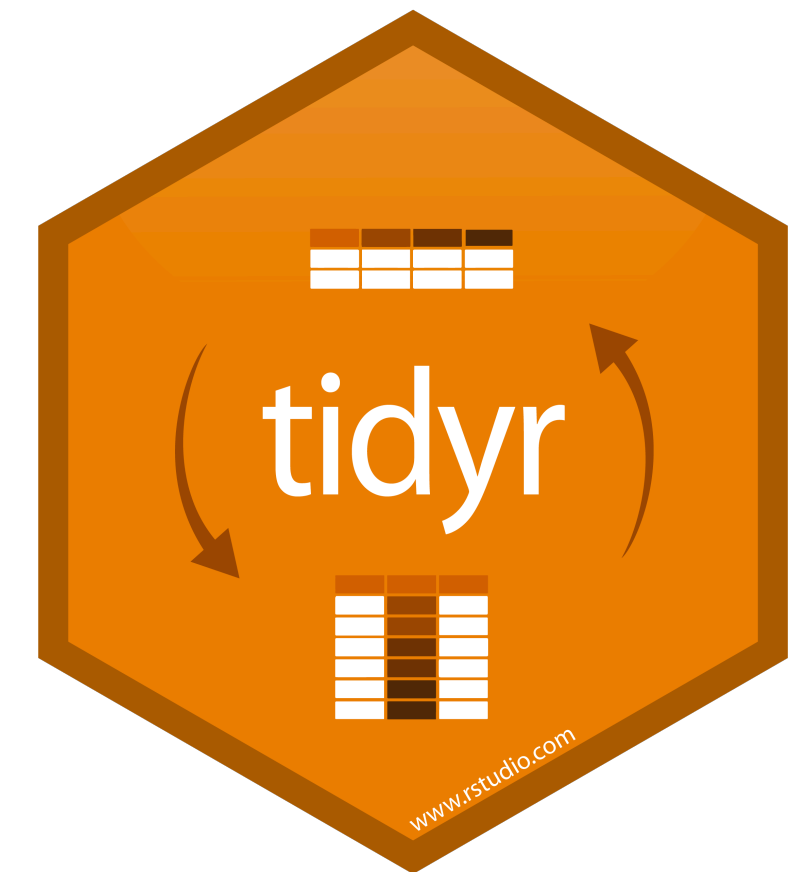
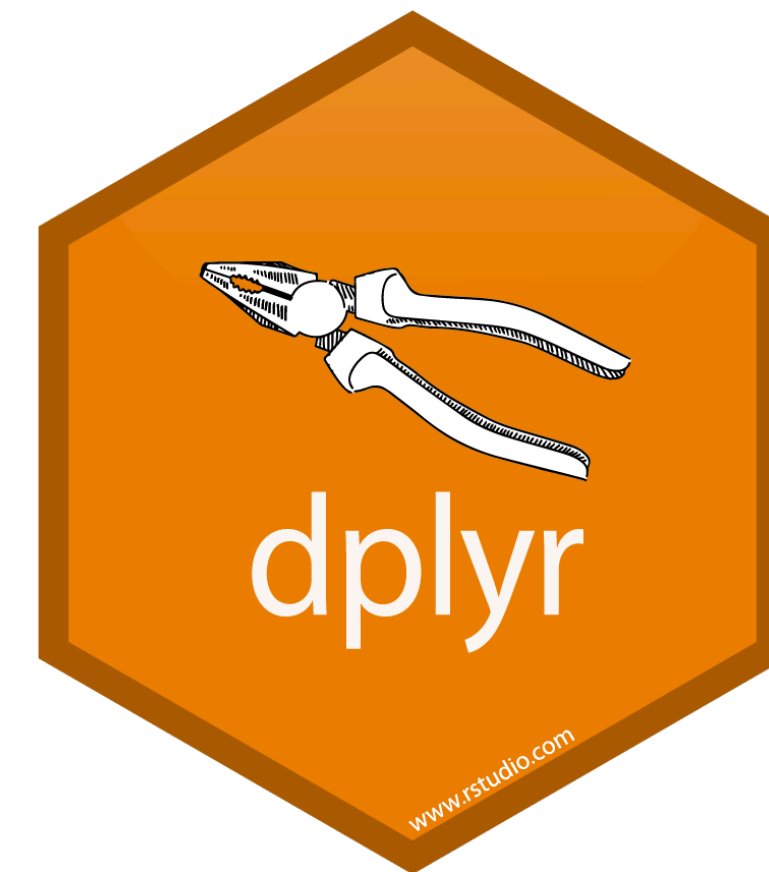


