Programming Assignment 3

R Programming

Introduction

Download the file ProgAssignment3-data.zip file containing the data for Programming Assignment 3 from the Coursera web site. Unzip the file in a directory that will serve as your working directory. When you start up R make sure to change your working directory to the directory where you unzipped the data.

The data for this assignment come from the Hospital Compare web site (http://hospitalcompare.hhs.gov) run by the U.S. Department of Health and Human Services. The purpose of the web site is to provide data and information about the quality of care at over 4,000 Medicare-certified hospitals in the U.S. This dataset essentially covers all major U.S. hospitals. This dataset is used for a variety of purposes, including determining whether hospitals should be fined for not providing high quality care to patients (see http://goo.gl/jAXFX for some background on this particular topic).

The Hospital Compare web site contains a lot of data and we will only look at a small subset for this assignment. The zip file for this assignment contains three files

- outcome-of-care-measures.csv: Contains information about 30-day mortality and readmission rates for heart attacks, heart failure, and pneumonia for over 4,000 hospitals.
- hospital-data.csv: Contains information about each hospital.
- Hospital_Revised_Flatfiles.pdf: Descriptions of the variables in each file (i.e the code book).

A description of the variables in each of the files is in the included PDF file named Hospital_Revised_Flatfiles.pdf. This document contains information about many other files that are not included with this programming assignment. You will want to focus on the variables for Number 19 ("Outcome of Care Measures.csv") and Number 11 ("Hospital_Data.csv"). You may find it useful to print out this document (at least the pages for Tables 19 and 11) to have next to you while you work on this assignment. In particular, the numbers of the variables for each table indicate column indices in each table (i.e. "Hospital Name" is column 2 in the outcome-of-care-measures.csv file).

1 Plot the 30-day mortality rates for heart attack

Read the outcome data into R via the read.csv function and look at the first few rows.

```
> outcome <- read.csv("outcome-of-care-measures.csv", colClasses = "character")
> head(outcome)
```

There are many columns in this dataset. You can see how many by typing ncol(outcome) (you can see the number of rows with the nrow function). In addition, you can see the names of each column by typing names(outcome) (the names are also in the PDF document.

To make a simple histogram of the 30-day death rates from heart attack (column 11 in the outcome dataset), run

```
> outcome[, 11] <- as.numeric(outcome[, 11])
> ## You may get a warning about NAs being introduced; that is okay
> hist(outcome[, 11])
```

Because we originally read the data in as character (by specifying colClasses = "character" we need to coerce the column to be numeric. You may get a warning about NAs being introduced but that is okay.

There is nothing to submit for this part.

2 Finding the best hospital in a state

Write a function called best that take two arguments: the 2-character abbreviated name of a state and an outcome name. The function reads the outcome-of-care-measures.csv file and returns a character vector with the name of the hospital that has the best (i.e. lowest) 30-day mortality for the specified outcome in that state. The hospital name is the name provided in the Hospital.Name variable. The outcomes can be one of "heart attack", "heart failure", or "pneumonia". Hospitals that do not have data on a particular outcome should be excluded from the set of hospitals when deciding the rankings.

Handling ties. If there is a tie for the best hospital for a given outcome, then the hospital names should be sorted in alphabetical order and the first hospital in that set should be chosen (i.e. if hospitals "b", "c", and "f" are tied for best, then hospital "b" should be returned).

The function should use the following template.

```
best <- function(state, outcome) {
    ## Read outcome data

    ## Check that state and outcome are valid

    ## Return hospital name in that state with lowest 30-day death
    ## rate
}</pre>
```

The function should check the validity of its arguments. If an invalid state value is passed to best, the function should throw an error via the stop function with the exact message "invalid state". If an invalid outcome value is passed to best, the function should throw an error via the stop function with the exact message "invalid outcome".

Here is some sample output from the function.

```
> source("best.R")
> best("TX", "heart attack")
[1] "CYPRESS FAIRBANKS MEDICAL CENTER"
> best("TX", "heart failure")
[1] "FORT DUNCAN MEDICAL CENTER"
> best("MD", "heart attack")
[1] "JOHNS HOPKINS HOSPITAL, THE"
> best("MD", "pneumonia")
[1] "GREATER BALTIMORE MEDICAL CENTER"
> best("BB", "heart attack")
Error in best("BB", "heart attack"): invalid state
> best("NY", "hert attack")
Error in best("NY", "hert attack"): invalid outcome
```

Save your code for this function to a file named best.R.

Use the submit script provided to submit your solution to this part. There are 3 tests that need to be passed for this part of the assignment.

3 Ranking hospitals by outcome in a state

Write a function called rankhospital that takes three arguments: the 2-character abbreviated name of a state (state), an outcome (outcome), and the ranking of a hospital in that state for that outcome (num). The function reads the outcome-of-care-measures.csv file and returns a character vector with the name of the hospital that has the ranking specified by the num argument. For example, the call

```
rankhospital("MD", "heart failure", 5)
```

would return a character vector containing the name of the hospital with the 5th lowest 30-day death rate for heart failure. The num argument can take values "best", "worst", or an integer indicating the ranking (smaller numbers are better). If the number given by num is larger than the number of hospitals in that state, then the function should return NA. Hospitals that do not have data on a particular outcome should be excluded from the set of hospitals when deciding the rankings.

Handling ties. It may occur that multiple hospitals have the same 30-day mortality rate for a given cause of death. In those cases ties should be broken by using the hospital name. For example, in Texas ("TX"), the hospitals with lowest 30-day mortality rate for heart failure are shown here.

