

Data Wrangling

Fundamentals of Computing and Data Display

Christoph Kern Ruben Bach¹

(c.kern, r.bach)@uni-mannheim.de

¹Slides draw on Golemund, G (2018). Master the Tidyverse: An Introduction to R for Data Science.
<https://github.com/rstudio/master-the-tidyverse>

Outline

- 1 Introduction
 - The tidyverse
- 2 Tidy data
 - tibble
 - tidyr
- 3 Split-Apply-Combine
 - apply
 - plyr
 - dplyr
- 4 joins
- 5 Resources

Introduction

Base R

- “R is a free software environment for statistical computing and graphics.”
- Focus on flexibility and the ability to program
- Leans towards developing rather than basic data analysis

```
> # Collapsing data the R way  
> cc_agg <- aggregate(cc_sub[, 2:5],  
>   by = list(cc_sub$district),  
>   FUN = sum)
```

The tidyverse

“A collection of modern R packages that share common philosophies, embed best practices, and are designed to work together”



<https://www.tidyverse.org/>

The tidyverse

<code>install.packages("tidyverse")</code>	<code>library("tidyverse")</code>
<code>install.packages("ggplot2")</code>	<code>library("ggplot2")</code>
<code>install.packages("dplyr")</code>	<code>library("dplyr")</code>
<code>install.packages("tidyr")</code>	<code>library("tidyr")</code>
<code>install.packages("readr")</code>	<code>library("readr")</code>
<code>install.packages("purrr")</code>	<code>library("purrr")</code>
<code>install.packages("tibble")</code>	<code>library("tibble")</code>
<code>install.packages("stringr")</code>	<code>library("stringr")</code>
<code>install.packages("forcats")</code>	<code>library("forcats")</code>
<code>install.packages("hms")</code>	
<code>install.packages("lubridate")</code>	
<code>install.packages("DBI")</code>	
<code>install.packages("haven")</code>	
<code>install.packages("httr")</code>	
<code>install.packages("jsonlite")</code>	
<code>install.packages("readxl")</code>	
<code>install.packages("rvest")</code>	
<code>install.packages("xml2")</code>	
<code>install.packages("modelr")</code>	
<code>install.packages("broom")</code>	

The tidyverse

Core tidyverse

- `ggplot2`, for data visualisation
- `dplyr`, for data manipulation
- `tidyr`, for data tidying
- `readr`, for data import
- `purrr`, for functional programming
- `tibble`, for tibbles
- `stringr`, for strings
- `forcats`, for factors

- 1 Introduction
 - The tidyverse
- 2 **Tidy data**
 - tibble
 - tidyr
- 3 Split-Apply-Combine
 - apply
 - plyr
 - dplyr
- 4 joins
- 5 Resources

Tidy data

Many ways to structure the same underlying data

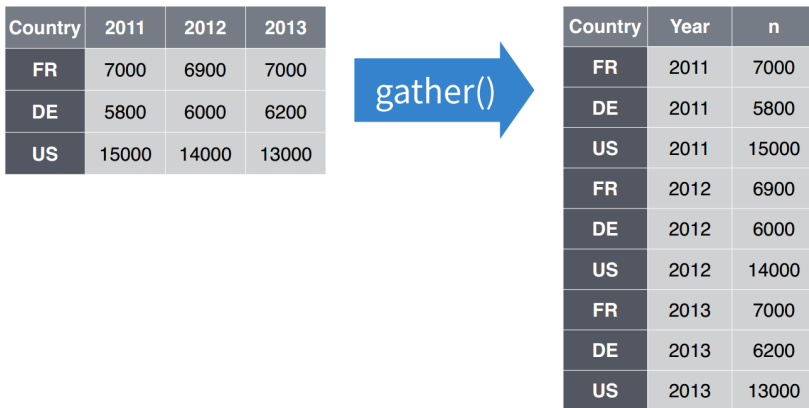
- R prefers **tidy data**
 - ① Each variable is in its own column
 - ② Each case is in its own row
 - ③ Each value is in its own cell
- Advantages of tidy data
 - ① Consistent way of storing data eases data transformations and reshaping
 - ② Utilizes vectorized nature of many (built-in) R functions
 - ③ Tidyverse designed to work with tidy data