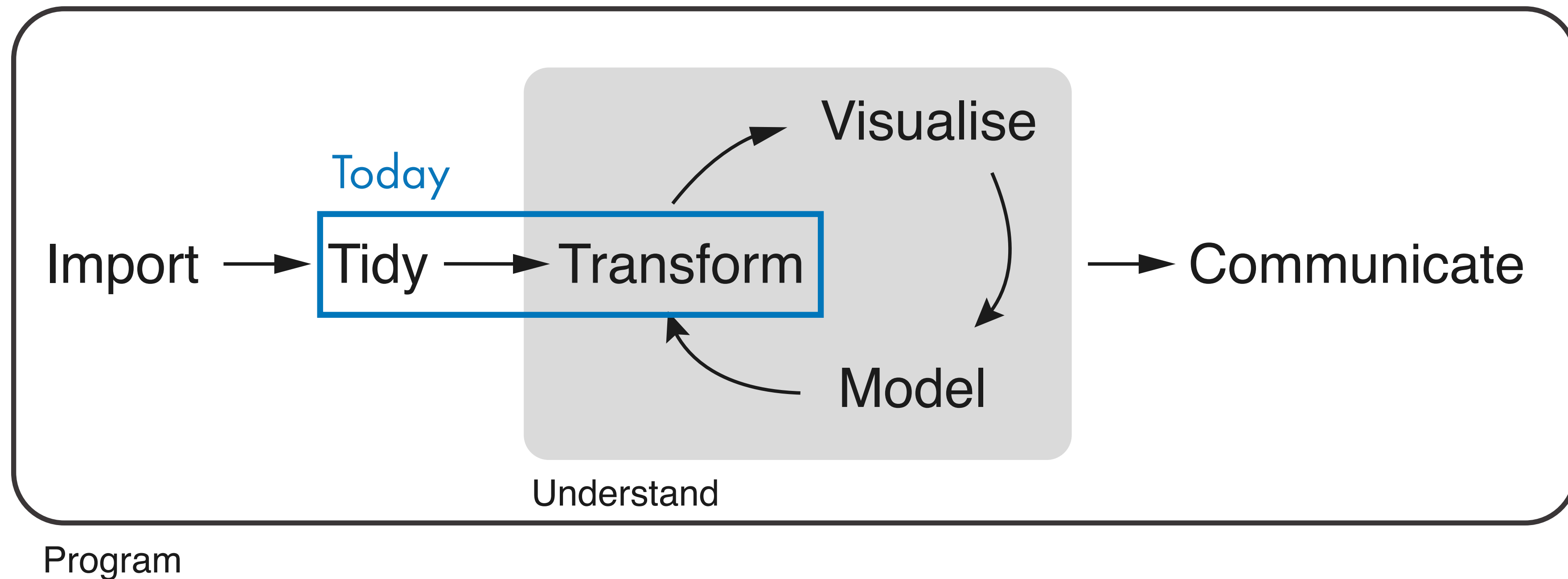
Programme for today

- Introductory lecture
- Hands-on tutorial on data wrangling
- Individual exercises



