Data structures in R

Abhijit Dasgupta

BIOF 339

A quick refresh

- R is a scripting language for data analysis and statistics
- R Markdown is a way of combining textual information and R code to produce reproducible documents
- RStudio is an integrated environment that makes it easier to work with R

You type commands (code) for R to run.

- objects like data (nouns)
- functions that do something to R objects (verbs)

Examples

airquality
diamonds
summary(airquality)

Let's start with the airquality data.

- It is an object
- of class class(airquality) = data.frame

How about each column?

Let's look at the Ozone and Wind columns

- We can access them using airquality\$Ozone and airquality\$Wind
 - o class(airquality\$0zone) = integer
 - o class(airquality\$Wind) = numeric

```
Ozone Solar.R Wind Temp Month Day
                         67
      36
              149 12.6
                          74
                                 5
5
      28
      19
               99 13.8
               19 20.1
      NA
                                    10
              194
                   8.6
                          74
               NA
                   6.9
      16
                                    12
             256
                                    13
                                    14
             274 10.9
                                    15
      18
                                    16
      34
                                    17
       6
                         57
                                    18
      30
                         68
                                    19
20
                         62
                                    20
                         59
                                    21
      11
                         73
                                    22
             320 16.6
23
                         61
                                    23
      32
                                    24
               66 16.6
                         57
                                    25
```

```
head(iris)
```

```
Sepal.Length Sepal.Width Petal.Length Petal.Width Species
        5.1
                   3.5
                                1.4
                                           0.2 setosa
                                            0.2 setosa
                    3.0
                                1.4
                   3.2
                                1.3
                                           0.2 setosa
        4.6
                   3.1
                                1.5
                                           0.2 setosa
        5.0
                   3.6
                                           0.2 setosa
                                1.4
                   3.9
                                1.7
                                            0.4 setosa
```

```
str(iris)
```

```
'data.frame': 150 obs. of 5 variables:
$ Sepal.Length: num 5.1 4.9 4.7 4.6 5 5.4 4.6 5 4.4 4.9 ...
$ Sepal.Width: num 3.5 3 3.2 3.1 3.6 3.9 3.4 3.4 2.9 3.1 ...
$ Petal.Length: num 1.4 1.4 1.3 1.5 1.4 1.7 1.4 1.5 1.4 1.5 ...
$ Petal.Width: num 0.2 0.2 0.2 0.2 0.2 0.4 0.3 0.2 0.2 0.1 ...
$ Species : Factor w/ 3 levels "setosa", "versicolor", ..: 1 1 1 1 1 1 1 1 1 1 ...
```

Now we see another type of object, a factor

```
library(ggplot2)
str(midwest)
```

```
tibble [437 \times 28] (S3: tbl_df/tbl/data.frame)
$ PID
                       : int [1:437] 561 562 563 564 565 566 567 568 569 570 ...
                       : chr [1:437] "ADAMS" "ALEXANDER" "BOND" "BOONE" ...
$ county
                       : chr [1:437] "IL" "IL" "IL" "IL" ...
$ state
                       : num [1:437] 0.052 0.014 0.022 0.017 0.018 0.05 0.017 0.027 0.024 0.058 ...
$ area
$ poptotal
                       : int [1:437] 66090 10626 14991 30806 5836 35688 5322 16805 13437 173025 ...
$ popdensity
                       : num [1:437] 1271 759 681 1812 324 ...
$ popwhite
                       : int [1:437] 63917 7054 14477 29344 5264 35157 5298 16519 13384 146506 ...
$ popblack
                       : int [1:437] 1702 3496 429 127 547 50 1 111 16 16559 ...
$ popamerindian
                       : int [1:437] 98 19 35 46 14 65 8 30 8 331 ...
$ popasian
                       : int [1:437] 249 48 16 150 5 195 15 61 23 8033 ...
$ popother
                       : int [1:437] 124 9 34 1139 6 221 0 84 6 1596 ...
$ percwhite
                       : num [1:437] 96.7 66.4 96.6 95.3 90.2 ...
$ percblack
                       : num [1:437] 2.575 32.9 2.862 0.412 9.373 ...
$ percamerindan
                       : num [1:437] 0.148 0.179 0.233 0.149 0.24 ...
$ percasian
                       : num [1:437] 0.3768 0.4517 0.1067 0.4869 0.0857 ...
$ percother
                       : num [1:437] 0.1876 0.0847 0.2268 3.6973 0.1028 ...
$ popadults
                       : int [1:437] 43298 6724 9669 19272 3979 23444 3583 11323 8825 95971 ....
$ perchsd
                       : num [1:437] 75.1 59.7 69.3 75.5 68.9 ...
$ percollege
                       : num [1:437] 19.6 11.2 17 17.3 14.5 ...
$ percprof
                       : num [1:437] 4.36 2.87 4.49 4.2 3.37 ...
$ poppovertyknown
                       : int [1:437] 63628 10529 14235 30337 4815 35107 5241 16455 13081 154934 ....
$ percpovertyknown
                       : num [1:437] 96.3 99.1 95 98.5 82.5 ...
                       : num [1:437] 13.15 32.24 12.07 7.21 13.52 ...
$ percbelowpoverty
$ percchildbelowpovert: num [1:437] 18 45.8 14 11.2 13 ...
```

