

W271 Live Session 9: ARIMA and SARIMA

Devesh Tiwari

7/11/2017

Main topics covered in Week 9

- Mixed Autoregressive Moving Average (ARMA) Models
 - Mathematical formulation and derivation of key properties
 - Comparing ARMA models and AR models using simulated series
 - Comparing ARMA models and AR models using an example
- An introduction to non-stationary time series model
- Random walk and integrated processes
- Autoregressive Integrated Moving Average (ARIMA) Models
 - Review the steps to build ARIMA time series model
 - Simulation
 - Modeling with simulated data using the Box-Jenkins approach
 - Estimation, model diagnostics, model identification, model selection, assumption testing, and statistical inference / forecasting, backtesting
- Seasonal ARIMA (SARIMA) Models
 - Mathematical formulation
 - An empirical example
- Putting everything together: ARIMA modeling

Readings

CM2009: Paul S.P. Cowpertwait and Andrew V. Metcalfe. *Introductory Time Series with R*. Springer. 2009.

- Ch. 4.3 – 4.7, 6, 7.1 – 7.3

SS2016: Robert H. Shumway and David S. Stoffer. *Time Series Analysis and Applications*. EZ Edition with R Examples:

- 3.7 – 3.10, review 3.1 – 3.6

HA: Rob J Hyndman and George Athanasopoulos. *Forecasting: Principles and Practice*:

- Ch. 8.5 – 8.9

Agenda for this week's live session:

1. Questions from last week
2. Recap
3. Breakout Session 1: ARIMA nomenclature
4. Breakout Session 2: SARIMA EDA

5. SARIMA modeling (in-class discussion / take-home)

Recap and overview

1. Last week, we were introduced to autoregressive (AR), moving average (MA), and autoregressive moving average (ARMA) models. These models are only appropriate for time-series that are weakly stationary (stationary in the mean and the variance).
2. We often are confronted with time-series that is not stationary in the mean and variance. Luckily for us, we can transform these series in order to make them stationary!
3. Here is an incomplete list of how/why time-series might not be stationary in the mean:
 - a. Data has a trend
 - b. Data contains a unit-root (next week)
 - c. Data contains seasonal elements
4. We can take care of these problems either by detrending the data or by differencing the data. Once the data are transformed into a weakly stationary series, we can model the resulting series with an ARMA model. We call these models ARIMA models if the data do not exhibit any seasonality. If the data are seasonal, then these models are called SARIMA models.
5. Remember, here are the steps to building an ARIMA model!
 - i. Conduct an EDA to determine if you need to transform the data in order to make it stationary.
 - ii. Transform the data if needed.
 - iii. Estimate several Arima(p,d,q) models. Remember, you set the value of d in the first step! So really, you are trying to find the appropriate values of p and q.
 - iv. Evaluate the residuals of models with the lowest AIC/BIC values and simpler models. Select the model where the residuals resemble white noise.
 - v. If you still have some candidate models remaining, then conduct an out of sample test and select the model with the lowest forecasting error.
 - vi. Answer your question / generate forecasts!

Breakout Session 1: Review Questions

1. Consider the following time-series process:

$$x_t = t + \alpha_1 x_{t-1} + \alpha_2 x_{t-2} + w_t$$

where t is time.

Hint:

$$x_{t-1} = t - 1 + \alpha_1 x_{t-2} + \alpha_2 x_{t-3} + w_{t-1}$$

and

$$\mathbf{B}^n(x_t) = x_{t-n}$$

and

$$\nabla = (1 - \mathbf{B})$$

- (a) What does this equation mean, in words?
- (b) Take the first difference of this equation and then re-write it using the backshift operator. What kind of model is it? What R commands would you use to model this?

- (c) How would you know if this model is stationary?
2. Consider a time series represented by the following equation:

$$(1 - \alpha_1 \mathbf{B} - \alpha_2 \mathbf{B}^2)(1 - \mathbf{B})x_t = \omega_t(1 + \beta_1 \mathbf{B})$$

- (a) Describe this time series using the ARIMA(p,d,q) notation.
- (b) Can you tell if the time series that is generated by this model is stationary? If so, is it? If you cannot determine whether it is stationary, what additional information do you need?