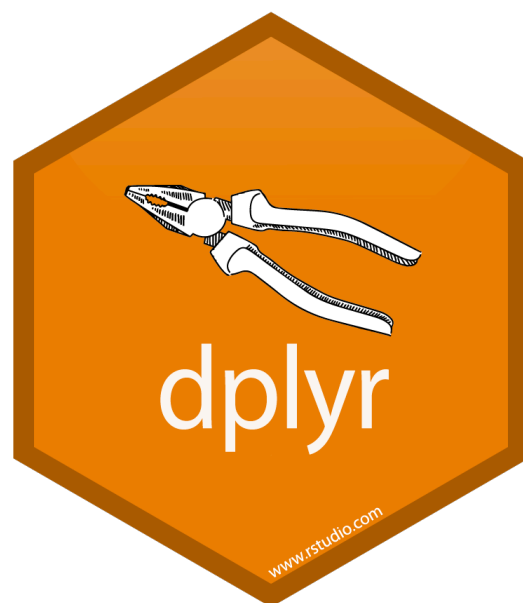
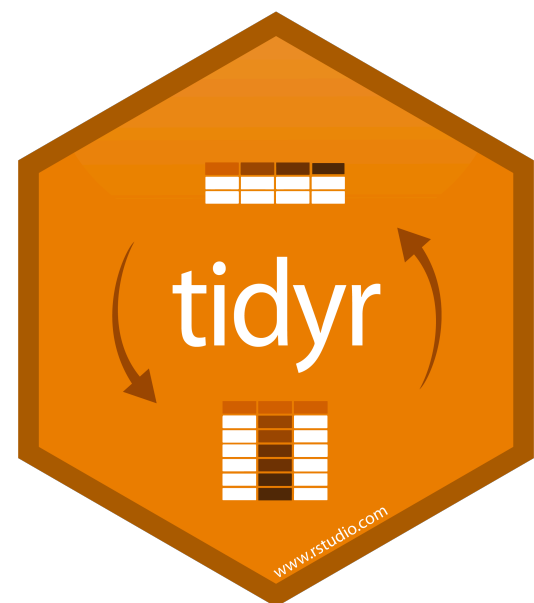
At the end of the day you should be familiar with transforming and tidying data in R using the dplyr and tidy packages

A data science project workflow



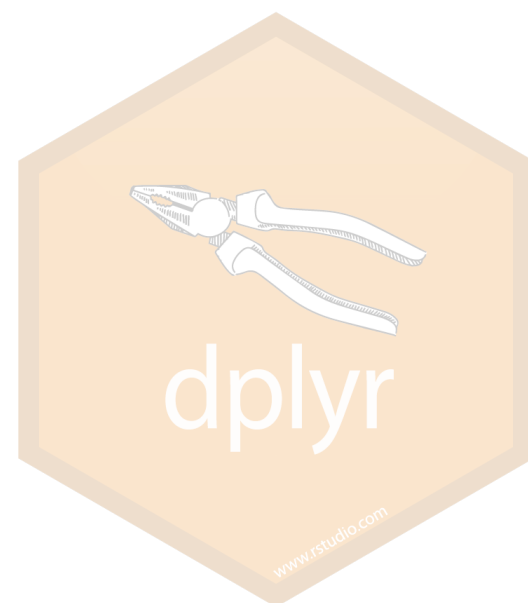
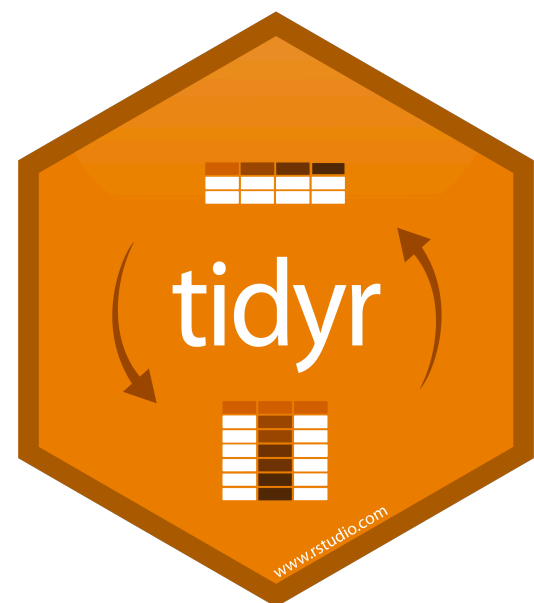
Data tidying