The most common types of data we see are numeric, character, factor. You can also see Date and logical

You can test to see if data is of a particular type, or convert from one data type to another

Data type	Test	Convert
numeric	is.numeric	as.numeric
character	is.character	as.character
factor	is.factor	as.factor

This conversion is important. Why?

05:00

Factors

Factors are uniquely an R thing.

They are meant to represent categorical data (gender, race, state, phenotype)

They look like character vectors, but internally act like integers, so you have to be a bit careful with them

Whenever you're in doubt, convert them to characters using as.character.

We'll see the utility of factors when we do data munging, summaries and modeling

Every object in R has a name

You give an object a name using the syntax name <- object

Naming conventions:

- 1. Snake_case or pothole_case
- 2. CamelCase
- 3. Some.people.use.periods

I'm partial to snake_case.

The point here is to make expressive names using English so you know what is stored in the name.

A silly exercise

From the iris dataset, save each column into a new object, giving it a name. Then see what kind of data that object contains.

05:00

Bigger objects

Scalar -> vector (array) -> matrix (2-d array)

- A scalar is a single number or word
- A vector is a bunch of scalars arranged in a row or column
- A matrix is a bunch of scalars arranged in rows and columns

Each of these must be of the same data type

Examples

```
[1] 2
c(2,3,4,5,6)
[1] 2 3 4 5 6
      c() is the concatenate function
matrix(c(1,2,3,4), byrow = T, nrow = 2)
     [,1] [,2]
1 2
3 4
```

Data comes in many flavors

Heterogeneous data

From Excel, we are familiar with keeping different kinds of data together in a spreadsheet

- Expression levels (numeric)
- Gene names (character)
- Date of experiment (Date)

In R, the objects that can hold heterogeneous data are data. frame and list

Data sets

Typical data structure

- Data is typically in a rectangular format
 - spreadsheet, database table
 - CSV (comma-separated values) or TSV (tab-separated values) files
- Characteristic
 - Rows are observations
 - Columns are variables
 - Each column has the same number of observations
 - Tidy data is a particularly amenable format for data analysis.

The data.frame

Dataframes are the primary mode of storing datasets in R

They were revolutionary in that they kept heterogeneous data together

They share properties of both a **matrix** and a **list**

class(airquality)

[1] "data.frame"

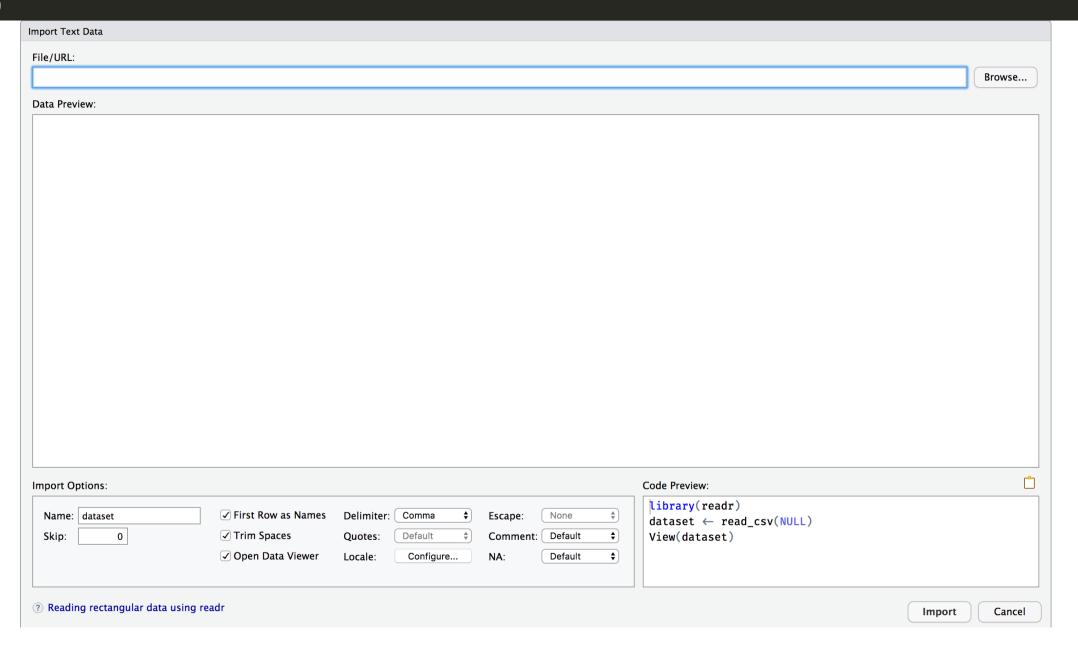
Technically, a data.frame is a list of vectors (or objects, generally) of the same length