Seasonality and EDA For SARIMA models

Many time-series data exhibit a clear seasonal pattern. Consider data that are measured every month (such as monthly sales or even unemployment). There are 12 months in a year and it is very common to see a strong relationship for a month across several years. In other words, sales in January are often correlated with sales in the prior Januarys etc. In addition, sales in January are also likely to be correlated with prior months' sales as well. We would say that such data are seasonal with a period of 12 and we can directly model the seasonal effect within the ARIMA framework.

You can detect seasonality during the EDA process. Sometimes, the seasonality is “obvious” by simply plotting the time-series. In addition to the plot, you can check to see if there is an especially high spike on the ACF and PACF charts at m ($m = 12$ if you have monthly data, $m = 52$ if you have weekly data), and you can use the monthplot (which you will do shortly).

Once you have identified seasonality, you need to model it. However, within the ARIMA framework, you need to make sure that your data are stationary. We know that we can difference the time-series to make it stationary in the mean, but we can also take the seasonal difference to make it stationary in the mean with respect to a seasonal effect also. After you have made the data stationary (by setting the d and D values), examine *that* data: plot it and look at the ACF/PACF charts and make your initial guesses about p, q, P , and Q .

Breakout session 2 Seasonality and EDA

Load the following three datasets and conduct an EDA on them. For each, come prepared to describe what you see in the initial plot of the data, what steps you took to make the data stationary in the mean, and what your initial ARIMA function would look like in R (`Arima(x, order = c(p,d,q), seasonal = list(order = c(P,D,Q)), method = “ML”)`).