tibble

The natural habitat of tidy data

- A Tibble: Tidyverse version of a `data.frame`
 - Print method shows only the first 10 rows
 - Subsetting with `[` returns a tibble, `[[`, `$` return vectors
- Converting into tibbles: `as_tibble()`
 - Prevents type conversion, renaming, `row.names`, error prone recycling
 - Tidyverse functions typically return tibbles
- Enhanced `str()`: `glimpse()`
- `vignette("tibble")`

Figure: Reshape from wide to long format



```
> # Reshape with gather()
> gather(cases, key = "year", value = "n", 2:4)
```

Figure: Reshape from wide to long format

Country	2011	2012	2013
FR	7000	6900	7000
DE	5800	6000	6200
US	15000	14000	13000

Country	Year	n
FR	2011	7000
DE	2011	5800
US	2011	15000
FR	2012	6900
DE	2012	6000
US	2012	14000
FR	2013	7000
DE	2013	6200
US	2013	13000

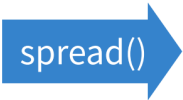
Figure: Reshape from wide to long format

Country	2011	2012	2013
FR	7000	6900	7000
DE	5800	6000	6200
US	15000	14000	13000

Country	Year	n
FR	2011	7000
DE	2011	5800
US	2011	15000
FR	2012	6900
DE	2012	6000
US	2012	14000
FR	2013	7000
DE	2013	6200
US	2013	13000

Figure: Reshape from long to wide format

city	size	amount
New York	large	23
New York	small	14
London	large	22
London	small	16
Beijing	large	121
Beijing	small	56



city	large	small
New York	23	14
London	22	16
Beijing	121	56

```
> # Reshape with spread()
> spread(pollution, key = size, value = amount)
```

Figure: Reshape from long to wide format

city	size	amount
New York	large	23
New York	small	14
London	large	22
London	small	16
Beijing	large	121
Beijing	small	56

city	large	small
New York	23	14
London	22	16
Beijing	121	56

Figure: Reshape from long to wide format

city	size	amount
New York	large	23
New York	small	14
London	large	22
London	small	16
Beijing	large	121
Beijing	small	56

city	large	small
New York	23	14
London	22	16
Beijing	121	56

Figure: Split a column by dividing values at a specific character


country <chr>	year <int>	codes <chr>	new <chr>	type <chr>	sexage <chr>	n <int>
Afghanistan	1980	new_sp_m014	new	sp	m014	NA
Afghanistan	1981	new_sp_m014	new	sp	m014	NA
Afghanistan	1982	new_sp_m014	new	sp	m014	NA
Afghanistan	1983	new_sp_m014	new	sp	m014	NA
Afghanistan	1984	new_sp_m014	new	sp	m014	NA
Afghanistan	1985	new_sp_m014	new	sp	m014	NA
Afghanistan	1986	new_sp_m014	new	sp	m014	NA
Afghanistan	1987	new_sp_m014	new	sp	m014	NA
Afghanistan	1988	new_sp_m014	new	sp	m014	NA
Afghanistan	1989	new_sp_m014	new	sp	m014	NA

```
> # Split column with separate()
> separate(who$codes, into = c("new", "type", "sexage"), sep = "_")
> # separate(who$codes, c("new", "type", "sexage"), sep = c(4, 7))
```


Figure: Drop rows that contain NA's (in the specified columns)