Data transformation



Data tidying

Data transformation

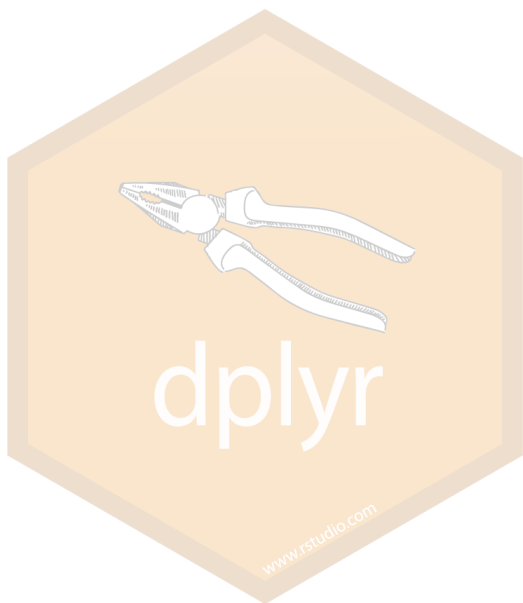
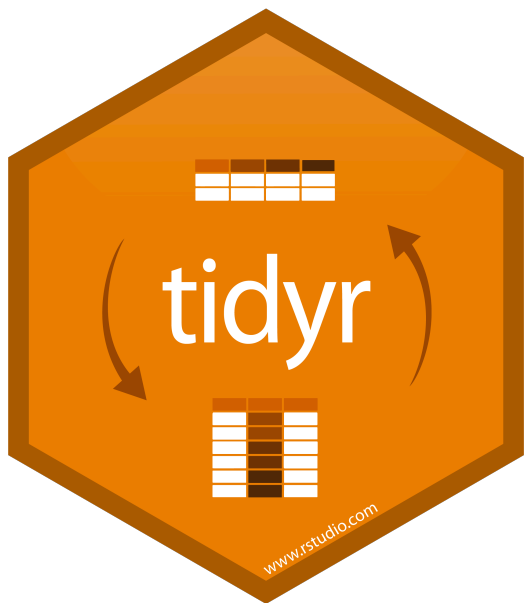


Data tidying

Data transformation

Untidy Data

species	habitat	weight	length	latitude/longitude	date
Alligator mississippiensis	swamp	431 lb	4 ft 2	29.531,-82.184	Sept 15, 2015
Puma concolor	forest	125 lb	2.2m	29.125,-81.682	08/10/2015
Ursus americanus	forest	88 kg	133 cm	N29°7'30"/W81°40'55.2"	07-13-2015



In which ways is this dataset “untidy”?

Data tidying

Data transformation

Untidy Data

species	habitat	weight	length	latitude/longitude	date
Alligator mississippiensis	swamp	431 lb	4 ft 2	29.531,-82.184	Sept 15, 2015
Puma concolor	forest	125 lb	2.2m	29.125,-81.682	08/10/2015
Ursus americanus	forest	88 kg	133 cm	N29°7'30"/W81°40'55.2"	07-13-2015

Tidy Data

meta-data

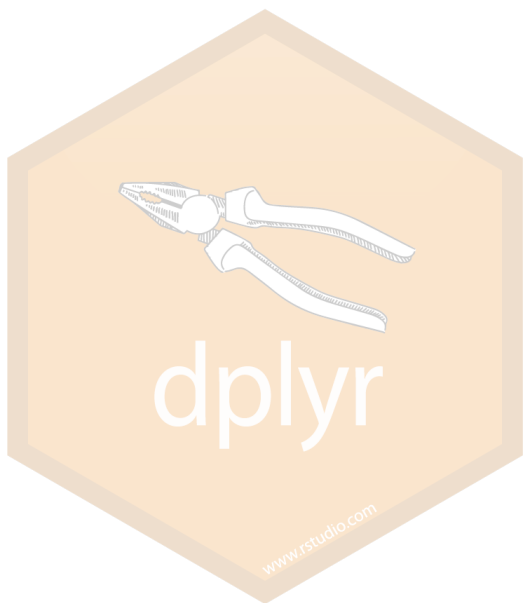
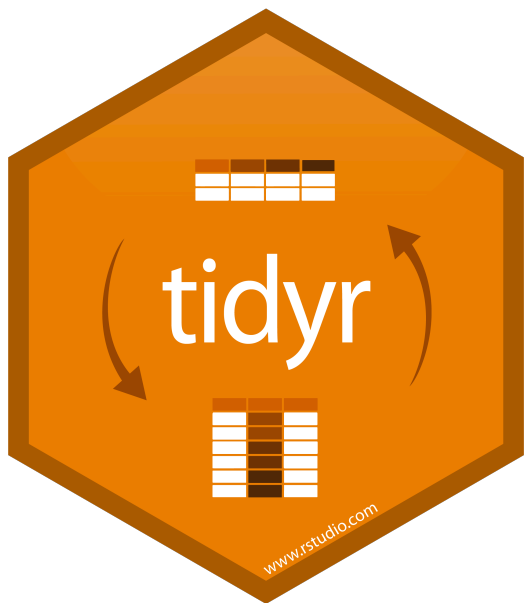
data

