following statements:

```
PROC ARIMA options ;
  BY variables ;
  IDENTIFY VAR=variable options ;
  ESTIMATE options ;
  OUTLIER options ;
  FORECAST options ;
```

The **PROC ARIMA** and **IDENTIFY** statements are required.

---

## Functional Summary

The statements and options that control the ARIMA procedure are summarized in [Table 7.3](#).

**Table 7.3** Functional Summary

Description	Statement	Option
<b>Data Set Options</b>		
specify the input data set	PROC ARIMA	DATA=
	IDENTIFY	DATA=
specify the output data set	PROC ARIMA	OUT=
	FORECAST	OUT=
include only forecasts in the output data set	FORECAST	NOOUTALL
write autocovariances to output data set	IDENTIFY	OUTCOV=
write parameter estimates to an output data set	ESTIMATE	OUTTEST=
write correlation of parameter estimates	ESTIMATE	OUTCORR
write covariance of parameter estimates	ESTIMATE	OUTCOV
write estimated model to an output data set	ESTIMATE	OUTMODEL=
write statistics of fit to an output data set	ESTIMATE	OUTSTAT=
<b>Options for Identifying the Series</b>		
difference time series and plot autocorrelations	IDENTIFY	
specify response series and differencing	IDENTIFY	VAR=
specify and cross-correlate input series	IDENTIFY	CROSSCORR=
center data by subtracting the mean	IDENTIFY	CENTER
exclude missing values	IDENTIFY	NOMISS
delete previous models and start	IDENTIFY	CLEAR
specify the significance level for tests	IDENTIFY	ALPHA=
perform tentative ARMA order identification	IDENTIFY	ESACF
by using the ESACF method		
perform tentative ARMA order identification	IDENTIFY	MINIC
by using the MINIC method		

**Table 7.3** *continued*

Description	Statement	Option
perform tentative ARMA order identification by using the SCAN method	IDENTIFY	SCAN
specify the range of autoregressive model orders for estimating the error series for the MINIC method	IDENTIFY	PERROR=
determine the AR dimension of the SCAN, ESACF, and MINIC tables	IDENTIFY	P=
determine the MA dimension of the SCAN, ESACF, and MINIC tables	IDENTIFY	Q=
perform stationarity tests	IDENTIFY	STATIONARITY=
selection of white noise test statistic in the presence of missing values	IDENTIFY	WHITENOISE=
<b>Options for Defining and Estimating the Model</b>		
specify and estimate ARIMA models	ESTIMATE	
specify autoregressive part of model	ESTIMATE	P=
specify moving-average part of model	ESTIMATE	Q=
specify input variables and transfer functions	ESTIMATE	INPUT=
drop mean term from the model	ESTIMATE	NOINT
specify the estimation method	ESTIMATE	METHOD=
use alternative form for transfer functions	ESTIMATE	ALTPARM
suppress degrees-of-freedom correction in variance estimates	ESTIMATE	NODF
selection of white noise test statistic in the presence of missing values	ESTIMATE	WHITENOISE=
<b>Options for Outlier Detection</b>		
specify the significance level for tests	OUTLIER	ALPHA=
identify detected outliers with variable	OUTLIER	ID=
limit the number of outliers	OUTLIER	MAXNUM=
limit the number of outliers to a percentage of the series	OUTLIER	MAXPCT=
specify the variance estimator used for testing	OUTLIER	SIGMA=
specify the type of level shifts	OUTLIER	TYPE=
<b>Printing Control Options</b>		
limit number of lags shown in correlation plots	IDENTIFY	NLAG=
suppress printed output for identification	IDENTIFY	NOPRINT
plot autocorrelation functions of the residuals	ESTIMATE	PLOT
print log-likelihood around the estimates	ESTIMATE	GRID
control spacing for GRID option	ESTIMATE	GRIDVAL=
print details of the iterative estimation process	ESTIMATE	PRINTALL
suppress printed output for estimation	ESTIMATE	NOPRINT

**Table 7.3** *continued*

Description	Statement	Option
suppress printing of the forecast values print the one-step forecasts and residuals	FORECAST FORECAST	NOPRINT PRINTALL
<b>Plotting Control Options</b> request plots associated with model identification, residual analysis, and forecasting	PROC ARIMA	PLOTS=
<b>Options to Specify Parameter Values</b> specify autoregressive starting values specify moving-average starting values specify a starting value for the mean parameter specify starting values for transfer functions	ESTIMATE ESTIMATE ESTIMATE ESTIMATE	AR= MA= MU= INITVAL=
<b>Options to Control the Iterative Estimation Process</b> specify convergence criterion specify the maximum number of iterations specify criterion for checking for singularity suppress the iterative estimation process omit initial observations from objective specify perturbation for numerical derivatives omit stationarity and invertibility checks use preliminary estimates as starting values for ML and ULS	ESTIMATE ESTIMATE ESTIMATE ESTIMATE ESTIMATE ESTIMATE ESTIMATE ESTIMATE	CONVERGE= MAXITER= SINGULAR= NOEST BACKLIM= DELTA= NOSTABLE NOLS
<b>Options for Forecasting</b> forecast the response series specify how many periods to forecast specify the ID variable specify the periodicity of the series specify size of forecast confidence limits start forecasting before end of the input data specify the variance term used to compute forecast standard errors and confidence limits control the alignment of SAS date values	FORECAST FORECAST FORECAST FORECAST FORECAST FORECAST FORECAST FORECAST	LEAD= ID= INTERVAL= ALPHA= BACK= SIGSQ= ALIGN=
<b>BY Groups</b> specify BY group processing	BY	

## PROC ARIMA Statement

**PROC ARIMA** *options* ;

The following options can be used in the PROC ARIMA statement.

**DATA=SAS-data-set**

specifies the name of the SAS data set that contains the time series. If different DATA= specifications appear in the PROC ARIMA and IDENTIFY statements, the one in the IDENTIFY statement is used. If the DATA= option is not specified in either the PROC ARIMA or IDENTIFY statement, the most recently created SAS data set is used.

**PLOTS<(global-plot-options)> <= plot-request <(options)>>**

**PLOTS<(global-plot-options)> <= (plot-request <(options)> <... plot-request <(options)>>)**

controls the plots produced through ODS Graphics. When you specify only one plot request, you can omit the parentheses around the plot request.

Here are some examples:

```
plots=none
plots=all
plots(unpack)=series(corr crosscorr)
plots(only)=(series(corr crosscorr) residual(normal smooth))
```

**Global Plot Options:**

The *global-plot-options* apply to all relevant plots generated by the ARIMA procedure. The following *global-plot-options* are supported:

**ONLY**

suppresses the default plots. Only the plots specifically requested are produced.

**UNPACK**

displays each graph separately. (By default, some graphs can appear together in a single panel.)

**Specific Plot Options**

The following list describes the specific plots and their options.

**ALL**

produces all plots appropriate for the particular analysis.

**NONE**

suppresses all plots.

**SERIES(<series-plot-options> )**

produces plots associated with the identification stage of the modeling. The panel plots corresponding to the CORR and CROSSCORR options are produced by default. The following *series-plot-options* are available:

**ACF**

produces the plot of autocorrelations.

**ALL**

produces all the plots associated with the identification stage.

**CORR**

produces a panel of plots that are useful in the trend and correlation analysis of the series.

The panel consists of the following:

- the time series plot
- the series-autocorrelation plot
- the series-partial-autocorrelation plot
- the series-inverse-autocorrelation plot

**CROSSCORR**

produces panels of cross-correlation plots.

**IACF**

produces the plot of inverse-autocorrelations.

**PACF**

produces the plot of partial-autocorrelations.

**RESIDUAL(<residual-plot-options> )**

produces the residuals plots. The residual correlation and normality diagnostic panels are produced by default. The following *residual-plot-options* are available:

**ACF**

produces the plot of residual autocorrelations.

**ALL**

produces all the residual diagnostics plots appropriate for the particular analysis.

**CORR**

produces a summary panel of the residual correlation diagnostics that consists of the following:

- the residual-autocorrelation plot
- the residual-partial-autocorrelation plot
- the residual-inverse-autocorrelation plot
- a plot of Ljung-Box white-noise test *p*-values at different lags

**HIST**

produces the histogram of the residuals.

**IACF**

produces the plot of residual inverse-autocorrelations.

**NORMAL**

produces a summary panel of the residual normality diagnostics that consists of the following:

- histogram of the residuals
- normal quantile plot of the residuals

**PACF**

produces the plot of residual partial-autocorrelations.

**QQ**

produces the normal quantile plot of the residuals.

**SMOOTH**

produces a scatter plot of the residuals against time, which has an overlaid smooth fit.

**WN**

produces the plot of Ljung-Box white-noise test  $p$ -values at different lags.

**FORECAST(<forecast-plot-options> )**

produces the forecast plots in the forecasting stage. The forecast-only plot that shows the multistep forecasts in the forecast region is produced by default.

The following *forecast-plot-options* are available:

**ALL**

produces the forecast-only plot as well as the forecast plot.

**FORECAST**

produces a plot that shows the one-step-ahead forecasts as well as the multistep-ahead forecasts.

**FORECASTONLY**

produces a plot that shows only the multistep-ahead forecasts in the forecast region.

**OUT=SAS-data-set**

specifies a SAS data set to which the forecasts are output. If different OUT= specifications appear in the PROC ARIMA and FORECAST statements, the one in the FORECAST statement is used.

## BY Statement

**BY** *variables* ;

A BY statement can be used in the ARIMA procedure to process a data set in groups of observations defined by the BY variables. Note that all IDENTIFY, ESTIMATE, and FORECAST statements specified are applied to all BY groups.

Because of the need to make data-based model selections, BY-group processing is not usually done with PROC ARIMA. You usually want to use different models for the different series contained in different BY groups, and the PROC ARIMA BY statement does not let you do this.

Using a BY statement imposes certain restrictions. The BY statement must appear before the first RUN statement. If a BY statement is used, the input data must come from the data set specified in the PROC statement; that is, no input data sets can be specified in IDENTIFY statements.

When a BY statement is used with PROC ARIMA, interactive processing applies only to the first BY group. Once the end of the PROC ARIMA step is reached, all ARIMA statements specified are executed again for each of the remaining BY groups in the input data set.

## **IDENTIFY Statement**

**IDENTIFY VAR=variable options ;**

The IDENTIFY statement specifies the time series to be modeled, differences the series if desired, and computes statistics to help identify models to fit. Use an IDENTIFY statement for each time series that you want to model.

If other time series are to be used as inputs in a subsequent ESTIMATE statement, they must be listed in a CROSSCORR= list in the IDENTIFY statement.

The following options are used in the IDENTIFY statement. The VAR= option is required.

**ALPHA=significance-level**

The ALPHA= option specifies the significance level for tests in the IDENTIFY statement. The default is 0.05.

### **CENTER**

centers each time series by subtracting its sample mean. The analysis is done on the centered data. Later, when forecasts are generated, the mean is added back. Note that centering is done after differencing. The CENTER option is normally used in conjunction with the NOCONSTANT option of the ESTIMATE statement.

### **CLEAR**

deletes all old models. This option is useful when you want to delete old models so that the input variables are not prewhitened. (See the section “[Prewhitening](#)” on page 238 for more information.)

**CROSSCORR=variable (d11, d12, ..., d1k )**

**CROSSCORR= (variable (d11, d12, ..., d1k )... variable (d21, d22, ..., d2k ))**

names the variables cross-correlated with the response variable given by the VAR= specification.

Each variable name can be followed by a list of differencing lags in parentheses, the same as for the VAR= specification. If differencing is specified for a variable in the CROSSCORR= list, the differenced series is cross-correlated with the VAR= option series, and the differenced series is used when the ESTIMATE statement INPUT= option refers to the variable.

**DATA=SAS-data-set**

specifies the input SAS data set that contains the time series. If the DATA= option is omitted, the DATA= data set specified in the PROC ARIMA statement is used; if the DATA= option is omitted from the PROC ARIMA statement as well, the most recently created data set is used.

**ESACF**

computes the extended sample autocorrelation function and uses these estimates to tentatively identify the autoregressive and moving-average orders of mixed models.

The ESACF option generates two tables. The first table displays extended sample autocorrelation estimates, and the second table displays probability values that can be used to test the significance of these estimates. The P=( $p_{min} : p_{max}$ ) and Q=( $q_{min} : q_{max}$ ) options determine the size of the table.

The autoregressive and moving-average orders are tentatively identified by finding a triangular pattern in which all values are insignificant. The ARIMA procedure finds these patterns based on the IDENTIFY statement ALPHA= option and displays possible recommendations for the orders.

The following code generates an ESACF table with dimensions of p=(0:7) and q=(0:8).

```
proc arima data=test;
  identify var=x esacf p=(0:7) q=(0:8);
run;
```

See the section “The ESACF Method” on page 233 for more information.

**MINIC**

uses information criteria or penalty functions to provide tentative ARMA order identification. The MINIC option generates a table that contains the computed information criterion associated with various ARMA model orders. The PERROR=( $p_{\epsilon,min} : p_{\epsilon,max}$ ) option determines the range of the autoregressive model orders used to estimate the error series. The P=( $p_{min} : p_{max}$ ) and Q=( $q_{min} : q_{max}$ ) options determine the size of the table. The ARMA orders are tentatively identified by those orders that minimize the information criterion.

The following statements generate a MINIC table with default dimensions of p=(0:5) and q=(0:5) and with the error series estimated by an autoregressive model with an order,  $p_\epsilon$ , that minimizes the AIC in the range from 8 to 11.

```
proc arima data=test;
  identify var=x minic perror=(8:11);
run;
```

See the section “The MINIC Method” on page 235 for more information.

**NLAG=number**

indicates the number of lags to consider in computing the autocorrelations and cross-correlations. To obtain preliminary estimates of an ARIMA( $p, d, q$ ) model, the NLAG= value must be at least  $p + q + d$ . The number of observations must be greater than or equal to the NLAG= value. The default value for NLAG= is 24 or one-fourth the number of observations, whichever is less. Even though the NLAG= value is specified, the NLAG= value can be changed according to the data set.

**NOMISS**

uses only the first continuous sequence of data with no missing values. By default, all observations are used.

**NOPRINT**

suppresses the normal printout (including the correlation plots) generated by the IDENTIFY statement.

**OUTCOV=SAS-data-set**

writes the autocovariances, autocorrelations, inverse autocorrelations, partial autocorrelations, and cross covariances to an output SAS data set. If the OUTCOV= option is not specified, no covariance output data set is created. See the section “[OUTCOV= Data Set](#)” on page 253 for more information.

**P=( $p_{min} : p_{max}$ )**

see the ESACF, MINIC, and SCAN options for details.

**PERROR=( $p_{\epsilon,min} : p_{\epsilon,max}$ )**

determines the range of the autoregressive model orders used to estimate the error series in MINIC, a tentative ARMA order identification method. See the section “[The MINIC Method](#)” on page 235 for more information. By default  $p_{\epsilon,min}$  is set to  $p_{max}$  and  $p_{\epsilon,max}$  is set to  $p_{max} + q_{max}$ , where  $p_{max}$  and  $q_{max}$  are the maximum settings of the P= and Q= options on the IDENTIFY statement.

**Q=( $q_{min} : q_{max}$ )**

see the ESACF, MINIC, and SCAN options for details.

**SCAN**

computes estimates of the squared canonical correlations and uses these estimates to tentatively identify the autoregressive and moving-average orders of mixed models.

The SCAN option generates two tables. The first table displays squared canonical correlation estimates, and the second table displays probability values that can be used to test the significance of these estimates. The P=( $p_{min} : p_{max}$ ) and Q=( $q_{min} : q_{max}$ ) options determine the size of each table.

The autoregressive and moving-average orders are tentatively identified by finding a rectangular pattern in which all values are insignificant. The ARIMA procedure finds these patterns based on the IDENTIFY statement ALPHA= option and displays possible recommendations for the orders.

The following code generates a SCAN table with default dimensions of p=(0:5) and q=(0:5). The recommended orders are based on a significance level of 0.1.

```
proc arima data=test;
  identify var=x scan alpha=0.1;
  run;
```

See the section “[The SCAN Method](#)” on page 236 for more information.

**STATIONARITY=**

performs stationarity tests. Stationarity tests can be used to determine whether differencing terms should be included in the model specification. In each stationarity test, the autoregressive orders can be specified by a range,  $test=ar_{max}$ , or as a list of values,  $test=(ar_1,..,ar_n)$ , where  $test$  is ADF, PP, or RW. The default is (0,1,2).

See the section “[Stationarity Tests](#)” on page 238 for more information.

**STATIONARITY=(ADF= AR orders DLAG= s )****STATIONARITY=(DICKEY= AR orders DLAG= s )**

performs augmented Dickey-Fuller tests. If the DLAG=s option is specified with s is greater than one, seasonal Dickey-Fuller tests are performed. The maximum allowable value of s is 12. The default value of s is 1. The following code performs augmented Dickey-Fuller tests with autoregressive orders 2 and 5.

```
proc arima data=test;
  identify var=x stationarity=(adf=(2,5));
run;
```

**STATIONARITY=(PP= AR orders )****STATIONARITY=(PHILLIPS= AR orders )**

performs Phillips-Perron tests. The following statements perform augmented Phillips-Perron tests with autoregressive orders ranging from 0 to 6.

```
proc arima data=test;
  identify var=x stationarity=(pp=6);
run;
```

**STATIONARITY=(RW=AR orders )****STATIONARITY=(RANDOMWALK=AR orders )**

performs random-walk-with-drift tests. The following statements perform random-walk-with-drift tests with autoregressive orders ranging from 0 to 2.

```
proc arima data=test;
  identify var=x stationarity=(rw);
run;
```

**VAR=variable****VAR= variable ( d1, d2, ..., dk )**

names the variable that contains the time series to analyze. The VAR= option is required.

A list of differencing lags can be placed in parentheses after the variable name to request that the series be differenced at these lags. For example, VAR=X(1) takes the first differences of X. VAR=X(1,1) requests that X be differenced twice, both times with lag 1, producing a second difference series, which is

$$(X_t - X_{t-1}) - (X_{t-1} - X_{t-2}) = X_t - 2X_{t-1} + X_{t-2}.$$

VAR=X(2) differences X once at lag two ( $X_t - X_{t-2}$ ).

If differencing is specified, it is the differenced series that is processed by any subsequent ESTIMATE statement.

**WHITENOISE=ST | IGNOREMISS**

specifies the type of test statistic that is used in the white noise test of the series when the series contains missing values. If WHITENOISE=IGNOREMISS, the standard Ljung-Box test statistic is used. If WHITENOISE=ST, a modification of this statistic suggested by Stoffer and Toloi (1992) is used. The default is WHITENOISE=ST.

---

## ESTIMATE Statement

*<label:>ESTIMATE options ;*

The ESTIMATE statement specifies an ARMA model or transfer function model for the response variable specified in the previous IDENTIFY statement, and produces estimates of its parameters. The ESTIMATE statement also prints diagnostic information by which to check the model. The label in the ESTIMATE statement is optional. Include an ESTIMATE statement for each model that you want to estimate.

Options used in the ESTIMATE statement are described in the following sections.

### Options for Defining the Model and Controlling Diagnostic Statistics

The following options are used to define the model to be estimated and to control the output that is printed.

#### ALTPARM

specifies the alternative parameterization of the overall scale of transfer functions in the model. See the section “[Alternative Model Parameterization](#)” on page 245 for details.

#### INPUT=variable

#### INPUT=( transfer-function variable ... )

specifies input variables and their transfer functions.

The variables used on the INPUT= option must be included in the CROSSCORR= list in the previous IDENTIFY statement. If any differencing is specified in the CROSSCORR= list, then the differenced series is used as the input to the transfer function.

The transfer function specification for an input variable is optional. If no transfer function is specified, the input variable enters the model as a simple regressor. If specified, the transfer function specification has the following syntax:

$S\$ (L_{1,1}, L_{1,2}, \dots) (L_{2,1}, \dots) \dots / (L_{j,1}, \dots) \dots$

Here,  $S$  is a shift or lag of the input variable, the terms before the slash (/) are numerator factors, and the terms after the slash (/) are denominator factors of the transfer function. All three parts are optional. See the section “[Specifying Inputs and Transfer Functions](#)” on page 244 for details.

#### METHOD=value

specifies the estimation method to use. METHOD=ML specifies the maximum likelihood method. METHOD=ULS specifies the unconditional least squares method. METHOD=CLS specifies the conditional least squares method. METHOD=CLS is the default. See the section “[Estimation Details](#)” on page 240 for more information.

#### NOCONSTANT

#### NOINT

suppresses the fitting of a constant (or intercept) parameter in the model. (That is, the parameter  $\mu$  is omitted.)

#### NODF

estimates the variance by dividing the error sum of squares (SSE) by the number of residuals. The default is to divide the SSE by the number of residuals minus the number of free parameters in the model.

**NOPRINT**

suppresses the normal printout generated by the ESTIMATE statement. If the NOPRINT option is specified for the ESTIMATE statement, then any error and warning messages are printed to the SAS log.

**P=order****P=(lag, ..., lag) ... (lag, ..., lag)**

specifies the autoregressive part of the model. By default, no autoregressive parameters are fit.

$P=(l_1, l_2, \dots, l_k)$  defines a model with autoregressive parameters at the specified lags.  $P=order$  is equivalent to  $P=(1, 2, \dots, order)$ .

A concatenation of parenthesized lists specifies a factored model. For example,  $P=(1,2,5)(6,12)$  specifies the autoregressive model

$$(1 - \phi_{1,1}B - \phi_{1,2}B^2 - \phi_{1,3}B^5)(1 - \phi_{2,1}B^6 - \phi_{2,2}B^{12})$$

**PLOT**

plots the residual autocorrelation functions. The sample autocorrelation, the sample inverse autocorrelation, and the sample partial autocorrelation functions of the model residuals are plotted.

**Q=order****Q=(lag, ..., lag) ... (lag, ..., lag)**

specifies the moving-average part of the model. By default, no moving-average part is included in the model.

$Q=(l_1, l_2, \dots, l_k)$  defines a model with moving-average parameters at the specified lags.  $Q=order$  is equivalent to  $Q=(1, 2, \dots, order)$ . A concatenation of parenthesized lists specifies a factored model. The interpretation of factors and lags is the same as for the P= option.

**WHITENOISE=ST | IGNOREMISS**

specifies the type of test statistic that is used in the white noise test of the series when the series contains missing values. If WHITENOISE=IGNOREMISS, the standard Ljung-Box test statistic is used. If WHITENOISE=ST, a modification of this statistic suggested by Stoffer and Toloi (1992) is used. The default is WHITENOISE=ST.

**Options for Output Data Sets**

The following options are used to store results in SAS data sets:

**OUTTEST=SAS-data-set**

writes the parameter estimates to an output data set. If the OUTCORR or OUTCOV option is used, the correlations or covariances of the estimates are also written to the OUTTEST= data set. See the section “[OUTTEST= Data Set](#)” on page 254 for a description of the OUTTEST= output data set.

**OUTCORR**

writes the correlations of the parameter estimates to the OUTTEST= data set.

**OUTCOV**

writes the covariances of the parameter estimates to the OUTTEST= data set.

**OUTMODEL=SAS-data-set**

writes the model and parameter estimates to an output data set. If OUTMODEL= is not specified, no model output data set is created. See the section “[OUTMODEL= SAS Data Set](#)” on page 257 for a description of the OUTMODEL= output data set.

**OUTSTAT=SAS-data-set**

writes the model diagnostic statistics to an output data set. If OUTSTAT= is not specified, no statistics output data set is created. See the section “[OUTSTAT= Data Set](#)” on page 258 for a description of the OUTSTAT= output data set.

## Options to Specify Parameter Values

The following options enable you to specify values for the model parameters. These options can provide starting values for the estimation process, or you can specify fixed parameters for use in the FORECAST stage and suppress the estimation process with the NOEST option. By default, the ARIMA procedure finds initial parameter estimates and uses these estimates as starting values in the iterative estimation process.

If values for any parameters are specified, values for all parameters should be given. The number of values given must agree with the model specifications.

**AR=value ...**

lists starting values for the autoregressive parameters. See the section “[Initial Values](#)” on page 245 for more information.

**INITVAL=(initializer-spec variable ...)**

specifies starting values for the parameters in the transfer function parts of the model. See the section “[Initial Values](#)” on page 245 for more information.

**MA=value ...**

lists starting values for the moving-average parameters. See the section “[Initial Values](#)” on page 245 for more information.

**MU=value**

specifies the MU parameter.

**NOEST**

uses the values specified with the AR=, MA=, INITVAL=, and MU= options as final parameter values. The estimation process is suppressed except for estimation of the residual variance. The specified parameter values are used directly by the next FORECAST statement. When NOEST is specified, standard errors, *t* values, and the correlations between estimates are displayed as 0 or missing. (The NOEST option is useful, for example, when you want to generate forecasts that correspond to a published model.)

## Options to Control the Iterative Estimation Process

The following options can be used to control the iterative process of minimizing the error sum of squares or maximizing the log-likelihood function. These tuning options are not usually needed but can be useful if convergence problems arise.

**BACKLIM=-n**

omits the specified number of initial residuals from the sum of squares or likelihood function. Omitting values can be useful for suppressing transients in transfer function models that are sensitive to start-up values.

**CONVERGE=value**

specifies the convergence criterion. Convergence is assumed when the largest change in the estimate for any parameter is less than the CONVERGE= option value. If the absolute value of the parameter estimate is greater than 0.01, the relative change is used; otherwise, the absolute change in the estimate is used. The default is CONVERGE=0.001.

**DELTA=value**

specifies the perturbation value for computing numerical derivatives. The default is DELTA=0.001.

**GRID**

prints the error sum of squares (SSE) or concentrated log-likelihood surface in a small grid of the parameter space around the final estimates. For each pair of parameters, the SSE is printed for the nine parameter-value combinations formed by the grid, with a center at the final estimates and with spacing given by the GRIDVAL= specification. The GRID option can help you judge whether the estimates are truly at the optimum, since the estimation process does not always converge. For models with a large number of parameters, the GRID option produces voluminous output.

**GRIDVAL=number**

controls the spacing in the grid printed by the GRID option. The default is GRIDVAL=0.005.

**MAXITER=n****MAXIT=n**

specifies the maximum number of iterations allowed. The default is MAXITER=50.

**NOLS**

begins the maximum likelihood or unconditional least squares iterations from the preliminary estimates rather than from the conditional least squares estimates that are produced after four iterations. See the section “[Estimation Details](#)” on page 240 for more information.

**NOSTABLE**

specifies that the autoregressive and moving-average parameter estimates for the noise part of the model not be restricted to the stationary and invertible regions, respectively. See the section “[Stationarity and Invertibility](#)” on page 246 for more information.

**PRINTALL**

prints preliminary estimation results and the iterations in the final estimation process.

**NOTFSTABLE**

specifies that the parameter estimates for the denominator polynomial of the transfer function part of the model not be restricted to the stability region. See the section “[Stationarity and Invertibility](#)” on page 246 for more information.

**SINGULAR=value**

specifies the criterion for checking singularity. If a pivot of a sweep operation is less than the SINGULAR= value, the matrix is deemed singular. Sweep operations are performed on the Jacobian matrix during final estimation and on the covariance matrix when preliminary estimates are obtained. The default is SINGULAR=1E-7.

---

## OUTLIER Statement

**OUTLIER** *options* ;

The OUTLIER statement can be used to detect shifts in the level of the response series that are not accounted for by the previously estimated model. An ESTIMATE statement must precede the OUTLIER statement. The following options are used in the OUTLIER statement:

**TYPE=ADDITIVE**

**TYPE=SHIFT**

**TYPE=TEMP (  $d_1, \dots, d_k$  )**

**TYPE=(*<ADDITIVE><SHIFT>* <TEMP (  $d_1, \dots, d_k$  ) )>**

specifies the types of level shifts to search for. The default is TYPE=(ADDITIVE SHIFT), which requests searching for additive outliers and permanent level shifts. The option TEMP(  $d_1, \dots, d_k$  ) requests searching for temporary changes in the level of durations  $d_1, \dots, d_k$ . These options can also be abbreviated as AO, LS, and TC.

**ALPHA=***significance-level*

specifies the significance level for tests in the OUTLIER statement. The default is 0.05.

**SIGMA=ROBUST | MSE**

specifies the type of error variance estimate to use in the statistical tests performed during the outlier detection. SIGMA=MSE corresponds to the usual mean squared error (MSE) estimate, and SIGMA=ROBUST corresponds to a robust estimate of the error variance. The default is SIGMA=ROBUST.

**MAXNUM=***number*

limits the number of outliers to search. The default is MAXNUM=5.

**MAXPCT=***number*

limits the number of outliers to search for according to a percentage of the series length. The default is MAXPCT=2. When both the MAXNUM= and MAXPCT= options are specified, the minimum of the two search numbers is used.

**ID=***Date-Time ID variable*

specifies a SAS date, time, or datetime identification variable to label the detected outliers. This variable must be present in the input data set.

