Hannah Roach

3/29/2019

## Lab Exercise 10: Time Series Analysis with ARIMA

| Purpose: | This lab is designed to investigate and practice Time Series Analysis with ARIMA models (Box-Jenkins-methodology). After completing the tasks in this lab you should able to:<br><br>• Use R functions for ARIMA models<br>• Apply the requirements for generating appropriate training data<br>• Validate the effectiveness of the ARIMA models |
|---|---|

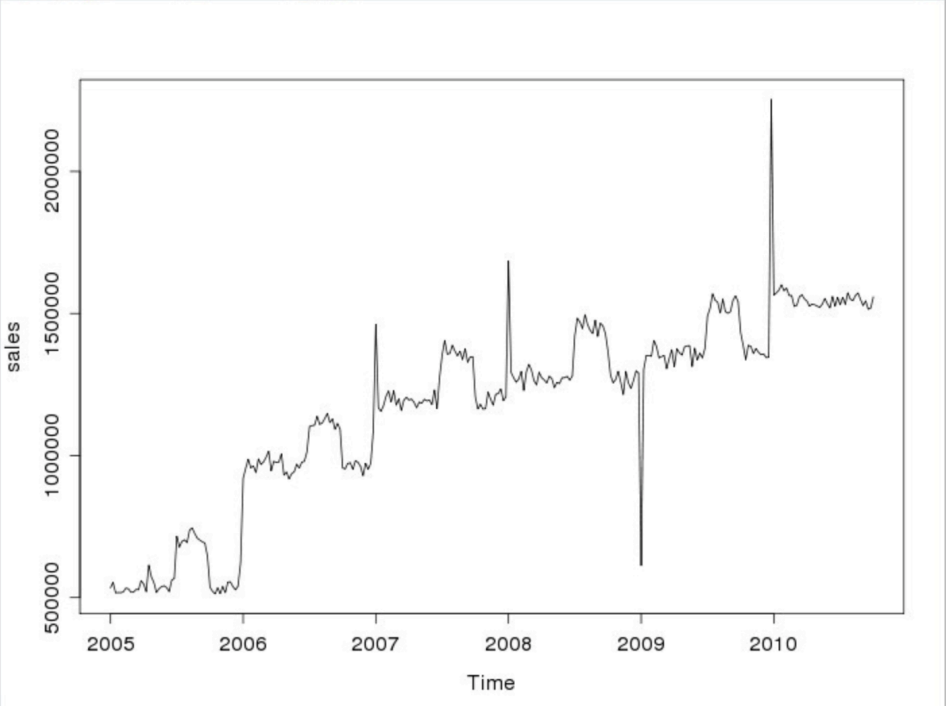| Tasks: | Tasks you will complete in this lab include:<br><br>• Use the R –Studio environment to code ARIMA models<br>• Use the ODBC connection to the database to create the weekly sales data from the retail database<br>• Prepare the data (sorting and rendering the data as a Time series)<br>• Generate a model and evaluate how well it predicts the results and compare the results with original data |
|---|---|

