

SESSION 2

COMPUTING AND R

R FOR SOCIAL DATA SCIENCE

JEFFREY ZIEGLER, PHD

ASSISTANT PROFESSOR IN POLITICAL SCIENCE & DATA SCIENCE
TRINITY COLLEGE DUBLIN

FALL 2022

ROAD MAP FOR TODAY

- Backstory of R
- R operators and objects
- Data structures and types

R BACKGROUND



- **S** (for **S**tatistics) is a programming language for statistical analysis developed in 1976 in AT&T Bell Labs
- Original S language and its extension S-PLUS were closed source
- In 1991 **R**oss Ihaka and **R**obert Gentleman began developing R, an open-source alternative to **S**

R BASICS

- R is an *interpreted* language (like Python and Stata)
- It is geared towards statistical analysis
- R is often used for interactive data analysis (one command at a time)
- But it also permits to execute entire scripts in *batch* mode

```
1 print("Hello World!")
```

```
[1] "Hello World!"
```

OPERATORS

Key operators ("infix" functions) in R are:

- Arithmetic ('+', '-', '*', '^', '/', '%/%', '%%', '%*%')
- Boolean ('&', '&&', '|', '||', 'i')
- Relational ('==', '!=', '>', '>=', '<', '<=')
- Assignment ('<-', '«-', '=')
- Membership ('%in%')

BASIC MATHEMATICAL OPERATIONS IN R

```
1 1+1
```

```
[1] 2
```

```
1 5-3
```

```
[1] 2
```

```
1 6/2
```

```
[1] 3
```

```
1 4*4
```

```
[1] 16
```

ADVANCED MATHEMATICAL OPERATIONS IN R

```
1 # Let's try exponentiation
2 # Note that 2 ** 4 also works, but is not recommended 2^4
3 2^4
```

```
[1] 16
```

```
1 # Integer division, equivalent to Python's '//'
2 7%/3
```

```
[1] 2
```

```
1 # Modulo operation (remainder of division), equivalent to
  Python's '%'
2 7%%3
```

```
[1] 1
```

BASIC LOGICAL OPERATIONS IN R

```
1 3!=1 # Not equal
```

```
[1] TRUE
```

```
1 3>3 # Greater than
```

```
[1] FALSE
```

```
1 FALSE|TRUE # True if either 1st or 2nd operand is T
```

```
[1] TRUE
```

```
1 F|T # R also treats F and T as Boolean, not recommended
```

```
[1] TRUE
```

```
1 3>3|3>=3 # Combining 3 Boolean expressions
```

```
[1] TRUE
```


R OBJECTS

Everything is an object.

- John Chambers

- Everything you are dealing with in R is an **object**
- That includes individual variables, datasets, functions, and many other classes of objects
- The key reference to an object is its *name*
- Typically, reference is established through assignment operation

ASSIGNMENT OPERATIONS

- '`<-`' is standard assignment operator in R
- While '`=`' is also supported, not recommended
- As it hides the difference between '`<-`' and '`«-`' (deep assignment)

Extra: [R Documentation on assignment](#)

```
1  
2 x <- 3
```

```
[1] 3
```

```
1 f <- function(){  
2   x<-1 # Modifies existing var in parent namespace  
3 }  
4 f()  
5 x
```

```
[1] 1
```

MEMBERSHIP OPERATIONS

Operator `'%in%'` returns `'TRUE'` if an object of the left side is in a sequence on the right

```
1 "a" %in% "abc" # Note that R strings are not sequences
```

```
[1] FALSE
```

```
1 3 %in% c(1,2,3) # c(1, 2, 3) is a vector
```

```
[1] TRUE
```

```
1 !(3 %in% c(1,2,3))
```

```
[1] FALSE
```

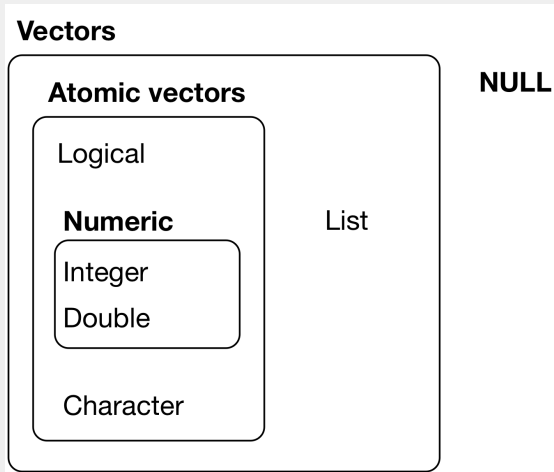
DATA STRUCTURES

- Base R data structures can be classified along their *dimensionality* and *homogeneity*
- 5 main built-in data structures in R:
 - ▶ Atomic vector ('vector')
 - ▶ Matrix ('matrix')
 - ▶ Array ('array')
 - ▶ List ('list')
 - ▶ Data frame ('data.frame')