The following examples illustrate a few possibilities for the OUTLIER statement.

The most basic usage, shown as follows, sets all the options to their default values.

```
outlier;
```

That is, it is equivalent to

```
outlier type=(ao ls) alpha=0.05 sigma=robust maxnum=5 maxpct=2;
```

The following statement requests a search for permanent level shifts and for temporary level changes of durations 6 and 12. The search is limited to at most three changes and the significance level of the underlying tests is 0.001. MSE is used as the estimate of error variance. It also requests labeling of the detected shifts using an ID variable *date*.

```
outlier type=(ls tc(6 12)) alpha=0.001 sigma=mse maxnum=3 ID=date;
```

## FORECAST Statement

**FORECAST** *options* ;

The FORECAST statement generates forecast values for a time series by using the parameter estimates produced by the previous ESTIMATE statement. See the section “[Forecasting Details](#)” on page 248 for more information about calculating forecasts.

The following options can be used in the FORECAST statement:

**ALIGN=***option*

controls the alignment of SAS dates used to identify output observations. The ALIGN= option allows the following values: BEGINNING|BEG|B, MIDDLE|MID|M, and ENDING|END|E. BEGINNING is the default.

**ALPHA=***n*

sets the size of the forecast confidence limits. The ALPHA= value must be between 0 and 1. When you specify ALPHA= $\alpha$ , the upper and lower confidence limits have a  $1 - \alpha$  confidence level. The default is ALPHA=0.05, which produces 95% confidence intervals. ALPHA values are rounded to the nearest hundredth.

**BACK=***n*

specifies the number of observations before the end of the data where the multistep forecasts are to begin. The BACK= option value must be less than or equal to the number of observations minus the number of parameters.

The default is BACK=0, which means that the forecast starts at the end of the available data. The end of the data is the last observation for which a noise value can be calculated. If there are no input series, the end of the data is the last nonmissing value of the response time series. If there are input series, this observation can precede the last nonmissing value of the response variable, since there may be missing values for some of the input series.

**ID=***variable*

names a variable in the input data set that identifies the time periods associated with the observations. The ID= variable is used in conjunction with the INTERVAL= option to extrapolate ID values from the end of the input data to identify forecast periods in the OUT= data set.

If the INTERVAL= option specifies an interval type, the ID variable must be a SAS date or datetime variable with the spacing between observations indicated by the INTERVAL= value. If the INTERVAL= option is not used, the last input value of the ID= variable is incremented by one for each forecast period to extrapolate the ID values for forecast observations.

**INTERVAL=***interval***INTERVAL=n**

specifies the time interval between observations. See Chapter 4, “[Date Intervals, Formats, and Functions](#),” for information about valid INTERVAL= values.

The value of the INTERVAL= option is used by PROC ARIMA to extrapolate the ID values for forecast observations and to check that the input data are in order with no missing periods. See the section “[Specifying Series Periodicity](#)” on page 250 for more details.

**LEAD=n**

specifies the number of multistep forecast values to compute. For example, if LEAD=10, PROC ARIMA forecasts for ten periods beginning with the end of the input series (or earlier if BACK= is specified). It is possible to obtain fewer than the requested number of forecasts if a transfer function model is specified and insufficient data are available to compute the forecast. The default is LEAD=24.

**NOOUTALL**

includes only the final forecast observations in the OUT= output data set, not the one-step forecasts for the data before the forecast period.

**NOPRINT**

suppresses the normal printout of the forecast and associated values.

**OUT=SAS-data-set**

writes the forecast (and other values) to an output data set. If OUT= is not specified, the OUT= data set specified in the PROC ARIMA statement is used. If OUT= is also not specified in the PROC ARIMA statement, no output data set is created. See the section “[OUT= Data Set](#)” on page 252 for more information.

**PRINTALL**

prints the FORECAST computation throughout the whole data set. The forecast values for the data before the forecast period (specified by the BACK= option) are one-step forecasts.

**SIGSQ=***value*

specifies the variance term used in the formula for computing forecast standard errors and confidence limits. The default value is the variance estimate computed by the preceding ESTIMATE statement. This option is useful when you wish to generate forecast standard errors and confidence limits based on a published model. It would often be used in conjunction with the NOEST option in the preceding ESTIMATE statement.

## Details: ARIMA Procedure

### The Inverse Autocorrelation Function

The sample inverse autocorrelation function (SIACF) plays much the same role in ARIMA modeling as the sample partial autocorrelation function (SPACF), but it generally indicates subset and seasonal autoregressive models better than the SPACF.

Additionally, the SIACF can be useful for detecting over-differencing. If the data come from a nonstationary or nearly nonstationary model, the SIACF has the characteristics of a noninvertible moving-average. Likewise, if the data come from a model with a noninvertible moving average, then the SIACF has nonstationary characteristics and therefore decays slowly. In particular, if the data have been over-differenced, the SIACF looks like a SACF from a nonstationary process.

The inverse autocorrelation function is not often discussed in textbooks, so a brief description is given here. For more complete discussions, see: Cleveland (1972); Chatfield (1980); Priestley (1981).

Let  $W_t$  be generated by the ARMA( $p, q$ ) process

$$\phi(B)W_t = \theta(B)a_t$$

where  $a_t$  is a white noise sequence. If  $\theta(B)$  is invertible (that is, if  $\theta$  considered as a polynomial in  $B$  has no roots less than or equal to 1 in magnitude), then the model

$$\theta(B)Z_t = \phi(B)a_t$$

is also a valid ARMA( $q, p$ ) model. This model is sometimes referred to as the dual model. The autocorrelation function (ACF) of this dual model is called the inverse autocorrelation function (IACF) of the original model.

Notice that if the original model is a pure autoregressive model, then the IACF is an ACF that corresponds to a pure moving-average model. Thus, it cuts off sharply when the lag is greater than  $p$ ; this behavior is similar to the behavior of the partial autocorrelation function (PACF).

The sample inverse autocorrelation function (SIACF) is estimated in the ARIMA procedure by the following steps. A high-order autoregressive model is fit to the data by means of the Yule-Walker equations. The order of the autoregressive model used to calculate the SIACF is the minimum of the NLAG= value and one-half the number of observations after differencing. The SIACF is then calculated as the autocorrelation function that corresponds to this autoregressive operator when treated as a moving-average operator. That is, the autoregressive coefficients are convolved with themselves and treated as autocovariances.

Under certain conditions, the sampling distribution of the SIACF can be approximated by the sampling distribution of the SACF of the dual model (Bhansali 1980). In the plots generated by ARIMA, the confidence limit marks (.) are located at  $\pm 2/\sqrt{n}$ . These limits bound an approximate 95% confidence interval for the hypothesis that the data are from a white noise process.

## The Partial Autocorrelation Function

The approximation for a standard error for the estimated partial autocorrelation function at lag  $k$  is based on a null hypothesis that a pure autoregressive Gaussian process of order  $k-1$  generated the time series. This standard error is  $1/\sqrt{n}$  and is used to produce the approximate 95% confidence intervals depicted by the dots in the plot.

## The Cross-Correlation Function

The autocorrelation and partial and inverse autocorrelation functions described in the preceding sections help when you want to model a series as a function of its past values and past random errors. When you want to

include the effects of past and current values of other series in the model, the correlations of the response series and the other series must be considered.

The CROSSCORR= option in the IDENTIFY statement computes cross-correlations of the VAR= series with other series and makes these series available for use as inputs in models specified by later ESTIMATE statements.

When the CROSSCORR= option is used, PROC ARIMA prints a plot of the cross-correlation function for each variable in the CROSSCORR= list. This plot is similar in format to the other correlation plots, but it shows the correlation between the two series at both lags and leads. For example,

```
identify var=y crosscorr=x ...;
```

plots the cross-correlation function of Y and X,  $\text{Cor}(y_t, x_{t-s})$ , for  $s = -L$  to  $L$ , where  $L$  is the value of the NLAG= option. Study of the cross-correlation functions can indicate the transfer functions through which the input series should enter the model for the response series.

The cross-correlation function is computed after any specified differencing has been done. If differencing is specified for the VAR= variable or for a variable in the CROSSCORR= list, it is the differenced series that is cross-correlated (and the differenced series is processed by any following ESTIMATE statement).

For example,

```
identify var=y(1) crosscorr=x(1);
```

computes the cross-correlations of the changes in Y with the changes in X. When differencing is specified, the subsequent ESTIMATE statement models changes in the variables rather than the variables themselves.

## The ESACF Method

The extended sample autocorrelation function (ESACF) method can tentatively identify the orders of a *stationary or nonstationary* ARMA process based on iterated least squares estimates of the autoregressive parameters. Tsay and Tiao (1984) proposed the technique, and Choi (1992) provides useful descriptions of the algorithm.

Given a stationary or nonstationary time series  $\{z_t : 1 \leq t \leq n\}$  with mean corrected form  $\tilde{z}_t = z_t - \mu_z$  with a true autoregressive order of  $p + d$  and with a true moving-average order of  $q$ , you can use the ESACF method to estimate the unknown orders  $p + d$  and  $q$  by analyzing the autocorrelation functions associated with filtered series of the form

$$w_t^{(m,j)} = \hat{\Phi}_{(m,j)}(B)\tilde{z}_t = \tilde{z}_t - \sum_{i=1}^m \hat{\phi}_i^{(m,j)}\tilde{z}_{t-i}$$

where  $B$  represents the backshift operator, where  $m = p_{min}, \dots, p_{max}$  are the autoregressive *test* orders, where  $j = q_{min} + 1, \dots, q_{max} + 1$  are the moving-average *test* orders, and where  $\hat{\phi}_i^{(m,j)}$  are the autoregressive parameter estimates under the assumption that the series is an ARMA( $m, j$ ) process.

For purely autoregressive models ( $j = 0$ ), ordinary least squares (OLS) is used to consistently estimate  $\hat{\phi}_i^{(m,0)}$ . For ARMA models, consistent estimates are obtained by the iterated least squares recursion formula, which is initiated by the pure autoregressive estimates:

$$\hat{\phi}_i^{(m,j)} = \hat{\phi}_i^{(m+1,j-1)} - \hat{\phi}_{i-1}^{(m,j-1)} \frac{\hat{\phi}_{m+1}^{(m+1,j-1)}}{\hat{\phi}_m^{(m,j-1)}}$$

The  $j$ th lag of the sample autocorrelation function of the filtered series  $w_t^{(m,j)}$  is the *extended sample autocorrelation function*, and it is denoted as  $r_{j(m)} = r_j(w^{(m,j)})$ .

The standard errors of  $r_{j(m)}$  are computed in the usual way by using Bartlett's approximation of the variance of the sample autocorrelation function,  $\text{var}(r_{j(m)}) \approx (1 + \sum_{t=1}^{j-1} r_j^2(w^{(m,j)}))$ .

If the true model is an ARMA  $(p+d, q)$  process, the filtered series  $w_t^{(m,j)}$  follows an MA( $q$ ) model for  $j \geq q$  so that

$$r_{j(p+d)} \approx 0 \quad j > q$$

$$r_{j(p+d)} \neq 0 \quad j = q$$

Additionally, Tsay and Tiao (1984) show that the extended sample autocorrelation satisfies

$$r_{j(m)} \approx 0 \quad j - q > m - p - d \leq 0$$

$$r_{j(m)} \neq c(m - p - d, j - q) \quad 0 \leq j - q \leq m - p - d$$

where  $c(m - p - d, j - q)$  is a nonzero constant or a continuous random variable bounded by  $-1$  and  $1$ .

An ESACF table is then constructed by using the  $r_{j(m)}$  for  $m = p_{\min}, \dots, p_{\max}$  and  $j = q_{\min} + 1, \dots, q_{\max} + 1$  to identify the ARMA orders (see Table 7.4). The orders are tentatively identified by finding a right (maximal) triangular pattern with vertices located at  $(p+d, q)$  and  $(p+d, q_{\max})$  and in which all elements are insignificant (based on asymptotic normality of the autocorrelation function). The vertex  $(p+d, q)$  identifies the order. Table 7.5 depicts the theoretical pattern associated with an ARMA(1,2) series.

**Table 7.4** ESACF Table

<b>AR</b>	<b>MA</b>					
	0	1	2	3	·	·
0	$r_{1(0)}$	$r_{2(0)}$	$r_{3(0)}$	$r_{4(0)}$	·	·
1	$r_{1(1)}$	$r_{2(1)}$	$r_{3(1)}$	$r_{4(1)}$	·	·
2	$r_{1(2)}$	$r_{2(2)}$	$r_{3(2)}$	$r_{4(2)}$	·	·
3	$r_{1(3)}$	$r_{2(3)}$	$r_{3(3)}$	$r_{4(3)}$	·	·
·	·	·	·	·	·	·
·	·	·	·	·	·	·

**Table 7.5** Theoretical ESACF Table for an ARMA(1,2) Series

AR	MA							
	0	1	2	3	4	5	6	7
0	*	X	X	X	X	X	X	X
1	*	X	0	0	0	0	0	0
2	*	X	X	0	0	0	0	0
3	*	X	X	X	0	0	0	0
4	*	X	X	X	X	0	0	0
	X = significant terms 0 = insignificant terms * = no pattern							

## The MINIC Method

The minimum information criterion (MINIC) method can tentatively identify the order of a *stationary and invertible* ARMA process. Note that Hannan and Rissanen (1982) proposed this method; for useful descriptions of the algorithm, see: Box, Jenkins, and Reinsel (1994); Choi (1992).

Given a stationary and invertible time series  $\{z_t : 1 \leq t \leq n\}$  with mean corrected form  $\tilde{z}_t = z_t - \mu_z$  with a true autoregressive order of  $p$  and with a true moving-average order of  $q$ , you can use the MINIC method to compute information criteria (or penalty functions) for various autoregressive and moving average orders. The following paragraphs provide a brief description of the algorithm.

If the series is a stationary and invertible ARMA( $p, q$ ) process of the form

$$\Phi_{(p,q)}(B)\tilde{z}_t = \Theta_{(p,q)}(B)\epsilon_t$$

the error series can be approximated by a high-order AR process

$$\hat{\epsilon}_t = \hat{\Phi}_{(p_\epsilon, q)}(B)\tilde{z}_t \approx \epsilon_t$$

where the parameter estimates  $\hat{\Phi}_{(p_\epsilon, q)}$  are obtained from the Yule-Walker estimates. The choice of the autoregressive order  $p_\epsilon$  is determined by the order that minimizes the Akaike information criterion (AIC) in the range  $p_{\epsilon, \min} \leq p_\epsilon \leq p_{\epsilon, \max}$

$$AIC(p_\epsilon, 0) = \ln(\tilde{\sigma}_{(p_\epsilon, 0)}^2) + 2(p_\epsilon + 0)/n$$

where

$$\tilde{\sigma}_{(p_\epsilon, 0)}^2 = \frac{1}{n} \sum_{t=p_\epsilon+1}^n \hat{\epsilon}_t^2$$

Note that Hannan and Rissanen (1982) use the Bayesian information criterion (BIC) to determine the autoregressive order used to estimate the error series while others recommend the AIC (Box, Jenkins, and Reinsel 1994; Choi 1992).

Once the error series has been estimated for autoregressive test order  $m = p_{min}, \dots, p_{max}$  and for moving-average test order  $j = q_{min}, \dots, q_{max}$ , the OLS estimates  $\hat{\Phi}_{(m,j)}$  and  $\hat{\Theta}_{(m,j)}$  are computed from the regression model

$$\tilde{z}_t = \sum_{i=1}^m \phi_i^{(m,j)} \tilde{z}_{t-i} + \sum_{k=1}^j \theta_k^{(m,j)} \hat{\epsilon}_{t-k} + \text{error}$$

From the preceding parameter estimates, the BIC is then computed

$$BIC(m, j) = \ln(\tilde{\sigma}_{(m,j)}^2) + 2(m + j)\ln(n)/n$$

where

$$\tilde{\sigma}_{(m,j)}^2 = \frac{1}{n} \sum_{t=t_0}^n \left( \tilde{z}_t - \sum_{i=1}^m \phi_i^{(m,j)} \tilde{z}_{t-i} - \sum_{k=1}^j \theta_k^{(m,j)} \hat{\epsilon}_{t-k} \right)^2$$

where  $t_0 = p_\epsilon + \max(m, j)$ .

A MINIC table is then constructed using  $BIC(m, j)$ ; see Table 7.6. If  $p_{max} > p_{\epsilon, min}$ , the preceding regression might fail due to linear dependence on the estimated error series and the mean-corrected series. Values of  $BIC(m, j)$  that cannot be computed are set to missing. For large autoregressive and moving-average test orders with relatively few observations, a nearly perfect fit can result. This condition can be identified by a large negative  $BIC(m, j)$  value.

**Table 7.6** MINIC Table

<b>AR</b>	<b>MA</b>					
	0	1	2	3	.	.
0	$BIC(0, 0)$	$BIC(0, 1)$	$BIC(0, 2)$	$BIC(0, 3)$	.	.
1	$BIC(1, 0)$	$BIC(1, 1)$	$BIC(1, 2)$	$BIC(1, 3)$	.	.
2	$BIC(2, 0)$	$BIC(2, 1)$	$BIC(2, 2)$	$BIC(2, 3)$	.	.
3	$BIC(3, 0)$	$BIC(3, 1)$	$BIC(3, 2)$	$BIC(3, 3)$	.	.
.	.	.	.	.	.	.
.	.	.	.	.	.	.

---

## The SCAN Method

The smallest canonical (SCAN) correlation method can tentatively identify the orders of a *stationary or nonstationary* ARMA process. Tsay and Tiao (1985) proposed the technique, and for useful descriptions of the algorithm, see: Box, Jenkins, and Reinsel (1994); Choi (1992).

Given a stationary or nonstationary time series  $\{z_t : 1 \leq t \leq n\}$  with mean corrected form  $\tilde{z}_t = z_t - \mu_z$  with a true autoregressive order of  $p + d$  and with a true moving-average order of  $q$ , you can use the SCAN method to analyze eigenvalues of the correlation matrix of the ARMA process. The following paragraphs provide a brief description of the algorithm.

For autoregressive test order  $m = p_{min}, \dots, p_{max}$  and for moving-average test order  $j = q_{min}, \dots, q_{max}$ , perform the following steps.

1. Let  $Y_{m,t} = (\tilde{z}_t, \tilde{z}_{t-1}, \dots, \tilde{z}_{t-m})'$ . Compute the following  $(m+1) \times (m+1)$  matrix

$$\begin{aligned}\hat{\beta}(m, j+1) &= \left( \sum_t Y_{m,t-j-1} Y'_{m,t-j-1} \right)^{-1} \left( \sum_t Y_{m,t-j-1} Y'_{m,t} \right) \\ \hat{\beta}^*(m, j+1) &= \left( \sum_t Y_{m,t} Y'_{m,t} \right)^{-1} \left( \sum_t Y_{m,t} Y'_{m,t-j-1} \right) \\ \hat{A}^*(m, j) &= \hat{\beta}^*(m, j+1) \hat{\beta}(m, j+1)\end{aligned}$$

where  $t$  ranges from  $j+m+2$  to  $n$ .

2. Find the *smallest* eigenvalue,  $\hat{\lambda}^*(m, j)$ , of  $\hat{A}^*(m, j)$  and its corresponding *normalized* eigenvector,  $\Phi_{m,j} = (1, -\phi_1^{(m,j)}, -\phi_2^{(m,j)}, \dots, -\phi_m^{(m,j)})$ . The squared canonical correlation estimate is  $\hat{\lambda}^*(m, j)$ .
3. Using the  $\Phi_{m,j}$  as AR( $m$ ) coefficients, obtain the residuals for  $t = j+m+1$  to  $n$ , by following the formula:  $w_t^{(m,j)} = \tilde{z}_t - \phi_1^{(m,j)} \tilde{z}_{t-1} - \phi_2^{(m,j)} \tilde{z}_{t-2} - \dots - \phi_m^{(m,j)} \tilde{z}_{t-m}$ .
4. From the sample autocorrelations of the residuals,  $r_k(w)$ , approximate the standard error of the squared canonical correlation estimate by

$$var(\hat{\lambda}^*(m, j)^{1/2}) \approx d(m, j)/(n - m - j)$$

where  $d(m, j) = (1 + 2 \sum_{i=1}^{j-1} r_k(w^{(m,j)}))$ .

The test statistic to be used as an identification criterion is

$$c(m, j) = -(n - m - j) \ln(1 - \hat{\lambda}^*(m, j)/d(m, j))$$

which is asymptotically  $\chi_1^2$  if  $m = p + d$  and  $j \geq q$  or if  $m \geq p + d$  and  $j = q$ . For  $m > p$  and  $j < q$ , there is more than one theoretical zero canonical correlation between  $Y_{m,t}$  and  $Y_{m,t-j-1}$ . Since the  $\hat{\lambda}^*(m, j)$  are the smallest canonical correlations for each  $(m, j)$ , the percentiles of  $c(m, j)$  are less than those of a  $\chi_1^2$ ; therefore, Tsay and Tiao (1985) state that it is safe to assume a  $\chi_1^2$ . For  $m < p$  and  $j < q$ , no conclusions about the distribution of  $c(m, j)$  are made.

A SCAN table is then constructed using  $c(m, j)$  to determine which of the  $\hat{\lambda}^*(m, j)$  are significantly different from zero (see Table 7.7). The ARMA orders are tentatively identified by finding a (maximal) rectangular pattern in which the  $\hat{\lambda}^*(m, j)$  are insignificant for all test orders  $m \geq p + d$  and  $j \geq q$ . There may be more than one pair of values  $(p + d, q)$  that permit such a rectangular pattern. In this case, parsimony and the number of insignificant items in the rectangular pattern should help determine the model order. Table 7.8 depicts the theoretical pattern associated with an ARMA(2,2) series.

**Table 7.7** SCAN Table

<b>AR</b>	<b>MA</b>					
	0	1	2	3	.	.
0	$c(0, 0)$	$c(0, 1)$	$c(0, 2)$	$c(0, 3)$	.	.
1	$c(1, 0)$	$c(1, 1)$	$c(1, 2)$	$c(1, 3)$	.	.
2	$c(2, 0)$	$c(2, 1)$	$c(2, 2)$	$c(2, 3)$	.	.
3	$c(3, 0)$	$c(3, 1)$	$c(3, 2)$	$c(3, 3)$	.	.
.	.	.	.	.	.	.
.	.	.	.	.	.	.

**Table 7.8** Theoretical SCAN Table for an ARMA(2,2) Series

AR	MA							
	0	1	2	3	4	5	6	7
0	*	X	X	X	X	X	X	X
1	*	X	X	X	X	X	X	X
2	*	X	0	0	0	0	0	0
3	*	X	0	0	0	0	0	0
4	*	X	0	0	0	0	0	0
	X = significant terms 0 = insignificant terms * = no pattern							

## Stationarity Tests

When a time series has a unit root, the series is nonstationary and the ordinary least squares (OLS) estimator is not normally distributed. Dickey and Fuller studied the limiting distribution of the OLS estimator of autoregressive models for time series with a simple unit root (Dickey 1976; Dickey and Fuller 1979). Dickey, Hasza, and Fuller obtained the limiting distribution for time series with seasonal unit roots (Dickey, Hasza, and Fuller 1984). Hamilton (1994) discusses the various types of unit root testing.

For a description of Dickey-Fuller tests, see the section “[PROBDF Function for Dickey-Fuller Tests](#)” on page 157 in [Chapter 5](#). See Chapter 8, “[The AUTOREG Procedure](#),” for a description of Phillips-Perron tests.

The random-walk-with-drift test recommends whether or not an integrated times series has a drift term. Hamilton (1994) discusses this test.

## Prewitening

If, as is usually the case, an input series is autocorrelated, the direct cross-correlation function between the input and response series gives a misleading indication of the relation between the input and response series.

One solution to this problem is called *prewhitening*. You first fit an ARIMA model for the input series sufficient to reduce the residuals to white noise; then, filter the input series with this model to get the white noise residual series. You then filter the response series with the same model and cross-correlate the filtered response with the filtered input series.

The ARIMA procedure performs this prewhitening process automatically when you precede the IDENTIFY statement for the response series with IDENTIFY and ESTIMATE statements to fit a model for the input series. If a model with no inputs was previously fit to a variable specified by the CROSSCORR= option, then that model is used to prewhiten both the input series and the response series before the cross-correlations are computed for the input series.

For example,

```
proc arima data=in;
  identify var=x;
```

```

estimate p=1 q=1;
identify var=y crosscorr=x;
run;

```

Both X and Y are filtered by the ARMA(1,1) model fit to X before the cross-correlations are computed.

Note that prewhitening is done to estimate the cross-correlation function; the unfiltered series are used in any subsequent ESTIMATE or FORECAST statements, and the correlation functions of Y with its own lags are computed from the unfiltered Y series. But initial values in the ESTIMATE statement are obtained with prewhitened data; therefore, the result with prewhitening can be different from the result without prewhitening.

To suppress prewhitening for all input variables, use the CLEAR option in the IDENTIFY statement to make PROC ARIMA disregard all previous models.

### Prewhitening and Differencing

If the VAR= and CROSSCORR= options specify differencing, the series are differenced before the prewhitening filter is applied. When the differencing lists specified in the VAR= option for an input and in the CROSSCORR= option for that input are not the same, PROC ARIMA combines the two lists so that the differencing operators used for prewhitening include all differences in either list (in the least common multiple sense).

## Identifying Transfer Function Models

When identifying a transfer function model with multiple input variables, the cross-correlation functions can be misleading if the input series are correlated with each other. Any dependencies among two or more input series will confound their cross-correlations with the response series.

The prewhitening technique assumes that the input variables do not depend on past values of the response variable. If there is feedback from the response variable to an input variable, as evidenced by significant cross-correlation at negative lags, both the input and the response variables need to be prewhitened before meaningful cross-correlations can be computed.

PROC ARIMA cannot handle feedback models. The STATESPACE and VARMAX procedures are more appropriate for models with feedback.

### Missing Values and Autocorrelations

To compute the sample autocorrelation function when missing values are present, PROC ARIMA uses only crossproducts that do not involve missing values and employs divisors that reflect the number of crossproducts used rather than the total length of the series. Sample partial autocorrelations and inverse autocorrelations are then computed by using the sample autocorrelation function. If necessary, a taper is employed to transform the sample autocorrelations into a positive definite sequence before calculating the partial autocorrelation and inverse correlation functions. The confidence intervals produced for these functions might not be valid when there are missing values. The distributional properties for sample correlation functions are not clear for finite samples. See Dunsmuir (1984) for some asymptotic properties of the sample correlation functions.

---

## Estimation Details

The ARIMA procedure primarily uses the computational methods outlined by Box and Jenkins. Marquardt's method is used for the nonlinear least squares iterations. Numerical approximations of the derivatives of the sum-of-squares function are taken by using a fixed delta (controlled by the `DELTA=` option).

The methods do not always converge successfully for a given set of data, particularly if the starting values for the parameters are not close to the least squares estimates.

### Back-Forecasting

The unconditional sum of squares is computed exactly; thus, back-forecasting is not performed. Early versions of SAS/ETS software used the back-forecasting approximation and allowed a positive value of the `BACKLIM=` option to control the extent of the back-forecasting. In the current version, requesting a positive number of back-forecasting steps with the `BACKLIM=` option has no effect.

### Preliminary Estimation

If an autoregressive or moving-average operator is specified with no missing lags, preliminary estimates of the parameters are computed by using the autocorrelations computed in the IDENTIFY stage. Otherwise, the preliminary estimates are arbitrarily set to values that produce stable polynomials.

When preliminary estimation is not performed by PROC ARIMA, then initial values of the coefficients for any given autoregressive or moving-average factor are set to 0.1 if the degree of the polynomial associated with the factor is 9 or less. Otherwise, the coefficients are determined by expanding the polynomial  $(1 - 0.1B)$  to an appropriate power by using a recursive algorithm.

These preliminary estimates are the starting values in an iterative algorithm to compute estimates of the parameters.

### Estimation Methods

#### Maximum Likelihood

The `METHOD= ML` option produces maximum likelihood estimates. The likelihood function is maximized via nonlinear least squares using Marquardt's method. Maximum likelihood estimates are more expensive to compute than the conditional least squares estimates; however, they may be preferable in some cases (Ansley and Newbold 1980; Davidson 1981).