_

# Workflow Overview

1. • Set the Working Directory

2. • Establish the ODBC Connection

3. • Open Connections to ODBC Database

4. • Get Data from the Database

5. • Review, Update, and Prepare DataFrame "msales" File for ARIMA  Modeling

6. • Convert "sales" into Time Series Type Data

7. • Plot the Time Series

8. • Analyze the ACF and PACF

9. • Difference the Data to Make it Stationary

10. • Plot ACF and PACF for the Differenced Data

11. • Fit the ARIMA Model

12. • Generate Predictions

13. • Compare predicted values with actual values

# LAB Instructions

| Step | Action |
|------|--------|
| 1 | Log in with GPADMIN credentials on to R-Studio. |
| 2 | **Set the working directory**<br>Set working directory to ~/LAB10/ , execute the command:<br>`setwd("~/LAB10")`<br><br>&bull; (Or using the "Tools" option in the tool bar in the RStudio environment.) |
| 3 | **Establish the ODBC Connection:**<br><br>Load the RODBC package using the following command:<br>`library('RODBC')`<br><br>`Minimize your screen` |
| 4 | **Open Connections to ODBC Database:**<br><br>1. Before connecting to the ODBC database make sure the file, /etc/odbc.ini, is properly set to point to database "training1".<br>2. If not, edit the line that starts with, "Database =" within the file /etc/odbc.ini, to point to, "training1".<br><br>OPEN PUTTY and connect to your BE IP<br>LOGIN: GPADMIN<br>Type cd/etc/<br>Type nano odbc.ini<br>  If Database does not say training1  changed to training1<br>Then hit CTRL X type y to save and enter<br>Type exit to close PUTTY<br><br>**OPEN your minimized SAFARI and continue**<br><br>3. Ensure the username(uid) and password (pwd) are provided correctly in the following command:<br><br>`ch <- odbcConnect("Greenplum",uid="gpadmin",`<br>`     case="postgresql",pwd="changeme")` |
| | |

| Step | Action |
|------|--------|
| 5 | **Get Data from the Database:**<br><br>1. Drop the table, weekly_sales, from the schema ddemo:<br><br>`sqlDrop(ch,"ddemo.weekly_sales")`<br>2. Execute an SQL query using the sqlQuery command, creating a table, weekly_sales, in which the weekly sales are grouped first by year and the week within a year:<br><br><pre>> sqlQuery(ch,<br>"CREATE TABLE<br>    ddemo.weekly_sales (<br>      sale INTEGER,<br>      Y1 INTEGER,<br>      W1 INTEGER )<br>      DISTRIBUTED BY (sale);<br>INSERT INTO ddemo.weekly_sales<br>SELECT<br>    SUM((o.item_price*o.item_quantity)) as sale ,<br>    EXTRACT(YEAR FROM o.order_datetime) as y1,<br>    CASE<br>      WHEN (EXTRACT(WEEK FROM o.order_datetime) = 53) THEN 52<br>      ELSE<br>      EXTRACT(WEEK FROM o.order_datetime)<br>    END w1<br>FROM<br>    ddemo.order_lineitems o<br>group<br>    by y1,w1<br>;<br>"<br>)</pre><br>- **Note:** Note the use of the **EXTRACT** function to obtain the year and the week within the year from the order_datetime field<br>- The sales number is accumulated for each week<br>  - The ISO Standard for numbering weeks within a year may lead to a year containing 52 or 53 weeks.  In order to work with Time Series data you need a consistent "periodicity" with the data.  You must accumulate the vales for week 53 with that of week 52 in the same year. Case statement is used to accumulate the week 53 (if exists) along with week 52 data<br>1. Get the results from the table into data frame msales. Execute the command:<br><br>`msales <- (sqlFetch(ch,"ddemo.weekly_sales"))`<br><br>4. Close the ODBC channel.<br>`odbcClose(ch)` |

| Step | Action |
|------|--------|
| 6 | **Review, Update and Prepare DataFrame "msales" File for ARIMA Modeling:**<br>1. Sort the data in the order of Year and Week:<br>   Use the R function "order" :<br><br>```<br>attach(msales)<br>msales <- msales[order(y1,w1),]<br>detach(msales)<br>```<br><br>2. Extract 300 values from "asales" for modeling and 12 values to compare with the predictions done by the model. Store them in two different vectors: "sales" and "csales". Use the following command:<br><br>```<br>>sales <- c(rep(0,300))<br>>csales <- c(rep(0,12))<br>>sales[1:300] <- msales[1:300, 1]<br>>csales[1:12] <- msales[301:312, 1]<br>``` |
| 7 | **Convert "sales" into Time Series Type Data:**<br><br>Convert the "sales" into a time series.<br>• This "transformation" is required for most of the time–series functions, since a time series contains more information than the values themself, namely information about dates and frequencies at which the time series has been recorded.<br><br>```<br>> sales <- ts(sales,start=2005,frequency=52)<br>``` |

