

# Problem Set 3

**Issued:** Wednesday, October 13, 2021

**Part 1 Due:** Friday, Oct 22nd, 11:59 PM EST   **Part 2 Due:** Monday, Nov 1st, 11:59 PM EST

---

## Part 1

### Problem 3.1: The Mauna Loa CO<sub>2</sub> concentration

In 1958, Charles David Keeling (1928-2005) from the Scripps Institution of Oceanography began recording carbon dioxide (CO<sub>2</sub>) concentrations in the atmosphere at an observatory located at about 3,400 m altitude on the Mauna Loa Volcano on Hawaii Island. The location was chosen because it is not influenced by changing CO<sub>2</sub> levels due to the local vegetation and because prevailing wind patterns on this tropical island tend to bring well-mixed air to the site. While the recordings are made near a volcano (which tends to produce CO<sub>2</sub>), wind patterns tend to blow the volcanic CO<sub>2</sub> away from the recording site. Air samples are taken several times a day, and concentrations have been observed using the same measuring method for over 60 years. In addition, samples are stored in flasks and periodically reanalyzed for calibration purposes. The observational study is now run by Ralph Keeling, Charles's son. The result is a data set with very few interruptions and very few inhomogeneities. It has been called the “most important data set in modern climate research.”

The data set for this problem can be found in `C02Data.csv`. It provides the concentration of CO<sub>2</sub> recorded at Mauna Loa for each month starting March 1958. More description is provided in the data set file. We will be considering only the CO<sub>2</sub> concentration given in **column 5**. The goal of the problem is to fit the data and understand its variations. You will encounter missing data points, part of the exercise is to deal with them appropriately.

Let  $C_i$  be the average CO<sub>2</sub> concentration in month  $i$  ( $i = 1, 2, \dots$ , counting from March 1958). We will look for a description of the form:

$$C_i = F(t_i) + P_i + R_i$$

where:

- $F : t \mapsto F(t)$  accounts for the long-term trend.
- $t_i$  is time at the middle of month  $i$ , measured in fractions of years after Jan 1, 1958.
- $P_i$  is periodic in  $i$  with period 12, accounting for the seasonal pattern.
- $R_i$  is the remaining residual that accounts for all other influences.

The decomposition is meaningful only if the range of  $F$  is much larger than the amplitude of the  $P_i$  and this amplitude in turn is substantially larger than that of  $R_i$ .

- Fit the data to a linear model  $F_1(t) \sim \alpha_1 + \alpha_2 t$ . Plot the data and the fit. What are the values of  $\hat{\alpha}_1$  and  $\hat{\alpha}_2$ ? Plot the residual error and comment about it. Report an appropriate goodness-of-fit metric of your choice and provide a justification for it.
- Fit the data to a quadratic model  $F_2(t) \sim \beta_1 + \beta_2 t + \beta_3 t^2$ . Plot the data and the fit. What are the values of  $\hat{\beta}_1$ ,  $\hat{\beta}_2$  and  $\hat{\beta}_3$ ? Similarly, plot the residual error and comment about it. Report an appropriate goodness-of-fit metric of your choice and provide a justification for it.

- (c) Fit the data to a quartic model  $F_3(t) \sim \gamma_1 + \gamma_2 t + \gamma_3 t^2 + \gamma_4 t^3 + \gamma_5 t^4$ . Which fit ( $F_1$ ,  $F_2$ ,  $F_3$ ) is better at capturing the trend in the data? Explain. What are the pros and cons of using a higher-order model? How would you go about selecting the order of your model in general? Describe your proposed approach.
- (d) Consider  $F_2(t)$ . We will now extract the periodic component which appears in the data. Average the residual  $C_i - F_2(t_i)$  over each month  $i$ . Namely, collect all the data for Jan (resp. Feb, Mar, etc) and average them to get one data point for Jan (resp. Feb, Mar, etc). The collection of those points can be interpolated to form a periodic signal  $P_i$ . Plot  $P_i$  over time. Compare what you get after adjusting for seasonality with the column in the data that is seasonally adjusted (column 6).
- (e) Plot the fit  $F_2(t_i) + P_i$  for each month  $i$  since 1958, i.e., over the overall time period considered in this problem. What can we conclude about the variation of CO<sub>2</sub> concentration since 1958? For instance, you can comment and provide potential reasons for the trend and seasonality observed.