The maximum likelihood estimates are computed as follows. Let the univariate ARMA model be

$$\phi(B)(W_t - \mu_t) = \theta(B)a_t$$

where  $a_t$  is an independent sequence of normally distributed innovations with mean 0 and variance  $\sigma^2$ . Here  $\mu_t$  is the mean parameter  $\mu$  plus the transfer function inputs. The log-likelihood function can be written as follows:

$$-\frac{1}{2\sigma^2} \mathbf{x}' \boldsymbol{\Omega}^{-1} \mathbf{x} - \frac{1}{2} \ln(|\boldsymbol{\Omega}|) - \frac{n}{2} \ln(\sigma^2)$$

In this equation,  $n$  is the number of observations,  $\sigma^2 \boldsymbol{\Omega}$  is the variance of  $\mathbf{x}$  as a function of the  $\phi$  and  $\theta$  parameters, and  $|\boldsymbol{\Omega}|$  denotes the determinant. The vector  $\mathbf{x}$  is the time series  $W_t$  minus the structural part of

the model  $\mu_t$ , written as a column vector, as follows:

$$\mathbf{x} = \begin{bmatrix} W_1 \\ W_2 \\ \vdots \\ W_n \end{bmatrix} - \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{bmatrix}$$

The maximum likelihood estimate (MLE) of  $\sigma^2$  is

$$s^2 = \frac{1}{n} \mathbf{x}' \boldsymbol{\Omega}^{-1} \mathbf{x}$$

Note that the default estimator of the variance divides by  $n - r$ , where  $r$  is the number of parameters in the model, instead of by  $n$ . Specifying the NODF option causes a divisor of  $n$  to be used.

The log-likelihood concentrated with respect to  $\sigma^2$  can be taken up to additive constants as

$$-\frac{n}{2} \ln(\mathbf{x}' \boldsymbol{\Omega}^{-1} \mathbf{x}) - \frac{1}{2} \ln(|\boldsymbol{\Omega}|)$$

Let  $\mathbf{H}$  be the lower triangular matrix with positive elements on the diagonal such that  $\mathbf{H}\mathbf{H}' = \boldsymbol{\Omega}$ . Let  $\mathbf{e}$  be the vector  $\mathbf{H}^{-1}\mathbf{x}$ . The concentrated log-likelihood with respect to  $\sigma^2$  can now be written as

$$-\frac{n}{2} \ln(\mathbf{e}' \mathbf{e}) - \ln(|\mathbf{H}|)$$

or

$$-\frac{n}{2} \ln(|\mathbf{H}|^{1/n} \mathbf{e}' \mathbf{e} |\mathbf{H}|^{1/n})$$

The MLE is produced by using a Marquardt algorithm to minimize the following sum of squares:

$$|\mathbf{H}|^{1/n} \mathbf{e}' \mathbf{e} |\mathbf{H}|^{1/n}$$

The subsequent analysis of the residuals is done by using  $\mathbf{e}$  as the vector of residuals.

### **Unconditional Least Squares**

The METHOD=ULS option produces unconditional least squares estimates. The ULS method is also referred to as the *exact least squares* (ELS) method. For METHOD=ULS, the estimates minimize

$$\sum_{t=1}^n \tilde{a}_t^2 = \sum_{t=1}^n (x_t - \mathbf{C}_t \mathbf{V}_t^{-1} (x_1, \dots, x_{t-1})')^2$$

where  $\mathbf{C}_t$  is the covariance matrix of  $x_t$  and  $(x_1, \dots, x_{t-1})$ , and  $\mathbf{V}_t$  is the variance matrix of  $(x_1, \dots, x_{t-1})$ . In fact,  $\sum_{t=1}^n \tilde{a}_t^2$  is the same as  $\mathbf{x}' \boldsymbol{\Omega}^{-1} \mathbf{x}$ , and hence  $\mathbf{e}' \mathbf{e}$ . Therefore, the unconditional least squares estimates are obtained by minimizing the sum of squared residuals rather than using the log-likelihood as the criterion function.

### Conditional Least Squares

The METHOD=CLS option produces conditional least squares estimates. The CLS estimates are conditional on the assumption that the past unobserved errors are equal to 0. The series  $x_t$  can be represented in terms of the previous observations, as follows:

$$x_t = a_t + \sum_{i=1}^{\infty} \pi_i x_{t-i}$$

The  $\pi$  weights are computed from the ratio of the  $\phi$  and  $\theta$  polynomials, as follows:

$$\frac{\phi(B)}{\theta(B)} = 1 - \sum_{i=1}^{\infty} \pi_i B^i$$

The CLS method produces estimates minimizing

$$\sum_{t=1}^n \hat{a}_t^2 = \sum_{t=1}^n (x_t - \sum_{i=1}^{\infty} \hat{\pi}_i x_{t-i})^2$$

where the unobserved past values of  $x_t$  are set to 0 and  $\hat{\pi}_i$  are computed from the estimates of  $\phi$  and  $\theta$  at each iteration.

For METHOD=ULS and METHOD=ML, initial estimates are computed using the METHOD=CLS algorithm.

### Start-up for Transfer Functions

When computing the noise series for transfer function and intervention models, the start-up for the transferred variable is done by assuming that past values of the input series are equal to the first value of the series. The estimates are then obtained by applying least squares or maximum likelihood to the noise series. Thus, for transfer function models, the ML option does not generate the full (multivariate ARMA) maximum likelihood estimates, but it uses only the univariate likelihood function applied to the noise series.

Because PROC ARIMA uses all of the available data for the input series to generate the noise series, other start-up options for the transferred series can be implemented by prefixing an observation to the beginning of the real data. For example, if you fit a transfer function model to the variable Y with the single input X, then you can employ a start-up using 0 for the past values by prefixing to the actual data an observation with a missing value for Y and a value of 0 for X.

### Information Criteria

PROC ARIMA computes and prints two information criteria, Akaike's information criterion (AIC) (Akaike 1974; Harvey 1981) and Schwarz's Bayesian criterion (SBC) (Schwarz 1978). The AIC and SBC are used to compare competing models fit to the same series. The model with the smaller information criteria is said to fit the data better. The AIC is computed as

$$-2\ln(L) + 2k$$

where  $L$  is the likelihood function and  $k$  is the number of free parameters. The SBC is computed as

$$-2\ln(L) + \ln(n)k$$

where  $n$  is the number of residuals that can be computed for the time series. Sometimes Schwarz's Bayesian criterion is called the Bayesian information criterion (BIC).

If METHOD=CLS is used to do the estimation, an approximation value of  $L$  is used, where  $L$  is based on the conditional sum of squares instead of the exact sum of squares, and a Jacobian factor is left out.

## Tests of Residuals

A table of test statistics for the hypothesis that the model residuals are white noise is printed as part of the ESTIMATE statement output. The chi-square statistics used in the test for lack of fit are computed using the Ljung-Box formula

$$\chi_m^2 = n(n + 2) \sum_{k=1}^m \frac{r_k^2}{(n - k)}$$

where

$$r_k = \frac{\sum_{t=1}^{n-k} a_t a_{t+k}}{\sum_{t=1}^n a_t^2}$$

and  $a_t$  is the residual series.

This formula has been suggested by Ljung and Box (1978) as yielding a better fit to the asymptotic chi-square distribution than the Box-Pierce Q statistic. Some simulation studies of the finite sample properties of this statistic are given by Davies, Triggs, and Newbold (1977); Ljung and Box (1978). When the time series has missing values, Stoffer and Toloi (1992) suggest a modification of this test statistic that has improved distributional properties over the standard Ljung-Box formula given above. When the series contains missing values, this modified test statistic is used by default.

Each chi-square statistic is computed for all lags up to the indicated lag value and is not independent of the preceding chi-square values. The null hypotheses tested is that the current set of autocorrelations is white noise.

## *t*-values

The  $t$  values reported in the table of parameter estimates are approximations whose accuracy depends on the validity of the model, the nature of the model, and the length of the observed series. When the length of the observed series is short and the number of estimated parameters is large with respect to the series length, the  $t$  approximation is usually poor. Probability values that correspond to a  $t$  distribution should be interpreted carefully because they may be misleading.

## Cautions during Estimation

The ARIMA procedure uses a general nonlinear least squares estimation method that can yield problematic results if your data do not fit the model. Output should be examined carefully. The GRID option can be used to ensure the validity and quality of the results. Problems you might encounter include the following:

- Preliminary moving-average estimates might not converge. If this occurs, preliminary estimates are derived as described previously in “Preliminary Estimation” on page 240. You can supply your own preliminary estimates with the ESTIMATE statement options.

- The estimates can lead to an unstable time series process, which can cause extreme forecast values or overflows in the forecast.
- The Jacobian matrix of partial derivatives might be singular; usually, this happens because not all the parameters are identifiable. Removing some of the parameters or using a longer time series might help.
- The iterative process might not converge. PROC ARIMA's estimation method stops after  $n$  iterations, where  $n$  is the value of the MAXITER= option. If an iteration does not improve the SSE, the Marquardt parameter is increased by a factor of ten until parameters that have a smaller SSE are obtained or until the limit value of the Marquardt parameter is exceeded.
- For METHOD=CLS, the estimates might converge but not to least squares estimates. The estimates might converge to a local minimum, the numerical calculations might be distorted by data whose sum-of-squares surface is not smooth, or the minimum might lie outside the region of invertibility or stationarity.
- If the data are differenced and a moving-average model is fit, the parameter estimates might try to converge exactly on the invertibility boundary. In this case, the standard error estimates that are based on derivatives might be inaccurate.

## Specifying Inputs and Transfer Functions

Input variables and transfer functions for them can be specified using the INPUT= option in the ESTIMATE statement. The variables used in the INPUT= option must be included in the CROSSCORR= list in the previous IDENTIFY statement. If any differencing is specified in the CROSSCORR= list, then the differenced variable is used as the input to the transfer function.

### General Syntax of the INPUT= Option

The general syntax of the INPUT= option is

**ESTIMATE ... INPUT=(*transfer-function variable* ...)**

The transfer function for an input variable is optional. The name of a variable by itself can be used to specify a pure regression term for the variable.

If specified, the syntax of the transfer function is

$S \$ (L_{1,1}, L_{1,2}, \dots)(L_{2,1}, \dots) \dots / (L_{i,1}, L_{i,2}, \dots)(L_{i+1,1}, \dots) \dots$

$S$  is the number of periods of time delay (lag) for this input series. Each term in parentheses specifies a polynomial factor with parameters at the lags specified by the  $L_{i,j}$  values. The terms before the slash (/) are numerator factors. The terms after the slash (/) are denominator factors. All three parts are optional.

Commas can optionally be used between input specifications to make the INPUT= option more readable. The \$ sign after the shift is also optional.

Except for the first numerator factor, each of the terms  $L_{i,1}, L_{i,2}, \dots, L_{i,k}$  indicates a factor of the form

$$(1 - \omega_{i,1} B^{L_{i,1}} - \omega_{i,2} B^{L_{i,2}} - \dots - \omega_{i,k} B^{L_{i,k}})$$

The form of the first numerator factor depends on the ALTPARM option. By default, the constant 1 in the first numerator factor is replaced with a free parameter  $\omega_0$ .

## Alternative Model Parameterization

When the ALTPARM option is specified, the  $\omega_0$  parameter is factored out so that it multiplies the entire transfer function, and the first numerator factor has the same form as the other factors.

The ALTPARM option does not materially affect the results; it just presents the results differently. Some people prefer to see the model written one way, while others prefer the alternative representation. Table 7.9 illustrates the effect of the ALTPARM option.

**Table 7.9** The ALTPARM Option

INPUT= Option	ALTPARM	Model
INPUT=((1 2)(12)/(1)X);	No	$(\omega_0 - \omega_1 B - \omega_2 B^2)(1 - \omega_3 B^{12})/(1 - \delta_1 B)X_t$
	Yes	$\omega_0(1 - \omega_1 B - \omega_2 B^2)(1 - \omega_3 B^{12})/(1 - \delta_1 B)X_t$

## Differencing and Input Variables

If you difference the response series and use input variables, take care that the differencing operations do not change the meaning of the model. For example, if you want to fit the model

$$Y_t = \frac{\omega_0}{(1 - \delta_1 B)} X_t + \frac{(1 - \theta_1 B)}{(1 - B)(1 - B^{12})} a_t$$

then the IDENTIFY statement must read

```
identify var=y(1,12) crosscorr=x(1,12);
estimate q=1 input=(/ (1)x) noconstant;
```

If instead you specify the differencing as

```
identify var=y(1,12) crosscorr=x;
estimate q=1 input=(/ (1)x) noconstant;
```

then the model being requested is

$$Y_t = \frac{\omega_0}{(1 - \delta_1 B)(1 - B)(1 - B^{12})} X_t + \frac{(1 - \theta_1 B)}{(1 - B)(1 - B^{12})} a_t$$

which is a very different model.

The point to remember is that a differencing operation requested for the response variable specified by the VAR= option is applied only to that variable and not to the noise term of the model.

## Initial Values

The syntax for giving initial values to transfer function parameters in the INITVAL= option parallels the syntax of the INPUT= option. For each transfer function in the INPUT= option, the INITVAL= option should

give an initialization specification followed by the input series name. The initialization specification for each transfer function has the form

$$C \$ (V_{1,1}, V_{1,2}, \dots) (V_{2,1}, \dots) \dots / (V_{i,1}, \dots) \dots$$

where  $C$  is the lag 0 term in the first numerator factor of the transfer function (or the overall scale factor if the ALTPARM option is specified) and  $V_{i,j}$  is the coefficient of the  $L_{i,j}$  element in the transfer function.

To illustrate, suppose you want to fit the model

$$Y_t = \mu + \frac{(\omega_0 - \omega_1 B - \omega_2 B^2)}{(1 - \delta_1 B - \delta_2 B^2 - \delta_3 B^3)} X_{t-3} + \frac{1}{(1 - \phi_1 B - \phi_2 B^3)} a_t$$

and start the estimation process with the initial values  $\mu=10$ ,  $\omega_0=1$ ,  $\omega_1=0.5$ ,  $\omega_2=0.03$ ,  $\delta_1=0.8$ ,  $\delta_2=-0.1$ ,  $\delta_3=0.002$ ,  $\phi_1=0.1$ ,  $\phi_2=0.01$ . (These are arbitrary values for illustration only.) You would use the following statements:

```
identify var=y crosscorr=x;
estimate p=(1, 3) input=(3$(1, 2) / (1, 2, 3)x)
        mu=10 ar=.1 .01
        initval=(1$(.5, .03) / (.8, -.1, .002)x);
```

Note that the lags specified for a particular factor are sorted, so initial values should be given in sorted order. For example, if the P= option had been entered as P=(3,1) instead of P=(1,3), the model would be the same and so would the AR= option. Sorting is done within all factors, including transfer function factors, so initial values should always be given in order of increasing lags.

Here is another illustration, showing initialization for a factored model with multiple inputs. The model is

$$Y_t = \mu + \frac{\omega_{1,0}}{(1 - \delta_{1,1} B)} W_t + (\omega_{2,0} - \omega_{2,1} B) X_{t-3} + \frac{1}{(1 - \phi_1 B)(1 - \phi_2 B^6 - \phi_3 B^{12})} a_t$$

and the initial values are  $\mu=10$ ,  $\omega_{1,0}=5$ ,  $\delta_{1,1}=0.8$ ,  $\omega_{2,0}=1$ ,  $\omega_{2,1}=0.5$ ,  $\phi_1=0.1$ ,  $\phi_2=0.05$ , and  $\phi_3=0.01$ . You would use the following statements:

```
identify var=y crosscorr=(w x);
estimate p=(1) (6,12) input=(/(1)w, 3$(1)x)
        mu=10 ar=.1 .05 .01
        initval=(5$/(.8)w 1$(.5)x);
```

## Stationarity and Invertibility

By default, PROC ARIMA requires that the parameter estimates for the AR and MA parts of the model always remain in the stationary and invertible regions, respectively. The NOSTABLE option removes this restriction and for high-order models can save some computer time. Note that using the NOSTABLE option does not necessarily result in an unstable model being fit, since the estimates can leave the stable region for some iterations but still ultimately converge to stable values. Similarly, by default, the parameter estimates for the denominator polynomial of the transfer function part of the model are also restricted to be stable. The NOTFSTABLE option can be used to remove this restriction.

---

## Naming of Model Parameters

In the table of parameter estimates produced by the ESTIMATE statement, model parameters are referred to by using the naming convention described in this section.

The parameters in the noise part of the model are named as  $AR_{i,j}$  or  $MA_{i,j}$ , where AR refers to autoregressive parameters and MA to moving-average parameters. The subscript  $i$  refers to the particular polynomial factor, and the subscript  $j$  refers to the  $j$ th term within the  $i$ th factor. These terms are sorted in order of increasing lag within factors, so the subscript  $j$  refers to the  $j$ th term after sorting.

When inputs are used in the model, the parameters of each transfer function are named  $NUM_{i,j}$  and  $DEN_{i,j}$ . The  $j$ th term in the  $i$ th factor of a numerator polynomial is named  $NUM_{i,j}$ . The  $j$ th term in the  $i$ th factor of a denominator polynomial is named  $DEN_{i,j}$ .

This naming process is repeated for each input variable, so if there are multiple inputs, parameters in transfer functions for different input series have the same name. The table of parameter estimates shows in the “Variable” column the input with which each parameter is associated. The parameter name shown in the “Parameter” column and the input variable name shown in the “Variable” column must be combined to fully identify transfer function parameters.

The lag 0 parameter in the first numerator factor for the first input variable is named NUM1. For subsequent input variables, the lag 0 parameter in the first numerator factor is named NUM $k$ , where  $k$  is the position of the input variable in the INPUT= option list. If the ALTPARM option is specified, the NUM $k$  parameter is replaced by an overall scale parameter named SCALE $k$ .

For the mean and noise process parameters, the response series name is shown in the “Variable” column. The lag and shift for each parameter are also shown in the table of parameter estimates when inputs are used.

---

## Missing Values and Estimation and Forecasting

Estimation and forecasting are carried out in the presence of missing values by forecasting the missing values with the current set of parameter estimates. The maximum likelihood algorithm employed was suggested by Jones (1980) and is used for both unconditional least squares (ULS) and maximum likelihood (ML) estimation.

The CLS algorithm simply fills in missing values with infinite memory forecast values, computed by forecasting ahead from the nonmissing past values as far as required by the structure of the missing values. These artificial values are then employed in the nonmissing value CLS algorithm. Artificial values are updated at each iteration along with parameter estimates.

For models with input variables, embedded missing values (that is, missing values other than at the beginning or end of the series) are not generally supported. Embedded missing values in input variables are supported for the special case of a multiple regression model that has ARIMA errors. A multiple regression model is specified by an INPUT= option that simply lists the input variables (possibly with lag shifts) without any numerator or denominator transfer function factors. One-step-ahead forecasts are not available for the response variable when one or more of the input variables have missing values.

When embedded missing values are present for a model with complex transfer functions, PROC ARIMA uses the first continuous nonmissing piece of each series to do the analysis. That is, PROC ARIMA skips observations at the beginning of each series until it encounters a nonmissing value and then uses the data from

there until it encounters another missing value or until the end of the data is reached. This makes the current version of PROC ARIMA compatible with earlier releases that did not allow embedded missing values.

## Forecasting Details

If the model has input variables, a forecast beyond the end of the data for the input variables is possible only if univariate ARIMA models have previously been fit to the input variables or future values for the input variables are included in the DATA= data set.

If input variables are used, the forecast standard errors and confidence limits of the response depend on the estimated forecast error variance of the predicted inputs. If several input series are used, the forecast errors for the inputs should be independent; otherwise, the standard errors and confidence limits for the response series will not be accurate. If future values for the input variables are included in the DATA= data set, the standard errors of the forecasts will be underestimated since these values are assumed to be known with certainty.

The forecasts are generated using forecasting equations consistent with the method used to estimate the model parameters. Thus, the estimation method specified in the ESTIMATE statement also controls the way forecasts are produced by the FORECAST statement. If METHOD=CLS is used, the forecasts are *infinite memory forecasts*, also called *conditional forecasts*. If METHOD=ULS or METHOD=ML, the forecasts are *finite memory forecasts*, also called *unconditional forecasts*. A complete description of the steps to produce the series forecasts and their standard errors by using either of these methods is quite involved, and only a brief explanation of the algorithm is given in the next two sections. Additional details about the finite and infinite memory forecasts can be found in Brockwell and Davis (1991). The prediction of stationary ARMA processes is explained in Chapter 5, and the prediction of nonstationary ARMA processes is given in Chapter 9 of Brockwell and Davis (1991).

### Infinite Memory Forecasts

If METHOD=CLS is used, the forecasts are *infinite memory forecasts*, also called *conditional forecasts*. The term *conditional* is used because the forecasts are computed by assuming that the unknown values of the response series before the start of the data are equal to the mean of the series. Thus, the forecasts are conditional on this assumption.

The series  $x_t$  can be represented as

$$x_t = a_t + \sum_{i=1}^{\infty} \pi_i x_{t-i}$$

where  $\phi(B)/\theta(B) = 1 - \sum_{i=1}^{\infty} \pi_i B^i$ .

The  $k$ -step forecast of  $x_{t+k}$  is computed as

$$\hat{x}_{t+k} = \sum_{i=1}^{k-1} \hat{\pi}_i \hat{x}_{t+k-i} + \sum_{i=k}^{\infty} \hat{\pi}_i x_{t+k-i}$$

where unobserved past values of  $x_t$  are set to zero and  $\hat{\pi}_i$  is obtained from the estimated parameters  $\hat{\phi}$  and  $\hat{\theta}$ .

## Finite Memory Forecasts

For METHOD=ULS or METHOD=ML, the forecasts are *finite memory forecasts*, also called *unconditional forecasts*. For finite memory forecasts, the covariance function of the ARMA model is used to derive the best linear prediction equation.

That is, the  $k$ -step forecast of  $x_{t+k}$ , given  $(x_1, \dots, x_{t-1})$ , is

$$\tilde{x}_{t+k} = \mathbf{C}_{k,t} \mathbf{V}_t^{-1} (x_1, \dots, x_{t-1})'$$

where  $\mathbf{C}_{k,t}$  is the covariance of  $x_{t+k}$  and  $(x_1, \dots, x_{t-1})$  and  $\mathbf{V}_t$  is the covariance matrix of the vector  $(x_1, \dots, x_{t-1})$ .  $\mathbf{C}_{k,t}$  and  $\mathbf{V}_t$  are derived from the estimated parameters.

Finite memory forecasts minimize the mean squared error of prediction if the parameters of the ARMA model are known exactly. (In most cases, the parameters of the ARMA model are estimated, so the predictors are not true best linear forecasts.)

If the response series is differenced, the final forecast is produced by summing the forecast of the differenced series. This summation and the forecast are conditional on the initial values of the series. Thus, when the response series is differenced, the final forecasts are not true finite memory forecasts because they are derived by assuming that the differenced series begins in a steady-state condition. Thus, they fall somewhere between finite memory and infinite memory forecasts. In practice, there is seldom any practical difference between these forecasts and true finite memory forecasts.

## Forecasting Log Transformed Data

The log transformation is often used to convert time series that are nonstationary with respect to the innovation variance into stationary time series. The usual approach is to take the log of the series in a DATA step and then apply PROC ARIMA to the transformed data. A DATA step is then used to transform the forecasts of the logs back to the original units of measurement. The confidence limits are also transformed by using the exponential function.

As one alternative, you can simply exponentiate the forecast series. This procedure gives a forecast for the median of the series, but the antilog of the forecast log series underpredicts the mean of the original series. If you want to predict the expected value of the series, you need to take into account the standard error of the forecast, as shown in the following example, which uses an AR(2) model to forecast the log of a series Y:

```

data in;
  set in;
  ylog = log( y );
run;

proc arima data=in;
  identify var=ylog;
  estimate p=2;
  forecast lead=10 out=out;
run;

data out;
  set out;
  y    = exp( ylog );

```

```

195 = exp( 195 );
u95 = exp( u95 );
forecast = exp( forecast + std*std/2 );
run;

```

## Specifying Series Periodicity

The INTERVAL= option is used together with the ID= variable to describe the observations that make up the time series. For example, INTERVAL=MONTH specifies a monthly time series in which each observation represents one month. See Chapter 4, “[Date Intervals, Formats, and Functions](#),” for details about the interval values supported.

The variable specified by the ID= option in the PROC ARIMA statement identifies the time periods associated with the observations. Usually, SAS date, time, or datetime values are used for this variable. PROC ARIMA uses the ID= variable in the following ways:

- to validate the data periodicity. When the INTERVAL= option is specified, PROC ARIMA uses the ID variable to check the data and verify that successive observations have valid ID values that correspond to successive time intervals. When the INTERVAL= option is not used, PROC ARIMA verifies that the ID values are nonmissing and in ascending order.
- to check for gaps in the input observations. For example, if INTERVAL=MONTH and an input observation for April 1970 follows an observation for January 1970, there is a gap in the input data with two omitted observations (namely February and March 1970). A warning message is printed when a gap in the input data is found.
- to label the forecast observations in the output data set. PROC ARIMA extrapolates the values of the ID variable for the forecast observations from the ID value at the end of the input data according to the frequency specifications of the INTERVAL= option. If the INTERVAL= option is not specified, PROC ARIMA extrapolates the ID variable by incrementing the ID variable value for the last observation in the input data by 1 for each forecast period. Values of the ID variable over the range of the input data are copied to the output data set.

The ALIGN= option is used to align the ID variable to the beginning, middle, or end of the time ID interval specified by the INTERVAL= option.

## Detecting Outliers

You can use the OUTLIER statement to detect changes in the level of the response series that are not accounted for by the estimated model. The types of changes considered are additive outliers (AO), level shifts (LS), and temporary changes (TC).

Let  $\eta_t$  be a regression variable that describes some type of change in the mean response. In time series literature  $\eta_t$  is called a shock signature. An additive outlier at some time point  $s$  corresponds to a shock signature  $\eta_t$  such that  $\eta_s = 1.0$  and  $\eta_t$  is 0.0 at all other points. Similarly a permanent level shift that

originates at time  $s$  has a shock signature such that  $\eta_t$  is 0.0 for  $t < s$  and 1.0 for  $t \geq s$ . A temporary level shift of duration  $d$  that originates at time  $s$  has  $\eta_t$  equal to 1.0 between  $s$  and  $s + d$  and 0.0 otherwise.

Suppose that you are estimating the ARIMA model

$$D(B)Y_t = \mu_t + \frac{\theta(B)}{\phi(B)}a_t$$

where  $Y_t$  is the response series,  $D(B)$  is the differencing polynomial in the backward shift operator  $B$  (possibly identity),  $\mu_t$  is the transfer function input,  $\phi(B)$  and  $\theta(B)$  are the AR and MA polynomials, respectively, and  $a_t$  is the Gaussian white noise series.

The problem of detection of level shifts in the OUTLIER statement is formulated as a problem of sequential selection of shock signatures that improve the model in the ESTIMATE statement. This is similar to the forward selection process in the stepwise regression procedure. The selection process starts with considering shock signatures of the type specified in the TYPE= option, originating at each nonmissing measurement. This involves testing  $H_0: \beta = 0$  versus  $H_a: \beta \neq 0$  in the model

$$D(B)(Y_t - \beta\eta_t) = \mu_t + \frac{\theta(B)}{\phi(B)}a_t$$

for each of these shock signatures. The most significant shock signature, if it also satisfies the significance criterion in ALPHA= option, is included in the model. If no significant shock signature is found, then the outlier detection process stops; otherwise this augmented model, which incorporates the selected shock signature in its transfer function input, becomes the null model for the subsequent selection process. This iterative process stops if at any stage no more significant shock signatures are found or if the number of iterations exceeds the maximum search number that results due to the MAXNUM= and MAXPCT= settings. In all these iterations, the parameters of the ARIMA model in the ESTIMATE statement are held fixed.

The precise details of the testing procedure for a given shock signature  $\eta_t$  are as follows:

The preceding testing problem is equivalent to testing  $H_0: \beta = 0$  versus  $H_a: \beta \neq 0$  in the following “regression with ARMA errors” model

$$N_t = \beta\zeta_t + \frac{\theta(B)}{\phi(B)}a_t$$

where  $N_t = (D(B)Y_t - \mu_t)$  is the “noise” process and  $\zeta_t = D(B)\eta_t$  is the “effective” shock signature.

In this setting, under  $H_0$ ,  $N = (N_1, N_2, \dots, N_n)^T$  is a mean zero Gaussian vector with variance covariance matrix  $\sigma^2\Omega$ . Here  $\sigma^2$  is the variance of the white noise process  $a_t$  and  $\Omega$  is the variance-covariance matrix associated with the ARMA model. Moreover, under  $H_a$ ,  $N$  has  $\beta\zeta$  as the mean vector where  $\zeta = (\zeta_1, \zeta_2, \dots, \zeta_n)^T$ . Additionally, the generalized least squares estimate of  $\beta$  and its variance is given by

$$\begin{aligned}\hat{\beta} &= \delta/\kappa \\ \text{Var}(\hat{\beta}) &= \sigma^2/\kappa\end{aligned}$$

where  $\delta = \zeta^T\Omega^{-1}N$  and  $\kappa = \zeta^T\Omega^{-1}\zeta$ . The test statistic  $\tau^2 = \delta^2/(\sigma^2\kappa)$  is used to test the significance of  $\beta$ , which has an approximate chi-squared distribution with 1 degree of freedom under  $H_0$ . The type of estimate of  $\sigma^2$  used in the calculation of  $\tau^2$  can be specified by the SIGMA= option. The default setting is SIGMA=ROBUST, which corresponds to a robust estimate suggested in an outlier detection procedure in

X-12-ARIMA, the Census Bureau's time series analysis program; see Findley et al. (1998) for additional information. The robust estimate of  $\sigma^2$  is computed by the formula

$$\hat{\sigma}^2 = (1.49 \times \text{Median}(|\hat{a}_t|))^2$$

where  $\hat{a}_t$  are the standardized residuals of the null ARIMA model. The setting SIGMA=MSE corresponds to the usual mean squared error estimate (MSE) computed the same way as in the ESTIMATE statement with the NODF option.