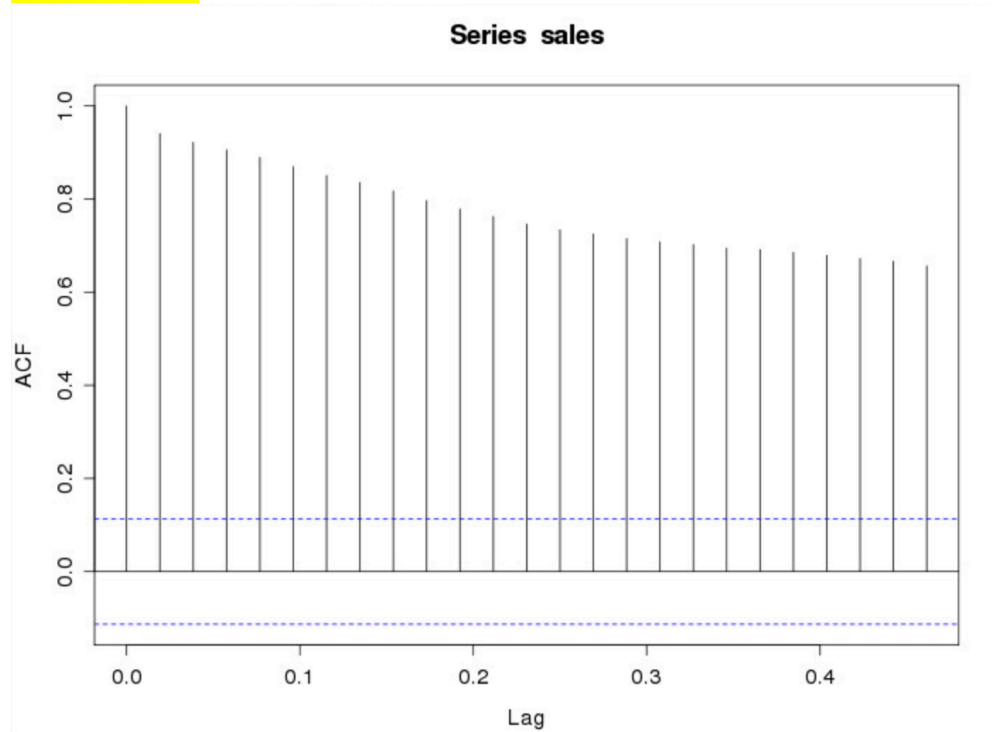| Step | Action |
|------|--------|
| 8 | **Plot the Time Series:**<br><br>1. Plot the Time Series using the following command:<br><br>`plot(sales,type="l")`<br>***SCREENSHOT<br><br><br><br>2. Review the plot of the Time Series.<br>3. Identify the seasonality features in the graph. – It appears that once a year there is a dip spike in sales from June to August and another spike in December.<br>4. Is the data Seasonal (Do you see patterns that repeat at a particular frequency)? Yes<br>5. Is the data stationary? The data does not appear to be stationary because there is a linear trend in the data and the data appears to be seasonal.<br>6. Is there a trend to the data? Yes, there is an upward trend in the data. |

