# SESSION 12 DATA WRANGLING 2

R FOR SOCIAL DATA SCIENCE

JEFFREY ZIEGLER, PHD

ASSISTANT PROFESSOR IN POLITICAL SCIENCE & DATA SCIENCE TRINITY COLLEGE DUBLIN

**FALL 2022** 

## **ROAD MAP FOR TODAY**

#### Last time:

- Data frames in base R
- Data input and output

#### This time:

- Alternatives to data frames
- 'tidyverse' packages
- Working with tabular data
- Summary statistics

#### **ALTERNATIVES TO DATA FRAME**

- Two major alternatives to data frames are:
  - 'tibble' from 'tibble' package (part of 'tidyverse' package ecosystem)
  - ▶ 'data.table' from 'data.table'
- 'tibble' provides features enhancing user experience (readability, ease of manipulation)
- 'data.table' provides speed

## DATA TABLE - FAST DATA FRAME

- As opposed to data frames, data tables are updated by reference
- Frees up a lot of RAM for big data!
- Provides low-level parallelism
- SQL-like operations for data manipulation
- Has no external dependencies (other than base R itself)

## 'TIDYVERSE' PACKAGES

- 'tidyverse' package ecosystem rich collection of data science packages
- Designed with consistent interfaces and generally higher usability than base R function
- Notable packages:
  - 'readr' data input/output (also 'readxl' for spreadsheets, 'haven' for SPSS/Stata)
  - 'dplyr' data manipulation (also 'tidyr' for pivoting)
  - ▶ 'ggplot2' data visualisation
  - ► 'lubridate' working with dates and time
  - ▶ 'tibble' enhanced data frame

install.packages("tidyverse")

#### TIBBLE: USER-FRIENDLY DATA FRAME

- Tibbles are designed to be backward compatible with base R data frames
- Console printing of tibbles is cleaner (prettified, only first 10 rows by default)
- Tibbles can have columns that themselves contain lists as elements
- Tibbles can be created with 'tibble::tibble()' function
- Or objects can be coerced into a tibble using 'tibble::as\_tibble()' function

## **TIBBLE: USER-FRIENDLY DATA FRAME**

```
tb <- tibble::tibble(
   X = 1:4.
          y = c("a", "b", "c", "d"),
          z = c(TRUE, FALSE, FALSE, TRUE)
t h
# A tibble: 4 × 3
    6v
<int> <chr> <lgl>
     1 a TRUE
2 2 b FALSE
3 3 C FALSE
   4 d
            TRUE
```

# TIBBLES WORK (MOSTLY) LIKE DATA FRAMES

```
str(tb)
tibble [4 × 3] (S3: tbl_df/tbl/data.frame)
$ x: int [1:4] 1 2 3 4
$ y: chr [1:4] "a" "b" "c" "d"
$ z: logi [1:4] TRUE FALSE FALSE TRUE
dim(tb)
[1] 4 3
tb[c("x", "z")]
# A tibble: 4 × 2
ΧZ
<int> <lgl>
      1 TRUE
2 2 FALSE
  3 FALSE
  4 TRUE
tb[tb$y=="b", ]
# A tibble: 1 × 3
<int> <chr> <lgl>
       2 b
                FALSE
```

## DATA MANIPULATION WITH 'DPLYR'

- 'dplyr' core package for data manipulation in 'tidyverse'
- Principal functions are:
  - ► 'filter()' subset rows from data
  - 'mutate()' add new/modify existing variables
  - ► 'rename()' rename existing variable
  - 'select()' subset columns from data
  - ► 'arrange()' order data by some variable
- For data summary:
  - ► 'group\_by()' aggregate data by some variable
  - 'summarise()' create a summary of aggregated variables

library("dplyr")

## SUBSETTING WITH 'DPLYR'

```
dplyr::filter(tb, y == 'b', z == FALSE)
# A tibble: 1 × 3
x y z
<int> <chr> <lgl>
     2 b FALSE
# Note that dplyr functions do not require enquoted variable names
dplyr::select(tb, x, z)
# A tibble: 4 × 2
<int> <lgl>
     1 TRUE
2
     2 FALSE
    3 FALSE
  4 TRUE
# We can also use helpful tidyselect functions for more complex rules
dplyr::select(tb, tidyselect::starts with('x'))
# A tibble: 4 × 1
<int>
      1
2
      2
3
```

## RENAMING/MODIFYING COLUMNS WITH 'DPLYR'

```
# Data is not modified in-place, you need to re-assign the results
tb <- dplyr::rename(tb, random = x)
dplyr::mutate(tb, random_8plus = ifelse(random >= 3, TRUE, FALSE))

# A tibble: 4 × 4
random y z random_8plus
<int> <chr> <lgl> <lgl> 1
1
2
2
3
3
5
ALSE
4
4
4
4
5
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
7
<
```