The quantities  $\delta$  and  $\kappa$  are efficiently computed by a method described in De Jong and Penzer (1998); see also Kohn and Ansley (1985).

## Modeling in the Presence of Outliers

In practice, modeling and forecasting time series data in the presence of outliers is a difficult problem for several reasons. The presence of outliers can adversely affect the model identification and estimation steps. Their presence close to the end of the observation period can have a serious impact on the forecasting performance of the model. In some cases, level shifts are associated with changes in the mechanism that drives the observation process, and separate models might be appropriate to different sections of the data. In view of all these difficulties, diagnostic tools such as outlier detection and residual analysis are essential in any modeling process.

The following modeling strategy, which incorporates level shift detection in the familiar Box-Jenkins modeling methodology, seems to work in many cases:

1. Proceed with model identification and estimation as usual. Suppose this results in a tentative ARIMA model, say M.
2. Check for additive and permanent level shifts unaccounted for by the model M by using the OUTLIER statement. In this step, unless there is evidence to justify it, the number of level shifts searched should be kept small.
3. Augment the original dataset with the regression variables that correspond to the detected outliers.
4. Include the first few of these regression variables in M, and call this model M1. Reestimate all the parameters of M1. It is important not to include too many of these outlier variables in the model in order to avoid the danger of over-fitting.
5. Check the adequacy of M1 by examining the parameter estimates, residual analysis, and outlier detection. Refine it more if necessary.

## OUT= Data Set

The output data set produced by the OUT= option of the PROC ARIMA or FORECAST statements contains the following:

- the BY variables
- the ID variable

- the variable specified by the VAR= option in the IDENTIFY statement, which contains the actual values of the response series
- FORECAST, a numeric variable that contains the one-step-ahead predicted values and the multistep forecasts
- STD, a numeric variable that contains the standard errors of the forecasts
- a numeric variable that contains the lower confidence limits of the forecast. This variable is named L95 by default but has a different name if the ALPHA= option specifies a different size for the confidence limits.
- RESIDUAL, a numeric variable that contains the differences between actual and forecast values
- a numeric variable that contains the upper confidence limits of the forecast. This variable is named U95 by default but has a different name if the ALPHA= option specifies a different size for the confidence limits.

The ID variable, the BY variables, and the response variable are the only ones copied from the input to the output data set. In particular, the input variables are not copied to the OUT= data set.

Unless the NOOUTALL option is specified, the data set contains the whole time series. The FORECAST variable has the one-step forecasts (predicted values) for the input periods, followed by  $n$  forecast values, where  $n$  is the LEAD= value. The actual and RESIDUAL values are missing beyond the end of the series.

If you specify the same OUT= data set in different FORECAST statements, the latter FORECAST statements overwrite the output from the previous FORECAST statements. If you want to combine the forecasts from different FORECAST statements in the same output data set, specify the OUT= option once in the PROC ARIMA statement and omit the OUT= option in the FORECAST statements.

When a global output data set is created by the OUT= option in the PROC ARIMA statement, the variables in the OUT= data set are defined by the first FORECAST statement that is executed. The results of subsequent FORECAST statements are vertically concatenated onto the OUT= data set. Thus, if no ID variable is specified in the first FORECAST statement that is executed, no ID variable appears in the output data set, even if one is specified in a later FORECAST statement. If an ID variable is specified in the first FORECAST statement that is executed but not in a later FORECAST statement, the value of the ID variable is the same as the last value processed for the ID variable for all observations created by the later FORECAST statement. Furthermore, even if the response variable changes in subsequent FORECAST statements, the response variable name in the output data set is that of the first response variable analyzed.

## OUTCOV= Data Set

The output data set produced by the OUTCOV= option of the IDENTIFY statement contains the following variables:

- LAG, a numeric variable that contains the lags that correspond to the values of the covariance variables. The values of LAG range from 0 to N for covariance functions and from -N to N for cross-covariance functions, where N is the value of the NLAG= option.
- VAR, a character variable that contains the name of the variable specified by the VAR= option.

- CROSSVAR, a character variable that contains the name of the variable specified in the CROSSCORR= option, which labels the different cross-covariance functions. The CROSSVAR variable is blank for the autocovariance observations. When there is no CROSSCORR= option, this variable is not created.
- N, a numeric variable that contains the number of observations used to calculate the current value of the covariance or cross-covariance function.
- COV, a numeric variable that contains the autocovariance or cross-covariance function values. COV contains the autocovariances of the VAR= variable when the value of the CROSSVAR variable is blank. Otherwise COV contains the cross covariances between the VAR= variable and the variable named by the CROSSVAR variable.
- CORR, a numeric variable that contains the autocorrelation or cross-correlation function values. CORR contains the autocorrelations of the VAR= variable when the value of the CROSSVAR variable is blank. Otherwise CORR contains the cross-correlations between the VAR= variable and the variable named by the CROSSVAR variable.
- STDERR, a numeric variable that contains the standard errors of the autocorrelations. The standard error estimate is based on the hypothesis that the process that generates the time series is a pure moving-average process of order LAG–1. For the cross-correlations, STDERR contains the value  $1/\sqrt{n}$ , which approximates the standard error under the hypothesis that the two series are uncorrelated.
- INVCORR, a numeric variable that contains the inverse autocorrelation function values of the VAR= variable. For cross-correlation observations (that is, when the value of the CROSSVAR variable is not blank), INVCORR contains missing values.
- PARTCORR, a numeric variable that contains the partial autocorrelation function values of the VAR= variable. For cross-correlation observations (that is, when the value of the CROSSVAR variable is not blank), PARTCORR contains missing values.

## OUTEST= Data Set

PROC ARIMA writes the parameter estimates for a model to an output data set when the OUTTEST= option is specified in the ESTIMATE statement. The OUTTEST= data set contains the following:

- the BY variables
- \_MODLABEL\_, a character variable that contains the model label, if it is provided by using the label option in the ESTIMATE statement (otherwise this variable is not created).
- \_NAME\_, a character variable that contains the name of the parameter for the covariance or correlation observations or is blank for the observations that contain the parameter estimates. (This variable is not created if neither OUTCOV nor OUTCORR is specified.)
- \_TYPE\_, a character variable that identifies the type of observation. A description of the \_TYPE\_ variable values is given below.
- variables for model parameters

The variables for the model parameters are named as follows:

ERRORVAR	This numeric variable contains the variance estimate. The _TYPE_=EST observation for this variable contains the estimated error variance, and the remaining observations are missing.
MU	This numeric variable contains values for the mean parameter for the model. (This variable is not created if NOCONSTANT is specified.)
MAj_k	These numeric variables contain values for the moving-average parameters. The variables for moving-average parameters are named MAj_k, where j is the factor-number and k is the index of the parameter within a factor.
ARj_k	These numeric variables contain values for the autoregressive parameters. The variables for autoregressive parameters are named ARj_k, where j is the factor number and k is the index of the parameter within a factor.
Ij_k	These variables contain values for the transfer function parameters. Variables for transfer function parameters are named Ij_k, where j is the number of the INPUT variable associated with the transfer function component and k is the number of the parameter for the particular INPUT variable. INPUT variables are numbered according to the order in which they appear in the INPUT= list.
_STATUS_	This variable describes the convergence status of the model. A value of 0_CONVERGED indicates that the model converged.

The value of the \_TYPE\_ variable for each observation indicates the kind of value contained in the variables for model parameters for the observation. The OUTTEST= data set contains observations with the following \_TYPE\_ values:

EST	The observation contains parameter estimates.
STD	The observation contains approximate standard errors of the estimates.
CORR	The observation contains correlations of the estimates. OUTCORR must be specified to get these observations.
COV	The observation contains covariances of the estimates. OUTCOV must be specified to get these observations.
FACTOR	The observation contains values that identify for each parameter the factor that contains it. Negative values indicate denominator factors in transfer function models.
LAG	The observation contains values that identify the lag associated with each parameter.
SHIFT	The observation contains values that identify the shift associated with the input series for the parameter.

The values given for \_TYPE\_=FACTOR, \_TYPE\_=LAG, or \_TYPE\_=SHIFT observations enable you to reconstruct the model employed when provided with only the OUTTEST= data set.

## OUTEST= Examples

This section clarifies how model parameters are stored in the OUTTEST= data set with two examples.

Consider the following example:

```

proc arima data=input;
  identify var=y cross=(x1 x2);
  estimate p=(1) (6) q=(1,3) (12) input=(x1 x2) outest=est;
run;

proc print data=est;
run;

```

The model specified by these statements is

$$Y_t = \mu + \omega_{1,0}X_{1,t} + \omega_{2,0}X_{2,t} + \frac{(1 - \theta_{11}B - \theta_{12}B^3)(1 - \theta_{21}B^{12})}{(1 - \phi_{11}B)(1 - \phi_{21}B^6)}a_t$$

The OUTEST= data set contains the values shown in [Table 7.10](#).

**Table 7.10** OUTEST= Data Set for First Example

Obs	_TYPE_	Y	MU	MA1_1	MA1_2	MA2_1	AR1_1	AR2_1	I1_1	I2_1
1	EST	$\sigma^2$	$\mu$	$\theta_{11}$	$\theta_{12}$	$\theta_{21}$	$\phi_{11}$	$\phi_{21}$	$\omega_{1,0}$	$\omega_{2,0}$
2	STD	.	se $\mu$	se $\theta_{11}$	se $\theta_{12}$	se $\theta_{21}$	se $\phi_{11}$	se $\phi_{21}$	se $\omega_{1,0}$	se $\omega_{2,0}$
3	FACTOR	.	0	1	1	2	1	2	1	1
4	LAG	.	0	1	3	12	1	6	0	0
5	SHIFT	.	0	0	0	0	0	0	0	0

Note that the symbols in the rows for \_TYPE\_=EST and \_TYPE\_=STD in [Table 7.10](#) would be numeric values in a real data set.

Next, consider the following example:

```

proc arima data=input;
  identify var=y cross=(x1 x2);
  estimate p=1 q=1 input=(2 $ (1)/(1,2)x1 1 $ /(1)x2) outest=est;
run;

proc print data=est;
run;

```

The model specified by these statements is

$$Y_t = \mu + \frac{\omega_{10} - \omega_{11}B}{1 - \delta_{11}B - \delta_{12}B^2}X_{1,t-2} + \frac{\omega_{20}}{1 - \delta_{21}B}X_{2,t-1} + \frac{(1 - \theta_1B)}{(1 - \phi_1B)}a_t$$

The OUTEST= data set contains the values shown in [Table 7.11](#).

**Table 7.11** OUTEST= Data Set for Second Example

Obs	_TYPE_	Y	MU	MA1_1	AR1_1	I1_1	I1_2	I1_3	I1_4	I2_1	I2_2
1	EST	$\sigma^2$	$\mu$	$\theta_1$	$\phi_1$	$\omega_{10}$	$\omega_{11}$	$\delta_{11}$	$\delta_{12}$	$\omega_{20}$	$\delta_{21}$
2	STD	.	se $\mu$	se $\theta_1$	se $\phi_1$	se $\omega_{10}$	se $\omega_{11}$	se $\delta_{11}$	se $\delta_{12}$	se $\omega_{20}$	se $\delta_{21}$
3	FACTOR	.	0	1	1	1	-1	-1	1	-1	
4	LAG	.	0	1	1	0	1	1	2	0	1
5	SHIFT	.	0	0	0	2	2	2	2	1	1

---

## OUTMODEL= SAS Data Set

The OUTMODEL= option in the ESTIMATE statement writes an output data set that enables you to reconstruct the model. The OUTMODEL= data set contains much the same information as the OUTEST= data set but in a transposed form that might be more useful for some purposes. In addition, the OUTMODEL= data set includes the differencing operators.

The OUTMODEL data set contains the following:

- the BY variables
- \_MODLABEL\_, a character variable that contains the model label, if it is provided by using the label option in the ESTIMATE statement (otherwise this variable is not created).
- \_NAME\_, a character variable that contains the name of the response or input variable for the observation.
- \_TYPE\_, a character variable that contains the estimation method that was employed. The value of \_TYPE\_ can be CLS, ULS, or ML.
- \_STATUS\_, a character variable that describes the convergence status of the model. A value of 0\_CONVERGED indicates that the model converged.
- \_PARM\_, a character variable that contains the name of the parameter given by the observation. \_PARM\_ takes on the values ERRORVAR, MU, AR, MA, NUM, DEN, and DIF.
- \_VALUE\_, a numeric variable that contains the value of the estimate defined by the \_PARM\_ variable.
- \_STD\_, a numeric variable that contains the standard error of the estimate.
- \_FACTOR\_, a numeric variable that indicates the number of the factor to which the parameter belongs.
- \_LAG\_, a numeric variable that contains the number of the term within the factor that contains the parameter.
- \_SHIFT\_, a numeric variable that contains the shift value for the input variable associated with the current parameter.

The values of \_FACTOR\_ and \_LAG\_ identify which particular MA, AR, NUM, or DEN parameter estimate is given by the \_VALUE\_ variable. The \_NAME\_ variable contains the response variable name for the MU, AR, or MA parameters. Otherwise, \_NAME\_ contains the input variable name associated with NUM or DEN parameter estimates. The \_NAME\_ variable contains the appropriate variable name associated with the current DIF observation as well. The \_VALUE\_ variable is 1 for all DIF observations, and the \_LAG\_ variable indicates the degree of differencing employed.

The observations contained in the OUTMODEL= data set are identified by the \_PARM\_ variable. A description of the values of the \_PARM\_ variable follows:

NUMRESID      \_VALUE\_ contains the number of residuals.

NPARMS      \_VALUE\_ contains the number of parameters in the model.

NDIFS	_VALUE_ contains the sum of the differencing lags employed for the response variable.
ERRORVAR	_VALUE_ contains the estimate of the innovation variance.
MU	_VALUE_ contains the estimate of the mean term.
AR	_VALUE_ contains the estimate of the autoregressive parameter indexed by the _FACTORTOR_ and _LAG_ variable values.
MA	_VALUE_ contains the estimate of a moving-average parameter indexed by the _FACTORTOR_ and _LAG_ variable values.
NUM	_VALUE_ contains the estimate of the parameter in the numerator factor of the transfer function of the input variable indexed by the _FACTOR_, _LAG_, and _SHIFT_ variable values.
DEN	_VALUE_ contains the estimate of the parameter in the denominator factor of the transfer function of the input variable indexed by the _FACTOR_, _LAG_, and _SHIFT_ variable values.
DIF	_VALUE_ contains the difference operator defined by the difference lag given by the value in the _LAG_ variable.

## OUTSTAT= Data Set

PROC ARIMA writes the diagnostic statistics for a model to an output data set when the OUTSTAT= option is specified in the ESTIMATE statement. The OUTSTAT data set contains the following:

- the BY variables.
- \_MODLABEL\_, a character variable that contains the model label, if it is provided by using the label option in the ESTIMATE statement (otherwise this variable is not created).
- \_TYPE\_, a character variable that contains the estimation method used. \_TYPE\_ can have the value CLS, ULS, or ML.
- \_STAT\_, a character variable that contains the name of the statistic given by the \_VALUE\_ variable in this observation. \_STAT\_ takes on the values AIC, SBC, LOGLIK, SSE, NUMRESID, NPARMS, NDIFS, ERRORVAR, MU, CONV, and NITER.
- \_VALUE\_, a numeric variable that contains the value of the statistic named by the \_STAT\_ variable.

The observations contained in the OUTSTAT= data set are identified by the \_STAT\_ variable. A description of the values of the \_STAT\_ variable follows:

AIC	Akaike's information criterion
SBC	Schwarz's Bayesian criterion
LOGLIK	the log-likelihood, if METHOD=ML or METHOD=ULS is specified
SSE	the sum of the squared residuals
NUMRESID	the number of residuals

NPARMS	the number of parameters in the model
NDIFS	the sum of the differencing lags employed for the response variable
ERRORVAR	the estimate of the innovation variance
MU	the estimate of the mean term
CONV	tells if the estimation converged. The value of 0 signifies that estimation converged. Nonzero values reflect convergence problems.
NITER	the number of iterations

**Remark.** CONV takes an integer value that corresponds to the error condition of the parameter estimation process. The value of 0 signifies that estimation process has converged. The higher values signify convergence problems of increasing severity. Specifically:

- CONV= 0 indicates that the estimation process has converged.
- CONV= 1 or 2 indicates that the estimation process has run into numerical problems (such as encountering an unstable model or a ridge) during the iterations.
- CONV= 3 or greater indicates that the estimation process has failed to converge.

## Printed Output

The ARIMA procedure produces printed output for each of the IDENTIFY, ESTIMATE, and FORECAST statements. The output produced by each ARIMA statement is described in the following sections.

### IDENTIFY Statement Printed Output

The printed output of the IDENTIFY statement consists of the following:

- a table of summary statistics, including the name of the response variable, any specified periods of differencing, the mean and standard deviation of the response series after differencing, and the number of observations after differencing
- a plot of the sample autocorrelation function for lags up to and including the NLAG= option value. Standard errors of the autocorrelations also appear to the right of the autocorrelation plot if the value of LINESIZE= option is sufficiently large. The standard errors are derived using Bartlett's approximation (Box and Jenkins 1976, p. 177). The approximation for a standard error for the estimated autocorrelation function at lag  $k$  is based on a null hypothesis that a pure moving-average Gaussian process of order  $k-1$  generated the time series. The relative position of an approximate 95% confidence interval under this null hypothesis is indicated by the dots in the plot, while the asterisks represent the relative magnitude of the autocorrelation value.
- a plot of the sample inverse autocorrelation function. See the section “[The Inverse Autocorrelation Function](#)” on page 231 for more information about the inverse autocorrelation function.
- a plot of the sample partial autocorrelation function

- a table of test statistics for the hypothesis that the series is white noise. These test statistics are the same as the tests for white noise residuals produced by the ESTIMATE statement and are described in the section “[Estimation Details](#)” on page 240.
- a plot of the sample cross-correlation function for each series specified in the CROSSCORR= option. If a model was previously estimated for a variable in the CROSSCORR= list, the cross-correlations for that series are computed for the prewhitened input and response series. For each input variable with a prewhitening filter, the cross-correlation report for the input series includes the following:
  - a table of test statistics for the hypothesis of no cross-correlation between the input and response series
  - the prewhitening filter used for the prewhitening transformation of the predictor and response variables
- ESACF tables if the ESACF option is used
- MINIC table if the MINIC option is used
- SCAN table if the SCAN option is used
- STATIONARITY test results if the STATIONARITY option is used

## **ESTIMATE Statement Printed Output**

The printed output of the ESTIMATE statement consists of the following:

- if the PRINTALL option is specified, the preliminary parameter estimates and an iteration history that shows the sequence of parameter estimates tried during the fitting process
- a table of parameter estimates that show the following for each parameter: the parameter name, the parameter estimate, the approximate standard error,  $t$  value, approximate probability ( $Pr > |t|$ ), the lag for the parameter, the input variable name for the parameter, and the lag or “Shift” for the input variable
- the estimates of the constant term, the innovation variance (variance estimate), the innovation standard deviation (Std Error Estimate), Akaike’s information criterion (AIC), Schwarz’s Bayesian criterion (SBC), and the number of residuals
- the correlation matrix of the parameter estimates
- a table of test statistics for hypothesis that the residuals of the model are white noise. The table is titled “Autocorrelation Check of Residuals.”
- if the PLOT option is specified, autocorrelation, inverse autocorrelation, and partial autocorrelation function plots of the residuals
- if an INPUT variable has been modeled in such a way that prewhitening is performed in the IDENTIFY step, a table of test statistics titled “Cross-correlation Check of Residuals.” The test statistic is based on the chi-square approximation suggested by Box and Jenkins (1976, pp. 395–396). The cross-correlation function is computed by using the residuals from the model as one series and the prewhitened input variable as the other series.

- if the GRID option is specified, the sum-of-squares or likelihood surface over a grid of parameter values near the final estimates
- a summary of the estimated model that shows the autoregressive factors, moving-average factors, and transfer function factors in backshift notation with the estimated parameter values.

## **OUTLIER Statement Printed Output**

The printed output of the OUTLIER statement consists of the following:

- a summary that contains the information about the maximum number of outliers searched, the number of outliers actually detected, and the significance level used in the outlier detection.
- a table that contains the results of the outlier detection process. The outliers are listed in the order in which they are found. This table contains the following columns:
  - The Obs column contains the observation number of the start of the level shift.
  - If an ID= option is specified, then the Time ID column contains the time identification labels of the start of the outlier.
  - The Type column lists the type of the outlier.
  - The Estimate column contains  $\hat{\beta}$ , the estimate of the regression coefficient of the shock signature.
  - The Chi-Square column lists the value of the test statistic  $\tau^2$ .
  - The Approx Prob > ChiSq column lists the approximate *p*-value of the test statistic.

## **FORECAST Statement Printed Output**

The printed output of the FORECAST statement consists of the following:

- a summary of the estimated model
- a table of forecasts with following columns:
  - The Obs column contains the observation number.
  - The Forecast column contains the forecast values.
  - The Std Error column contains the forecast standard errors.
  - The Lower and Uppers columns contain the approximate 95% confidence limits. The ALPHA= option can be used to change the confidence interval for forecasts.
  - If the PRINTALL option is specified, the forecast table also includes columns for the actual values of the response series (Actual) and the residual values (Residual).

---

## ODS Table Names

PROC ARIMA assigns a name to each table it creates. You can use these names to reference the table when you use the Output Delivery System (ODS) to select tables and create output data sets. These names are listed in [Table 7.12](#).

**Table 7.12** ODS Tables Produced by PROC ARIMA

ODS Table Name	Description	Statement	Option
ChiSqAuto	chi-square statistics table for autocorrelation	IDENTIFY	
ChiSqCross	chi-square statistics table for cross-correlations	IDENTIFY	CROSSCORR
AutoCorrGraph	Correlations graph	IDENTIFY	
CrossCorrGraph	Cross-Correlations graph	IDENTIFY	
DescStats	Descriptive statistics	IDENTIFY	
ESACF	Extended sample autocorrelation function	IDENTIFY	ESACF
ESACFPValues	ESACF probability values	IDENTIFY	ESACF
IACFGraph	Inverse autocorrelations graph	IDENTIFY	
InputDescStats	Input descriptive statistics	IDENTIFY	
MINIC	Minimum information criterion	IDENTIFY	MINIC
PACFGraph	Partial autocorrelations graph	IDENTIFY	
SCAN	Squared canonical correlation estimates	IDENTIFY	SCAN
SCANPValues	SCAN chi-square probability values	IDENTIFY	SCAN
StationarityTests	Stationarity tests	IDENTIFY	STATIONARITY
TentativeOrders	Tentative order selection	IDENTIFY	MINIC, ESACF, or SCAN
ARPolynomial	Filter equations	ESTIMATE	
ChiSqAuto	chi-square statistics table for autocorrelation	ESTIMATE	
ChiSqCross	chi-square statistics table for cross-correlations	ESTIMATE	
CorrB	Correlations of the estimates	ESTIMATE	
DenPolynomial	Filter equations	ESTIMATE	
FitStatistics	Fit statistics	ESTIMATE	
IterHistory	Iteration history	ESTIMATE	PRINTALL
InitialAREstimates	Initial autoregressive parameter estimates	ESTIMATE	
InitialMAEstimates	Initial moving-average parameter estimates	ESTIMATE	
InputDescription	Input description	ESTIMATE	
MAPolynomial	Filter equations	ESTIMATE	

**Table 7.12** *continued*

<b>ODS Table Name</b>	<b>Description</b>	<b>Statement</b>	<b>Option</b>
ModelDescription	Model description	ESTIMATE	
NumPolynomial	Filter equations	ESTIMATE	
ParameterEstimates	Parameter estimates	ESTIMATE	
PrelimEstimates	Preliminary estimates	ESTIMATE	
ObjectiveGrid	Objective function grid matrix	ESTIMATE	GRID
OptSummary	ARIMA estimation optimization	ESTIMATE	PRINTALL
OutlierDetails	Detected outliers	OUTLIER	
Forecasts	Forecast	FORECAST	

## Statistical Graphics

Statistical procedures use ODS Graphics to create graphs as part of their output. ODS Graphics is described in detail in Chapter 21, “Statistical Graphics Using ODS” (*SAS/STAT User’s Guide*).

Before you create graphs, ODS Graphics must be enabled (for example, with the ODS GRAPHICS ON statement). For more information about enabling and disabling ODS Graphics, see the section “Enabling and Disabling ODS Graphics” in that chapter.

The overall appearance of graphs is controlled by ODS styles. Styles and other aspects of using ODS Graphics are discussed in the section “A Primer on ODS Statistical Graphics” in that chapter.

This section provides information about the graphics produced by the ARIMA procedure. (See Chapter 21, “Statistical Graphics Using ODS” (*SAS/STAT User’s Guide*), for more information about ODS statistical graphics.) The main types of plots available are as follows:

- plots useful in the trend and correlation analysis of the dependent and input series
- plots useful for the residual analysis of an estimated model
- forecast plots

You can obtain most plots relevant to the specified model by default. For finer control of the graphics, you can use the **PLOTS=** option in the PROC ARIMA statement. The following example is a simple illustration of how to use the **PLOTS=** option.

### Airline Series: Illustration of ODS Graphics

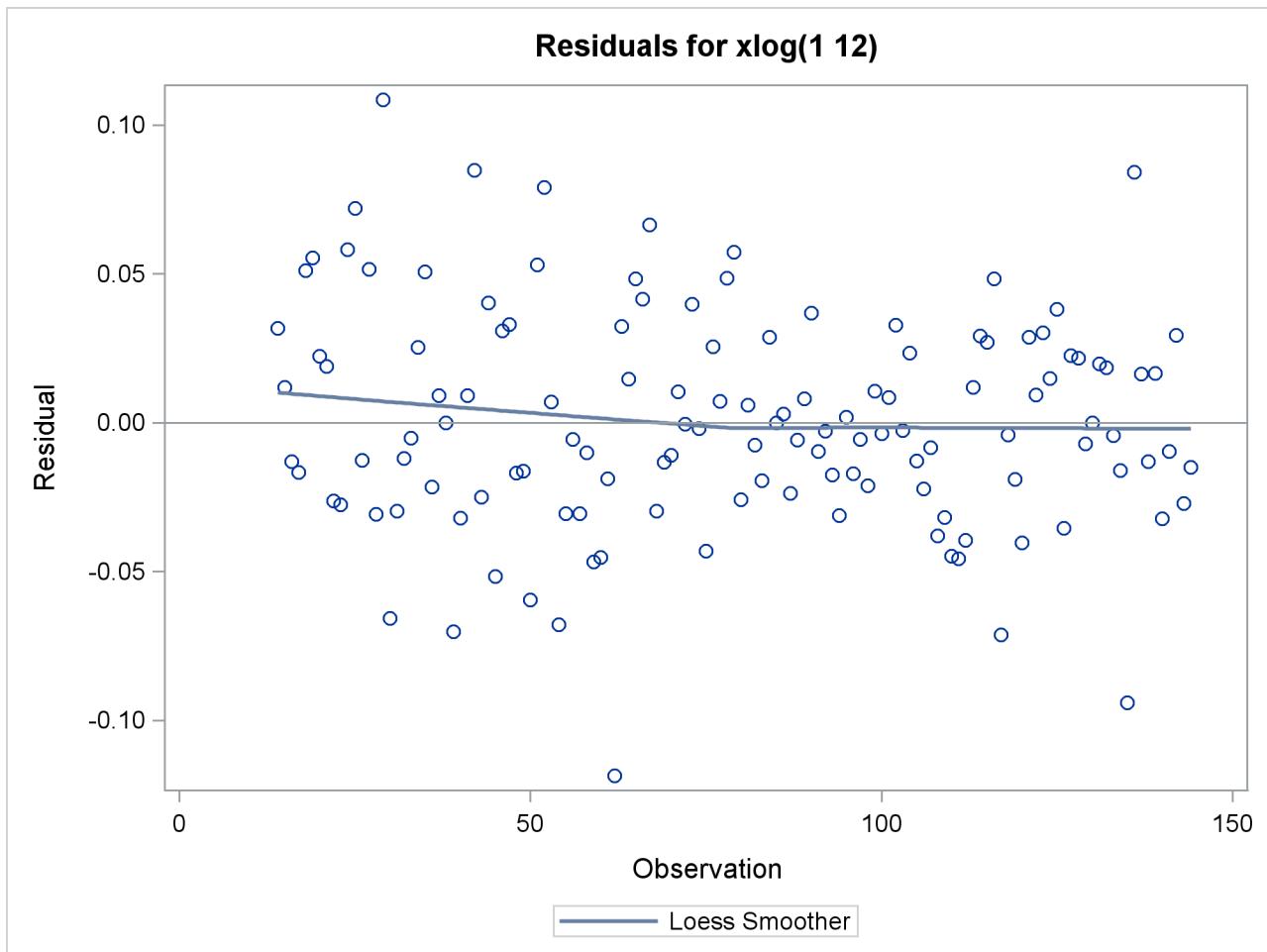
The series in this example, the monthly airline passenger series, is also discussed later, in Example 7.2.

