

# Tell me which packages you use and I'll tell you who you are

tidyverse vs data.table vs base R

Anna Skrzydło & Bartosz Kowalski

## WhyR? 2018

04 | 07 | 2018



What happens when you use packages from 10<sup>th</sup> page of google results...

ywność koloru oznacza wielkość inwestycji me

Error: could not find function "ggplotify"



# If you know only base R, tidyverse might surprise you...

```
1 library(tidyverse)
2 library(MBSWarsawEconometricTools)
3
4 # Read dictionaries
5 country.dict.df <- read_csv2('02_Working/country_of_residence_dictionary.csv')
6 type.dict.df <- read_csv2('02_Working/type_dictionary.csv')
7
8 # Prepare monthly data
9 tourism.df <- read_csv('01_Input/DCSC_TUR_28062018103715289.csv') %>%
10   filter('Country of residence of guests' %in% c('Foreign countries', 'Italy', 'All countries
11   filter(TIPO_ALLOGGIO2 == 'ALL') %>%
12   select(Indicators, 'Country of residence of guests', TIME, Value) %>%
13   mutate('Code var type' = 'TU') %>%
14   left_join(country.dict.df, by = "Country of residence of guests") %>%
15   left_join(type.dict.df, by = "Indicators") %>%
16   mutate(Variable = paste('Code var type', 'Code type', 'Code country', sep = '_')) %>%
17   select(TIME, Variable, Value) %>%
18   spread(Variable, Value) %>%
19   mutate(Start.date = as.Date(paste0(TIME, '-01'))) %>%
20   mutate(End.date = as.Date(lead(Start.date) - 1)) %>%
21   select(Start.date, End.date, everything(), -TIME)
22
23 tourism.df$End.date[nrow(tourism.df)] <- as.Date('2018-03-31')
```

```
1 library(MBSWarsawEconometricTools)
2
3 # Read dictionaries
4 country.dict.df <- read_csv2('02_Working/country_of_residence_dictionary.csv', stringsAsFactors = FALSE)
5 type.dict.df <- read_csv2('02_Working/type_dictionary.csv', stringsAsFactors = FALSE)
6
7 tourism.raw.df <- read_csv('01_Input/DCSC_TUR_28062018103715289.csv', stringsAsFactors = FALSE)
8
9 tourism.df <- tourism.raw.df[tourism.raw.df$Country.of.residence.of.guests %in% c('Foreign countries', 'Italy', 'All countries of the world'), ]
10 tourism.df <- tourism.raw.df[tourism.raw.df$TIPO_ALLOGGIO2 == 'ALL', ]
11 tourism.df <- tourism.df[, c('Indicators', 'Country.of.residence.of.guests', 'TIME', 'Value')]
12 tourism.df$Code.var.type <- 'TU'
13 tourism.df <- merge(tourism.df, country.dict.df, by = "Country.of.residence.of.guests")
14 tourism.df <- merge(tourism.df, type.dict.df, by = "Indicators")
15 tourism.df$Variable <- paste(tourism.df$Code.var.type, tourism.df$Code.type, tourism.df$Code.country, sep = '_')
16 tourism.df <- tourism.df[, c('TIME', 'Variable', 'Value')]
17
18 tourism.monthly.df <- data.frame(TIME = unique(tourism.df$TIME))
19 tourism.monthly.df <- tourism.monthly.df[order(tourism.monthly.df$TIME), , drop = FALSE]
20 tourism.monthly.df$TIME <- as.character(tourism.monthly.df$TIME)
21
22 for (var in unique(tourism.df$Variable)) {
23   current.var.df <- tourism.df[tourism.df$Variable == var, c('TIME', 'Value')]
24   colnames(current.var.df) <- c('TIME', var)
25   tourism.monthly.df <- merge(tourism.monthly.df, current.var.df, by = 'TIME')
26 }
27
28 tourism.monthly.df$Start.date <- as.Date(paste0(tourism.monthly.df$TIME, '-01'))
29 tourism.monthly.df$End.date <- c(as.Date(tourism.monthly.df$Start.date[2:nrow(tourism.monthly.df)]) - 1, as.Date('2018-03-31'))
30
31 tourism.monthly.df <- tourism.monthly.df[, c('Start.date', 'End.date', setdiff(colnames(tourism.monthly.df), c('Start.date', 'End.date', 'TIME')))]
32 |
```

# Tidyverse is faster than base R

## Reading a data.frame with 468 observations and 5 575 variables

```
library(microbenchmark)
file <- "speed_case/Sales_Nielsen.csv"

microbenchmark(base=read.csv2(file),
               tidy=suppressMessages(readr::read_csv2(file)), times = 25L)
```

Unit: seconds

expr	min	1q	mean	median	uq	max	neval
<b>base</b>	5.06	5.26	5.56	<b>5.47</b>	5.60	7.29	25
<b>tidy</b>	2.25	2.65	2.99	<b>3.07</b>	3.23	4.10	25

