

Programming

John M. Drake & Andrew W. Park

Introduction

This exercise is about writing *scripts*, purpose-built computer programs for performing some analysis. Our scripts will be written using the *RStudio Editor* and compiled using Rstudio. In this exercise we review many basic numerical operations and R functions, programming style, the development of custom functions, flow control and...

Data

West Nile virus (WNV) is a positive-sense single-stranded RNA virus transmitted by mosquitoes to a range of vertebrate hosts. WNV was first identified in Uganda in 1937 and found in parts of Europe, Asia, and Australia during the 1950s and 1960s. In 1999, WNV was first reported in the Americas in association with dieoffs of captive and wild birds. This outbreak initiated widespread epidemic that swept across North America and is now spreading in Central and South America. Humans are “dead end” hosts (humans do not achieve sufficiently high viremia to be infectious to mosquitoes). The majority of human cases are asymptomatic, but a small fraction of cases result in meningitis, encephalitis, and/or death. State-level data on the number of reported cases, meningitis/encephalitis, and fatalities are compiled and reported by the CDC and US Geological survey at <https://diseasemaps.usgs.gov/>. The file `wnv.csv` contains tabular data on the number of reported cases (mostly febrile cases), neuroinvasive cases (meningitis/encephalitis), and fatalities for all continental US states from 1999-2007. Additional data are the latitude and longitude of the centroid of each state.

Scripts

Exercise. Write a script to load the West Nile virus data and use `ggplot` to create a histogram for the total number of cases in each state in each year. Follow the format of the *prototypical script* advocated in the presentation: Header, Load Packages, Declare Functions, Load Data, Perform Analysis.

With each of the following exercises, extend your script so that at the end of the unit you have one script that performs the entire analysis.

Exercise. The state-level and case burden is evidently highly skewed. Plot a histogram for the logarithm of the number of cases. Do this two different ways.

Exercise. Use arithmetic operators to calculate the raw case fatality rate (CFR) in each state in each year. Plot a histogram of the calcated CFRs.

Exercise. Use arithmetic operators, logical operators, and the function `sum` to verify that the variable `Total` is simply the sum of the number of febrile cases, neuroinvasive cases, and other cases.

Exercise. Use modular arithmetic to provide an annual case count for each state rounded (down) to the nearest dozen. Use modular arithmetic to extract the rounding errors associated with this calculate, then add the errors to obtain the total error.

Functions

Exercise. Let us call the ratio of meningitis/encephalitis cases to the total number of cases the *neuroinvasive disease rate*. Write a function to calculate the mean and standard error (standard deviation divided by the square root of the sample size) of the neuroinvasive disease rate for all the states in a given list and given set of years. Follow the Google R style and remember to place the function near the top of your script. Use your function to calculate the average severe disease rate in California, Colorado, and New York.

Exercise. Use ggplot to show the neuroinvasive disease rate for these states as a bar graph with error bars to show the standard deviation.

Exercise. Use your function and ggplot to show the neuroinvasive disease rate for all states.

Pipes

Exercise. Use pipes to produce the same plots without using your function.

Flow control

Using help

Exercise. Use the help for function `lm`.