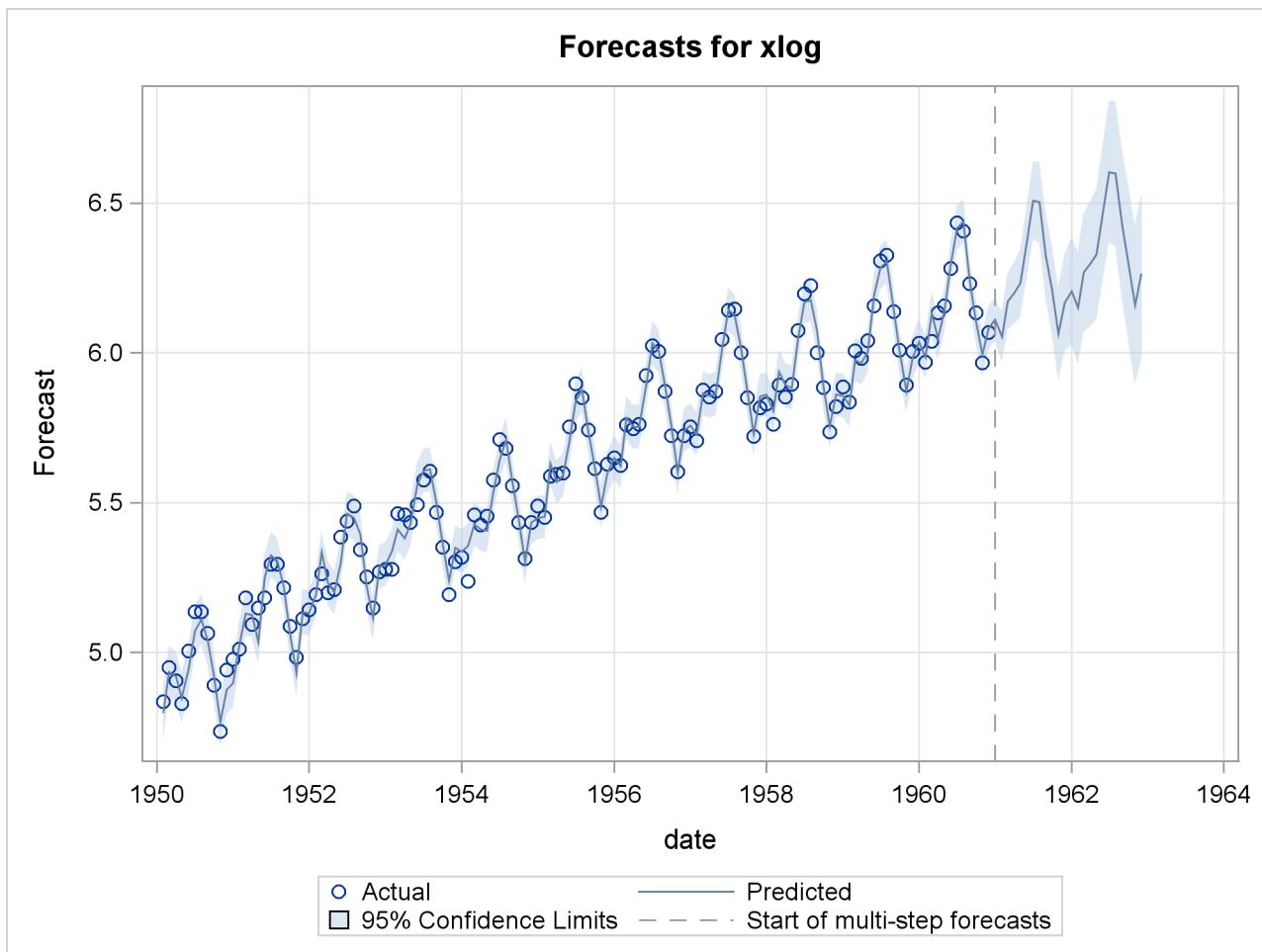
The following statements specify an ARIMA(0,1,1)×(0,1,1)<sub>12</sub> model without a mean term to the logarithms of the airline passengers series, xlog. Notice the use of the global plot option **ONLY** in the **PLOTS=** option

of the PROC ARIMA statement. It suppresses the production of default graphics and produces only the plots specified by the subsequent RESIDUAL and FORECAST plot options. The **RESIDUAL(SMOOTH)** plot specification produces a time series plot of residuals that has an overlaid loess fit; see Figure 7.21. The **FORECAST(FORECAST)** option produces a plot that shows the one-step-ahead forecasts, as well as the multistep-ahead forecasts; see Figure 7.22.

```
proc arima data=seriesg;
  plots(only)=(residual(smooth) forecast(forecasts));
  identify var=xlog(1,12);
  estimate q=(1) (12) noint method=ml;
  forecast id=date interval=month;
run;
```

**Figure 7.21** Residual Plot of the Airline Model



**Figure 7.22** Forecast Plot of the Airline Model

### ODS Graph Names

PROC ARIMA assigns a name to each graph it creates by using ODS. You can use these names to reference the graphs when you use ODS. The names are listed in [Table 7.13](#).

**Table 7.13** ODS Graphics Produced by PROC ARIMA

ODS Graph Name	Plot Description	Option
SeriesPlot	Time series plot of the dependent series	PLOTS(UNPACK)
SeriesACFPlot	Autocorrelation plot of the dependent series	PLOTS(UNPACK)
SeriesPACFPlot	Partial-autocorrelation plot of the dependent series	PLOTS(UNPACK)
SeriesIACFPlot	Inverse-autocorrelation plot of the dependent series	PLOTS(UNPACK)
SeriesCorrPanel	Series trend and correlation analysis panel	Default

**Table 7.13** *continued*

<b>ODS Graph Name</b>	<b>Plot Description</b>	<b>Option</b>
CrossCorrPanel	Cross-correlation plots, either individual or paneled. They are numbered 1, 2, and so on as needed.	Default
ResidualACFPlot	Residual-autocorrelation plot	PLOTS(UNPACK)
ResidualPACFPlot	Residual-partial-autocorrelation plot	PLOTS(UNPACK)
ResidualIACFPlot	Residual-inverse-autocorrelation plot	PLOTS(UNPACK)
ResidualWNPlot	Residual-white-noise-probability plot	PLOTS(UNPACK)
ResidualHistogram	Residual histogram	PLOTS(UNPACK)
ResidualQQPlot	Residual normal Q-Q Plot	PLOTS(UNPACK)
ResidualPlot	Time series plot of residuals with a superimposed smoother	PLOTS=RESIDUAL(SMOOTH)
ForecastsOnlyPlot	Time series plot of multistep forecasts	Default
ForecastsPlot	Time series plot of one-step-ahead as well as multistep forecasts	PLOTS=FORECAST(FORECAST)

---

## Examples: ARIMA Procedure

---

### Example 7.1: Simulated IMA Model

This example illustrates the ARIMA procedure results for a case where the true model is known. An integrated moving-average model is used for this illustration.

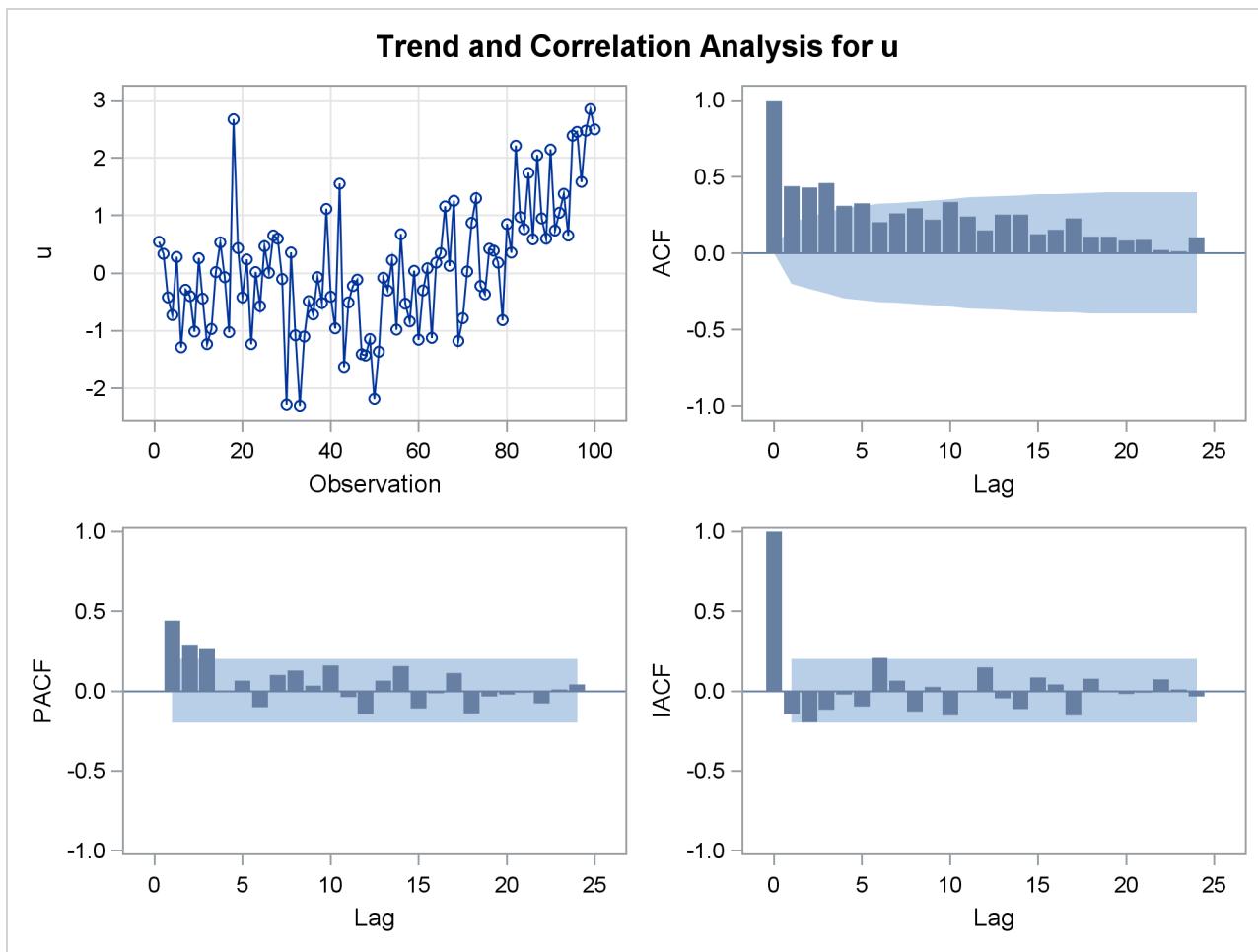
The following DATA step generates a pseudo-random sample of 100 periods from the ARIMA(0,1,1) process  $u_t = u_{t-1} + a_t - 0.8a_{t-1}$ ,  $a_t$  iid  $N(0, 1)$ :

```
title1 'Simulated IMA(1,1) Series';
data a;
  u1 = 0.9; a1 = 0;
  do i = -50 to 100;
    a = rannor( 32565 );
    u = u1 + a - .8 * a1;
    if i > 0 then output;
    a1 = a;
    u1 = u;
  end;
run;
```

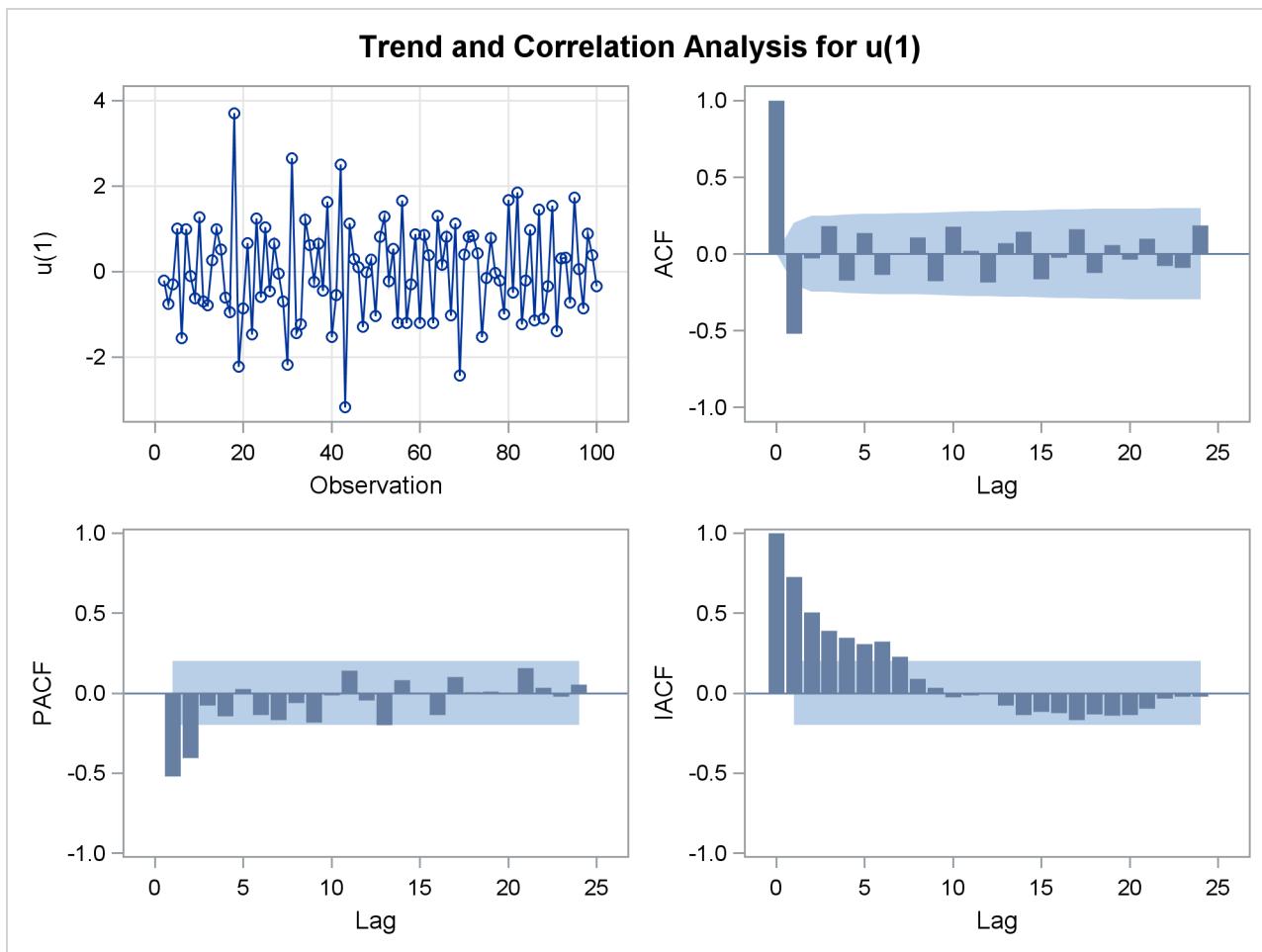
The following ARIMA procedure statements identify and estimate the model:

```
/*--- Simulated IMA Model ---*/
proc arima data=a;
  identify var=u;
  run;
  identify var=u(1);
  run;
  estimate q=1 ;
  run;
quit;
```

The graphical series correlation analysis output of the first IDENTIFY statement is shown in [Output 7.1.1](#). The output shows the behavior of the sample autocorrelation function when the process is nonstationary. Note that in this case the estimated autocorrelations are not very high, even at small lags. Nonstationarity is reflected in a pattern of significant autocorrelations that do not decline quickly with increasing lag, not in the size of the autocorrelations.

**Output 7.1.1** Correlation Analysis from the First IDENTIFY Statement


The second IDENTIFY statement differences the series. The results of the second IDENTIFY statement are shown in [Output 7.1.2](#). This output shows autocorrelation, inverse autocorrelation, and partial autocorrelation functions typical of MA(1) processes.

**Output 7.1.2** Correlation Analysis from the Second IDENTIFY Statement

The ESTIMATE statement fits an ARIMA(0,1,1) model to the simulated data. Note that in this case the parameter estimates are reasonably close to the values used to generate the simulated data. ( $\mu = 0$ ,  $\hat{\mu} = 0.02$ ;  $\theta_1 = 0.8$ ,  $\hat{\theta}_1 = 0.79$ ;  $\sigma^2 = 1$ ,  $\hat{\sigma}^2 = 0.82$ .) Moreover, the graphical analysis of the residuals shows no model inadequacies (see [Output 7.1.4](#) and [Output 7.1.5](#)).

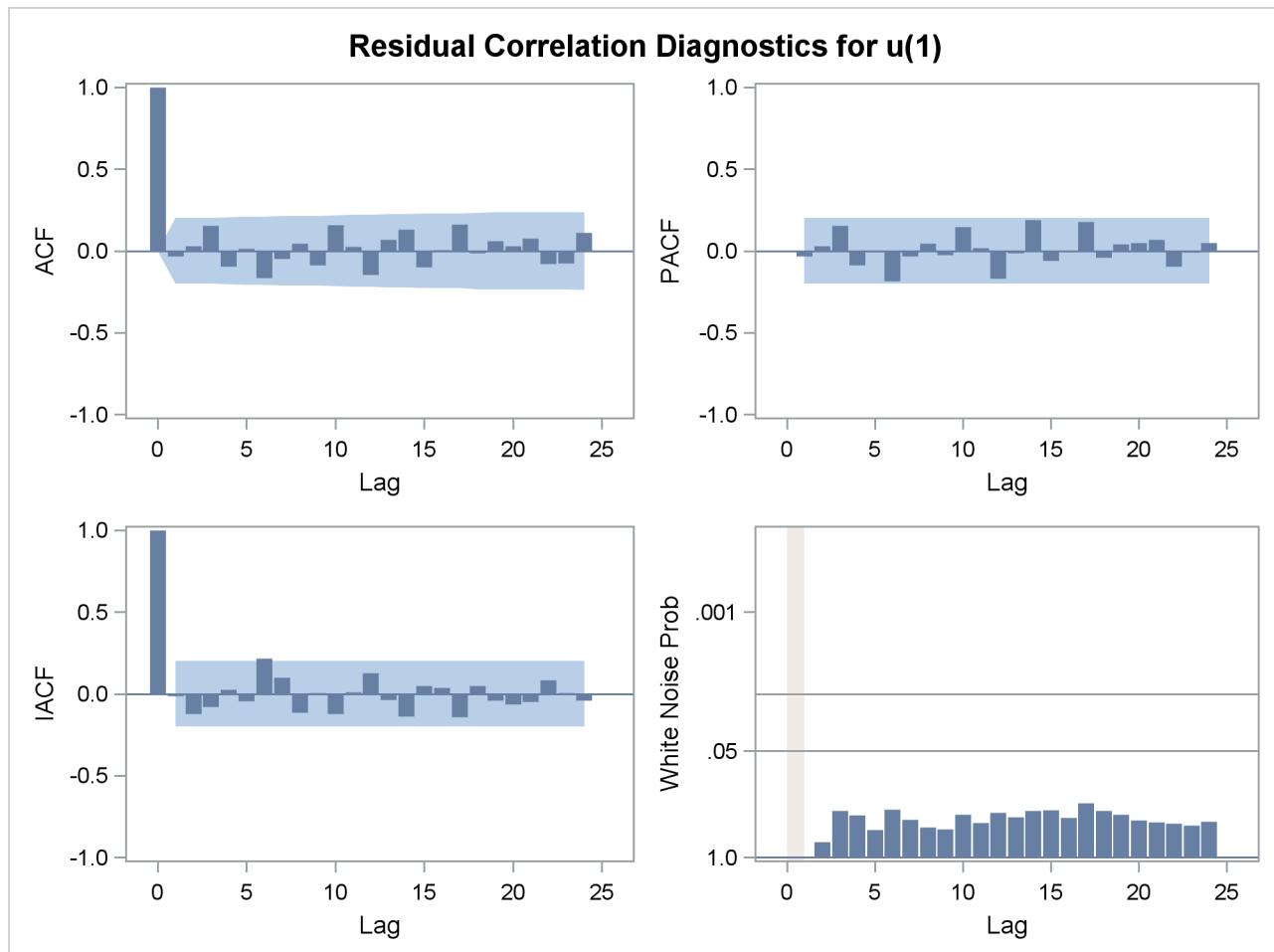
The ESTIMATE statement results are shown in [Output 7.1.3](#).

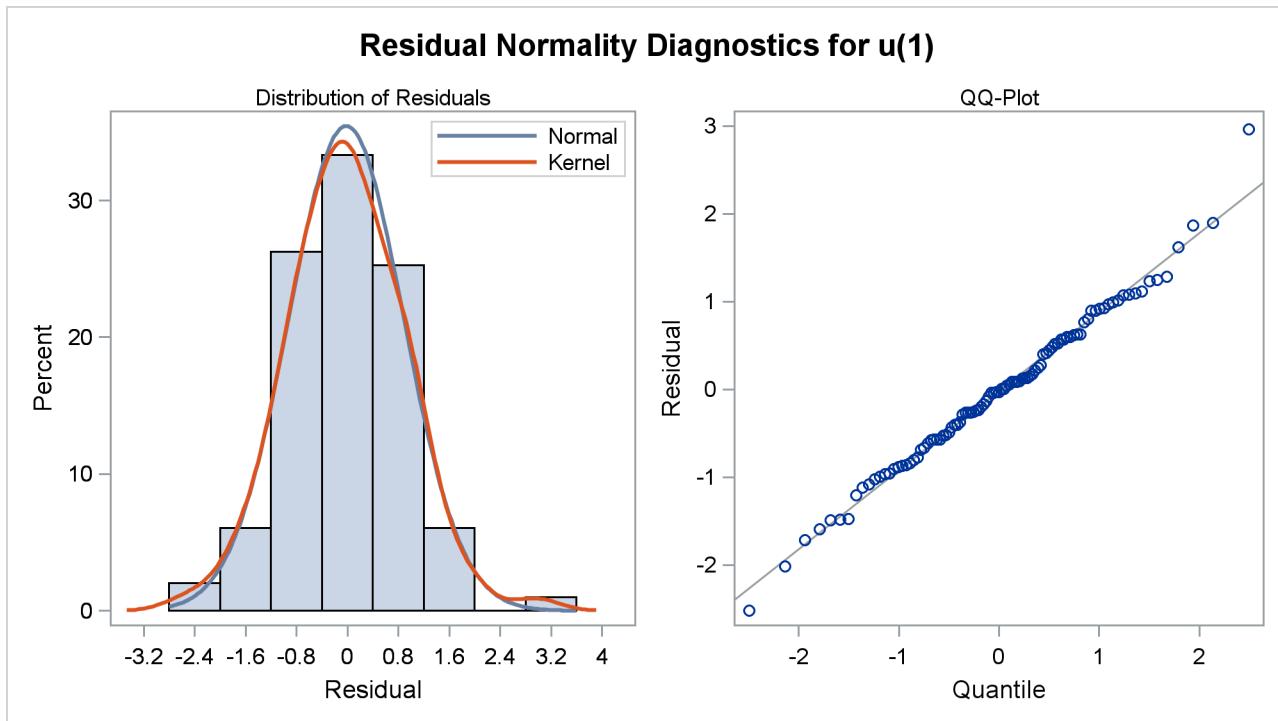
**Output 7.1.3** Output from Fitting ARIMA(0,1,1) Model

Parameter	Conditional Least Squares Estimation				
	Estimate	Standard Error	t Value	Approx	
				Pr >  t	Lag
MU	0.02056	0.01972	1.04	0.2997	0
MA1,1	0.79142	0.06474	12.22	<.0001	1

**Output 7.1.3** *continued*

<b>Constant Estimate</b>	0.020558
<b>Variance Estimate</b>	0.819807
<b>Std Error Estimate</b>	0.905432
<b>AIC</b>	263.2594
<b>SBC</b>	268.4497
<b>Number of Residuals</b>	99
<hr/>	
<b>Model for variable u</b>	
<b>Estimated Mean</b>	0.020558
<b>Period(s) of Differencing</b>	1
<hr/>	
<b>Moving Average Factors</b>	
<b>Factor 1:</b>	$1 - 0.79142 B^{**}(1)$

**Output 7.1.4** Residual Correlation Analysis of the ARIMA(0,1,1) Model

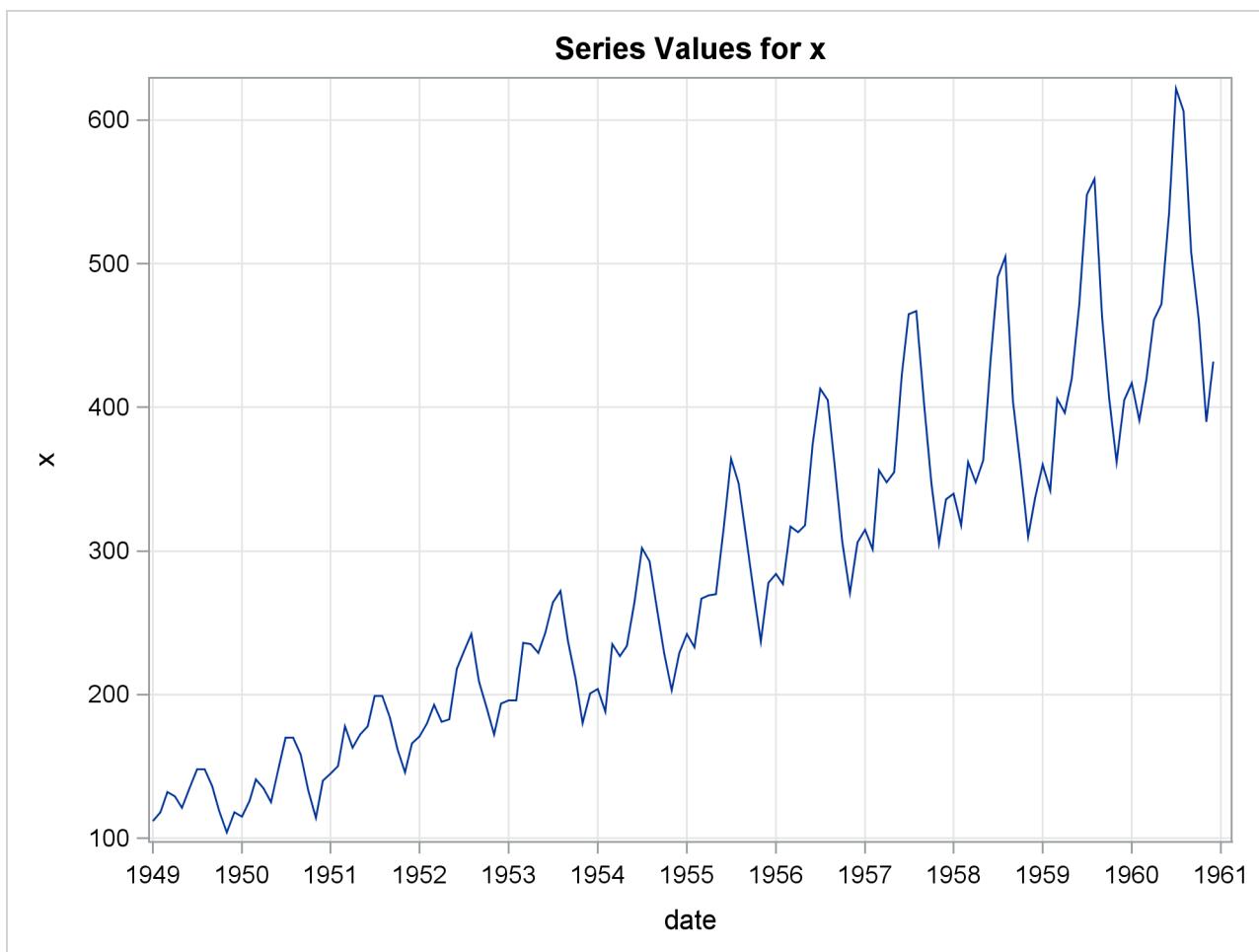
**Output 7.1.5** Residual Normality Analysis of the ARIMA(0,1,1) Model**Example 7.2: Seasonal Model for the Airline Series**