```
rm(list = ls())
library(knitr)
opts_chunk$set(tidy.opts=list(width.cutoff=60),tidy=TRUE)

library(forecast)
```

```
## Loading required package: zoo
##
## Attaching package: 'zoo'
##
## The following objects are masked from 'package:base':
##
##      as.Date, as.Date.numeric
## Loading required package: timeDate
```

```
## This is forecast 7.3
library(astsa)

##
## Attaching package: 'astsa'
## The following object is masked from 'package:forecast':
##
##      gas
path <- "~/Documents/Projects/MIDS/"
setwd(path)
df <- read.csv("Summer 2016/materials for ISVC/Lab 3/correlate-flight_prices.csv")

# We can actually analyze the relative seach activity for many phrases!
df <- df[,1:2] #focus on the phrase flight prices
str(df)

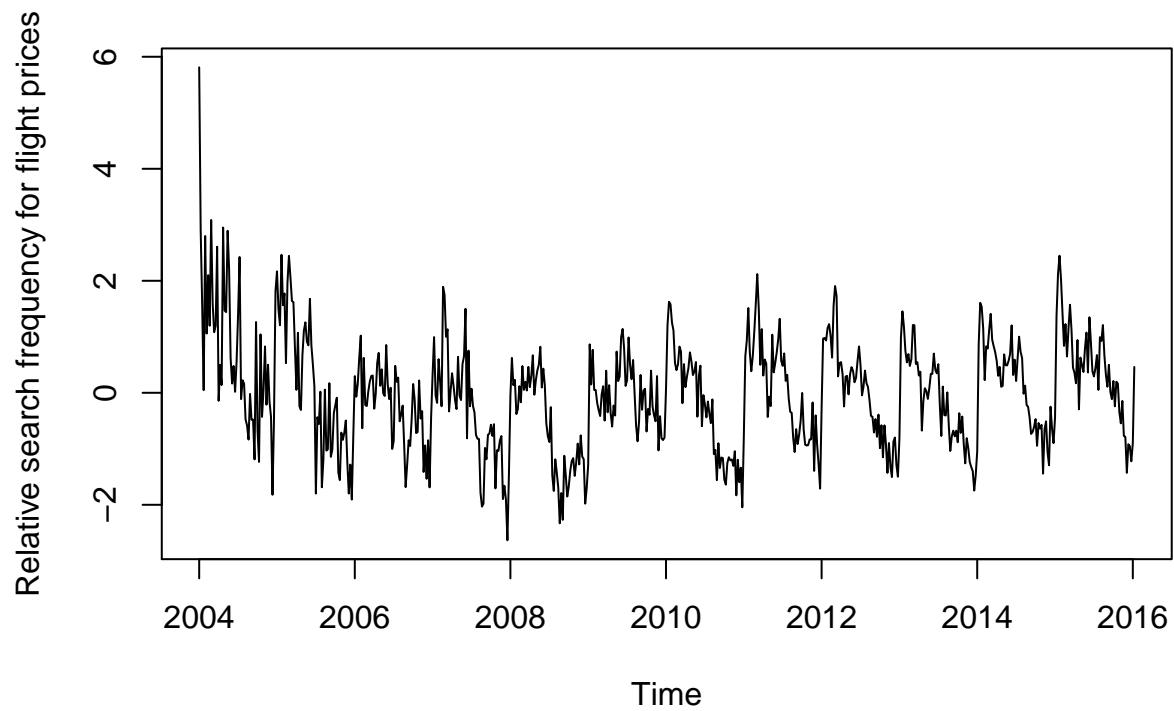
## 'data.frame':    626 obs. of  2 variables:
## $ Date          : Factor w/ 626 levels "1/1/06","1/1/12",...: 43 4 16 30 211 256 222 235 247 307 ...
## $ flight.prices: num  5.811 2.901 1.575 0.047 2.8 ...
cbind(head(df), tail(df))

##      Date flight.prices      Date flight.prices
## 1  1/4/04      5.811 11/22/15      -1.430
## 2 1/11/04      2.901 11/29/15      -0.918
## 3 1/18/04      1.575  12/6/15      -0.952
## 4 1/25/04      0.047 12/13/15      -1.224
## 5  2/1/04      2.800 12/20/15      -0.897
## 6  2/8/04      1.058 12/27/15       0.462

summary(df)

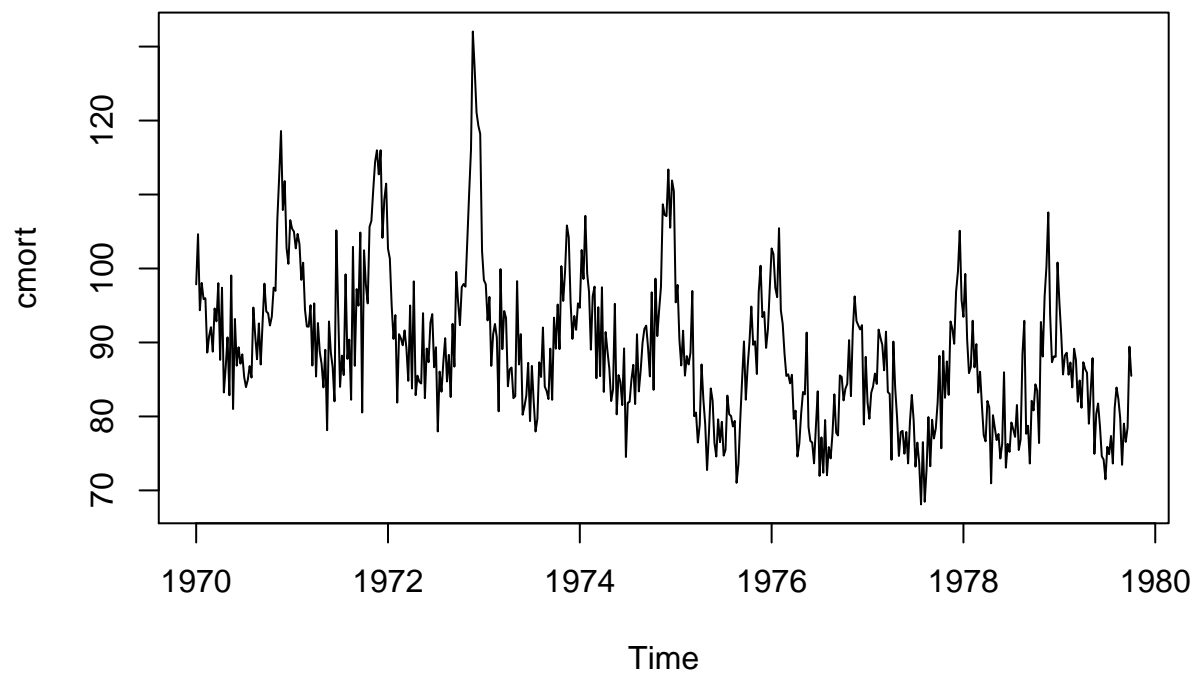
##      Date      flight.prices
## 1/1/06 : 1   Min.      :-2.6330
## 1/1/12 : 1   1st Qu.: -0.7480
## 1/10/10: 1   Median : 0.0175
## 1/11/04: 1   Mean      :-0.0123
## 1/11/09: 1   3rd Qu.: 0.5727
## 1/11/15: 1   Max.      : 5.8110
## (Other):620

# Create a time-series object
fp <- ts(df$flight.prices, frequency = 52, start = c(2004,1))
plot(fp, ylab = "Relative search frequency for flight prices")
```



