

**Name:** Your name here

**Due:** 2024/11/06

# Day 16 - Lab: Regression harmonics, GLS, and state-space models

## Introduction

In this assignment, you will simulate a time series using some of the regression techniques discussed in class. Then, you will fit an OLS, GLS, and state-space representation of the time series and compare and contrast the results. For this assignment, assume that we are simulating daily measurements over a 25 week period.

---

### **i** Regression harmonics

Sometimes, it can be difficult to put your finger on exactly what is causing a seasonal effect. In such case, it is useful to model the seasonal component as a sum of harmonic functions, like **sin** and **cos**. The **harmonic seasonal model** for time series  $\{x_t\}$  with frequency  $f$  is defined by:

$$x_t = m_t + \sum_{i=1}^{\lfloor f/2 \rfloor} (s_i \sin(2\pi i t / f) + c_i \cos(2\pi i t / f)) + z_t$$

where  $\lfloor f/2 \rfloor$  denotes the integer component of  $f/2$  (sometimes called the floor function).

For example, a harmonic seasonal model for daily observations over the course of weeks might look like:

$$x_t = \beta_0 + \beta_1 t + \beta_2 \sin\left(\frac{2\pi t}{7}\right) + \beta_3 \cos\left(\frac{2\pi t}{7}\right) + \beta_4 \sin\left(\frac{4\pi t}{7}\right) + \beta_5 \cos\left(\frac{4\pi t}{7}\right) + \beta_6 \sin\left(\frac{6\pi t}{7}\right) + \beta_7 \cos\left(\frac{6\pi t}{7}\right) + z_t$$

where  $z_t$  is the error series. By writing the model in this way, we can use regression to estimate all model parameters. In the following lab, we will simulate a time series using this method, then compare and contrast parameter estimates from OLS, GLS, and state-space models.

## Simulating a time series

1. [3 pt] Let us start by generating the mean function for our series. Simulate the mean function for daily measurements from a 25-week long series with a linear trend and a harmonic seasonal component. You may choose the values used to generate the mean function, but the resulting function should have clear linear trend and seasonality. Note that your seasonal effects should have 3 harmonic cycles. Plot the mean function.
2. [2 pt] Simulate a (stationary) AR(2) process to add to the mean function created in question 1. You may choose the values of the parameters  $\alpha_1$ ,  $\alpha_2$ , and  $\sigma$ , but ensure that the simulated error series is stationary, zero-mean, and exhibits meaningful serial autocorrelation (i.e. do not set the  $\alpha$ 's to be super small). Plot the error series.
3. [2 pt] Combine the mean function with the error series to create a complete time series and plot the result.
4. [1 pt] Finally, create two data sets: `obs_df` and `forecast_df`, comprised of the first 140 and last 35 time points, respectively. Henceforth, the `obs_df` data set will represent the data that we observe, and the `forecast_df` data set will represent unobserved values that we hope to forecast.

## OLS

5. [2 pt] Fit an ordinary least squares regression model to the `obs_df` time series, including `t` and all 6 harmonic terms, and create an PACF plot of the residuals. Comment on what you see.

### **i** Confidence vs prediction intervals

When forecasting future values, we make the distinction between estimating the *mean response* and predicting a future *individual* response; the former is called a confidence interval and the latter is called a prediction interval. The estimated value remains the same between the two intervals, but the standard error of the value changes (which affects the width of the confidence interval).

As an example, consider the `AirPassengers` data. When forecasting one time point ahead, a confidence interval would provide a range of plausible values for the average number of air passengers on all planes one time point ahead, while a prediction interval would provide a range of plausible values for the number of air passengers on a single plane. There is much more uncertainty in the latter than the former.

The formulas to calculate confidence intervals and predictions intervals for regression models are quite complicated. Fortunately, the `predict` function in R will do it for you by specifying the `interval` type in the call to `predict`.

6. [2 pt] Forecast the values of the series in the `forecast_df` data frame, and plot the observed series, fitted series, and forecasted series on a single plot, including 95% **confidence** intervals for the forecast. (see `?predict.lm` for examples)
7. [2 pt] Do the same thing again, but instead including a **prediction** interval. Comment on which intervals are wider: the confidence or the prediction.

**GLS**

8. [2 pt] Fit a GLS model with an AR(2) correlation structure to the `obs_df` time series and create a pacf plot of the normalized residuals. What are the estimated values of  $\alpha_1$  and  $\alpha_2$ ? Does the GLS model solve the issues with autocorrelation in the residuals?

No more questions on GLS in the lab - tune in for more in the homework. :)

## State-space

9. [2 pt] Fit a state-space model with an AR(2) correlation structure to `obs_df` time series and create a pacf plot of the residuals. What are the estimated values of  $\alpha_1$  and  $\alpha_2$ ? Does the SS model solve the issues with autocorrelation in the residuals?

One of the main advantages of the state-space representation of the GLS model is that we are able to obtain standard errors for predictions. Sadly, the `predict` function cannot generate confidence intervals and prediction intervals for forecasts from an `arima`, as there is no notion of confidence vs prediction with the state-space approach. Instead, we can generate an approximate 95% prediction interval by using the standard errors and t-distribution.

10. [4 pt] Forecast the last 35 observations and construct approximate 95% prediction intervals for the forecasted values. Plot the observed, fitted, and forecasted values on a single plot. To create the intervals, you should use a multiplier obtained from `qt(.975, 140-8)` (the degrees of freedom are determined by the sample size minus the number of estimated regression coefficients).
11. [4 pt] The state-space representation of the model does not provide t-statistics and p-values by default. However, we can calculate approximate tests using the same t-distribution with 132 degrees of freedom. Create a matrix that returns the estimated regression coefficients, standard errors, t-statistics, and two-sided p-values for the estimates from the state-space model.