The airline passenger data, given as Series G in Box and Jenkins (1976), have been used in time series analysis literature as an example of a nonstationary seasonal time series. This example uses PROC ARIMA to fit the airline model, ARIMA(0,1,1) $\times$ (0,1,1)<sub>12</sub>, to Box and Jenkins' Series G. The following statements read the data and log-transform the series:

```
title1 'International Airline Passengers';
title2 '(Box and Jenkins Series-G)';
data seriesg;
  input x @@;
  xlog = log( x );
  date = intnx( 'month', '31dec1948'd, _n_ );
  format date monyy.;
datalines;
112 118 132 129 121 135 148 148 136 119 104 118
... more lines ...
```

The following PROC TIMESERIES step plots the series, as shown in Output 7.2.1:

```
proc timeseries data=seriesg plot=series;
  id date interval=month;
  var x;
run;
```

**Output 7.2.1** Time Series Plot of the Airline Passenger Series

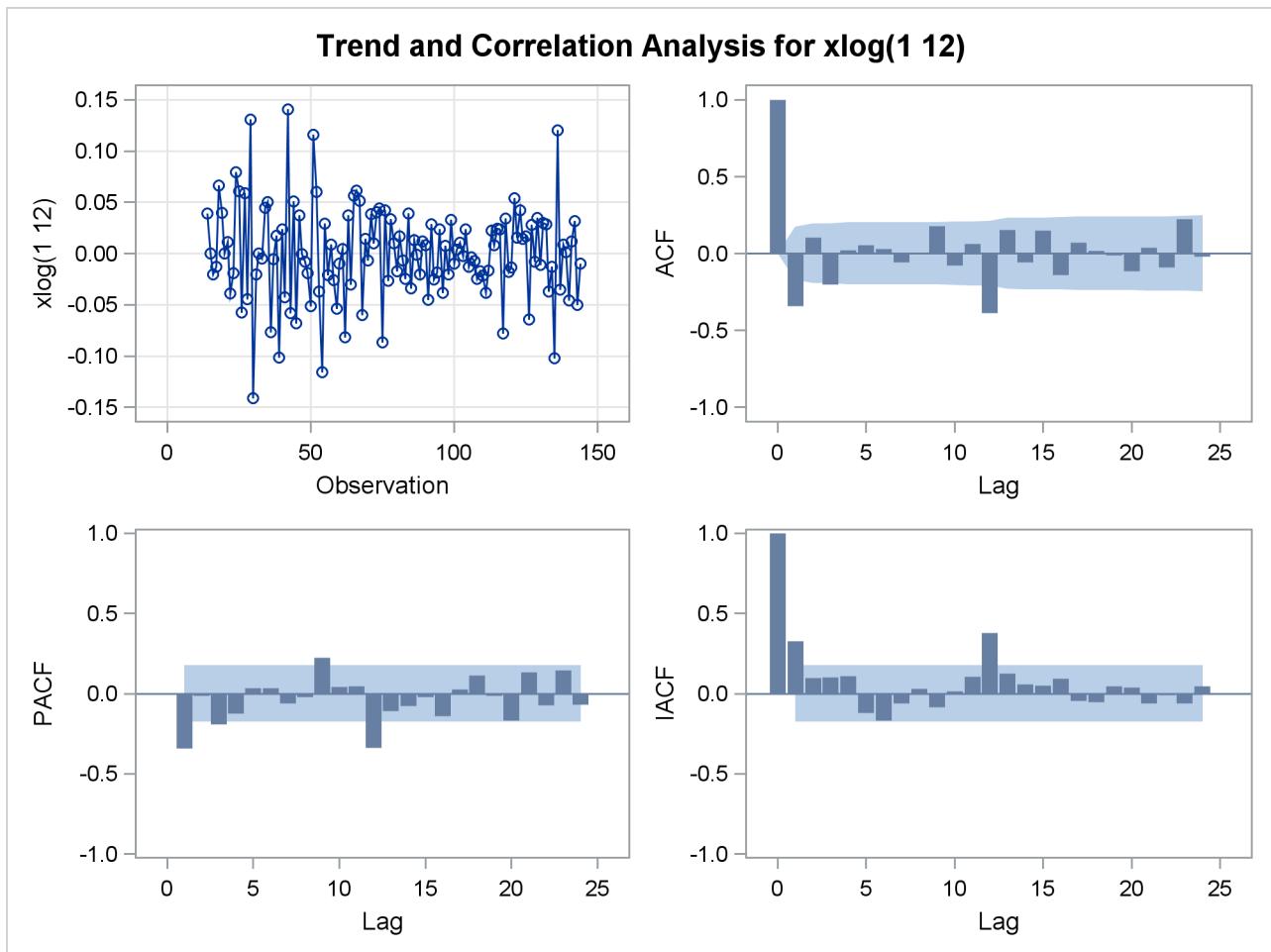
The following statements specify an ARIMA(0,1,1) $\times$ (0,1,1)<sub>12</sub> model without a mean term to the logarithms of the airline passengers series, xlog. The model is forecast, and the results are stored in the data set B.

```
/*--- Seasonal Model for the Airline Series ---*/
proc arima data=seriesg;
  identify var=xlog(1,12);
  estimate q=(1) (12) noint method=ml;
  forecast id=date interval=month printall out=b;
run;
```

The output from the IDENTIFY statement is shown in [Output 7.2.2](#). The autocorrelation plots shown are for the twice differenced series  $(1 - B)(1 - B^{12})XLOG$ . Note that the autocorrelation functions have the pattern characteristic of a first-order moving-average process combined with a seasonal moving-average process with lag 12.

**Output 7.2.2 IDENTIFY Statement Output****International Airline Passengers  
(Box and Jenkins Series-G)****The ARIMA Procedure**

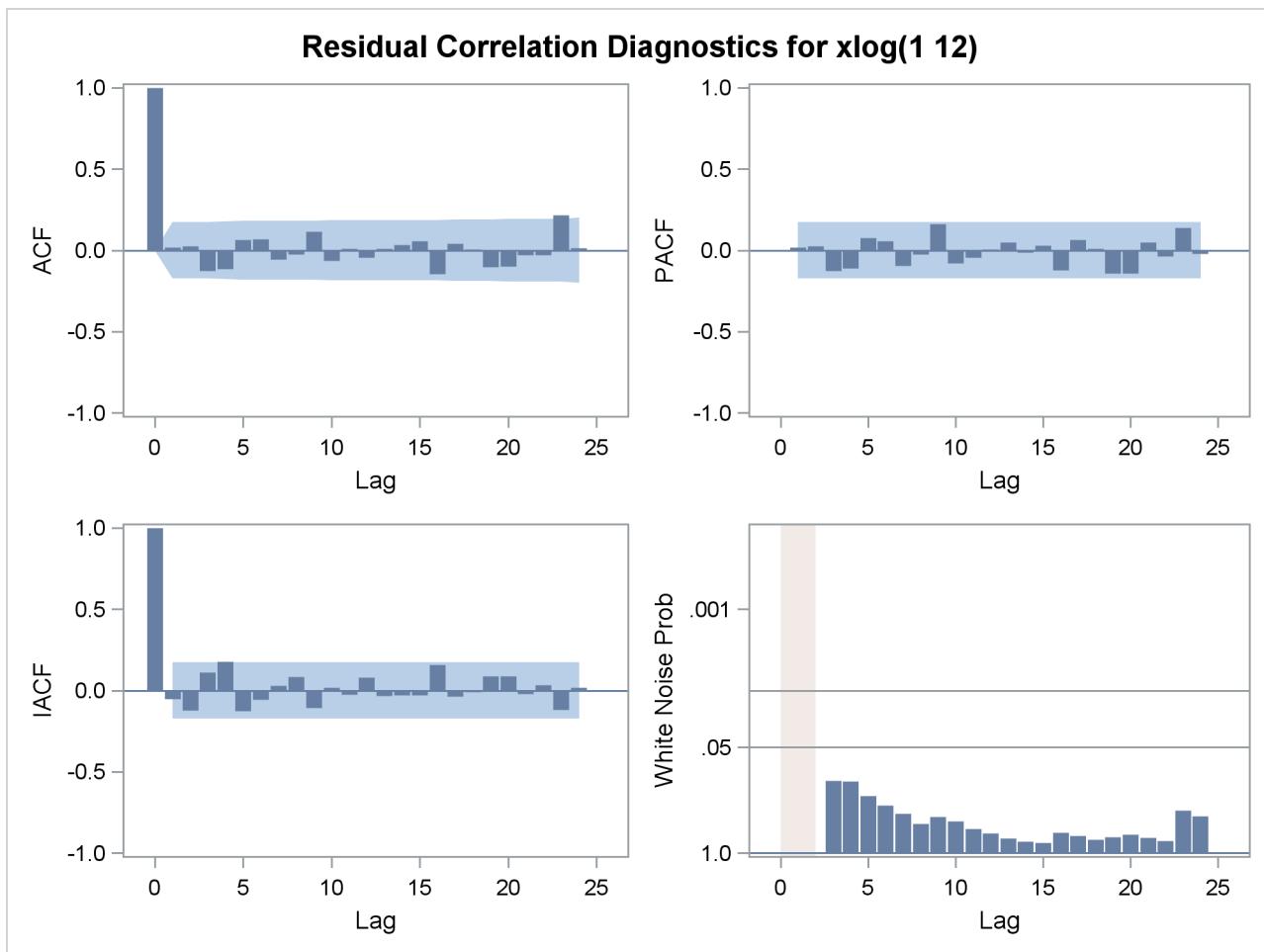
Name of Variable = xlog	
Period(s) of Differencing	1,12
Mean of Working Series	0.000291
Standard Deviation	0.045673
Number of Observations	131
Observation(s) eliminated by differencing	13

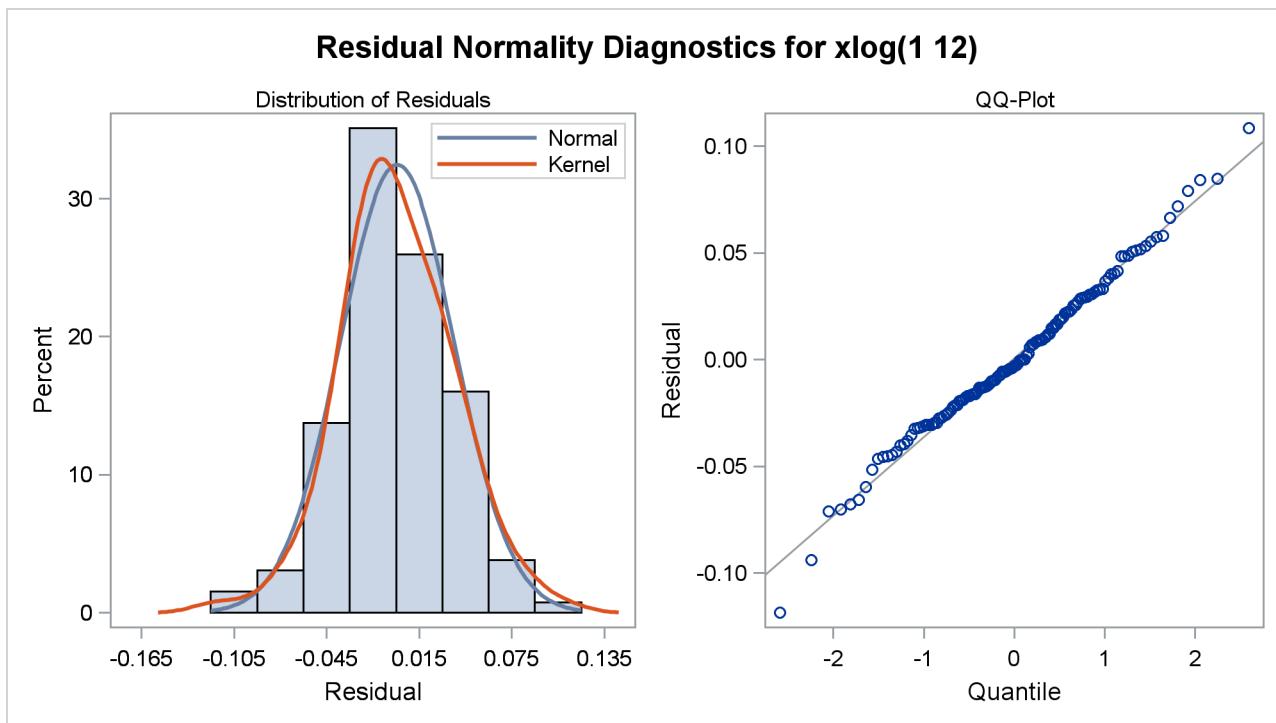
**Output 7.2.3 Trend and Correlation Analysis for the Twice Differenced Series**

The results of the ESTIMATE statement are shown in Output 7.2.4, Output 7.2.5, and Output 7.2.6. The model appears to fit the data quite well.

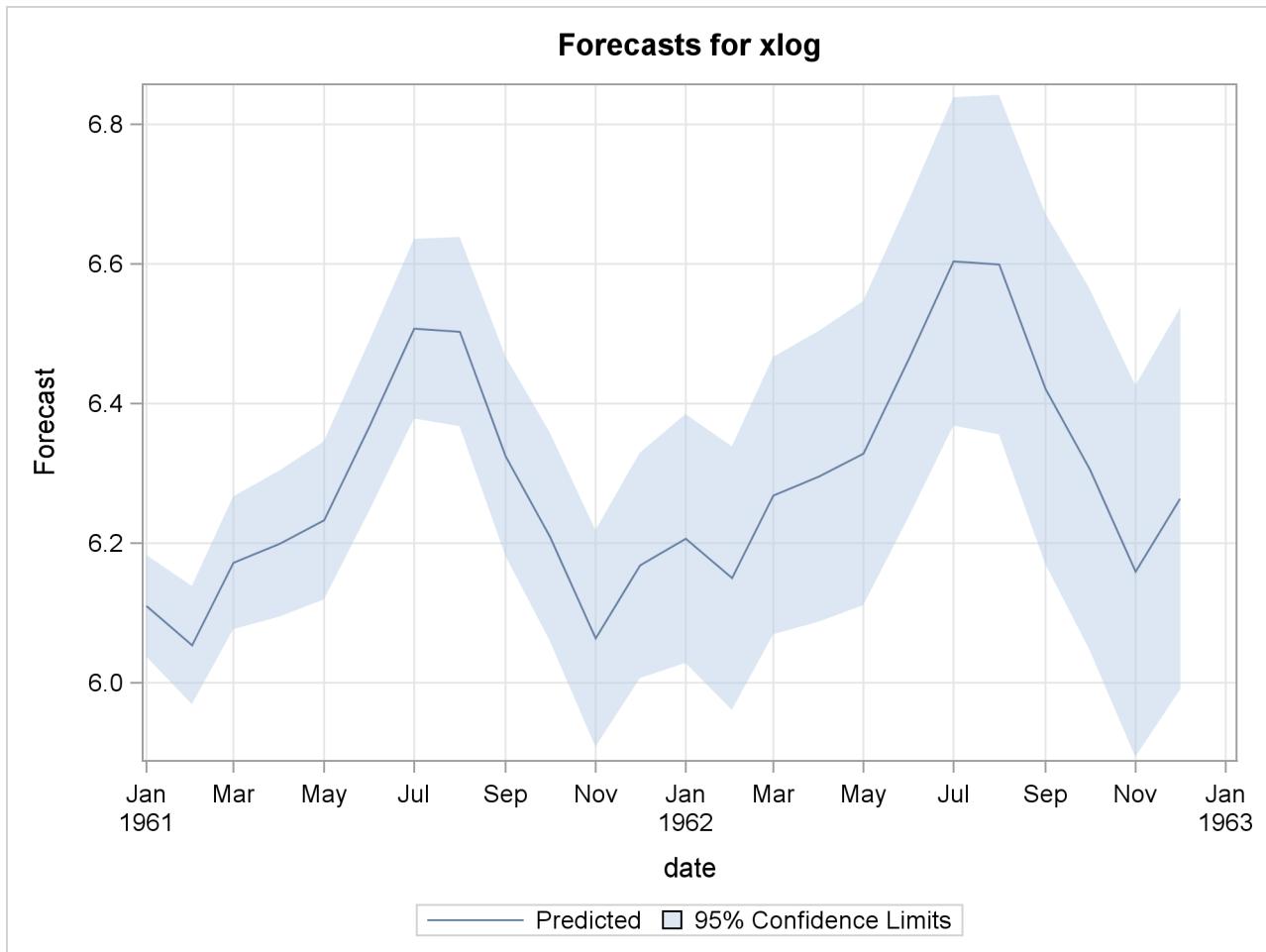
**Output 7.2.4** ESTIMATE Statement Output

Maximum Likelihood Estimation					
Parameter	Estimate	Standard Error	t Value	Approx Pr >  t	Lag
MA1,1	0.40194	0.07988	5.03	<.0001	1
MA2,1	0.55686	0.08403	6.63	<.0001	12
<hr/>					
Variance Estimate	0.001369				
Std Error Estimate		0.037			
AIC		-485.393			
SBC		-479.643			
Number of Residuals	131				
<hr/>					
<b>Model for variable xlog</b>					
<b>Period(s) of Differencing</b> 1,12					
<hr/>					
<b>Moving Average Factors</b>					
<b>Factor 1:</b> 1 - 0.40194 B**(1)					
<b>Factor 2:</b> 1 - 0.55686 B**(12)					

**Output 7.2.5** Residual Analysis of the Airline Model: Correlation

**Output 7.2.6** Residual Analysis of the Airline Model: Normality

The forecasts and their confidence limits for the transformed series are shown in Output 7.2.7.

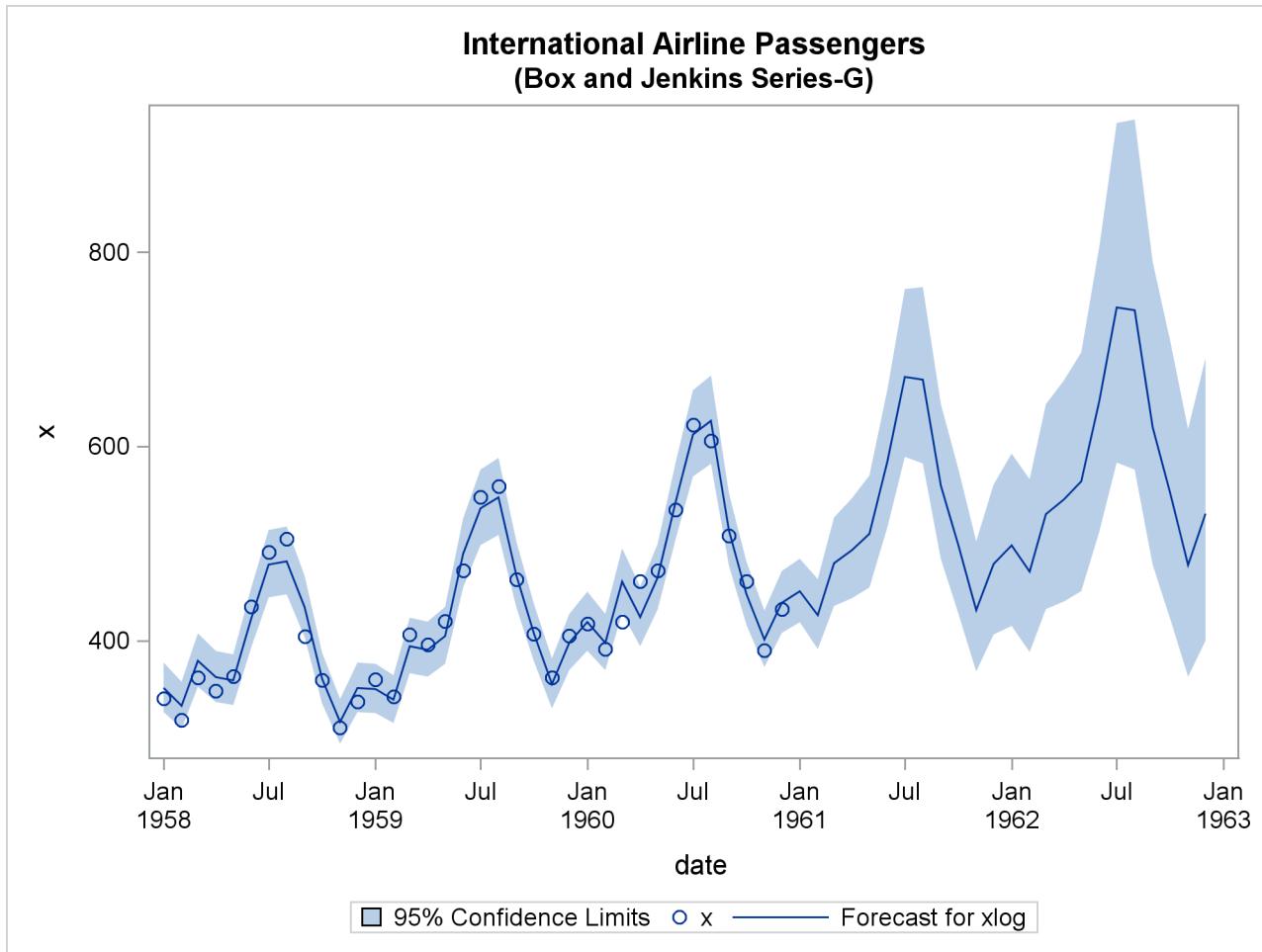
**Output 7.2.7** Forecast Plot for the Transformed Series

The following statements retransform the forecast values to get forecasts in the original scales. See the section “[Forecasting Log Transformed Data](#)” on page 249 for more information.

```
data c;
  set b;
  x      = exp( xlog );
  forecast = exp( forecast + std*std/2 );
  195    = exp( 195 );
  u95    = exp( u95 );
run;
```

The forecasts and their confidence limits are plotted by using the following PROC SGPLOT step. The plot is shown in [Output 7.2.8](#).

```
proc sgplot data=c;
  where date >= '1jan58'd;
  band Upper=u95 Lower=195 x=date
    / LegendLabel="95% Confidence Limits";
  scatter x=date y=x;
  series x=date y=forecast;
run;
```

**Output 7.2.8** Plot of the Forecast for the Original Series


---

### Example 7.3: Model for Series J Data from Box and Jenkins

This example uses the Series J data from Box and Jenkins (1976). First, the input series X is modeled with a univariate ARMA model. Next, the dependent series Y is cross-correlated with the input series. Since a model has been fit to X, both Y and X are prewhitened by this model before the sample cross-correlations are computed. Next, a transfer function model is fit with no structure on the noise term. The residuals from this model are analyzed; then, the full model, transfer function and noise, is fit to the data.

The following statements read 'Input Gas Rate' and 'Output CO<sub>2</sub>' from a gas furnace. (Data values are not shown. The full example including data is in the SAS/ETS sample library.)

```

title1 'Gas Furnace Data';
title2 '(Box and Jenkins, Series J)';
data seriesj;
  input x y @@;
  label x = 'Input Gas Rate'
        y = 'Output CO2';
  datalines;

```

```
-0.109 53.8 0.000 53.6 0.178 53.5 0.339 53.5
... more lines ...
```

The following statements produce Output 7.3.1 through Output 7.3.11:

```
proc arima data=seriesj;

/*---- Look at the input process -----*/
identify var=x;
run;

/*---- Fit a model for the input -----*/
estimate p=3 plot;
run;

/*---- Cross-correlation of prewhitened series -----*/
identify var=y crosscorr=(x) nlag=12;
run;

/*---- Fit a simple transfer function - look at residuals ---*/
estimate input=( 3 $ (1,2)/(1) x );
run;

/*---- Final Model - look at residuals -----*/
estimate p=2 input=( 3 $ (1,2)/(1) x );
run;

quit;
```

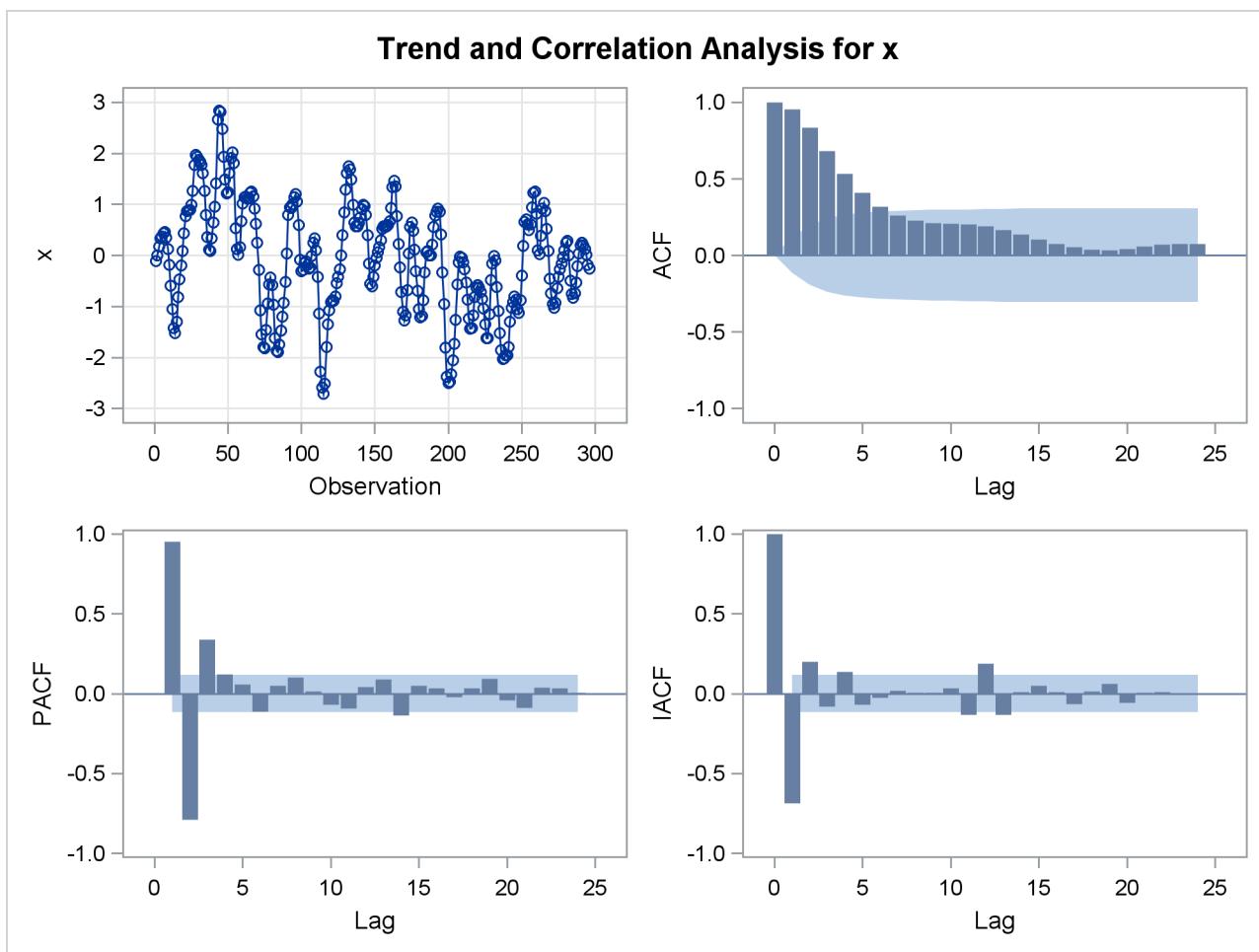
The results of the first IDENTIFY statement for the input series X are shown in Output 7.3.1. The correlation analysis suggests an AR(3) model.

#### Output 7.3.1 IDENTIFY Statement Results for X

##### Gas Furnace Data (Box and Jenkins, Series J)

##### The ARIMA Procedure

Name of Variable = x	
Mean of Working Series	-0.05683
Standard Deviation	1.070952
Number of Observations	296

**Output 7.3.2** IDENTIFY Statement Results for X: Trend and Correlation

The ESTIMATE statement results for the AR(3) model for the input series X are shown in Output 7.3.3.