\bar{x}

x1	x2
A	1
B	NA
C	NA
D	3
E	NA



x1	x2
A	1
D	3

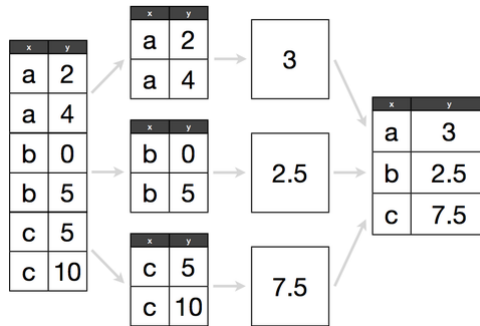
```
> # Listwise deletion with drop_na()
> drop_na(x, x2)
> # drop_na(x)
```

- 1 Introduction
 - The tidyverse
- 2 Tidy data
 - tibble
 - tidyr
- 3 Split-Apply-Combine**
 - apply
 - plyr
 - dplyr
- 4 joins
- 5 Resources

Split-Apply-Combine

- ① “split” problem into smaller pieces
- ② “apply” function to each piece separately
- ③ “combine” results back together

- In R, loops are not ideal for this problem
 - Slow, hard to read, much to code
- Vectorized functions, `plyr`, `dplyr`
 - Faster, easier to read, draw e.g. on C++



apply

The apply family (Base R)

- `apply(object, n, function)`
 - Apply a function to the rows ($n = 1$) or columns ($n = 2$) of an object
 - Input: array, matrix
 - Output: array, list
- “List apply”: `lapply(object, function)`
 - Apply a function to the columns of an object
 - Input: list, data.frame
 - Output: list
- “Simplified apply”: `sapply(object, function)`
 - Apply a function to the columns of an object
 - Input: list, data.frame
 - Output: vector, matrix
- `vapply()`, `mapply()`, `rapply()`

apply

Apply by groups (Base R)

- `tapply(object, factor, function)`
 - Apply a function to a vector grouped by factor levels
 - Input: vector
 - Output: array
- `by(object, factor, function)`
 - Apply a function to the columns of an object by factor levels
 - Input: data.frame, matrix
 - Output: list, array

```
> # Apply example from Advanced R book
> fix_missing <- function(x) {
>   x[x == -99] <- NA
>   x
> }
> df[] <- lapply(df, fix_missing)
```

plyr

Split-apply-combine with `plyr`

- Consistent naming schema, parallel processing
- `d*ply(.data, .variables, .fun, ...)`
 - First argument is data to be processed
 - Second argument describes how to split up input
 - Third argument is processing function

Table: The plyr family

Input	Output			
	Array	Data frame	List	Discarded
Array	<code>aaply</code>	<code>adply</code>	<code>alply</code>	<code>a_ply</code>
Data frame	<code>daply</code>	<code>ddply</code>	<code>dlply</code>	<code>d_ply</code>
List	<code>laply</code>	<code>ldply</code>	<code>llply</code>	<code>l_ply</code>

- 1 Introduction
 - The tidyverse
- 2 Tidy data
 - tibble
 - tidyr
- 3 Split-Apply-Combine
 - apply
 - plyr
 - **dplyr**
- 4 joins
- 5 Resources

dplyr

A fast `plyr` for data frames: **dplyr**

- Optimized for speed
- Consistent function names and interfaces
- Results are returned as tibbles
- Powerful in connection with piping



<https://dplyr.tidyverse.org/>

dplyr

Base R data transformation workflows

- ① Save each intermediate step as a new object
- ② Overwrite the original object many times
- ③ Compose functions

dplyr

- Provides the pipe operator `%>%` to pass results to next function
 - Run a sequence of transformations in one code block
 - Readable structure of data operations
 - Easy to modify, add and remove steps
- Left-hand side is piped in as first argument of right-hand side

```
> who %>%  
>   gather("codes", "n", 5:60) %>%  
>   separate(codes, into = c("new", "type", "sexage"), sep = "_")
```

dplyr

Base R data transformation workflows

- ① Save each intermediate step as a new object
- ② Overwrite the original object many times
- ③ Compose functions