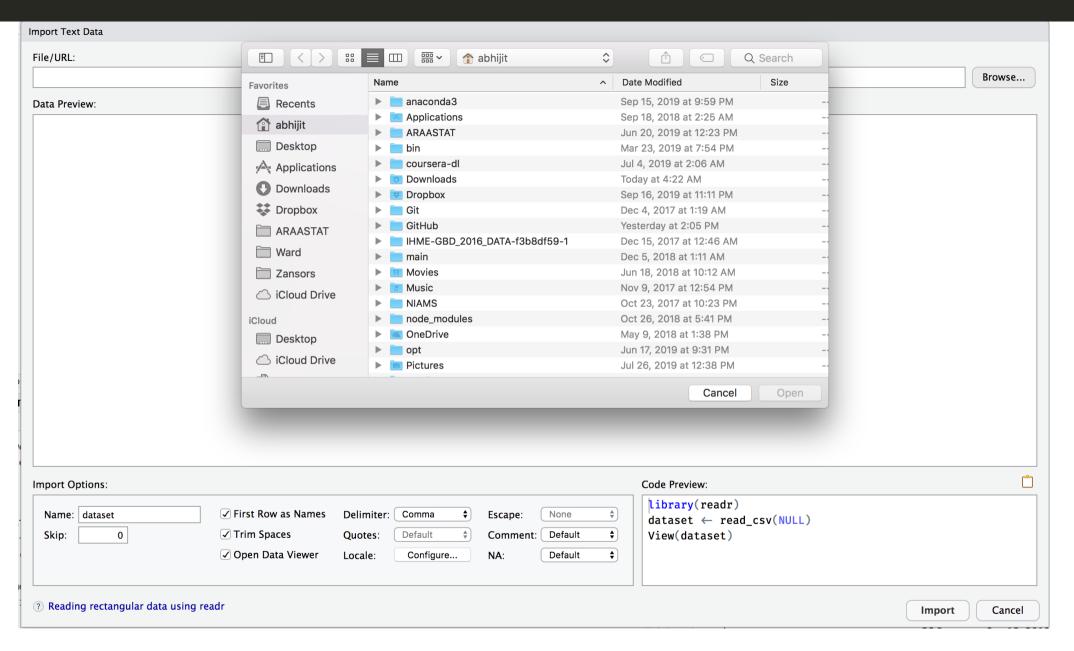
Load some data

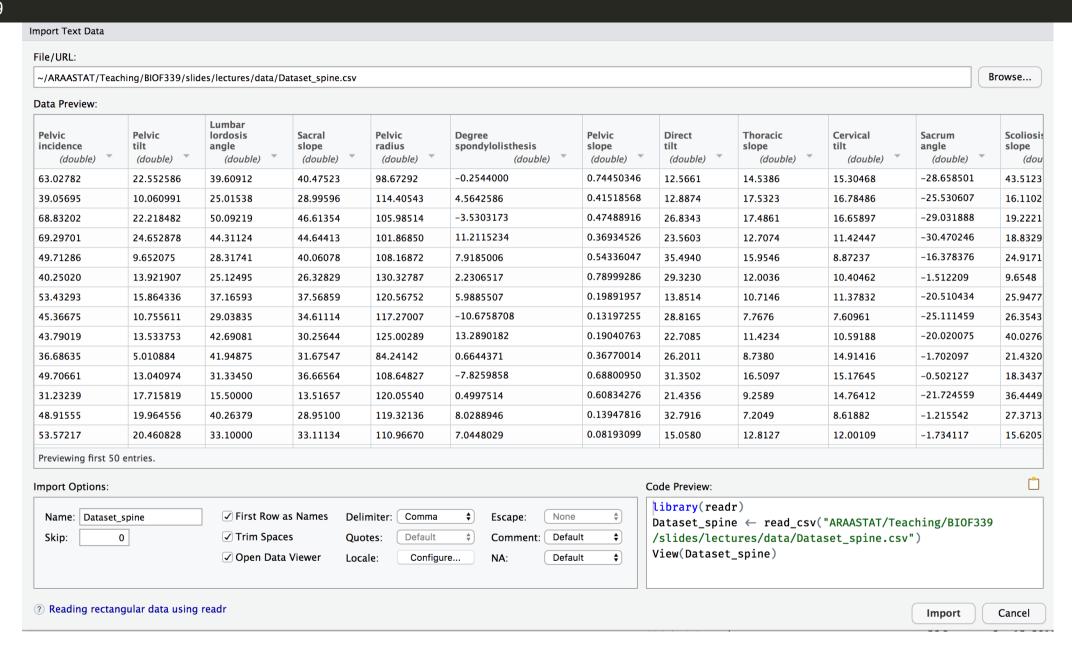
We'll load the spine dataset into R.

