

# Homework 2: Fun with Trends and Detrending

36-467, Fall 2020

Due at 6 pm on Thursday, 17 September 2020

AGENDA: Getting a feel for how linear smoothers work; understand properties of detrended data.

1. *What do splines like?* Refer back to the Kyoto cherry blossom data from Homework 1.

The function `smooth.spline` fits a smooth curve through the data points, balancing smoothness against coming close to the data. The command `lines(smooth.spline(x, y))` adds the fitted curve to the current plot.

- (a) (4) Plot the data (as in the solutions to HW 1). Use `smooth.spline()` to find the trend of flowering day vs. year, and then add it to the plot. Comment on the plot. How is it similar to the plots with the moving averages, and how does it differ? *Hint:* `smooth.spline()` is picky about NAs, and it may help to create a copy of the data without an missing values.
- (b) (4) The fitted values of the spline are stored in the `$y` component of its return value. Create a vector of residuals which stores, for each year, the difference between the actual day of flowering and the fitted value from the spline. Plot this against year. Describe the patterns in the plot, if any. *Hint:* Be careful about NAs when calculating the residuals.
- (c) (1) The value of  $\lambda$  (the penalty on curvature) selected by cross-validation is stored in the `$lambda` element of `kyoto.spline`. What is it?
- (d) (4) `smooth.spline` unfortunately does not store the smoother matrix `w`, though it does store its diagonal in the `$lev` component. We can recover the whole of `w` re-fitting the spline on artificial data. The following code will do it:

```
smoother.matrix <- function(a.spline, x) {  
  n <- length(x)  
  w <- matrix(0, nrow=n, ncol=n)  
  for (i in 1:n) {  
    y <- rep_len(0, n) # Equivalent to rep(0, length.out=n) but faster
```

```

        y[i] <- 1
        w[,i] <- fitted(smooth.spline(x, y, lambda=a.spline$lambda))
    }
    return(w)
}

```

(This code is also in a file on the class homepage.)

Run this function on the spline you fit from the data, and a suitable choice  $\mathbf{x}$ . Check that it isn't doing something obvious wrong by checking the dimensions of the resulting matrix (what should they be?), and by checking that the diagonal of this matrix matches `kyoto.spline$lev`. (For that latter, the `all.equal()` function is helpful.) In this problem, be explicit about the code you are using.

- (e) (3) Fix any matrix  $\mathbf{w}$ , and let  $\mathbf{e}_j$  be the  $n \times 1$  matrix which is 0 everywhere, except in row  $j$ , where it is 1. Explain why  $\mathbf{w}\mathbf{e}_j$  gives the  $j^{\text{th}}$  column of  $\mathbf{w}$ . Explain how this relates to the code in the previous problem.
  - (f) (4) Make a plot of all of the eigenvalues (not vectors) of the smoother matrix, in order. How many of them are  $> 0.95$ ? How many are  $> 0.1$ ? How many are  $> 0.01$ ? Are any of them exactly zero?
  - (g) (4) Following the example in the notes for lectures 3 and 4, make a plot of the first eight eigenvectors of the smoother matrix. Describe the kinds of patterns in the data these eigenvectors capture.
  - (h) (4) Make a similar plot for the last eight eigenvectors. Describe the kinds of patterns they capture.
  - (i) (5) What sorts of patterns will show up in the fitted values of the spline (= the estimate of the trend)? What sorts of patterns will show up in its residuals (= the estimate of the fluctuations)?
  - (j) (2) The example in the notes for the previous problem plots the entries in the eigenvectors against their position ("index") in the vector. Here, different positions correspond to different calendar years. Redo the plot from problem 1g so the eigenvectors are plotted against the `Year.AD` variable. Does this change your description of any of the patterns?
  - (k) (3) Use the `contour` function to make a contour plot of the smoother matrix. (You will probably want to adjust some of the default settings of `contour`.) Explain how this plot relates to the idea that the spline is doing a kind of local averaging.
2. *Degrees of freedom and moving averages* Assume, for simplicity, no missing values in the data. Also assume that we are only interested in points which are in the "interior" of the region where we gathered data, so they have  $k$  neighbors, and there are  $n$  such points.