dplyr

- Provides the pipe operator `%>%` to pass results to next function
 - Run a sequence of transformations in one code block
 - Readable structure of data operations
 - Easy to modify, add and remove steps
- Left-hand side is piped in as first argument of right-hand side

```
> who %>%  
>   gather("codes", "n", 5:60) %>%  
>   separate(codes, into = c("new", "type", "sexage"), sep = "_")
```

dplyr

Other pipe operators from **magrittr**

- **%T>%**
 - Returns the left-hand side value (not right-hand side result)
- **%%\$%**
 - Allows to refer to individual vectors (w/o data.frame context)
- **%<>%**
 - Assign result of pipeline to object

```
# equivalent to who %<>% gather...  
> who <- who %>%  
>   gather("codes", "n", 5:60) %>%  
>   separate(codes, into = c("new", "type", "sexage"), sep = "_")
```


Figure: Select columns

year	sex	name	n	prop		name	prop
1880	M	John	9655	0.0815	→	John	0.0815
1880	M	William	9532	0.0805		William	0.0805
1880	M	James	5927	0.0501		James	0.0501
1880	M	Charles	5348	0.0451		Charles	0.0451
1880	M	Garrett	13	0.0001		Garrett	0.0001
1881	M	John	8769	0.081		John	0.081

```
> # Select columns with select()
> select(babynames, name, prop)
> # select range of cols with :, exclude cols with -
```

Figure: Filter rows

year	sex	name	n	prop
1880	M	John	9655	0.0815
1880	M	William	9532	0.0805
1880	M	James	5927	0.0501
1880	M	Charles	5348	0.0451
1880	M	Garrett	13	0.0001
1881	M	John	8769	0.081



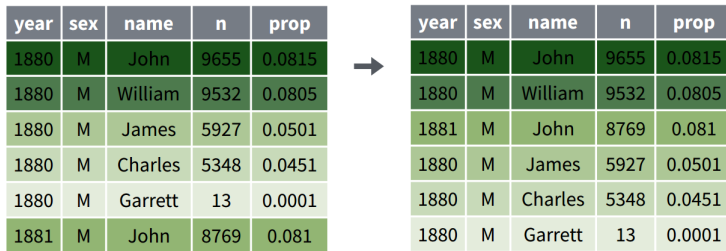
year	sex	name	n	prop
1880	M	Garrett	13	0.0001

```
> # Filter rows with filter()
> filter(babynames, name == "Garrett" & year == 1880)
```

Table: Operators

?comparison	?Logic
$x < y$	$a \& b$
$x > y$	$a \mid b$
$x == y$	$\text{xor}(a,b)$
$x \leq y$	$!a$
$x \geq y$	
$x \neq y$	
$x \%in\% y$	
$\text{is.na}(x)$	
$! \text{is.na}(x)$	

Figure: Order rows



The diagram illustrates the effect of the `arrange()` function with `desc(n)`. It shows two data tables. The left table is the original data, and the right table is the result of sorting by the `n` column in descending order. An arrow points from the original table to the sorted table.

year	sex	name	n	prop
1880	M	John	9655	0.0815
1880	M	William	9532	0.0805
1880	M	James	5927	0.0501
1880	M	Charles	5348	0.0451
1880	M	Garrett	13	0.0001
1881	M	John	8769	0.081

year	sex	name	n	prop
1880	M	John	9655	0.0815
1880	M	William	9532	0.0805
1881	M	John	8769	0.081
1880	M	James	5927	0.0501
1880	M	Charles	5348	0.0451
1880	M	Garrett	13	0.0001

