

ENV 790.30 - Time Series Analysis for Energy Data | Spring 2021

Assignment 7 - Due date 04/07/21

Student Name

Directions

You should open the .rmd file corresponding to this assignment on RStudio. The file is available on our class repository on Github. And to do so you will need to fork our repository and link it to your RStudio.

Once you have the project open the first thing you will do is change “Student Name” on line 3 with your name. Then you will start working through the assignment by **creating code and output** that answer each question. Be sure to use this assignment document. Your report should contain the answer to each question and any plots/tables you obtained (when applicable).

When you have completed the assignment, **Knit** the text and code into a single PDF file. Rename the pdf file such that it includes your first and last name (e.g., “LuanaLima_TSA_A07_Sp21.Rmd”). Submit this pdf using Sakai.

Set up

Some packages needed for this assignment: `forecast`, `tseries`, `smooth`. Do not forget to load them before running your script, since they are NOT default packages.

```
#Load/install required package here
```

```
library(forecast)
```

```
## Registered S3 method overwritten by 'quantmod':
```

```
##   method             from
```

```
##   as.zoo.data.frame zoo
```

```
library(tseries)
```

```
library(Kendall)
```

```
library(readxl)
```

```
library(dplyr)
```

```
##
```

```
## Attaching package: 'dplyr'
```

```
## The following objects are masked from 'package:stats':
```

```
##
```

```
##   filter, lag
```

```
## The following objects are masked from 'package:base':
```

```
##
```

```
##   intersect, setdiff, setequal, union
```

```
library(lubridate)
```

```
##
```

```
## Attaching package: 'lubridate'
```

```
## The following objects are masked from 'package:base':
##
##   date, intersect, setdiff, union

library(ggplot2)
library(tidyverse)

## -- Attaching packages ----- tidyverse 1.3.0 --

## v tibble  3.0.6      v purrr  0.3.4
## v tidyr   1.1.2      v stringr 1.4.0
## v readr   1.4.0      v forcats 0.5.1

## -- Conflicts ----- tidyverse_conflicts() --
## x lubridate::as.difftime() masks base::as.difftime()
## x lubridate::date()        masks base::date()
## x dplyr::filter()          masks stats::filter()
## x lubridate::intersect()   masks base::intersect()
## x dplyr::lag()              masks stats::lag()
## x lubridate::setdiff()     masks base::setdiff()
## x lubridate::union()       masks base::union()

library(readr)
```

Importing and processing the data set

Consider the data from the file “inflowtimeseries.txt”. The data corresponds to the monthly inflow in m^3/s for some hydro power plants in Brazil. You will only use the last column of the data set which represents one hydro plant in the Amazon river basin. The data span the period from January 1931 to August 2011 and is provided by the Brazilian ISO.

For all parts of the assignment prepare the data set such that the model consider only the data from January 2000 up to December 2009. Leave the year 2010 of data (January 2010 to December 2010) for the out-of-sample analysis. Do **NOT** use data from 2010 and 2011 for model fitting. You will only use it to compute forecast accuracy of your model.

Part I: Preparing the data sets

Q1

Read the file into a data frame. Prepare your time series data vector such that observations start in January 2000 and end in December 2009. Make you sure you specify the **start=** and **frequency=** arguments. Plot the time series over time, ACF and PACF.

```
setwd('/Users/rajatkhandelwal/Documents/GitHub/ENV790_30_TSA_S2021/Data')
data <- read_table("inflowtimeseries.txt")
```

```
## Warning: Duplicated column names deduplicated: '452' => '452_1' [15]

##
## -- Column specification -----
## cols(
##   Jan = col_character(),
##   `1931` = col_double(),
##   `4782` = col_double(),
##   `4076` = col_double(),
##   `2518` = col_double(),
##   `2450` = col_double(),
```

```
## `2649` = col_double(),
## `1462` = col_double(),
## `450` = col_double(),
## `968` = col_double(),
## `246` = col_double(),
## `2636` = col_double(),
## `452` = col_double(),
## `4870` = col_double(),
## `452_1` = col_double(),
## `17342` = col_double(),
## `31270` = col_double()
## )
```

Q2

Using the *decompose()* or *stl()* and the *seasadj()* functions create a series without the seasonal component, i.e., a deseasonalized inflow series. Plot the deseasonalized series and original series together using ggplot, make sure your plot includes a legend. Plot ACF and PACF for the deseasonalized series. Compare with the plots obtained in Q1.

Part II: Forecasting with ARIMA models and its variations

Q3

Fit a non-seasonal ARIMA(p, d, q) model using the *auto.arima()* function to the non-seasonal data. Forecast 12 months ahead of time using the *forecast()* function. Plot your forecasting results and further include on the plot the last year of non-seasonal data to compare with forecasted values (similar to the plot on the lesson file for M10).

Q4

Put the seasonality back on your forecasted values and compare with the original seasonal data values. *Hint* : One way to do it is by summing the last year of the seasonal component from your *decompose* object to the forecasted series.

Q5

Repeat Q3 for the original data, but now fit a seasonal ARIMA(p, d, q) $x(P, D, Q)_{12}$ also using the *auto.arima()*.

Q6

Compare the plots from Q4 and Q5 using the *autoplot()* function.

Part III: Forecasting with Other Models

Q7

Fit an exponential smooth model to the original time series using the function *es()* from package *smooth*. Note that this function automatically do the forecast. Do not forget to set the arguments: *silent=FALSE* and *holdout=FALSE*, so that the plot is produced and the forecast is for the year of 2010.

Q8

Fit a state space model to the original time series using the function *StructTS()* from package *stats*. Which one of the tree model we learned should you try: “local”, “trend”, or “BSM”. Why? Play with argument *fixed* a bit to try to understand how the different variances can affect the model. If you can’t seem to find a variance that leads to a good fit here is a hint: try *fixed = c(0.1, 0.001, NA, NA)*. Since *StructTS()* fits

a state space model to the data, you need to use *forecast()* to generate the forecasts. Like you do for the ARIMA fit.

Part IV: Checking Forecast Accuracy

Q9

Make one plot with the complete original seasonal historical data (Jan 2000 to Dec 2010). Now add the forecasts from each of the developed models in parts Q4, Q5, Q7 and Q8. You can do it using the *autoplot()* combined with *autolayer()*. If everything is correct in terms of time line, the forecasted lines should appear only in the final year. If you decide to use *ggplot()* you will need to create a data frame with all the series will need to plot. Remember to use a different color for each model and add a legend in the end to tell which forecast lines corresponds to each model.

Q10

From the plot in Q9 which model or model(s) are leading to the better forecasts? Explain your answer. Hint: Think about which models are doing a better job forecasting the high and low inflow months for example.

Q11

Now compute the following forecast metrics we learned in class: RMSE and MAPE, for all the models you plotted in part Q9. You can do this by hand since you have forecasted and observed values for the year of 2010. Or you can use R function *accuracy()* from package “forecast” to do it. Build a table with the results and highlight the model with the lowest MAPE. Does the lowest MAPE corresponds match your answer for part Q10?