SUMMARY OF DATA STRUCTURES IN R

Structure	Description	Dimensionality	Data Type
vector	Atomic vector (scalar)	1d	homogenous
matrix	Matrix	2d	homogenous
array	One-, two or n-dimensional array	1d/2d/nd	homogenous
list	List	1d	heterogeneous
data.frame	Rectangular data	2s	heterogeneous

VECTORS IN R



Source: R for Data Science

ATOMIC VECTORS

- *Vector* is core building block of R
- R has no scalars (they are just vectors of length 1)
- Vectors can be created with 'c()' function (short for concatenate)

```
1 v <- c(8,10,12)
2 v
```

```
[1] FALSE
```

```
1 v <- c(v,14) # Always flattened (even when nested)
2 v
```

```
[1] FALSE
```

DATA TYPES

4 common data types that are contained in R structures:

- Character ('character')
- Integer ('integer')
- Double/numeric ('double'/'numeric')
- Logical/boolean ('logical')

CHARACTER VECTOR

```
1 char_vec <- c("apple", "banana", "watermelon")  
2 char_vec
```

```
[1] "apple"      "banana"     "watermelon"
```

```
1 length(char_vec)
```

```
[1] 3
```

```
1 is.character(char_vec)
```

```
[1] TRUE
```

INTEGER VECTOR

```
1 # Note 'L' suffix to get an integer, not double
2 int_vec <- c(300L,200L,4L)
```

```
[1] 300 200 4
```

```
1 typeof(int_vec)
```

```
[1] "integer"
```

```
1 is.integer(int_vec)
```

```
[1] TRUE
```

DOUBLE VECTOR

```
1 # Note that even without decimal part R treats these  
   numbers as doubled
```

```
2 dbl_vec <- c(300,200,4)  
3 dbl_vec
```

```
[1] 300 200 4
```

```
1 typeof(dbl_vec)
```

```
[1] "double"
```

```
1 is.double(dbl_vec)
```

```
[1] TRUE
```

```
1 is.numeric(int_vec)
```

```
[1] TRUE
```

INTEGER VS DOUBLE

- Integers are used to store whole numbers (e.g. counts)
- Double is a floating-point numbers with double precisions
- 32-bit integer: $2^{32} = 4,294,967,296$
- Signed 32-bit integer: $[-2,147,483,648 \dots 2,147,483,648]$



Extra: [More on integer overflow on YouTube](#)

LOGICAL VECTOR

```
1 log_vec <- c(FALSE, FALSE, TRUE)
2 log_vec
```

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

```
1 # While more concise, using T/F instead of TRUE/FALSE can
  # be confusing
2 log_vec2 <- c(F, F, T)
3 log_vec2
```

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

```
1 typeof(log_vec)
```

```
[1] "logical"
```

TYPE COERCION IN VECTORS

- All elements of a vector must be of the same type
- If you try to combine vectors of different types, their elements will be *coerced* to most flexible type

```
1 # Note that logical vector get coerced to 0/1 for F/T
2 c(dbl_vec, log_vec)
```

```
[1] 300 200 4 0 0 1
```

```
1 c(char_vec, int_vec)
```

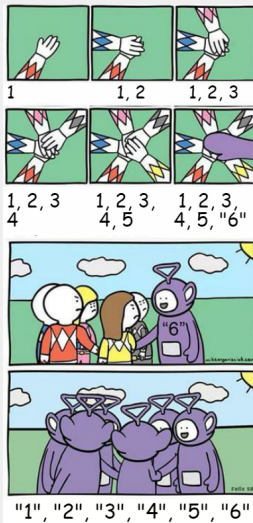
```
[1] "apple" "banana" "watermelon" "300" "200" "4"
```

```
1 # If no way of conversion exists, NAs are introduced
2 as.numeric(char_vec)
```

```
[1] NA NA NA
```

Warning message: NAs introduced by coercion

IMPLICIT TYPE COERCION



Extra: [Twitter](#)

NA AND NULL VALUES

- In Python we encountered 'None' value
- R makes a distinction between:
 - ▶ 'NA' - value exists, but is unknown (e.g. survey non-response)
 - ▶ 'NULL' - object does not exist
- 'NA's are defined for each data type (integer, character, numeric, etc.))