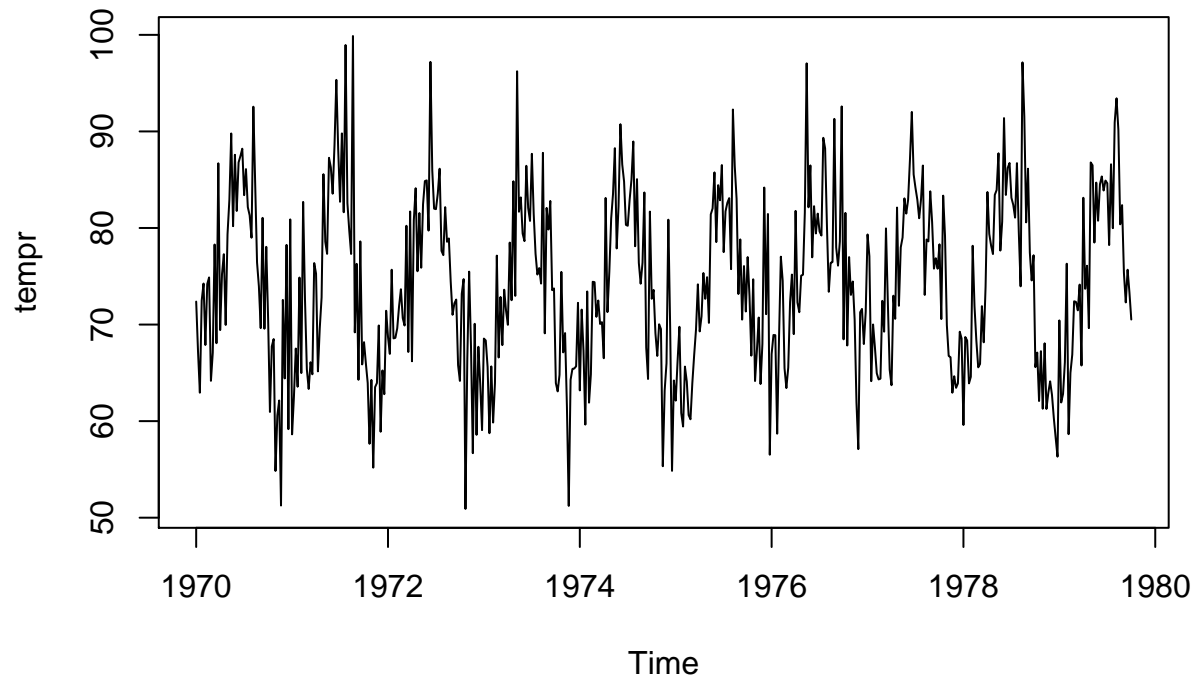
```
### INSERT OTHER EDA
```

```
plot(cmort)
```



```
### INSERT OTHER EDA
```

```
plot(tempr)
```



```
### INSERT OTHER EDA
```

Putting it all together: SARIMA Flight Prices

In this live session, we are going to build a SARIMA model on the relative search activity for the phrase, “flight prices.” These data are provided by google correlate and they are weekly. For the sake of simplicity, we will focus on 2010 onward.

Remember that we can express a SARIMA model as: $SARIMA(p,d,q) \times (P,D,Q)_m$.

Step 1: Conduct EDA and determine values for d and D

Step 2: Find values for P and Q (seasonal components)

Step 3: Find values for p and q

Step 4: Fine - tuning (play around with different values of p , q , P , and Q)

Step 5: Plot fitted values against time-series

Step 6: Conduct out of sample tests

Step 7: Select model and forecast

Note: Sometimes you might not be sure about what transformations you should make. In that case, DO BOTH and compare their forecasting accuracy.

```
fp.training <- window(fp, start = c(2005, 1), end = c(2014, 12))
fp.test <- window(fp, start = c(2015, 1))

plot(fp.training)
```