**Output 7.3.3** Estimates of the AR(3) Model for X

Conditional Least Squares Estimation					
Parameter	Estimate	Standard Error	t Value	Pr >  t	Approx Lag
MU	-0.12280	0.10902	-1.13	0.2609	0
AR1,1	1.97607	0.05499	35.94	<.0001	1
AR1,2	-1.37499	0.09967	-13.80	<.0001	2
AR1,3	0.34336	0.05502	6.24	<.0001	3

Constant Estimate	-0.00682
Variance Estimate	0.035797
Std Error Estimate	0.1892
AIC	-141.667
SBC	-126.906
Number of Residuals	296

**Output 7.3.3** *continued*

Model for variable x	
Estimated Mean	-0.1228
Autoregressive Factors	
Factor 1: 1 - 1.97607 B**(1) + 1.37499 B**(2) - 0.34336 B**(3)	

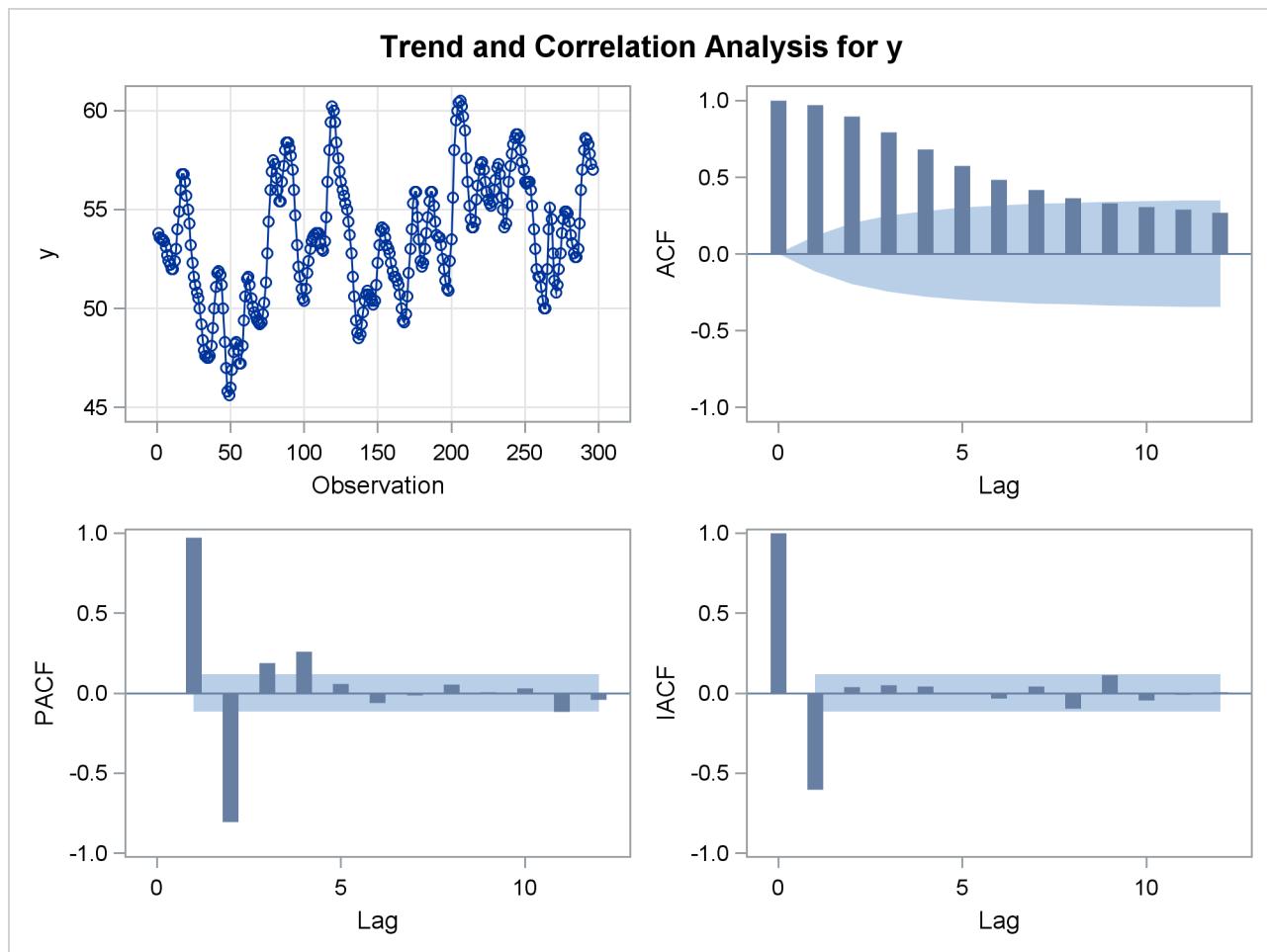
The IDENTIFY statement results for the dependent series Y cross-correlated with the input series X are shown in [Output 7.3.4](#), [Output 7.3.5](#), [Output 7.3.6](#), and [Output 7.3.7](#). Since a model has been fit to X, both Y and X are prewhitened by this model before the sample cross-correlations are computed.

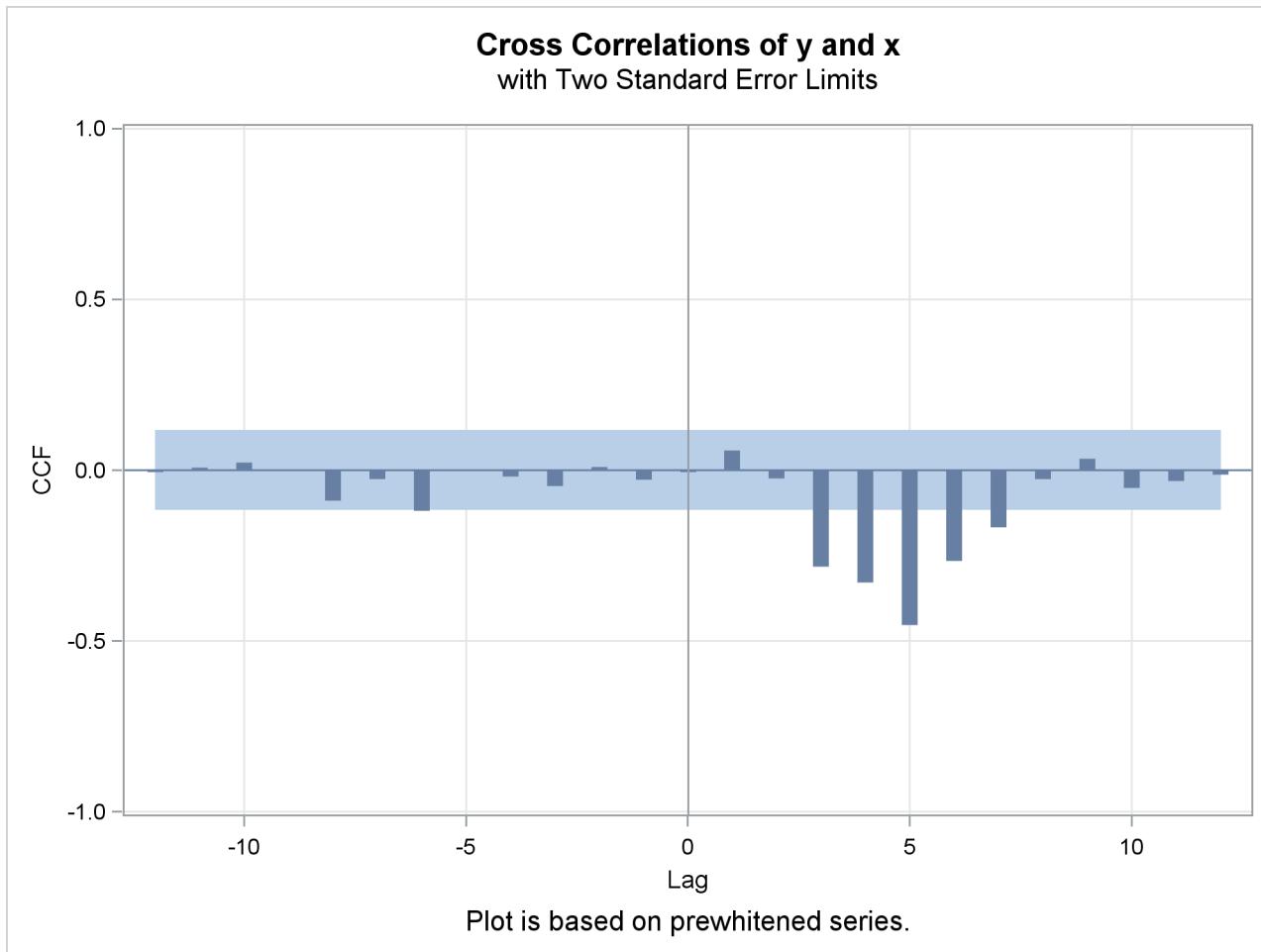
**Output 7.3.4** Summary Table: Y Cross-Correlated with X

Correlation of y and x	
Number of Observations	296
Variance of transformed series y	0.131438
Variance of transformed series x	0.035357
Both series have been prewhitened.	

**Output 7.3.5** Prewitening Filter

Autoregressive Factors	
Factor 1: 1 - 1.97607 B**(1) + 1.37499 B**(2) - 0.34336 B**(3)	

**Output 7.3.6** IDENTIFY Statement Results for Y: Trend and Correlation

**Output 7.3.7** IDENTIFY Statement for Y Cross-Correlated with X

The ESTIMATE statement results for the transfer function model with no structure on the noise term are shown in [Output 7.3.8](#), [Output 7.3.9](#), and [Output 7.3.10](#).

**Output 7.3.8** Estimation Output of the First Transfer Function Model

Conditional Least Squares Estimation						
Parameter	Estimate	Standard Error	t Value	Pr >  t	Lag	Variable Shift
MU	53.32256	0.04926	1082.51	<.0001	0	y
NUM1	-0.56467	0.22405	-2.52	0.0123	0	x
NUM1,1	0.42623	0.46472	0.92	0.3598	1	x
NUM1,2	0.29914	0.35506	0.84	0.4002	2	x
DEN1,1	0.60073	0.04101	14.65	<.0001	1	x

Constant Estimate	53.32256
Variance Estimate	0.702625
Std Error Estimate	0.838227
AIC	728.0754
SBC	746.442
Number of Residuals	291

**Output 7.3.9** Model Summary: First Transfer Function Model

---

**Model for variable y**

---

**Estimated Intercept** 53.32256

---

**Input Number 1**

**Input Variable** x  
**Shift** 3

---

**Numerator Factors**


---

**Factor 1:** -0.5647 - 0.42623 B\*\*(1) - 0.29914 B\*\*2)

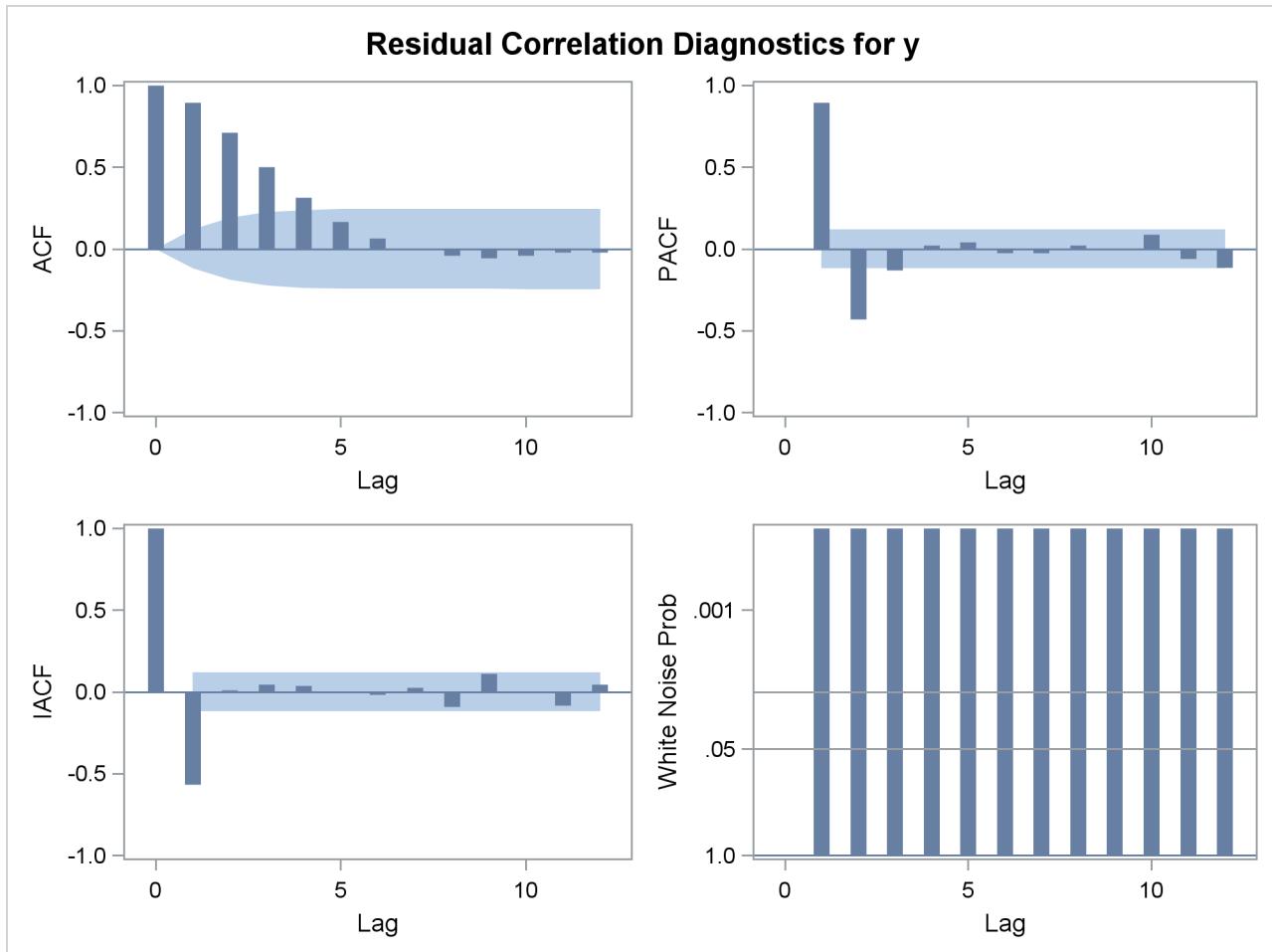
---

**Denominator Factors**


---

**Factor 1:** 1 - 0.60073 B\*\*1)

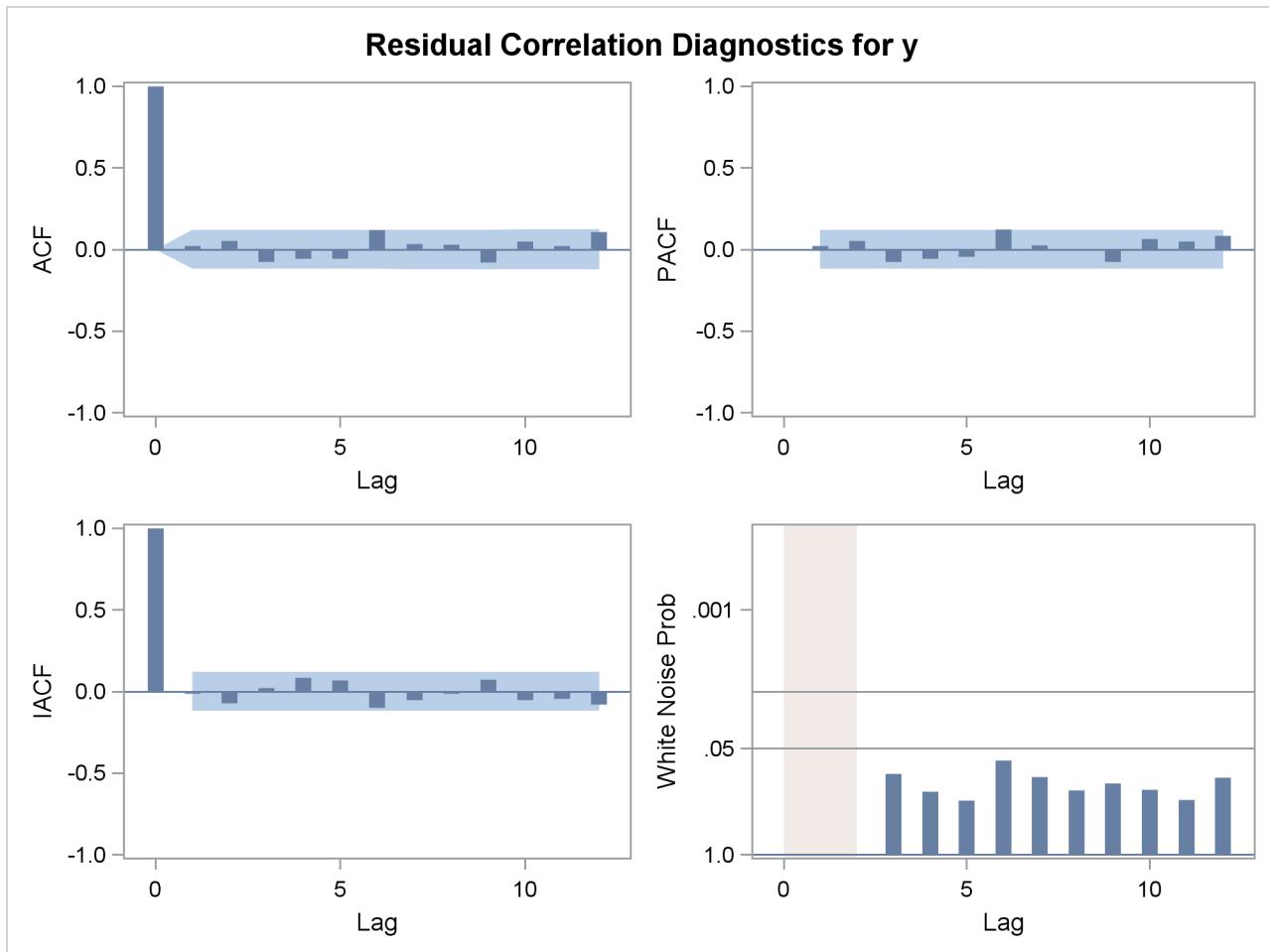
---

**Output 7.3.10** Residual Analysis: First Transfer Function Model

The residual correlation analysis suggests an AR(2) model for the noise part of the model. The ESTIMATE statement results for the final transfer function model with AR(2) noise are shown in [Output 7.3.11](#).

**Output 7.3.11** Estimation Output of the Final Model

Conditional Least Squares Estimation						
Parameter	Estimate	Standard Error	t Value	Pr >  t	Lag	Approx Variable Shift
MU	53.26304	0.11929	446.48	<.0001	0	y
AR1,1	1.53291	0.04754	32.25	<.0001	1	y
AR1,2	-0.63297	0.05006	-12.64	<.0001	2	y
NUM1	-0.53522	0.07482	-7.15	<.0001	0	x
NUM1,1	0.37603	0.10287	3.66	0.0003	1	x
NUM1,2	0.51895	0.10783	4.81	<.0001	2	x
DEN1,1	0.54841	0.03822	14.35	<.0001	1	x
<hr/>						
<b>Constant Estimate</b>		5.329425				
<b>Variance Estimate</b>		0.058828				
<b>Std Error Estimate</b>		0.242544				
<b>AIC</b>		8.292809				
<b>SBC</b>		34.00607				
<b>Number of Residuals</b>		291				

**Output 7.3.12** Residual Analysis of the Final Model

**Output 7.3.13** Model Summary of the Final Model

Model for variable y	
Estimated Intercept	53.26304
Autoregressive Factors	
Factor 1:	1 - 1.53291 B**(1) + 0.63297 B**(2)
Input Number 1	
Input Variable	x
Shift	3
Numerator Factors	
Factor 1:	-0.5352 - 0.37603 B**(1) - 0.51895 B**(2)
Denominator Factors	
Factor 1:	1 - 0.54841 B**(1)

**Example 7.4: An Intervention Model for Ozone Data**

This example fits an intervention model to ozone data as suggested by Box and Tiao (1975). Notice that the response variable, OZONE, and the innovation, X1, are seasonally differenced. The final model for the differenced data is a multiple regression model with a moving-average structure assumed for the residuals.

The model is fit by maximum likelihood. The seasonal moving-average parameter and its standard error are fairly sensitive to which method is chosen to fit the model (Ansley and Newbold 1980; Davidson 1981); thus, fitting the model by the unconditional or conditional least squares method produces somewhat different estimates for these parameters.

Some missing values are appended to the end of the input data to generate additional values for the independent variables. Since the independent variables are not modeled, values for them must be available for any times at which predicted values are desired. In this case, predicted values are requested for 12 periods beyond the end of the data. Thus, values for X1, WINTER, and SUMMER must be given for 12 periods ahead.

The following statements read in the data and compute dummy variables for use as intervention inputs:

```

title1 'Intervention Data for Ozone Concentration';
title2 '(Box and Tiao, JASA 1975 P.70)';
data air;
  input ozone @@;
  label ozone = 'Ozone Concentration'
    x1      = 'Intervention for post 1960 period'
    summer = 'Summer Months Intervention'
    winter = 'Winter Months Intervention';
  date = intnx( 'month', '31dec1954'd, _n_ );
  format date monyy.;
  month = month( date );
  year = year( date );
  x1 = year >= 1960;
  summer = ( 5 < month < 11 ) * ( year > 1965 );
  winter = ( year > 1965 ) - summer;
  datalines;

```

```

2.7  2.0  3.6  5.0  6.5  6.1  5.9  5.0  6.4  7.4  8.2  3.9
4.1  4.5  5.5  3.8  4.8  5.6  6.3  5.9  8.7  5.3  5.7  5.7
3.0  3.4  4.9  4.5  4.0  5.7  6.3  7.1  8.0  5.2  5.0  4.7
3.7  3.1  2.5  4.0  4.1  4.6  4.4  4.2  5.1  4.6  4.4  4.0

```

... more lines ...

The following statements produce [Output 7.4.1](#) through [Output 7.4.3](#):

```

proc arima data=air;

/* Identify and seasonally difference ozone series */
identify var=ozone(12)
            crosscorr=( x1(12) summer winter ) noprint;

/* Fit a multiple regression with a seasonal MA model */
/* by the maximum likelihood method */
estimate q=(1)(12) input=( x1 summer winter )
            noconstant method=ml;

/* Forecast */
forecast lead=12 id=date interval=month;

run;

```

The ESTIMATE statement results are shown in [Output 7.4.1](#) and [Output 7.4.2](#).

### Output 7.4.1 Parameter Estimates Intervention Data for Ozone Concentration (Box and Tiao, JASA 1975 P.70)

#### The ARIMA Procedure

Maximum Likelihood Estimation						
Parameter	Estimate	Standard Error	t Value	Pr >  t	Lag	Variable Shift
MA1,1	-0.26684	0.06710	-3.98	<.0001	1	ozone
MA2,1	0.76665	0.05973	12.83	<.0001	12	ozone
NUM1	-1.33062	0.19236	-6.92	<.0001	0	x1
NUM2	-0.23936	0.05952	-4.02	<.0001	0	summer
NUM3	-0.08021	0.04978	-1.61	0.1071	0	winter

Variance Estimate	0.634506
Std Error Estimate	0.796559
AIC	501.7696
SBC	518.3602
Number of Residuals	204

### Output 7.4.2 Model Summary

<b>Model for variable ozone</b>	
<b>Period(s) of Differencing</b>	12
<b>Moving Average Factors</b>	
<b>Factor 1:</b>	$1 + 0.26684 B^{**}(1)$
<b>Factor 2:</b>	$1 - 0.76665 B^{**}(12)$
<b>Input Number 1</b>	
<b>Input Variable</b>	x1
<b>Period(s) of Differencing</b>	12
<b>Overall Regression Factor</b>	-1.33062

The FORECAST statement results are shown in Output 7.4.3.

### Output 7.4.3 Forecasts

<b>Forecasts for variable ozone</b>				
<b>Obs</b>	<b>Forecast</b>	<b>Std Error</b>	<b>95% Confidence Limits</b>	
217	1.4205	0.7966	-0.1407	2.9817
218	1.8446	0.8244	0.2287	3.4604
219	2.4567	0.8244	0.8408	4.0725
220	2.8590	0.8244	1.2431	4.4748
221	3.1501	0.8244	1.5342	4.7659
222	2.7211	0.8244	1.1053	4.3370
223	3.3147	0.8244	1.6989	4.9306
224	3.4787	0.8244	1.8629	5.0946
225	2.9405	0.8244	1.3247	4.5564
226	2.3587	0.8244	0.7429	3.9746
227	1.8588	0.8244	0.2429	3.4746
228	1.2898	0.8244	-0.3260	2.9057

## Example 7.5: Using Diagnostics to Identify ARIMA Models

Fitting ARIMA models is as much an art as it is a science. The ARIMA procedure has diagnostic options to help tentatively identify the orders of both stationary and nonstationary ARIMA processes.

Consider the Series A in Box, Jenkins, and Reinsel (1994), which consists of 197 concentration readings taken every two hours from a chemical process. Let Series A be a data set that contains these readings in a variable named X. The following SAS statements use the SCAN option of the IDENTIFY statement to generate Output 7.5.1 and Output 7.5.2. See “The SCAN Method” on page 236 for details of the SCAN method.

```
/--- Order Identification Diagnostic with SCAN Method ---/
proc arima data=SeriesA;
  identify var=x scan;
run;
```

**Output 7.5.1** Example of SCAN Tables**SERIES A: Chemical Process Concentration Readings****The ARIMA Procedure**

Squared Canonical Correlation Estimates						
Lags	MA 0	MA 1	MA 2	MA 3	MA 4	MA 5
AR 0	0.3263	0.2479	0.1654	0.1387	0.1183	0.1417
AR 1	0.0643	0.0012	0.0028	<.0001	0.0051	0.0002
AR 2	0.0061	0.0027	0.0021	0.0011	0.0017	0.0079
AR 3	0.0072	<.0001	0.0007	0.0005	0.0019	0.0021
AR 4	0.0049	0.0010	0.0014	0.0014	0.0039	0.0145
AR 5	0.0202	0.0009	0.0016	<.0001	0.0126	0.0001

SCAN Chi-Square[1] Probability Values						
Lags	MA 0	MA 1	MA 2	MA 3	MA 4	MA 5
AR 0	<.0001	<.0001	<.0001	0.0007	0.0037	0.0024
AR 1	0.0003	0.6649	0.5194	0.9235	0.3993	0.8528
AR 2	0.2754	0.5106	0.5860	0.7346	0.6782	0.2766
AR 3	0.2349	0.9812	0.7667	0.7861	0.6810	0.6546
AR 4	0.3297	0.7154	0.7113	0.6995	0.5807	0.2205
AR 5	0.0477	0.7254	0.6652	0.9576	0.2660	0.9168

In Output 7.5.1, there is one (maximal) rectangular region in which all the elements are insignificant with 95% confidence. This region has a vertex at (1,1). Output 7.5.2 gives recommendations based on the significance level specified by the ALPHA=siglevel option.

**Output 7.5.2** Example of SCAN Option Tentative Order Selection

ARMA(p+d,q)
Tentative
Order
Selection
Tests
SCAN
<u>p+d</u>
1
<u>q</u>
1
(5% Significance Level)

Another order identification diagnostic is the extended sample autocorrelation function or ESACF method. See “The ESACF Method” on page 233 for details of the ESACF method.

The following statements generate Output 7.5.3 and Output 7.5.4:

```
/*-- Order Identification Diagnostic with ESACF Method --*/
proc arima data=SeriesA;
  identify var=x esacf;
run;
```

**Output 7.5.3** Example of ESACF Tables**SERIES A: Chemical Process Concentration Readings****The ARIMA Procedure**

Extended Sample Autocorrelation Function						
Lags	MA 0	MA 1	MA 2	MA 3	MA 4	MA 5
AR 0	0.5702	0.4951	0.3980	0.3557	0.3269	0.3498
AR 1	-0.3907	0.0425	-0.0605	-0.0083	-0.0651	-0.0127
AR 2	-0.2859	-0.2699	-0.0449	0.0089	-0.0509	-0.0140
AR 3	-0.5030	-0.0106	0.0946	-0.0137	-0.0148	-0.0302
AR 4	-0.4785	-0.0176	0.0827	-0.0244	-0.0149	-0.0421
AR 5	-0.3878	-0.4101	-0.1651	0.0103	-0.1741	-0.0231

ESACF Probability Values						
Lags	MA 0	MA 1	MA 2	MA 3	MA 4	MA 5
AR 0	<.0001	<.0001	0.0001	0.0014	0.0053	0.0041
AR 1	<.0001	0.5974	0.4622	0.9198	0.4292	0.8768
AR 2	<.0001	0.0002	0.6106	0.9182	0.5683	0.8592
AR 3	<.0001	0.9022	0.2400	0.8713	0.8930	0.7372
AR 4	<.0001	0.8380	0.3180	0.7737	0.8913	0.6213
AR 5	<.0001	<.0001	0.0765	0.9142	0.1038	0.8103

In Output 7.5.3, there are three right-triangular regions in which all elements are insignificant at the 5% level. The triangles have vertices (1,1), (3,1), and (4,1). Since the triangle at (1,1) covers more insignificant terms, it is recommended first. Similarly, the remaining recommendations are ordered by the number of insignificant terms contained in the triangle. Output 7.5.4 gives recommendations based on the significance level specified by the ALPHA=siglevel option.

**Output 7.5.4** Example of ESACF Option Tentative Order Selection

ARMA(p+d,q)	
Tentative	
Order	
Selection	
Tests	
SCAN	
p+d	q
1	1
(5% Significance Level)	

If you also specify the SCAN option in the same IDENTIFY statement, the two recommendations are printed side by side:

```
/*--- Combination of SCAN and ESACF Methods ---*/
proc arima data=SeriesA;
  identify var=x scan esacf;
run;
```

Output 7.5.5 shows the results.

