

# Programming Assignment 1: Air Pollution

**Introduction** For this first programming assignment you will write three functions that are meant to interact with dataset that accompanies this assignment. The dataset is contained in a zip file `specdata.zip` that you can download from the Coursera web site.

Although this is a programming assignment, you will be assessed using a separate quiz.

**Data** The zip file containing the data can be downloaded here:

`specdata.zip` [2.4MB]

The zip file contains 332 comma-separated-value (CSV) files containing pollution monitoring data for fine particulate matter (PM) air pollution at 332 locations in the United States. Each file contains data from a single monitor and the ID number for each monitor is contained in the file name. For example, data for monitor 200 is contained in the file “200.csv”. Each file contains three variables:

**Date:** the date of the observation in YYYY-MM-DD format (year-month-day)

**sulfate:** the level of sulfate PM in the air on that date (measured in micrograms per cubic meter)

**nitrate:** the level of nitrate PM in the air on that date (measured in micrograms per cubic meter)

For this programming assignment you will need to unzip this file and create the directory ‘specdata’. Once you have unzipped the zip file, do not make any modifications to the files in the ‘specdata’ directory. In each file you’ll notice that there are many days where either sulfate or nitrate (or both) are missing (coded as NA). This is common with air pollution monitoring data in the United States.

**Part 1** Write a function named ‘pollutantmean’ that calculates the mean of a pollutant (sulfate or nitrate) across a specified list of monitors. The function ‘pollutantmean’ takes three arguments: ‘directory’, ‘pollutant’, and ‘id’. Given a vector monitor ID numbers, ‘pollutantmean’ reads that monitors’ particulate matter data from the directory specified in the ‘directory’ argument and returns the mean of the pollutant across all of the monitors, ignoring any missing values coded as NA. A prototype of the function is as follows

You can see some example output from this function below. The function that you write should be able to match this output. Please save your code to a file named `pollutantmean.R`.

```
pollutantmean <- function(directory, pollutant, id = 1:332) {  
  # Initialize a vector to hold the pollutant data  
  pollutant_data <- c()  
  
  # Loop over each monitor ID in the id vector  
  for (i in id) {  
    # Create the filename using the directory and monitor ID  
    filename <- sprintf("%s/%03d.csv", directory, i)  
  
    # Read the data from the file  
    data <- read.csv(filename)  
  
    # Extract the pollutant data and remove NA values  
    pollutant_values <- data[[pollutant]]  
    pollutant_data <- c(pollutant_data, pollutant_values)  
  }  
}
```

```

# Calculate and return the mean of the pollutant data, ignoring NA values
mean_value <- mean(pollutant_data, na.rm = TRUE)
return(mean_value)
}

```

*# Example usage:*

```
pollutantmean("specdata", "sulfate", 1:10)
```

```
## [1] 4.064128
```

```
pollutantmean("specdata", "nitrate", 70:72)
```

```
## [1] 1.706047
```

```
pollutantmean("specdata", "nitrate", 23)
```

```
## [1] 1.280833
```

Part 2 Write a function that reads a directory full of files and reports the number of completely observed cases in each data file. The function should return a data frame where the first column is the name of the file and the second column is the number of complete cases. A prototype of this function follows

You can see some example output from this function below. The function that you write should be able to match this output. Please save your code to a file named complete.R. To run the submit script for this part, make sure your working directory has the file complete.R in it.

```

setwd(getwd())
source("complete.R")
complete("specdata", 1)

```

```
##   id nobs
## 1  1  117
```

```
complete("specdata", c(2, 4, 8, 10, 12))
```

```
##   id nobs
## 1  2 1041
## 2  4  474
## 3  8  192
## 4 10  148
## 5 12   96
```

```
complete("specdata", 30:25)
```

```
##   id nobs
## 1 30  932
## 2 29  711
## 3 28  475
## 4 27  338
## 5 26  586
## 6 25  463
```

```
complete("specdata", 3)
```

```
##    id nobs  
## 1   3   243
```

Part 3 Write a function that takes a directory of data files and a threshold for complete cases and calculates the correlation between sulfate and nitrate for monitor locations where the number of completely observed cases (on all variables) is greater than the threshold. The function should return a vector of correlations for the monitors that meet the threshold requirement. If no monitors meet the threshold requirement, then the function should return a numeric vector of length 0. A prototype of this function follows:

```
source("corr.R")  
source("complete.R")  
cr <- corr("specdata", 150)  
  
head(cr)
```

```
## [1] -0.01895754 -0.14051254 -0.04389737 -0.06815956 -0.12350667 -0.07588814
```

```
summary(cr)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.  
## -0.21057 -0.04999  0.09463  0.12525  0.26844  0.76313
```

```
cr <- corr("specdata", 400)  
head(cr)
```

```
## [1] -0.01895754 -0.04389737 -0.06815956 -0.07588814  0.76312884 -0.15782860
```

```
summary(cr)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.  
## -0.17623 -0.03109  0.10021  0.13969  0.26849  0.76313
```

```
cr <- corr("specdata", 5000)  
summary(cr)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.  
##
```

```
length(cr)
```

```
## [1] 0
```

```
cr <- corr("specdata")  
summary(cr)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.  
## -1.00000 -0.05282  0.10718  0.13684  0.27831  1.00000
```

```
length(cr)
```

```
## [1] 323
```