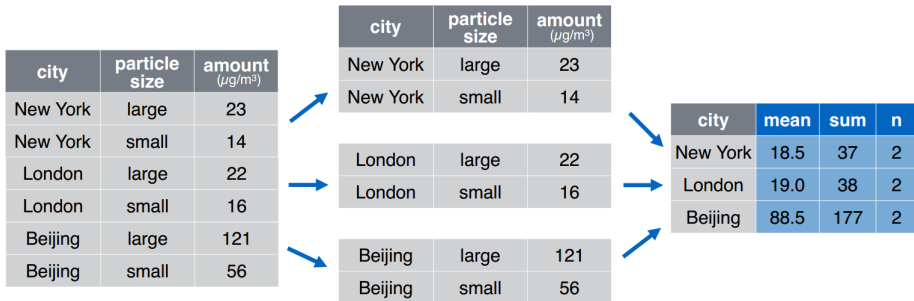
```
> # Order rows with arrange()  
> arrange(babynames, desc(n))
```

Figure: Summarize data

year	sex	name	n	prop		
1880	M	John	9655	0.0815	→	
1880	M	William	9532	0.0805		
1880	M	James	5927	0.0501		
1880	M	Charles	5348	0.0451		
1880	M	Garrett	13	0.0001		
1881	M	John	8769	0.081		
					total	max
					127538	99680

```
> # Summarize data with summarise()  
> babynames %>%  
>   summarise(total = sum(n), max = max(n))
```



Figure: Split-apply-combine with dplyr



```
> # Split data by groups with group_by()
> pollution %>%
>   group_by(city) %>%
>   summarise(mean = mean(amount), sum = sum(amount), n = n())
```

Figure: Create new variables

year	sex	name	n	prop	
1880	M	John	9655	0.0815	
1880	M	William	9532	0.0805	
1880	M	James	5927	0.0501	
1880	M	Charles	5348	0.0451	
1880	M	Garrett	13	0.0001	
1881	M	John	8769	0.081	



year	sex	name	n	prop	percent
1880	M	John	9655	0.0815	8.15
1880	M	William	9532	0.0805	8.05
1880	M	James	5927	0.0501	5.01
1880	M	Charles	5348	0.0451	4.51
1880	M	Garrett	13	0.0001	0.01
1881	M	John	8769	0.081	8.1

```
> # Create new variables with mutate()  
> babynames %>%  
>   mutate(percent = round(prop*100, 2))
```

dplyr

Data analysis workflows

```
> # Temporal objects
> boys_2015 <- filter(babynames, year == 2015, sex == "M")
> boys_2015 <- select(boys_2015, name, n)
> boys_2015 <- arrange(boys_2015, desc(n))
> boys_2015
```

```
> # Nesting functions
> arrange(select(filter(babynames, year == 2015, sex == "M"), name, n), desc(n))
```

```
> # Piping
> babynames %>%
>   filter(year == 2015, sex == "M") %>%
>   select(name, n) %>%
>   arrange(desc(n))
```

- 1 Introduction
 - The tidyverse
- 2 Tidy data
 - tibble
 - tidyr
- 3 Split-Apply-Combine
 - apply
 - plyr
 - dplyr
- 4 joins**
- 5 Resources

joins

Figure: Merge data I

band				instrument				
name	band			name	plays			
Mick	Stones			John	guitar			
John	Beatles			Paul	bass			
Paul	Beatles			Keith	guitar			

+ =

name	band	plays
Mick	Stones	<NA>
John	Beatles	guitar
Paul	Beatles	bass

```
> # Merge data with left_join()
> band %>% left_join(instrument, by = "name")
```

joins

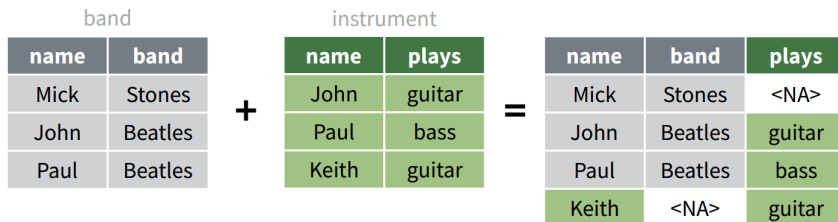
Figure: Merge data II

band				instrument				
name	band			name	plays	name	band	plays
Mick	Stones			John	guitar	John	Beatles	guitar
John	Beatles			Paul	bass	Paul	Beatles	bass
Paul	Beatles			Keith	guitar	Keith	<NA>	guitar

```
> # Merge data with right_join()
> band %>% right_join(instrument, by = "name")
```

