

Homework 4

Stats 32: Introduction to R for Undergraduates

Due: Friday, May 3, 2024, 11:59 pm PT

Instructions

You may (in fact, are encouraged to) use the Internet (including AI assistants, like ChatGPT) to search up any information to help you with this assignment, though you must cite any external (i.e. non-course related) resources that you use. Similarly, *after attempting this assignment by yourself*, you may collaborate with other students in the course, but you must each write your own code and acknowledge all students with whom you collaborated *for each problem* (you don't need to cite by subpart). However, you may not post on Internet forums (e.g. Stack Exchange) for help with this assignment; doing so is considered an Honor Code violation. You also may not copy verbatim any significant amount of code from the Internet (including AI assistants, like ChatGPT), even with citation. Feeding in the problems directly into AI assistants (or substantively paraphrased version) is also not permitted.

Please provide your code responses to each problem in the `.Rmd` file in the R code chunks directly below each subpart, inserting additional R code chunks if needed. Any text response can go right underneath the corresponding question.

On Gradescope, please submit a single `.pdf` file created by knitting the document with your responses. Problem 0 will provide guidance on how to do this.

Credit is given based on the approach and code, not necessarily the final answer.

All visualizations should have an appropriate title and axis labels.

Problem 1: Coffee and mental math

In November 2021, I ran an experiment to see whether coffee would improve my mental math performance. I took a mental math challenge 6 times a day for 12 days. On each day I randomly assigned myself a coffee dosage of 0 mL, 125 mL, or 250 mL at 12 noon. Then I took the test twice at 3 pm, 6 pm, and 9 pm. `coffee.csv` contains all of my scores.

- (2 points) Read in `coffee.csv`. You will note that the data is not in “tidy” format; there is 1 row for each date, but to get “tidy” data we want 1 row for each observation (i.e. each time I took a test). Use `pivot_longer()` to convert the data into tidy format. The end result, which you should store in a variable called `coffee`, should have a total of 5 columns: `Date`, `Dosage`, `Time` (3pm, 6pm, or 9pm), `Repetition` (R1 or R2), and `Score` (the test score). Hint: You may want to use the `names_sep` argument.
- (3 points) We consider `Dosage` as a categorical variable. Convert that column to a factor. Then create an appropriate visualization that shows the distribution of test scores for different dosages.
- (2 points) Fit a simple linear regression model to predict test score based on `Dosage`.
- (2 points) Compute the predictions for each of the 3 possible values of `Dosage`, under the model in (b). Verify, using some `dplyr` verbs, that these predictions are equivalent to the average score across all test scores from days with the particular coffee dosage.
- (3 points) Now, fit a multiple regression model using all three predictors `Dosage`, `Repetition`, and `Time`, all categorical. Interpret each of the resulting coefficient estimates.

Problem 2: Temperatures in San Jose

Daily high and low Fahrenheit temperatures in San Jose for 8,179 dates between 1998 and 2020 can be found in `san_jose.csv` (they are taken from Climate Data Online, at the station USW00023293). Your goal will be to see if you can predict the 2020 temperatures based on the data from 1998 to 2019.

- (a) (2 points) Read the data into a tibble and create a new column called `day_of_year` that converts the date to the day of the year. For example, January 4 would be the 4th day of the year, March 1 would be the 60th day of the year in non-leap years and the 61st day in leap years. Hint: The `yday()` function in the `lubridate` package may be helpful.
- (b) (2 points) Fit a simple linear regression model to predict the high temperature based on `day_of_year` **using only the data from 1998 through 2019**. Provide graphical evidence that suggests the linear model is a poor fit.
- (c) (3 points) Fit a better model for predicting the high temperature based on `day_of_year` **using only the data from 1998 through 2019**. Explain why your model choice is better than the simple linear regression in part (b) for this data, based on your visualizations in part (b).
- (d) (2 points) Repeat part (c) for the low temperature. What is the predicted low temperature on May 4, 2023?
- (e) (4 points) Create a scatterplot with date on the horizontal axis and (actual) high and low temperatures on the vertical axis. Only include data from 2020 (i.e. the year for which we did not train the model). Color the points corresponding to high temperatures in red, and the points corresponding to low temperatures in blue. Then, compute high and low temperature predictions for all dates in 2020 using the models you fit in parts (c) and (d). Plot those predictions as lines on the same axes (red for high temperatures, blue for low temperatures). Make sure there is a legend for these colors.