Extra: [R Documentation on NA](#)

NA AND NULL EXAMPLE

```
1 na_vec <- c(NA, NA, NA)
2 na_vec
```

```
[1] NA NA NA
```

```
1 length(na_vec)
```

```
[1] 3
```

```
1 null_vec <- c(NULL, NULL, NULL)
2 null_vec
```

```
[1] NULL
```

```
1 length(null_vec)
```

```
[1] 0
```

WORKING WITH NAs

```
1 v_na<-c(1,2,3,NA,5)
2 mean(v_na)
```

```
[1] NA
```

```
1 # NAs should be treated specially
2 mean(v_na, na.rm=TRUE)
```

```
[1] 2.75
```

```
1 # Remember NAs are missing values
2 NA==NA
```

```
[1] NA
```

WORKING WITH NAs

```
1 is.na(v_na)
```

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

```
1 # We can use such logical vectors for subsetting (more  
  below)
```

```
2 v_na[!is.na(v_na)]
```

```
[1] 1 2 3 5
```

"TUTORIAL": R SCRIPT

- Usually you want to have a record of what analysis was done and how you did it
- So, instead of writing all your R commands in the interactive console
- You can create an R script, write them there and run them together or one at a time
- R script is a file with .R extension and contains a collection of valid R commands

NAMING CONVENTIONS

- Even while allowed in R, do not use “.” in variable names (it works as an object attribute in Python)
- Do not name give objects the names of existing functions and variables (e.g. ‘c’, ‘T’, ‘list’, ‘mean’)
- Use **UPPER_CASE_WITH_UNDERSCORE** for named constants (e.g. variables that remain fixed and unmodified)
- Use **lower_case_with_underscores** for function and variable names

Extra: <http://adv-r.had.co.nz/Style.html>

CODE LAYOUT

- Limit all lines to a maximum of 79 characters
- Break up longer lines

```
my_long_vector <- c(
  1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22,
  23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41,
  42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60
)

long_function_name <- function(a = "a long argument",
  b = "another argument",
  c = "another long argument") {
  # As usual code is indented by two spaces
}
```

RESERVED WORDS

There are 14 (plus some variations of them) reserved words in R that cannot be used as identifiers

- | | | |
|-----------|--------------|-----------|
| ■ 'break' | ■ 'next' | ■ 'if' |
| ■ 'NA' | ■ 'for' | ■ 'TRUE' |
| ■ 'else' | ■ 'NULL' | ■ 'Inf' |
| ■ 'NaN' | ■ 'function' | ■ 'while' |
| ■ 'FALSE' | ■ 'repeat' | |

EXERCISE 1: VECTOR MANIPULATION

- Load built-in R object 'letters' (lower-case letters of the Roman alphabet)
- Calculate its length
- Generate a vector of integers that starts from 1 and has the same length as 'letters'
- Assign each integer corresponding lower-case letter as name
- Use these names to subset all vowels
- Now, repeat subsetting, but using indices not names

TABULATION AND CROSSTABULATION IN R

- R function 'table()' provides an easy way of summarizing categorical variables
- Note that implicitly variables represented as character vectors are converted to factors"

```
1 settlements <- c("Dublin", "Belfast", "Cork", "Limerick", "Derry", "Galway", "Newtownabbey", "Bangor", "Waterford", "Lisburn")
2 provinces <- c("Leinster", "Ulster", "Munster", "Munster", "Ulster", "Connacht", "Ulster", "Ulster", "Munster", "Ulster")
3 # Given that each town appears only once, crosstabs might not be informative
4 table(top_10_settlements, provinces)
```

Source: Top 10 most populous settlements on the island of Ireland

EXERCISE 2: WORKING WITH ATTRIBUTES AND FACTORS

- As you note the output of `'table(provinces)'` is sorted alphabetically
- Change this to reflect the actual counts
- First, let's store the result of tabulation for later re-use
- Start from exploring the structure of this object with `'str()'`
- What are the 2 main parts of this object? How are they stored?

EXERCISE 2: WORKING WITH ATTRIBUTES AND FACTORS

- Extract the relevant parts from the stored object
- Save them as a named vector with provinces as names and counts as values
- Use 'sort()' function to sort the vector in a decreasing order (from largest to smallest)
- Convert the original 'provinces' vector into a factor with the levels ordered accordingly
- Re-run 'table(provinces)'

WEEK 2 "TAKEHOME" EXERCISE

- Save a 'letters' object under a different name
- Convert saved object into a matrix of 13 rows and 2 columns
- Subset letter 'f' using indices
- Concatenate 3 copies of 'letters' object together in a single character vector
- Convert it into a 3-dimensional array, where each dimension appears as a matrix above
- Subset all letters 'f' across all 3 dimensions "

CLASS BUSINESS

Today, we talked about...

- Backstory
- R operators and objects
- Data structures and types

Next week...

- Indexing and subsetting
- Attributes