joins

Figure: Merge data III



```
> # Merge data with full_join()
> band %>% full_join(instrument, by = "name")
```

joins

Figure: Merge data IV

band				instrument						
name	band			name	plays			name	band	plays
Mick	Stones			John	guitar			John	Beatles	guitar
John	Beatles			Paul	bass			Paul	Beatles	bass
Paul	Beatles			Keith	guitar					

```

> # Merge data with inner_join()
> band %>% inner_join(instrument, by = "name")
> # Adjust by if names do not match
> # band %>% inner_join(instrument2, by = c("name" = "artist"))

```


joins

Figure: Filter data I

band		instrument					
name	band		name	plays		name	band
Mick	Stones	+	John	guitar	=	John	Beatles
John	Beatles		Paul	bass		Paul	Beatles
Paul	Beatles		Keith	guitar			

```
> # Filter data with semi_join()
> band %>% semi_join(instrument, by = "name")
```

joins

Figure: Filter data II

band				instrument			
name	band			name	plays	name	band
Mick	Stones	+		John	guitar	=	Mick Stones
John	Beatles			Paul	bass		
Paul	Beatles			Keith	guitar		

```
> # Filter data with anti_join()
> band %>% anti_join(instrument, by = "name")
```

- 1 Introduction
 - The tidyverse
- 2 Tidy data
 - tibble
 - tidyr
- 3 Split-Apply-Combine
 - apply
 - plyr
 - dplyr
- 4 joins
- 5 Resources**

APIs

Google Trends

- Jun, S.-P., Yoo, H. S., Choi, S. (2018). Ten years of research change using Google Trends: From the perspective of big data utilizations and applications. *Technological Forecasting and Social Change* 130, 69–87.
- Lazer, D., Kennedy, R., King, G., Vespignani, A. (2014). The Parable of Google Flu: Traps in Big Data Analysis. *Science* 343(6176), 1203–1205.
- Mellon, J. (2013). Where and When Can We Use Google Trends to Measure Issue Salience? *PS: Political Science & Politics*, 46(2), 280–290.
- Stephens-Davidowitz, S. (2014). The cost of racial animus on a black candidate: Evidence using Google search data. *Journal of Public Economics* 118, 26–40.
- Stephens-Davidowitz, S. and Varian, H (2015). A hands-on guide to Google data. <http://people.ischool.berkeley.edu/~hal/Papers/2015/primer.pdf>
- Vosen, S. and Schmidt, T. (2011). Forecasting Private Consumption: Survey-Based Indicators vs. Google Trends. *Journal of Forecasting* 30, 565–578.

APIs

Census API

- Overview: <https://www.census.gov/data/developers/data-sets.html>
- ACS 5-Year data: <https://www.census.gov/data/developers/data-sets/acs-5year.html>
- R package: <https://cran.r-project.org/web/packages/censusapi/vignettes/getting-started.html>

Resources

- Textbook
 - Boehmke, B. C. (2016). *Data Wrangling in R*. New York, NY: Springer.
- Paper
 - Wickham, H. (2014). Tidy Data. *Journal of Statistical Software* 59(10), 1–23.
 - Wickham, H. (2011). The Split-Apply-Combine Strategy for Data Analysis. *Journal of Statistical Software* 40(1), 1–29.
- Online Learning
 - <https://www.datacamp.com/community/tutorials/r-tutorial-apply-family>
 - <https://cran.r-project.org/web/packages/magrittr/vignettes/magrittr.html>
- Cheatsheets
 - <https://github.com/rstudio/cheatsheets/raw/master/data-import.pdf>
 - <https://github.com/rstudio/cheatsheets/raw/master/data-transformation.pdf>