Problem Set 2

Due Friday Sep. 20, 10 am

Comments

- This covers material in Units 3 and 4.
- It's due at 10 am (Pacific) on September 20, both submitted as a PDF to Gradescope as well as committed to your GitHub repository.
- Please see PS1 for formatting and attribution requirements.
- Note that using chunks of bash code in Qmd may be troublesome.
 - You will need to add engine: knitr to the YAML preface of your qmd document. The default jupyter engine won't run both bash and Python chunks in the same document because the Jupyter notebooks are associated with a single 'kernel' (i.e., a single language for the code chunks).
 - For the knitr engine, you'll need to have R installed on your computer, including the knitr package. Quarto will then process the code chunks through knitr (which will use R's reticulate package to handle Python chunks).
 - If you have trouble on your own computer, you can always render your solution for this
 problem set on an SCF machine (we won't generally use bash chunks in future problem
 sets). (In particular I'm not quite sure what will happen if you render on Windows.)
 - We can help troubleshoot and feel free to post on Ed.
 - Test things out well before the due date with a dummy qmd file with both a bash chunk and a Python chunk and make sure things work.
- You will probably need to use sed in a basic way as shown in the bash tutorial. You should not
 need to use more advanced functionality nor should you need to use awk, but you may if you
 want to.

Problems

1. A friend of mine is planning to get married in Death Valley National Park in March (this problem is based on real events...). She wants to hold it as late in March as possible but without having a high chance of a very hot day. This problem will automate the task of generating information about what day of March to hold the wedding using data from the Global Historical Climatology Network. All of your operations should be done using the bash shell except part (c). Also, ALL of your work should be done using shell commands that you save in your solution file. So you can't say "I downloaded the data from such-and-such website" or "I unzipped the file"; you need to give us the bash code that we could run to repeat what you did. This is partly for practice in writing shell code and partly to enforce the idea that your work should be reproducible and

documented.

- a. Download yearly climate data for a set of years of interest into a temporary directory. Do not download all the years and feel free to focus on a small number of years to reduce the amount of data you need to download. Note that data for Death Valley is only present in the last few decades. As you are processing the files, report the number of observations in each year by printing the information to the screen (i.e., stdout), including if there are no observations for that year.
- b. Subset to the station corresponding to Death Valley, to the TMAX (maximum daily temperature) variable, and to March, and put all the data into a single file. In subsetting to Death Valley, get the information programmatically from the ghcnd-stations.txt file one level up in the website. Do NOT type in the station ID code when you retrieve the Death Valley data from the yearly files.
- c. Create a Python chunk (or R would be fine too) that takes as input your single file from (b) and makes a single plot showing side-by-side boxplots containing the maximum daily temperatures on each calendar day in March. (If you somehow really have trouble mixing Python and bash chunks, it's ok to insert this figure manually, after running the Python code separately. In this case you could use the jupyter engine provided that a bash kernel is available for Jupyter.)
- d. Now generalize your code from parts (a) and (b). Write a shell function that takes as arguments a string for identifying the location, the weather variable of interest, and the time period (i.e., the years of interest and the month of interest), and returns the results. Your function should detect if the user provides the wrong number of arguments or a string that doesn't allow one to identify a single weather station and return a useful error message. It should also give useful help information if the user invokes the function as: get_weather -h. Finally the function should remove the raw downloaded data files (or you should download into your operating system's temporary file location).

Hint: to check for equality in an if statement, you generally need syntax like:

- 2. Add documentation, error-trapping and testing for your code from Problem 4, parts (b) and (c) of PS1. You may use a modified version of your PS1 solution, perhaps because you found errors in what you did or wanted to make changes based on Chris' solutions (to be distributed in class on Friday Sep. 13) or your discussions with other students. These topics will be covered in Lab 2 (Sep. 13) and are also discussed in Unit 4.
 - a. Add informative doc strings to your functions.
 - b. Add exceptions for handling run-time errors. You should try to catch the various incorrect inputs a user could provide and anything else that could go wrong (e.g., what happens if the server refuses the request or if one is not online?). In some cases you will want to raise an error, but in others you may want to catch an error with try-except and return None.
 - c. Use the pytest package to set up a thoughtful set of unit tests of your functions.