## Part 2

### Problem 3.2: BPP Data Analysis

The goal of this problem is to analyze the PriceStats data from the MIT Billion Prices Project, provided by Professor Rigobon. The file `PriceStats_CPI.csv` also contains a column “CPI”. CPI stands for consumer price index. It refers to the price of a “market basket of consumer goods and services”, which is a proxy for inflation. It is released monthly by the Bureau of Labor Statistics. The file `T10YIE.csv` provides the break-even rate (BER), i.e., the difference in yield between a fixed rate and inflation adjusted 10 year treasury note. Of note, the BER and CPI data are captured during the same time period. This difference can be interpreted as what the market views the inflation rate will be for the next 10 years, on average. The file `T5YIE.csv` lists the corresponding measure for the 5-year rather than the 10-year timescale. There are 122 months of data in `PriceStats_CPI.csv`. In the following questions, you may want to either work on a log scale or operate on inflation rates. Of note, CPI is an index and not a rate. If you would like to work with rates, make sure you appropriately transform the raw CPI data into rates. In part (a), you can use data from all prior months to make one-month ahead forecasts. For example, to predict the CPI in May 2015, you can use all the data before May 2015 (but not May 2015). For parts (b), (c), and (d), you should perform all of your model fitting on the months prior to September 2013, and use the remaining months for evaluation. For your results, report the mean squared prediction error (MSPE) for one-month ahead forecasts starting September 2013.

- (a) First, we will try to predict the monthly CPI without using the BER or PriceStats data. Fit an AR model to the CPI data. For this, take the first CPI value of each month as that month’s CPI. You may or may not want to work in log scale to make the model comparable to models you fit in part (c) and report the MSPE for one-month ahead forecasts. Which order model gives the best predictions? **Hint:** one way to determine the proper order to use is to examine the auto-correlation (ACF) and partial auto-correlation functions (PACF) of the residuals. Start with a single AR term and add other terms as necessary. Provide an intuitive explanation for why an AR model may be a more suitable starting point for our problem than other types of models (e.g., such as a MA model).

- (b) How might you calculate monthly inflation rates from the CPI data and your one-month ahead predictions? How about from PriceStats data? And from BER data, both 5-year and 10-year? In the latter scenario, what dates would you use? Or would you use an average of many dates? Overlay your estimates of monthly inflation rates over time (months from September 2013 onwards).
- (c) Next, we will include external regressors to try to improve the predictions. Include as external regressors monthly average PriceStats data and BER data to fit a new AR model to the CPI. Report your MSPE. Try instead using PriceStats data and BER data from the first day of each month as your external regressors. Fit another AR model. Which model performs better in prediction? Explain the difference in the resulting models and error that arises from using an average value rather than a single value. **Hint:** Again, in order to match the units of your predictors and responses, you want to either work on a log scale or work with inflation rates. Please justify your choice regarding the type of values you decide to work with.
- (d) Try to improve your model from part (c). What is the smallest MSPE you can obtain? You might consider including MA terms, adding a seasonal AR term, or adding multiple daily values (or values from different months) of PriceStats and BER data as external regressors.
- (e) Consider the MA(1) model,  $X_t = W_t + \theta W_{t-1}$ , where  $\{W_t\} \sim WN(0, \sigma^2)$ . Find the autocovariance function of  $\{X_t\}$ . Together with the mean and variance of the MA(1) model, what does it imply?
- (f) Consider the AR(1) model,  $X_t = \phi X_{t-1} + W_t$ , where  $\{W_t\} \sim WN(0, \sigma^2)$ . Suppose  $|\phi| < 1$ . Find the autocovariance function of  $\{X_t\}$ . Note that you may use, without proving it, the fact that  $\{X_t\}$  is stationary if  $|\phi| < 1$ . What order MA model is the AR(1) model equivalent to?