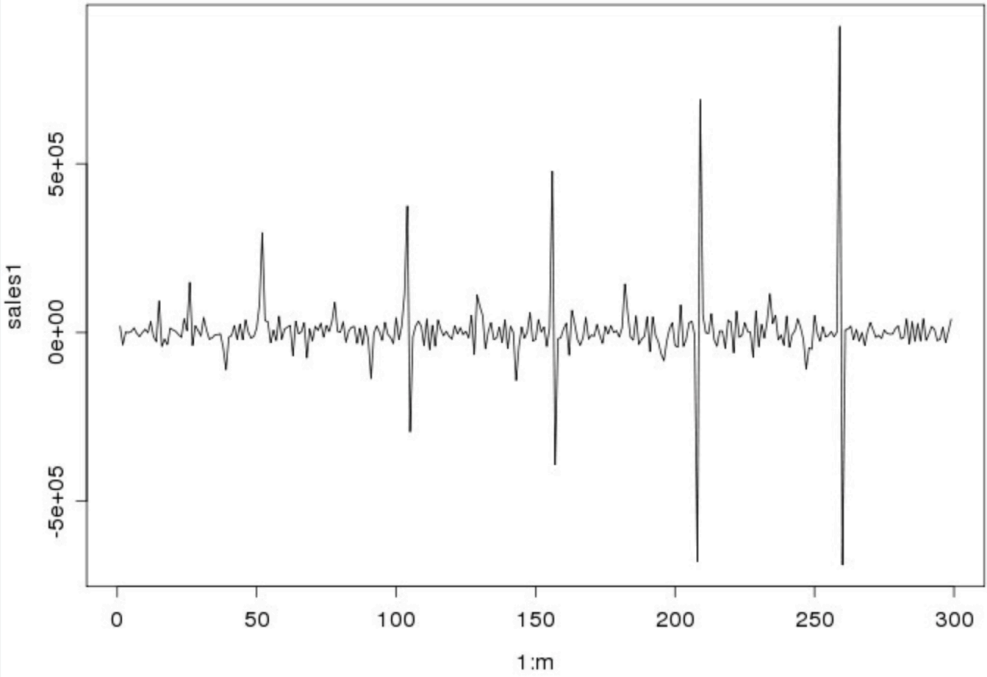| 9 | **Analyze the ACF and PACF :** |
|---|---|
| | The next step in analyzing time series is to examine the **autocorrelations (ACF)** and **partial autocorrelations (PACF).** R provides the functions acf( ) and pacf( ) for computing and plotting of ACF and PACF. |

1.  Use the **parameter function "par"** to set the plot window to display both the ACF and PACF plots. Plot ACF and PACF on the same graph using the command:
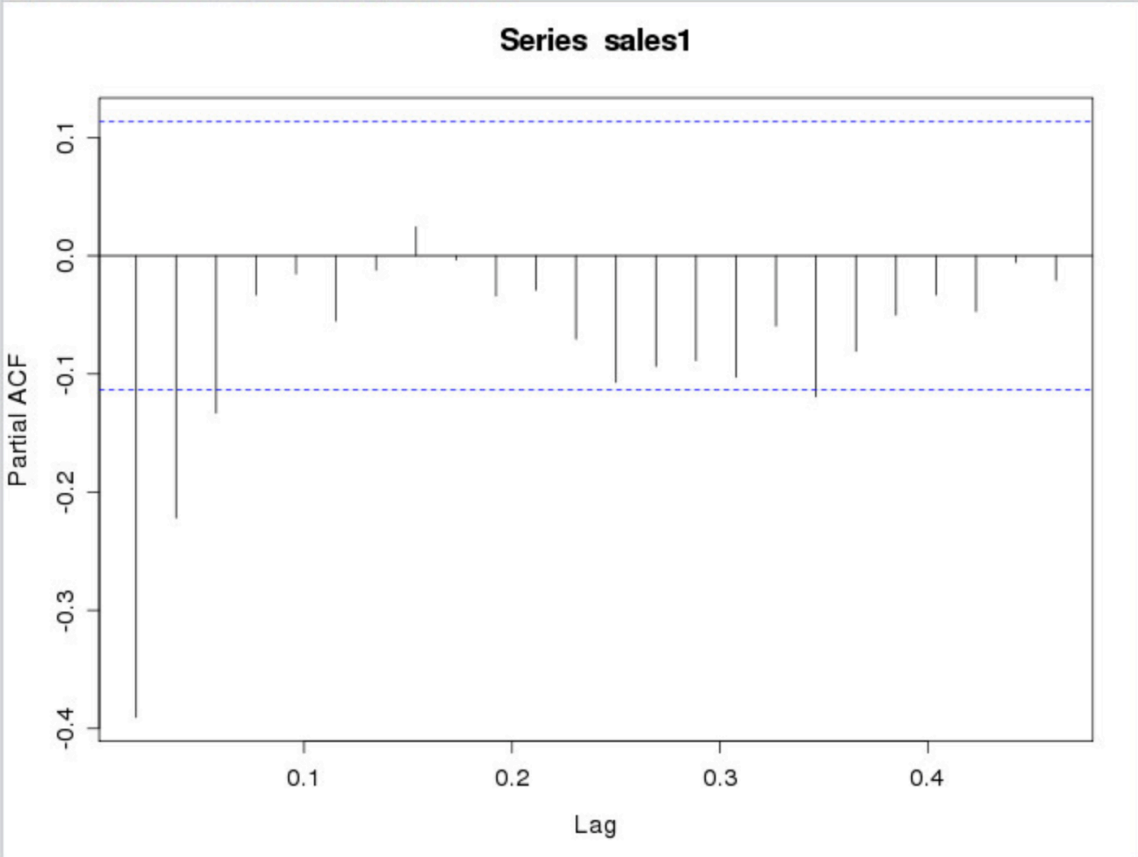
```
> par(mfrow=c(2,1))
> acf(sales)
> pacf(sales)
```
**SCREENSHOT**