**Output 7.5.5** Example of SCAN and ESACF Option Combined  
**SERIES A: Chemical Process Concentration Readings**

**The ARIMA Procedure**

ARMA(p+d,q)			
Tentative Order			
Selection Tests			
SCAN		ESACF	
p+d	q	p+d	q
1	1	1	1
		3	1
		4	1

**(5% Significance Level)**

From [Output 7.5.5](#), the autoregressive and moving-average orders are tentatively identified by both SCAN and ESACF tables to be  $(p + d, q) = (1, 1)$ . Because both the SCAN and ESACF indicate a  $p + d$  term of 1, a unit root test should be used to determine whether this autoregressive term is a unit root. Since a moving-average term appears to be present, a large autoregressive term is appropriate for the augmented Dickey-Fuller test for a unit root.

Submitting the following statements generates [Output 7.5.6](#):

```
/*** Augmented Dickey-Fuller Unit Root Tests ***/
proc arima data=SeriesA;
  identify var=x stationarity=(adf=(5, 6, 7, 8));
run;
```

**Output 7.5.6** Example of STATIONARITY Option Output  
**SERIES A: Chemical Process Concentration Readings**

**The ARIMA Procedure**

Augmented Dickey-Fuller Unit Root Tests							
Type	Lags	Rho	Pr < Rho	Tau	Pr < Tau	F	Pr > F
<b>Zero Mean</b>	5	0.0403	0.6913	0.42	0.8024		
	6	0.0479	0.6931	0.63	0.8508		
	7	0.0376	0.6907	0.49	0.8200		
	8	0.0354	0.6901	0.48	0.8175		
<b>Single Mean</b>	5	-18.4550	0.0150	-2.67	0.0821	3.67	0.1367
	6	-10.8939	0.1043	-2.02	0.2767	2.27	0.4931
	7	-10.9224	0.1035	-1.93	0.3172	2.00	0.5605
	8	-10.2992	0.1208	-1.83	0.3650	1.81	0.6108
<b>Trend</b>	5	-18.4360	0.0871	-2.66	0.2561	3.54	0.4703
	6	-10.8436	0.3710	-2.01	0.5939	2.04	0.7694
	7	-10.7427	0.3773	-1.90	0.6519	1.91	0.7956
	8	-10.0370	0.4236	-1.79	0.7081	1.74	0.8293

The preceding test results show that a unit root is very likely given that none of the  $p$ -values are small enough to cause you to reject the null hypothesis that the series has a unit root. Based on this test and the previous

results, the series should be differenced, and an ARIMA(0,1,1) would be a good choice for a tentative model for Series A.

Using the recommendation that the series be differenced, the following statements generate [Output 7.5.7](#):

```
/*-- Minimum Information Criterion --*/
proc arima data=SeriesA;
  identify var=x(1) minic;
run;
```

#### **Output 7.5.7 Example of MINIC Table**

#### **SERIES A: Chemical Process Concentration Readings**

##### **The ARIMA Procedure**

Minimum Information Criterion						
Lags	MA 0	MA 1	MA 2	MA 3	MA 4	MA 5
AR 0	-2.05761	-2.3497	-2.32358	-2.31298	-2.30967	-2.28528
AR 1	-2.23291	-2.32345	-2.29665	-2.28644	-2.28356	-2.26011
AR 2	-2.23947	-2.30313	-2.28084	-2.26065	-2.25685	-2.23458
AR 3	-2.25092	-2.28088	-2.25567	-2.23455	-2.22997	-2.20769
AR 4	-2.25934	-2.2778	-2.25363	-2.22983	-2.20312	-2.19531
AR 5	-2.2751	-2.26805	-2.24249	-2.21789	-2.19667	-2.17426

The error series is estimated by using an AR(7) model, and the minimum of this MINIC table is  $BIC(0, 1)$ . This diagnostic confirms the previous result which indicates that an ARIMA(0,1,1) is a tentative model for Series A.

If you also specify the SCAN or MINIC option in the same IDENTIFY statement as follows, the BIC associated with the SCAN table and ESACF table recommendations is listed. [Output 7.5.8](#) shows the results.

```
/*-- Combination of MINIC, SCAN and ESACF Options --*/
proc arima data=SeriesA;
  identify var=x(1) minic scan esacf;
run;
```

#### **Output 7.5.8 Example of SCAN, ESACF, MINIC Options Combined**

#### **SERIES A: Chemical Process Concentration Readings**

##### **The ARIMA Procedure**

ARMA(p+d,q) Tentative Order Selection Tests						
SCAN			ESACF			
p+d	q	BIC	p+d	q	BIC	
0	1	-2.3497	0	1	-2.3497	
			1	1	-2.32345	

**(5% Significance Level)**

---

## Example 7.6: Detection of Level Changes in the Nile River Data

This example shows how to use the OUTLIER statement to detect changes in the dynamics of the time series being modeled. The time series used here is discussed in De Jong and Penzer (1998). The data consist of readings of the annual flow volume of the Nile River at Aswan from 1871 to 1970. These data have also been studied by Cobb (1978). These studies indicate that river flow levels in the years 1877 and 1913 are strong candidates for additive outliers and that there was a shift in the flow levels starting from the year 1899. This shift in 1899 is attributed partly to the weather changes and partly to the start of construction work for a new dam at Aswan. The following DATA step statements create the input data set.

```

data nile;
  input level @@;
  year = intnx( 'year', '1jan1871'd, _n_-1 );
  format year year4.;
datalines;
1120 1160 963 1210 1160 1160 813 1230 1370 1140
995 935 1110 994 1020 960 1180 799 958 1140
1100 1210 1150 1250 1260 1220 1030 1100 774 840
...
  ... more lines ...

```

The following program fits an ARIMA model, ARIMA(0,1,1), similar to the structural model suggested in De Jong and Penzer (1998). This model is also suggested by the usual correlation analysis of the series. By default, the OUTLIER statement requests detection of additive outliers and level shifts, assuming that the series follows the estimated model.

```

/*--- ARIMA(0, 1, 1) Model ---*/
proc arima data=nile;
  identify var=level(1);
  estimate q=1 noint method=ml;
  outlier maxnum= 5 id=year;
run;

```

The outlier detection output is shown in Output 7.6.1.

### Output 7.6.1 ARIMA(0, 1, 1) Model

#### The ARIMA Procedure

Outlier Detection Summary	
Maximum number searched	5
Number found	5
Significance used	0.05

**Output 7.6.1** *continued*

Outlier Details					
Obs	Time ID	Type	Estimate	Chi-Square	Approx Prob>ChiSq
29	1899	Shift	-315.75346	13.13	0.0003
43	1913	Additive	-403.97105	11.83	0.0006
7	1877	Additive	-335.49351	7.69	0.0055
94	1964	Additive	305.03568	6.16	0.0131
18	1888	Additive	-287.81484	6.00	0.0143

Note that the first three outliers detected are indeed the ones discussed earlier. You can include the shock signatures that correspond to these three outliers in the Nile data set as follows:

```
data nile;
  set nile;
  AO1877 = ( year = '1jan1877'd );
  AO1913 = ( year = '1jan1913'd );
  LS1899 = ( year >= '1jan1899'd );
run;
```

Now you can refine the earlier model by including these outliers. After examining the parameter estimates and residuals (not shown) of the ARIMA(0,1,1) model with these regressors, the following stationary MA1 model (with regressors) appears to fit the data well:

```
/*--- MA1 Model with Outliers ---*/
proc arima data=nile;
  identify var=level
    crosscorr=( AO1877 AO1913 LS1899 );
  estimate q=1
    input=( AO1877 AO1913 LS1899 )
    method=ml;
  outlier maxnum=5 alpha=0.01 id=year;
run;
```

The relevant outlier detection process output is shown in [Output 7.6.2](#). No outliers, at significance level 0.01, were detected.

**Output 7.6.2** MA1 Model with Outliers**The ARIMA Procedure**

Outlier Detection Summary	
Maximum number searched	5
Number found	0
Significance used	0.01

**Example 7.7: Iterative Outlier Detection**

This example illustrates the iterative nature of the outlier detection process. This is done by using a simple test example where an additive outlier at observation number 50 and a level shift at observation number 100 are artificially introduced in the international airline passenger data used in [Example 7.2](#). The following

DATA step shows the modifications introduced in the data set:

```
data airline;
  set sashelp.air;
  logair = log(air);
  if _n_ = 50 then logair = logair - 0.25;
  if _n_ >= 100 then logair = logair + 0.5;
run;
```

In Example 7.2 the airline model, ARIMA(0, 1, 1)  $\times$  (0, 1, 1)<sub>12</sub>, was seen to be a good fit to the unmodified log-transformed airline passenger series. The preliminary identification steps (not shown) again suggest the airline model as a suitable initial model for the modified data. The following statements specify the airline model and request an outlier search.

```
/*--- Outlier Detection ---*/
proc arima data=airline;
  identify var=logair( 1, 12 )  noprint;
  estimate q= (1)(12) noint method= ml;
  outlier maxnum=3 alpha=0.01;
run;
```

The outlier detection output is shown in Output 7.7.1.

### Output 7.7.1 Initial Model

#### The ARIMA Procedure

Outlier Detection Summary		
Maximum number searched	3	
Number found	3	
Significance used	0.01	

Outlier Details				
Obs	Type	Estimate	Chi-Square	Approx Prob>ChiSq
100	Shift	0.49325	199.36	<.0001
50	Additive	-0.27508	104.78	<.0001
135	Additive	-0.10488	13.08	0.0003

Clearly the level shift at observation number 100 and the additive outlier at observation number 50 are the dominant outliers. Moreover, the corresponding regression coefficients seem to correctly estimate the size and sign of the change. You can augment the airline data with these two regressors, as follows:

```
data airline;
  set airline;
  if _n_ = 50 then AO = 1;
  else AO = 0.0;
  if _n_ >= 100 then LS = 1;
  else LS = 0.0;
run;
```

You can now refine the previous model by including these regressors, as follows. Note that the differencing order of the dependent series is matched to the differencing orders of the outlier regressors to get the correct “effective” outlier signatures.

```

/*--- Airline Model with Outliers ---*/
proc arima data=airline;
  identify var=logair(1, 12)
    crosscorr=( AO(1, 12) LS(1, 12) )
    noprint;
  estimate q= (1)(12) noint
    input=( AO LS )
    method=ml plot;
  outlier maxnum=3 alpha=0.01;
run;

```

The outlier detection results are shown in Output 7.7.2.

### Output 7.7.2 Airline Model with Outliers

#### The ARIMA Procedure

Outlier Detection Summary		
Maximum number searched	3	
Number found	3	
Significance used	0.01	

Outlier Details				
Obs	Type	Estimate	Chi-Square	Approx Prob>ChiSq
135	Additive	-0.10310	12.63	0.0004
62	Additive	-0.08872	12.33	0.0004
29	Additive	0.08686	11.66	0.0006

The output shows that a few outliers still remain to be accounted for and that the model could be refined further.

## References

- Akaike, H. (1974), “A New Look at the Statistical Model Identification,” *IEEE Transactions on Automatic Control*, AC-19, 716–723.
- Anderson, T. W. (1971), *The Statistical Analysis of Time Series*, New York: John Wiley & Sons.
- Andrews, D. F. and Herzberg, A. M. (1985), *A Collection of Problems from Many Fields for the Student and Research Worker*, New York: Springer-Verlag.
- Ansley, C. F. (1979), “An Algorithm for the Exact Likelihood of a Mixed Autoregressive–Moving Average Process,” *Biometrika*, 66, 59–65.
- Ansley, C. F. and Newbold, P. (1980), “Finite Sample Properties of Estimators for Autoregressive Moving-Average Models,” *Journal of Econometrics*, 13, 159–183.
- Bhansali, R. J. (1980), “Autoregressive and Window Estimates of the Inverse Correlation Function,” *Biometrika*, 67, 551–566.

- Box, G. E. P. and Jenkins, G. M. (1976), *Time Series Analysis: Forecasting and Control*, Rev. Edition, San Francisco: Holden-Day.
- Box, G. E. P., Jenkins, G. M., and Reinsel, G. C. (1994), *Time Series Analysis: Forecasting and Control*, 3rd Edition, Englewood Cliffs, NJ: Prentice-Hall.
- Box, G. E. P. and Tiao, G. C. (1975), "Intervention Analysis with Applications to Economic and Environmental Problems," *Journal of the American Statistical Association*, 70, 70–79.
- Brocklebank, J. C. and Dickey, D. A. (2003), *SAS for Forecasting Time Series*, 2nd Edition, Cary, NC: SAS Institute Inc.
- Brockwell, P. J. and Davis, R. A. (1991), *Time Series: Theory and Methods*, 2nd Edition, New York: Springer-Verlag.
- Chatfield, C. (1980), "Inverse Autocorrelations," *Journal of the Royal Statistical Society, Series A*, 142, 363–377.
- Choi, B. (1992), *ARMA Model Identification*, New York: Springer-Verlag.
- Cleveland, W. S. (1972), "The Inverse Autocorrelations of a Time Series and Their Applications," *Technometrics*, 14, 277.
- Cobb, G. W. (1978), "The Problem of the Nile: Conditional Solution to a Change Point Problem," *Biometrika*, 65, 243–251.
- Davidson, J. (1981), "Problems with the Estimation of Moving Average Models," *Journal of Econometrics*, 16, 295.
- Davies, N., Triggs, C. M., and Newbold, P. (1977), "Significance Levels of the Box-Pierce Portmanteau Statistic in Finite Samples," *Biometrika*, 64, 517–522.
- De Jong, P. and Penzer, J. (1998), "Diagnosing Shocks in Time Series," *Journal of the American Statistical Association*, 93, 796–806.
- Dickey, D. A. (1976), *Estimation and Testing of Nonstationary Time Series*, Ph.D. diss., Iowa State University.
- Dickey, D. A. and Fuller, W. A. (1979), "Distribution of the Estimators for Autoregressive Time Series with a Unit Root," *Journal of the American Statistical Association*, 74, 427–431.
- Dickey, D. A., Hasza, D. P., and Fuller, W. A. (1984), "Testing for Unit Roots in Seasonal Time Series," *Journal of the American Statistical Association*, 79, 355–367.
- Dunsmuir, W. (1984), "Large Sample Properties of Estimation in Time Series Observed at Unequally Spaced Times," in E. Parzen, ed., *Time Series Analysis of Irregularly Observed Data*, New York: Springer-Verlag.
- Findley, D. F., Monsell, B. C., Bell, W. R., Otto, M. C., and Chen, B. C. (1998), "New Capabilities and Methods of the X-12-ARIMA Seasonal Adjustment Program," *Journal of Business and Economic Statistics*, 16, 127–176.
- Fuller, W. A. (1976), *Introduction to Statistical Time Series*, New York: John Wiley & Sons.
- Hamilton, J. D. (1994), *Time Series Analysis*, Princeton, NJ: Princeton University Press.

- Hannan, E. J. and Rissanen, J. (1982), “Recursive Estimation of Mixed Autoregressive Moving Average Order,” *Biometrika*, 69, 81–94.
- Harvey, A. C. (1981), *Time Series Models*, New York: John Wiley & Sons.
- Jones, R. H. (1980), “Maximum Likelihood Fitting of ARMA Models to Time Series with Missing Observations,” *Technometrics*, 22, 389–396.
- Kohn, R. and Ansley, C. F. (1985), “Efficient Estimation and Prediction in Time Series Regression Models,” *Biometrika*, 72, 694–697.
- Ljung, G. M. and Box, G. E. P. (1978), “On a Measure of Lack of Fit in Time Series Models,” *Biometrika*, 65, 297–303.
- Montgomery, D. C. and Johnson, L. A. (1976), *Forecasting and Time Series Analysis*, New York: McGraw-Hill.
- Morf, M., Sidhu, G. S., and Kailath, T. (1974), “Some New Algorithms for Recursive Estimation on Constant Linear Discrete Time Systems,” *IEEE Transactions on Automatic Control*, 19, 315–323.
- Nelson, C. R. (1973), *Applied Time Series for Managerial Forecasting*, San Francisco: Holden-Day.
- Newbold, P. (1981), “Some Recent Developments in Time Series Analysis,” *International Statistical Review*, 49, 53–66.
- Newton, H. J. and Pagano, M. (1983), “The Finite Memory Prediction of Covariance Stationary Time Series,” *SIAM Journal on Scientific and Statistical Computing*, 4, 330–339.
- Pankratz, A. (1983), *Forecasting with Univariate Box-Jenkins Models: Concepts and Cases*, New York: John Wiley & Sons.
- Pankratz, A. (1991), *Forecasting with Dynamic Regression Models*, New York: John Wiley & Sons.
- Pearlman, J. G. (1980), “An Algorithm for the Exact Likelihood of a High-Order Autoregressive–Moving Average Process,” *Biometrika*, 67, 232–233.
- Priestley, M. B. (1981), *Spectral Analysis and Time Series*, London: Academic Press.
- Schwarz, G. (1978), “Estimating the Dimension of a Model,” *Annals of Statistics*, 6, 461–464.
- Stoffer, D. S. and Toloi, C. M. C. (1992), “A Note on the Ljung-Box-Pierce Portmanteau Statistic with Missing Data,” *Statistics and Probability Letters*, 13, 391–396.
- Tsay, R. S. and Tiao, G. C. (1984), “Consistent Estimates of Autoregressive Parameters and Extended Sample Autocorrelation Function for Stationary and Nonstationary ARMA Models,” *Journal of the American Statistical Association*, 79, 84–96.
- Tsay, R. S. and Tiao, G. C. (1985), “Use of Canonical Analysis in Time Series Model Identification,” *Biometrika*, 72, 299–315.
- Woodfield, T. J. (1987), “Time Series Intervention Analysis Using SAS Software,” in *Proceedings of the Twelfth Annual SAS Users Group International Conference*, Cary, NC: SAS Institute Inc.

# Subject Index

- additive model
  - ARIMA model, 206
- AIC, *see* Akaike information criterion
- Akaike information criterion
  - AIC, 242
  - ARIMA procedure, 242
- ARIMA model
  - additive model, 206
  - ARIMA procedure, 186
  - autoregressive integrated moving-average model, 186
  - Box-Jenkins model, 186
  - factored model, 206
  - multiplicative model, 206
  - notation for, 201
  - seasonal model, 206
  - subset model, 206
- ARIMA procedure
  - Akaike information criterion, 242
  - ARIMA model, 186
  - ARIMAX model, 186, 207
  - ARMA model, 186
  - autocorrelations, 189
  - autoregressive parameters, 247
  - BY groups, 220
  - conditional forecasts, 248
  - confidence limits, 248
  - correlation plots, 189
  - cross-correlation function, 233
  - data requirements, 214
  - differencing, 204, 239, 245
  - factored model, 206
  - finite memory forecasts, 249
  - forecasting, 248, 250
  - Gauss-Marquardt method, 241
  - ID variables, 250
  - infinite memory forecasts, 248
  - input series, 207
  - interaction effects, 211
  - intervention model, 207, 210, 212, 286
  - inverse autocorrelation function, 231
  - invertibility, 246
  - Iterative Outlier Detection, 294
  - log transformations, 249
  - Marquardt method, 241
  - Model Identification, 288
  - moving-average parameters, 247
  - naming model parameters, 247
  - ODS graph names, 265
  - ODS Graphics, 218
  - Outlier Detection, 293
  - output data sets, 252–254, 257, 258
  - output table names, 262
  - predicted values, 248
  - prewhitening, 238, 239
  - printed output, 259
  - rational transfer functions, 213
  - regression model with ARMA errors, 207, 208
  - residuals, 248
  - Schwarz Bayesian criterion, 242
  - seasonal model, 206
  - stationarity, 190
  - subset model, 206
  - syntax, 215
  - time intervals, 250
  - transfer function model, 207, 211, 242
  - unconditional forecasts, 249
- ARIMAX model
  - ARIMA procedure, 186, 207
- ARIMAX models and
  - design matrix, 211
- ARMA model
  - ARIMA procedure, 186
  - autoregressive moving-average model, 186
  - notation for, 202
  - augmented Dickey-Fuller tests, 224, 238
  - autocorrelations
    - ARIMA procedure, 189
    - plotting, 189
  - autoregressive integrated moving-average model, *see* ARIMA model
  - autoregressive moving-average model, *see* ARMA model
  - autoregressive parameters
    - ARIMA procedure, 247
- Box-Jenkins model, *see* ARIMA model
- BY groups
  - ARIMA procedure, 220
- conditional forecasts
  - ARIMA procedure, 248
- confidence limits
  - ARIMA procedure, 248
- correlation plots
  - ARIMA procedure, 189
- cross-correlation function

ARIMA procedure, 233

data requirements  
ARIMA procedure, 214

denominator factors  
transfer function model, 212

design matrix  
ARIMAX models and, 211

Dickey-Fuller tests, 224

differencing  
ARIMA procedure, 204, 239, 245

dynamic regression, 186, 207

ESACF (Extended Sample Autocorrelation Function method), 233

Extended Sample Autocorrelation Function (ESACF) method, 233

factored model  
ARIMA model, 206

ARIMA procedure, 206

finite memory forecasts  
ARIMA procedure, 249

forecasting  
ARIMA procedure, 248, 250

Gauss-Marquardt method  
ARIMA procedure, 241

ID variables  
ARIMA procedure, 250

impulse function  
intervention model and, 210

infinite memory forecasts  
ARIMA procedure, 248

input series  
ARIMA procedure, 207

interaction effects  
ARIMA procedure, 211

interrupted time series analysis, *see* intervention model

interrupted time series model, *see* intervention model

intervention analysis, *see* intervention model

intervention model  
ARIMA procedure, 207, 210, 212, 286

interrupted time series analysis, 210

interrupted time series model, 207

intervention analysis, 210

intervention model and  
impulse function, 210

step function, 211

inverse autocorrelation function  
ARIMA procedure, 231

invertibility  
ARIMA procedure, 246

Iterative Outlier Detection

ARIMA procedure, 294

log transformations  
ARIMA procedure, 249

Marquardt method  
ARIMA procedure, 241

MINIC (Minimum Information Criterion) method, 235

Minimum Information Criterion (MINIC) method, 235

Model Identification  
ARIMA procedure, 288

moving-average parameters  
ARIMA procedure, 247

multiplicative model  
ARIMA model, 206

naming model parameters  
ARIMA procedure, 247

nonstationarity, *see* stationarity

notation for  
ARIMA model, 201

ARMA model, 202

numerator factors  
transfer function model, 212

ODS graph names  
ARIMA procedure, 265

ODS Graphics  
ARIMA procedure, 218

of time series  
stationarity, 204

Outlier Detection  
ARIMA procedure, 293

output data sets  
ARIMA procedure, 252–254, 257, 258

output table names  
ARIMA procedure, 262

Phillips-Perron tests, 224

plotting  
autocorrelations, 189

predicted values  
ARIMA procedure, 248

prewhitening  
ARIMA procedure, 238, 239

printed output  
ARIMA procedure, 259

random-walk with drift tests, 224

rational transfer functions  
ARIMA procedure, 213

regression model with ARMA errors  
ARIMA procedure, 207, 208

residuals  
ARIMA procedure, 248

SBC, *see* Schwarz Bayesian criterion  
SCAN (Smallest Canonical) correlation method, 236  
Schwarz Bayesian criterion  
    ARIMA procedure, 242  
    SBC, 242  
seasonal model  
    ARIMA model, 206  
    ARIMA procedure, 206  
seasonal unit root test, 238  
Smallest Canonical (SCAN) correlation method, 236  
stationarity  
    ARIMA procedure, 190  
    nonstationarity, 190  
    of time series, 204  
stationarity tests, 223, 238  
step function  
    intervention model and, 211  
subset model  
    ARIMA model, 206  
    ARIMA procedure, 206  
time intervals  
    ARIMA procedure, 250  
transfer function model  
    ARIMA procedure, 207, 211, 242  
    denominator factors, 212  
    numerator factors, 212  
unconditional forecasts  
    ARIMA procedure, 249  
white noise test of the residuals, 226  
white noise test of the series, 224



# Syntax Index

- ALIGN= option
  - FORECAST statement (ARIMA), 230
- ALPHA= option
  - FORECAST statement (ARIMA), 230
  - IDENTIFY statement (ARIMA), 221
  - OUTLIER statement (ARIMA), 229
- ALTPARM option
  - ESTIMATE statement (ARIMA), 225, 245
- AR= option
  - ESTIMATE statement (ARIMA), 227
- ARIMA procedure, 215
  - syntax, 215
- ARIMA procedure, PROC ARIMA statement
  - PLOT option, 218
- BACK= option
  - FORECAST statement (ARIMA), 230
- BACKLIM= option
  - ESTIMATE statement (ARIMA), 228
- BY statement
  - ARIMA procedure, 220
- CENTER option
  - IDENTIFY statement (ARIMA), 221
- CLEAR option
  - IDENTIFY statement (ARIMA), 221
- CONVERGE= option
  - ESTIMATE statement (ARIMA), 228
- CROSSCORR= option
  - IDENTIFY statement (ARIMA), 221
- DATA= option
  - IDENTIFY statement (ARIMA), 222
  - PROC ARIMA statement, 218
- DELTA= option
  - ESTIMATE statement (ARIMA), 228
- ESACF option
  - IDENTIFY statement (ARIMA), 222
- ESTIMATE statement
  - ARIMA procedure, 225
- FORECAST statement
  - ARIMA procedure, 230
- GRID option
  - ESTIMATE statement (ARIMA), 228
- GRIDVAL= option
  - ESTIMATE statement (ARIMA), 228
- ID= option
  - FORECAST statement (ARIMA), 230
  - OUTLIER statement (ARIMA), 229
- IDENTIFY statement
  - ARIMA procedure, 221, 229
- INITVAL= option
  - ESTIMATE statement (ARIMA), 227
- INPUT= option
  - ESTIMATE statement (ARIMA), 225, 244
- INTERVAL= option
  - FORECAST statement (ARIMA), 231, 250
- LEAD= option
  - FORECAST statement (ARIMA), 231
- MA= option
  - ESTIMATE statement (ARIMA), 227
- MAXIT= option
  - ESTIMATE statement (ARIMA), 228
- MAXITER= option
  - ESTIMATE statement (ARIMA), 228
- MAXNUM= option
  - OUTLIER statement (ARIMA), 229
- MAXPCT= option
  - OUTLIER statement (ARIMA), 229
- METHOD= option
  - ESTIMATE statement (ARIMA), 225
- MINIC option
  - IDENTIFY statement (ARIMA), 222
- MU= option
  - ESTIMATE statement (ARIMA), 227
- NLAG= option
  - IDENTIFY statement (ARIMA), 222
- NOCONSTANT option
  - ESTIMATE statement (ARIMA), 225
- NODF option
  - ESTIMATE statement (ARIMA), 225
- NOEST option
  - ESTIMATE statement (ARIMA), 227
- NOINT option
  - ESTIMATE statement (ARIMA), 225
- NOLS option
  - ESTIMATE statement (ARIMA), 228
- NOMISS option
  - IDENTIFY statement (ARIMA), 223
- NOOUTALL option
  - FORECAST statement (ARIMA), 231
- NOPRINT option

ESTIMATE statement (ARIMA), 226  
FORECAST statement (ARIMA), 231  
IDENTIFY statement (ARIMA), 223  
NOSTABLE option  
    ESTIMATE statement (ARIMA), 228  
NOTFSTABLE option  
    ESTIMATE statement (ARIMA), 228

OUT= option  
    FORECAST statement (ARIMA), 231, 252  
    PROC ARIMA statement, 220

OUTCORR option  
    ESTIMATE statement (ARIMA), 226

OUTCOV option  
    ESTIMATE statement (ARIMA), 227

OUTCOV= option  
    IDENTIFY statement (ARIMA), 223, 253

OUTTEST= option  
    ESTIMATE statement (ARIMA), 226, 254

OUTMODEL= option  
    ESTIMATE statement (ARIMA), 227, 257

OUTSTAT= option  
    ESTIMATE statement (ARIMA), 227, 258

P= option  
    ESTIMATE statement (ARIMA), 226  
    IDENTIFY statement (ARIMA), 223

PERROR= option  
    IDENTIFY statement (ARIMA), 223

PLOT option  
    ESTIMATE statement (ARIMA), 226  
    PROC ARIMA statement, 218

PLOTS option  
    PROC ARIMA statement, 218

PRINTALL option  
    ESTIMATE statement (ARIMA), 228  
    FORECAST statement (ARIMA), 231

PROC ARIMA statement, 218

Q= option  
    ESTIMATE statement (ARIMA), 226  
    IDENTIFY statement (ARIMA), 223

SCAN option  
    IDENTIFY statement (ARIMA), 223

SIGMA= option  
    OUTLIER statement (ARIMA), 229

SIGSQ= option  
    FORECAST statement (ARIMA), 231

SINGULAR= option  
    ESTIMATE statement (ARIMA), 228

STATIONARITY= option  
    IDENTIFY statement (ARIMA), 223, 224

TYPE= option

OUTLIER statement (ARIMA), 229

VAR= option  
    IDENTIFY statement (ARIMA), 224

WHITENOISE= option  
    ESTIMATE statement (ARIMA), 226  
    IDENTIFY statement (ARIMA), 224