## **'%>%' OPERATOR**

- Users of 'tidyverse' packages are encouraged to use pipe operator '%>%'
- Allows to chain data transformations without creating intermidate variables
- Passes result of previous operation as a first first argument to next
- Base R now also includes its own pipe operator '|>' but it's still relatively uncommon

```
<result> <- <input> %>%
<function_name>(., arg_1, arg_2, ..., arg_n)
<result> <- <input> %>%
<function_name>(arg_1, arg_2, ..., arg_n)
```

## **'%>%' OPERATOR EXAMPLE**

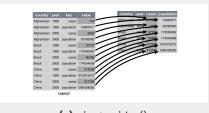
2 3 c FALSE 1.68

# '%>%' OPERATOR VS BUILT-IN '|>' OPERATOR

■ Since R version 4.1.0 (mid-2021), there is a built-in '|>' pipe operator

## **PIVOTING DATA**

- Sometimes you want to pivot you data by:
  - Spreading some variable across columns ('tidyr::pivot\_wider()')
  - Gathering some columns in a variable pair ('tidyr::pivot\_longer()')







(b) pivot\_longer()

Source: R for Data Science

## **PIVOTING DATA EXAMPLES**

```
th2 <- tibble::tibble(
  country = c("Afghanistan", "Brazil"),
'1999' = c(745, 2666),
 (2000) = (37737, 80488)
th2
# A tibble: 2 × 3
country '1999' '2000'
<chr> <dbl> <dbl>
1 Afghanistan 745 37737
2 Brazil 2666 80488
# Note that pivoting functions come 'tidyr' package
th2 <- th2 %>%
tidyr::pivot_longer(cols = c("1999", "2000"), names_to = "year", values_to = "cases")
tb2
# A tibble: 4 × 3
country vear cases
<chr> <chr> <chr> <chr> <dbl>
1 Afghanistan 1999 745
2 Afghanistan 2000 37737
3 Brazil 1999 2666
4 Brazil 2000 80488
```

## **PIVOTING DATA EXAMPLES**

```
tb2 <-tb2 %%

tidyr::pivot_wider(names_from="year", values_from="cases")

# A tibble: 4 × 3

country year cases

<chr> <chr> <chr> <dbl>
1 Afghanistan 1999 745

2 Afghanistan 2000 37737

3 Brazil 1999 2666

4 Brazil 2000 80488
```

## **SUMMARIZING NUMERIC VARIABLES**

## SUMMARIZING CATEGORICAL VARIABLES

```
# table() function is flexible, can tabulate a single variable and do crosstabs
table(kaggle2021[3])
Man
     Nonbinary
                     Prefer not to say Prefer to self-describe
                                                                 Woman
20598
          88
                                 355
                                        42
                                                     4890
# Wrapping it inside prop.table() gives proportions of each category
prop.table(table(kaggle2021[3]))
     Nonbinary Prefer not to say Prefer to self-describe
Man
                                                             Woman
0.793054326 0.003388134 0.013668040 0.001617064
                                                              0.188272437
# Wrapping it inside sort() gives value sorting, as opposed to alphabetic (or facto
     levels)
sort(table(kaggle2021[3]), decreasing = TRUE)[1]
Man
20598
```

## **TUTORIAL - TABULATION AND SUMMARY STATS**

- Load 'kaggle\_survey\_2021\_responses.csv' dataframe from GitHub repository to global environment
- Do so from the url, and from the local file
- Consider country of residence reported by respondents (question Q3)
- Make sure you can select the column both using both it name and index
- Calculate the percentages of top 3 countries of residence in the sample

#### **TUTORIAL - PIVOTING TABLES**

- To simplify working with the dataset, let's create a unique id for each respondent (you can use 'seq\_along()' function in combination with any other variable to do so)
- Finally, use 'pivot\_wider' function from 'tidyr' package to create a separate column for each age group
- If original pivoting produced columns that are populated by values of the categorical variable and 'NA"s, use 'mutate' function to replace them with o's and 1's
- Finally, use 'pivot\_longer' function to convert this representation of the dataset back into its original form
- You might also need to use 'dplyr::filter()' function to remove redundant rows

## **OVERVIEW**

#### This week:

- Data frames in base R
- Data input and output
- Alternatives to data frames
- 'tidyverse' packages
- Working with tabular data
- Summary statistics

#### Next week:

■ Data visualization