Series sales

| Step | Action |
|------|--------|
|  | 

Series sales

2. Using the plot generated in the plot window:
   - <mark>Does the ACF tail off quickly? Shows that the AR model handles seasonality in the data after 0.6. ACF 0.2, 0.4, and 0.6 have significance.</mark>
   - <mark>What does the ACF indicate with respect to stationarity of the data? It shows that the data is not stationary.</mark>
   - <mark>What will you do to make the data stationary? Take the difference of the linear trend.</mark> |

| Step | Action |
|------|--------|
| 10 | **Difference the Data to Make it Stationary:**<br><br>To difference the data and make it stationary you should use the "diff" function in R. The "diff" function takes the pair of each observation and differences it from the one previous to it.<br><br>1. Run the following code:<br>```r
> #Difference the series and plot it
> sales1 <- diff(sales)
> m <- length(sales1)
> par(mfrow=c(1,1))
> plot(1:m,sales1,type="l")
```<br>***SCREENSHOT<br><br><br><br>1.  Do you see any trend to the data now? No<br>2.  Is seasonality still shown in the data? Yes, in the y-axis.<br>3.  How does the trend of oscillating spikes change as you move from left to right on this plot? It increases |

| Step | Action |
|------|--------|
| 11 | **Plot ACF and PACF for the Differenced Data:**<br><br>1. Run the following code:<br><br>```<br>> #Plot ACF and PACF on the same graph<br>> par(mfrow=c(2,1))<br>> acf(sales1)<br>> pacf(sales1)<br>```<br>***SCREENSHOT<br><br><br><br>2. Do you see the ACF tailing off quickly? Yes |

| 12 | **Fit the ARIMA Model:**

Once you configure the data and review the seasonality elements, you are ready to fit an ARIMA model. Selecting the correct model to fit an ARIMA model is a bit of an art. Algorithms are used to select the correct parameters.

- Use the data configuration and the principle of "parsimony" to fit a basic model
- Use the statement ARIMA
- Use the time series
- Use "sales" and not the "diff(sales)" that we computed in step 10. ARIMA will automatically invoke the "diff" function based on the parameters specified.
- You need to specify the order (p,d,q) where "p" is the order of AR , "q" order of MA and "d" is the number of differences.
- Fit (1,1,0) so the model will difference once, use AR and MA parameters as 1 and 0.
- Use a seasonal statement since the data seems to have a seasonal component to it as you see spikes every 52 weeks.
- In the seasonal statement, you have a "list" where you put in the "order" and a (time) "period" which will be "52".
- Use "include.mean = false" - R will force a mean and continue any trend seen in the past for the model. This automatically defaults to "True". You turn the automatic default off with this statement.
- **Note:** You can also experiment with not using this statement during this lab.

1. Type in the following code:
```
sales.fit <- arima (sales,
                    order=c(1,1,0),
                    seasonal = list(order=c(1,1,0),period=52),
                    include.mean=FALSE)
> sales.fit
**SCREENSHOT
```

```
> sales.fit

Call:
arima(x = sales, order = c(1, 1, 0), seasonal = list(order = c(1, 1, 0), period = 52),
      include.mean = FALSE)

Coefficients:
         ar1      sar1
      -0.3769   -0.3770
s.e.   0.0608    0.0657