- (a) (3) Suppose we have one-dimensional data, and estimate the trend by averaging each observation with its  $k$  nearest neighbors. How many degrees of freedom does the smoother matrix  $\mathbf{w}$  have?
  - (b) (2) Assume we have two-dimensional data, laid out in a regular grid, and estimate the trend by averaging each observation with its  $k$  nearest neighbors in every direction. How many degrees of freedom does  $\mathbf{w}$  have?
  - (c) (1) Assume we have four-dimensional data (3D space + time), with measurements taken at irregular points and times, and we estimate the trend by averaging each observation with its  $k$  nearest neighbors in 4D. How many degrees of freedom does  $\mathbf{w}$  have?
3. *The Yule-Slutsky effect*<sup>1</sup> In this problem, assume that  $X(t) = \mu(t) + \epsilon(t)$ , that  $\hat{\mu}(t) = \frac{1}{3} \sum_{i=t-1}^{t+1} X_i$ , and that  $\hat{\epsilon}(t) = X(t) - \hat{\mu}(t)$ .
- (a) (5) Find an expression for  $\mathbb{E}[\hat{\epsilon}(t)]$  in terms of the  $\mu$ 's.
  - (b) (5) Find an expression for  $\text{Var}[\hat{\epsilon}(t)]$  in terms of the variances and covariances of the  $\epsilon$ 's (and possibly the  $\mu$ 's). *Hint:* Remember that for any random variables  $U, V, W$  and constants  $a, b$ ,  $\text{Cov}[aU + bV, W] = a\text{Cov}[U, W] + b\text{Cov}[V, W]$ , and that  $\text{Cov}[U, U] = \text{Var}[U]$ .
  - (c) (5) Find an expression for  $\text{Cov}[\hat{\epsilon}(t), \hat{\epsilon}(t+1)]$ .
  - (d) (5) Now further assume that  $\text{Var}[\epsilon(t)] = \sigma^2$ , and that  $\text{Cov}[\epsilon(t), \epsilon(s)] = 0$  when  $t \neq s$ . Find  $\text{Cov}[\hat{\epsilon}(t), \hat{\epsilon}(t+1)]$ . Explain why the de-trended residuals are correlated, even though the true fluctuations are not.
4. *Detrending by differencing* We have focused on detrending by smoothing, where we first estimate the trend, and then subtract it off. An alternative procedure is to remove trends by taking *differences* between nearby observations. This problem explores how this works, in a situation where we have one coordinate, so our data can be written  $X(t) = \mu(t) + \epsilon(t)$ . Assume  $t$  is discrete, so it can be  $1, 2, \dots$ , and that there are no missing values.
- Define  $\Delta(t)$  as  $X(t) - X(t-1)$ .
- (a) (5) Write  $\Delta(t)$  in terms of the  $\mu$ 's and  $\epsilon$ 's.
  - (b) (5) Find the expectation and variance of  $\Delta(t)$  in terms of  $\mu$ , and the variance and covariance of  $\epsilon$ .
  - (c) (5) Explain why  $\Delta(t)$  can be said to be “detrended”, if  $\mu$  changes slowly.
  - (d) (5) Find an expression for  $\text{Cov}[\Delta(t), \Delta(t+1)]$ .

---

<sup>1</sup>The fact that fitted values [= estimated trend] and residuals [= de-trended values] are correlated, even when there aren't correlations in the original data, was independently discovered in the 1920s by the British statistician G. Udny Yule and the Soviet statistician E. E. Slutsky; today it's known as the “Yule-Slutsky effect”.

- (e) (5) *Yule-Slutsky again* Assume  $\text{Var}[\epsilon(t)] = \sigma^2$ ,  $\text{Cov}[\epsilon(t), \epsilon(s)] = 0$  (unless  $t = s$ ). What is  $\text{Cov}[\Delta(t), \Delta(t+1)]$ ?

5. (1) How long did you spend on this problem set?

RUBRIC (10): The text is laid out cleanly, with clear divisions between problems and sub-problems. The writing itself is well-organized, free of grammatical and other mechanical errors, and easy to follow. Questions which ask for a plot or table are answered with both the figure itself and the command (or commands) use to make the plot. Plots are carefully labeled, with informative and legible titles, axis labels, and (if called for) sub-titles and legends; they are placed near the text of the corresponding problem. All quantitative and mathematical claims are supported by appropriate derivations, included in the text, or calculations in code. Numerical results are reported to appropriate precision. Code is properly integrated with a tool like R Markdown or knitr, and both the knitted file and the source file are submitted. The code is indented, commented, and uses meaningful names. All code is relevant, without dangling or useless commands. All parts of all problems are answered with coherent sentences, and raw computer code or output are only shown when explicitly asked for.