> head(texas)

```
Hospital.Name Rate Rank
3935
           FORT DUNCAN MEDICAL CENTER
4085
      TOMBALL REGIONAL MEDICAL CENTER
                                               2
4103 CYPRESS FAIRBANKS MEDICAL CENTER
                                               3
3954
               DETAR HOSPITAL NAVARRO
                                               4
                                        8.7
4010
               METHODIST HOSPITAL, THE
                                               5
3962
     MISSION REGIONAL MEDICAL CENTER
                                        8.8
                                               6
```

Note that Cypress Fairbanks Medical Center and Detar Hospital Navarro both have the same 30-day rate (8.7). However, because Cypress comes before Detar alphabetically, Cypress is ranked number 3 in this scheme and Detar is ranked number 4. One can use the order function to sort multiple vectors in this manner (i.e. where one vector is used to break ties in another vector).

The function should use the following template.

```
rankhospital <- function(state, outcome, num = "best") {
    ## Read outcome data

    ## Check that state and outcome are valid

    ## Return hospital name in that state with the given rank
    ## 30-day death rate
}</pre>
```

The function should check the validity of its arguments. If an invalid state value is passed to best, the function should throw an error via the stop function with the exact message "invalid state". If an invalid outcome value is passed to best, the function should throw an error via the stop function with the exact message "invalid outcome".

Here is some sample output from the function.

```
> source("rankhospital.R")
> rankhospital("TX", "heart failure", 4)

[1] "DETAR HOSPITAL NAVARRO"
> rankhospital("MD", "heart attack", "worst")

[1] "HARFORD MEMORIAL HOSPITAL"
> rankhospital("MN", "heart attack", 5000)

[1] NA
```

Save your code for this function to a file named rankhospital.R.

Use the submit script provided to submit your solution to this part. There are 4 tests that need to be passed for this part of the assignment.

4 Ranking hospitals in all states

Write a function called rankall that takes two arguments: an outcome name (outcome) and a hospital ranking (num). The function reads the outcome-of-care-measures.csv file and returns a 2-column data frame containing the hospital in each state that has the ranking specified in num. For example the function call rankall("heart attack", "best") would return a data frame containing the names of the hospitals that are the best in their respective states for 30-day heart attack death rates. The function should return a value for every state (some may be NA). The first column in the data frame is named hospital, which contains the hospital name, and the second column is named state, which contains the 2-character abbreviation for the state name. Hospitals that do not have data on a particular outcome should be excluded from the set of hospitals when deciding the rankings.

Handling ties. The rankall function should handle ties in the 30-day mortality rates in the same way that the rankhospital function handles ties.

The function should use the following template.

```
rankall <- function(outcome, num = "best") {
    ## Read outcome data

    ## Check that state and outcome are valid

    ## For each state, find the hospital of the given rank

    ## Return a data frame with the hospital names and the
    ## (abbreviated) state name
}</pre>
```

NOTE: For the purpose of this part of the assignment (and for efficiency), your function should NOT call the rankhospital function from the previous section.

The function should check the validity of its arguments. If an invalid outcome value is passed to rankall, the function should throw an error via the stop function with the exact message "invalid outcome". The num variable can take values "best", "worst", or an integer indicating the ranking (smaller numbers are better). If the number given by num is larger than the number of hospitals in that state, then the function should return NA.

Here is some sample output from the function.

```
> source("rankall.R")
> head(rankall("heart attack", 20), 10)
                               hospital state
AK
                                   <NA>
                                           AK
        D W MCMILLAN MEMORIAL HOSPITAL
AL
                                           AL
AR
     ARKANSAS METHODIST MEDICAL CENTER
                                           AR
AZ JOHN C LINCOLN DEER VALLEY HOSPITAL
                                           ΑZ
CA
                 SHERMAN OAKS HOSPITAL
                                           CA
              SKY RIDGE MEDICAL CENTER
CO
                                           CO
CT
               MIDSTATE MEDICAL CENTER
                                           CT
DC
                                   <NA>
                                           DC
DE
                                   <NA>
                                           DE
FL
        SOUTH FLORIDA BAPTIST HOSPITAL
                                           FL
> tail(rankall("pneumonia", "worst"), 3)
                                      hospital state
WI MAYO CLINIC HEALTH SYSTEM - NORTHLAND, INC
WV
                        PLATEAU MEDICAL CENTER
                                                   WV
WY
             NORTH BIG HORN HOSPITAL DISTRICT
                                                   WY
> tail(rankall("heart failure"), 10)
                                                              hospital state
TN
                            WELLMONT HAWKINS COUNTY MEMORIAL HOSPITAL
                                           FORT DUNCAN MEDICAL CENTER
UT VA SALT LAKE CITY HEALTHCARE - GEORGE E. WAHLEN VA MEDICAL CENTER
VA
                                              SENTARA POTOMAC HOSPITAL
VI
                               GOV JUAN F LUIS HOSPITAL & MEDICAL CTR
```

Save your code for this function to a file named rankall.R.

VT

WA

WI

WV

WY

Use the submit script provided to submit your solution to this part. There are 3 tests that need to be passed for this part of the assignment.

TN

TX

UT

VA

VI

VT

WA

WI

WV

WY

SPRINGFIELD HOSPITAL

HARBORVIEW MEDICAL CENTER

FAIRMONT GENERAL HOSPITAL

CHEYENNE VA MEDICAL CENTER

AURORA ST LUKES MEDICAL CENTER