

Lecture 19 - Interacting with data using R

Review and extension

In this short walkthrough, we will work with two tabular datasets to review and extend a few concepts from Monday and the Patient Data exercise from Software Carpentry.

First we will load a simple data table using `read.table()`:

```
data=read.table(file="test.dat",header=FALSE,sep=" ")
data
```

```
##   V1 V2 V3 V4 V5
## 1  1  5  9 13 17
## 2  2  6 10 14 18
## 3  3  7 11 15 19
## 4  4  8 12 16 20
```

Remember we can use logic tests and square brackets to index our data. This can be really useful when we have a large dataset that we want to access a subset of based on characteristics of the data itself.

```
# we can test for equality using double equal signs
data[,1]==1
```

```
## [1] TRUE FALSE FALSE FALSE
```

```
# we can test for greater than or less than as well
data[,1]>2
```

```
## [1] FALSE FALSE TRUE TRUE
```

```
# the logical values returned by a logic test can be used just like numbers to index a data structure
data[data[,1]>2,]
```

```
##   V1 V2 V3 V4 V5
## 3  3  7 11 15 19
## 4  4  8 12 16 20
```

A nice thing about `dataframes` in R is that they can hold more than one data mode (e.g. numbers and characters), like what we saw in `wages.csv` from our exercise last week.

```
#Load wages.csv; the stringsAsFactors argument prevents strings from being treated as factors
wages=read.csv(file="wages.csv",header=TRUE,stringsAsFactors=FALSE)
class(wages)
```

```
## [1] "data.frame"
```

```
dim(wages)
```

```
## [1] 3294    4
```

```
head(wages)
```

```
##   gender yearsExperience yearsSchool    wage
## 1 female                9         13 6.315296
```

```
## 2 female          12          12 5.479770
## 3 female          11          11 3.642170
## 4 female           9          14 4.593337
## 5 female           8          14 2.418157
## 6 female           9          14 2.094058
```

We can use square brackets to index a portion of a `dataframe`, but we can also use dollarsign notation. This is because a `dataframe` also behaves like a `list` in R.

```
# we can extract all of the female wage data using square brackets
females=wages[wages[,1]=="female",]
dim(females)
```

```
## [1] 1569    4
```

```
unique(females[,1])
```

```
## [1] "female"
```

```
# or a mix of square brackets and dollarsign notation
females2=wages[wages$gender=="female",]
dim(females2)
```

```
## [1] 1569    4
```

```
unique(females2$gender)
```

```
## [1] "female"
```

One thing we did in `bash` that we haven't covered in R is sorting. We can use the `sort()` function to sort a vector. However, to sort the rows of a matrix or dataframe by one or more columns we use a different function - `order()`. `order()` doesn't sort for us, but it does return a vector that essentially gives the "instructions" for sorting. We can then use those "instructions" returned by `order()` along with square brackets to index the rows in the order we want to generate a sorted matrix or dataframe.

```
# let's make a simple data frame to see how this works
df<-data.frame(name=c("abby","sam","clara","joe"),age=c(24,27,12,82))
```

```
# we can sort a single column with sort()
sort(df$age)
```

```
## [1] 12 24 27 82
```

```
sort(df$name)
```

```
## [1] "abby" "clara" "joe"  "sam"
```

```
# but to sort rows by a given column and keep the rows intact we must use order()
df_sorted<-df[order(df$age),]
df_sorted
```

```
##   name age
## 1 abby  24
## 3 clara 12
## 4  joe  82
## 2  sam  27
```

```
# this works because order() returns a vector of the row indexes and then we use those indexes in square brackets
order(df$age)
```

```
## [1] 1 3 4 2
```

Challenge

Let's redo one of our bash exercises in R!

1. Write a file containing the unique gender-yearsExperience combinations contained in the file "wages.csv". The file you create should contain gender in the first column and yearsExperience in a second column with a space separating the two columns. The rows should be sorted first by gender and then by yearsExperience, but remember to keep the pairings in a given row intact. Don't worry about column names in the output file.
2. Return the following information to the R console when the script is executed: the gender, yearsExperience, and wage for the highest earner, the gender, yearsExperience, and wage for the lowest earner, and the number of females in the top ten earners in this data set. Be sure to indicate, which output is which when returning them to the console.
3. Return one more piece of information to the console: the effect of graduating college (12 vs. 16 years of school) on the minimum wage for earners in this dataset. In other words, take the difference of the minimum wages for those with 12 versus 16 years of schooling