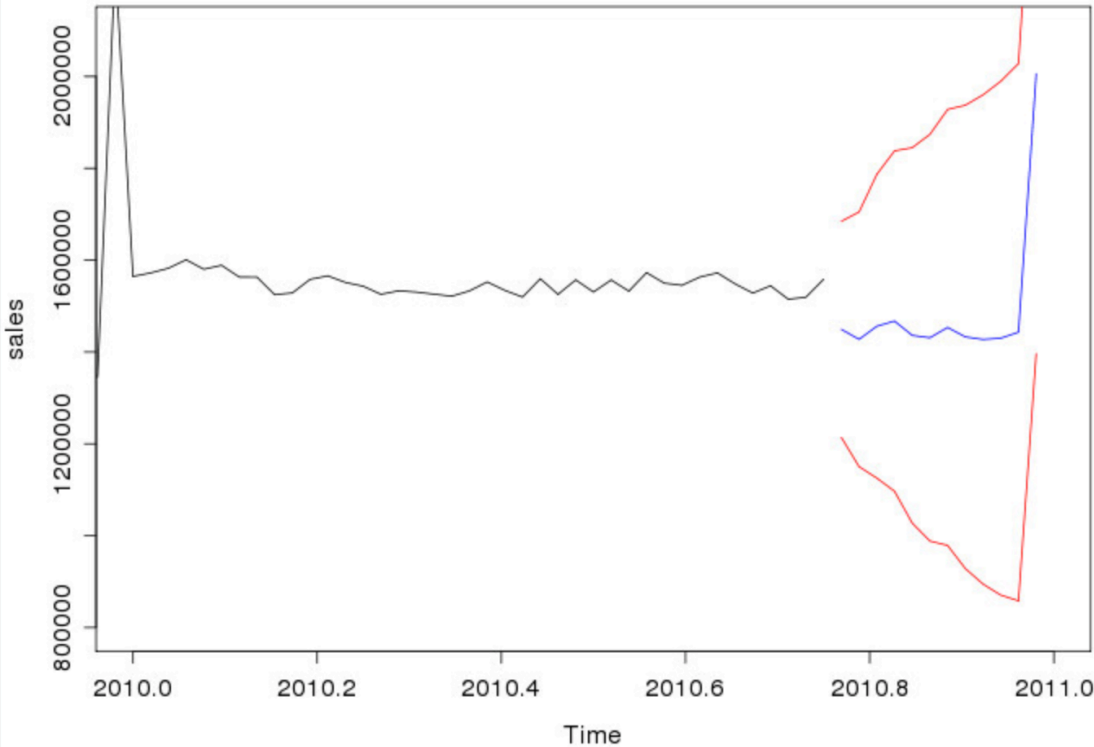sigma^2 estimated as 1.384e+10:  log likelihood = -3238.4,  aic = 6482.8
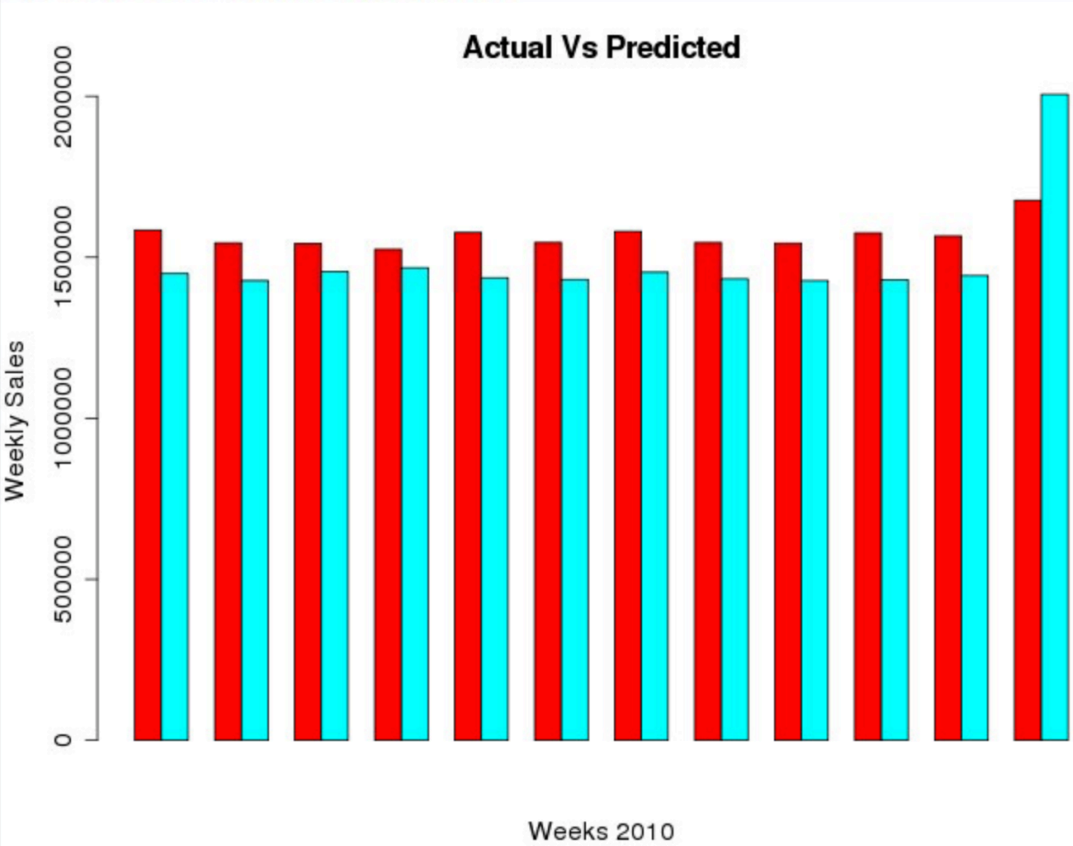```