## ■ If speed is a criterion... Go for data.table!

Reading a data.frame with 468 observations and 5 575 variables

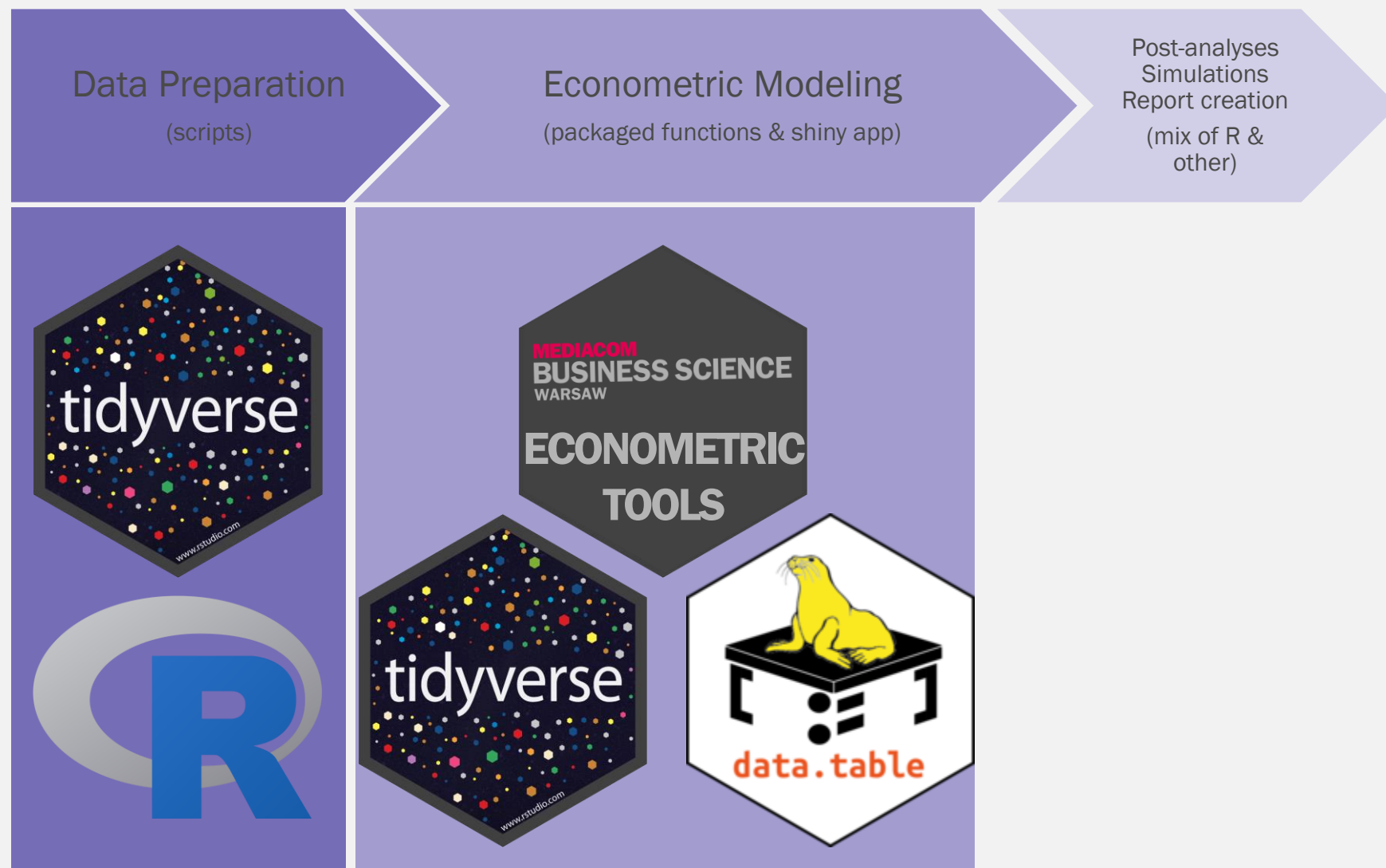
```
microbenchmark(base=read.csv2(file),  
               tidy=suppressMessages(readr::read_csv2(file)),  
               dt=data.table::fread(file, sep = ";", dec = ",", showProgress=FALSE), times = 25L)
```

Unit: milliseconds

expr	min	1q	mean	median	uq	max	neval
<b>base</b>	5049	5175	5487	<b>5315</b>	5569	6705	25
<b>tidy</b>	2324	2460	2894	<b>2795</b>	3050	4035	25
<b>dt</b>	181	186	228	<b>196</b>	269	393	25

# Where we are today

## R framework for Marketing Mix Modeling at MediaCom Business Science



## Our take on base vs Tidyverse vs data.table



- Everyone in the community knows it
- Stable, no sudden changes to functions structure

- Structures and orders thinking by providing a limited number of options (functions)
- When combined with meaningful naming convention for R objects, the code is extremely easy to understand
- Quick to learn by a non-coder

- It's FAST\*
- Memory efficiency
- Compact syntax;



- Outdated default arguments
- Slow to run


- The syntax is quite verbose. Thus, it takes longer to run
- Even though speed is getting better over time, it still underperforms in some cases, even on smaller datasets vs data.table. This includes reading data (read\_csv vs fread)

- The syntax is not self-explaining
- We find data.table more timely to learn



A person is rappelling down a dark, craggy rock face. The scene is set against a dramatic sunset sky with orange, yellow, and blue hues. In the background, a large, dark mountain peak rises from a body of water. The overall mood is adventurous and inspiring.

**MEDIACOM**  
**BUSINESS SCIENCE**  
WARSAW

 @MBSWarsaw  
anna.skrzydlo@mediacom.com  
bartosz.kowalski@mediacom.com

“

***An investment in knowledge  
Always pays the best interest***

Benjamin Franklin