To do this, download the data from the web, and store it in the main folder in your project.

Then, in the Environment pane, import it using the **Import Dataset** button. You will use the From text (readr) option







A digression: Lists and Matrices

Matrices

A matrix is a rectangular array of data of the same type

```
matrix(0, nrow=2, ncol=4)
     [,1] [,2] [,3] [,4]
matrix(letters, nrow=2)
matrix(letters, nrow=2, byrow=T)
```

Matrices

You can create a matrix from a set of vectors of the same length

```
x <- c(1,2,3,4)

y <- c(10,20,30,40)
```

Put columns together

```
cbind(c(1,2,3,4), c(10,20,30,40)) ## Column bind
```

```
[,1] [,2]
[1,] 1 10
[2,] 2 20
[3,] 3 30
[4,] 4 40
```

Matrices

You can create a matrix from a set of vectors of the same length

```
x <- c(1,2,3,4)

y <- c(10,20,30,40)
```

Put rows together

```
example_matrix <- rbind(c(1,2,3,4), c(10,20,30,40)) ## Row bind example_matrix
```

```
[,1] [,2] [,3] [,4]
[1,] 1 2 3 4
[2,] 10 20 30 40
```

Extracting elements

```
example_matrix
     [,1] [,2] [,3] [,4]
[1,]
[2,]
example_matrix[1,] ## Extracts 1st row
[1] 1 2 3 4
example_matrix[,2:3] ## extracts 2nd & 3rd columns
     [,1][,2]
example_matrix[1,4]
[1] 4
```

Matrix properties

```
example_matrix
     [,1] [,2] [,3] [,4]
[1,]
[2,]
nrow(example_matrix) ## Number of rows
[1] 2
ncol(example_matrix) ## Number of columns
[1] 4
dim(example_matrix) ## shortcut for above
[1] 2 4
```

Matrix arithmetic

```
example_matrix
    [,1] [,2] [,3] [,4]
[1,]
[2,]
example_matrix + 5 ## Add 5 to each element
     [,1] [,2] [,3] [,4]
[1,]
example_matrix * 2 ## Multiply each element by 2
     [,1] [,2] [,3] [,4]
       20 40 60 80
```

```
example_matrix
     [,1] [,2] [,3] [,4]
[1,]
[2,]
example_matrix2 <- rbind(3:6, 9:12)</pre>
example_matrix2
     [,1] [,2] [,3] [,4]
[1,]
[2,]
example_matrix + example_matrix2
     [,1] [,2] [,3] [,4]
       19 30 41 52
```

example_matrix

```
[,1] [,2] [,3] [,4]
[1,] 1 2 3 4
[2,] 10 20 30 40
```

example_matrix2

```
[,1] [,2] [,3] [,4]
[1,] 3 4 5 6
[2,] 9 10 11 12
```

example_matrix * example_matrix2 ## Not matrix multiplication, but element-wise multiplication

```
[,1] [,2] [,3] [,4]
[1,] 3 8 15 24
[2,] 90 200 330 480
```

```
rbind(example_matrix, example_matrix2)
```

```
[,1] [,2] [,3] [,4]
[1,] 1 2 3 4
[2,] 10 20 30 40
[3,] 3 4 5 6
[4,] 9 10 11 12
```

```
cbind(example_matrix, example_matrix2)
```

```
[,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8]
[1,] 1 2 3 4 3 4 5 6
[2,] 10 20 30 40 9 10 11 12
```

```
dim(example_matrix2)
[1] 2 4
t(example_matrix2) ## Transpose of a matrix
     [,1] [,2]
3 9
[1,]
[2,]
[3,]
[4,]
       4 10
5 11
6 12
example_matrix %*% t(example_matrix2) ## Matrix multiplication
     [,1][,2]
       50 110
     500 1100
```