2. Review the output displayed.
3. Record the coefficients for the AR term and the seasonal AR terms and the standard errors. AR term: -0.3769, seasonal AR terms: -0.3770, standard errors: 0.0608, 0.0657
4. What is your observation on the standard errors compared to the coefficients? The standard errors are positive whereas the coefficients are negative. The standard error is much smaller than coefficients.

- **Note:** The model gives you the "log likelihoods" that provide important input on model selection. |

Wait, the table structure — left column shows "12". Let me keep it simple.

| Step | Action |
|------|--------|
|      |        |

| Step | Action |
|------|--------|
| 13 | **Generate Predictions:**<br><br>1.  Use the fitted model "sales.fit" and the "predict" statement for 12  periods ahead:<br><br>`sales.predict <- predict (sales.fit, n.ahead=12)`<br><br>• **Note:** You can see the predictions by typing "sales.fit".<br><br>2.  Plot the predictions using plot statements:<br><br>```
> par(mfrow=c(1,1))
> plot (sales,xlim=c(2009.50,2011))
> lines (sales.predict$pred,col="blue")
> lines (sales.predict$pred+2*sales.predict$se,col="red")
> lines (sales.predict$pred-2*sales.predict$se,col="red")
```<br>**SCREENSHOT**<br><br><br><br>• **Note:** The first plot statement plots the original "sales" data<br>• EMC has restricted the "xlim" value to zoom in on the values just prior to the predicted values<br>• The "blue" line indicates the mean of the predictions<br>• The "red" lines denote the upper and lower bounds of the prediction |

| Step | Action |
|------|--------|
| 14 | **Compare predicted values with actual values**<br><br>1.  Compare the predicted values with the actual values and plot the actual Values Predicted as a barplot.  Use the following code:<br><br>```
> #comparing the predictions with the actual values
> par(mfrow=c(1,1))
> x <- c(rep(0,24))
> x[1:12] <- msales[1:12]
> x[13:24] <- as.numeric(sales.predict$pred)
> forbar <- matrix(x,ncol=12,byrow=TRUE)
> colnames(forbar)<- asales[1:12,3]
> rownames(forbar)<- c("Actual","Predicted")
> barplot(forbar,beside=TRUE,
+            main="Actual Vs Predicted",
+            ylab="Weekly Sales",
+            col=rainbow(2),
+            xlab="Weeks 2010")
```<br><br>2.  **SCREENSHOT**<br><br> |

*End of Lab Exercise*