Lists

Lists are collections of arbitrary objects in R

Extracting elements from lists

```
example_list[[3]]

[1] TRUE TRUE FALSE

example_list[1:2]

[[1]] "Andy" "Brian" "Harry"

[[2]] [1] 12 16 16
```

Extracting elements from lists

```
example_list[[4]]
     [,1][,2][,3]
[1,]
[2,]
class(example_list[[4]])
[1] "matrix" "array"
example_list[[4]][1,]
[1] 1 1 1
```

Named lists

```
example_named_list[['Names']]
[1] "Andy" "Brian" "Harry"
example_named_list$Names
[1] "Andy" "Brian" "Harry"
example_named_list$Names[3]
[1] "Harry"
```

Back to a Data Frame

A data.frame object is a **named list** where each element is of the same length

You can use both *matrix* and *list* functions to operate on data.frame objects!!

head(data_spine)

```
Pelvic.incidence Pelvic.tilt Lumbar.lordosis.angle Sacral.slope Pelvic.radius
        63.02782
                   22.552586
                                          39.60912
                                                       40.47523
                                                                      98.67292
                                                       28.99596
        39.05695
                   10.060991
                                          25.01538
                                                                     114.40543
        68.83202
                   22.218482
                                          50.09219
                                                       46.61354
                                                                     105.98514
                                                       44.64413
                                                                     101.86850
        69.29701
                   24.652878
                                          44.31124
                                                       40.06078
                                                                     108.16872
        49.71286
                   9.652075
                                          28.31741
        40.25020
                   13.921907
                                          25.12495
                                                       26.32829
                                                                     130.32787
Degree.spondylolisthesis Pelvic.slope Direct.tilt Thoracic.slope
               -0.254400
                                          12.5661
                            0.7445035
                                                         14.5386
                4.564259
                            0.4151857
                                          12.8874
                                                         17.5323
                                          26.8343
               -3.530317
                            0.4748892
                                                         17.4861
               11.211523
                            0.3693453
                                          23.5603
                                                         12.7074
                            0.5433605
                                          35.4940
               7.918501
                                                         15.9546
               2.230652
                            0.7899929
                                          29.3230
                                                         12.0036
Cervical.tilt Sacrum.angle Scoliosis.slope Class.attribute
                                   43.5123
     15.30468
                -28.658501
                                                  Abnormal
    16.78486
                -25.530607
                                   16.1102
                                                  Abnormal
    16.65897
                -29.031888
                                   19.2221
                                                  Abnormal
    11.42447
                -30.470246
                                   18.8329
                                                  Abnormal
     8.87237
                -16.378376
                                   24.9171
                                                  Abnormal
     10.40462
                 -1.512209
                                    9.6548
                                                  Abnormal
```

```
dim(data_spine)

[1] 310    13

nrow(data_spine)

[1] 310

data_spine_small <- data_spine[1:4,] ## Matrix operation</pre>
```

```
data_spine_small[,2] ## Matrix extraction by position
```

```
[1] 22.55259 10.06099 22.21848 24.65288
```

data_spine_small[[2]] ## List extraction by position

[1] 22.55259 10.06099 22.21848 24.65288

```
data_spine_small[['Pelvic.tilt']] ## Named list extraction

[1] 22.55259 10.06099 22.21848 24.65288

data_spine_small[,'Pelvic.tilt'] ## Data frame named column extraction

[1] 22.55259 10.06099 22.21848 24.65288

data_spine_small$Pelvic.tilt ## Dollar sign extraction

[1] 22.55259 10.06099 22.21848 24.65288
```

My preference is for

- 1. data frame named column extraction data_spine_small[,'Pelvic.tilt'],
- 2. named list extraction data_spine_small[['Pelvic.tilt']]
- 3. Dollar-based extraction data_spine_small\$Pelvic.tilt

```
names(data_spine_small)
```

```
data_spine_small[,c('Pelvic.tilt', 'Pelvic.slope','Class.attribute')]
```

```
Pelvic.tilt Pelvic.slope Class.attribute
1 22.55259 0.7445035 Abnormal
2 10.06099 0.4151857 Abnormal
3 22.21848 0.4748892 Abnormal
4 24.65288 0.3693453 Abnormal
```