species_code	date	station_code	weight_kg	length_cm
TSN 551771	2015-09-15	1	196	127
TSN 55247	2015-08-10	2	57	220
TSN 180544	2015-07-13	2	88	133

station_code	habitat	latitude	longitude
1	swamp	29.531	-82.184
2	forest	29.125	-81.682

species_code	class	genus	species
TSN 551771	Reptilia	Alligator	mississippiensis
TSN 55247	Mammalia	Puma	concolor
TSN 180544	Mammalia	Ursus	americanus

Example of a tidy version of dataset



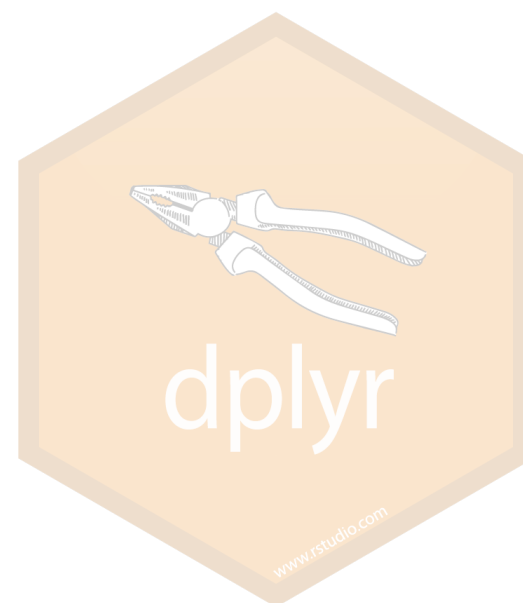
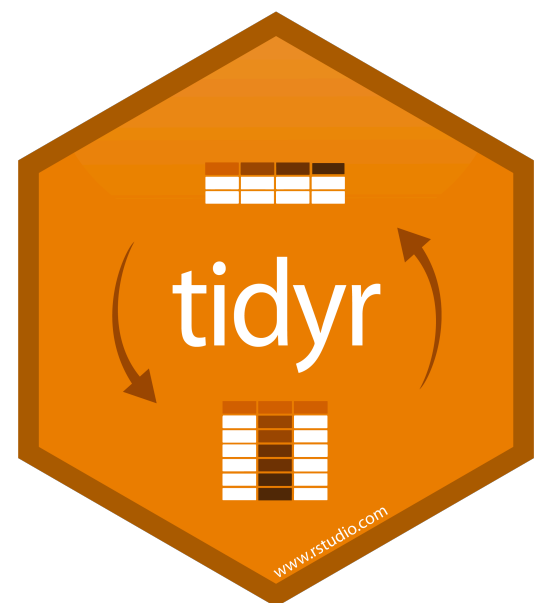
Data tidying

Data transformation

country	year	cases	population
Afghanistan	1999	1845	19957071
Afghanistan	2000	2666	20095360
Brazil	1999	31737	172006362
Brazil	2000	80488	174004898
China	1999	212258	1272015272
China	2000	210766	1280028583

variables

1. Each variable has its own column



Rules for tidy data

Data tidying

Data transformation

country	year	cases	population
Afghanistan	1999	745	19997071
Afghanistan	2000	2666	20095360
Brazil	1999	37737	17206362
Brazil	2000	80488	17404898
China	1999	212258	1272015272
China	2000	210766	128048583

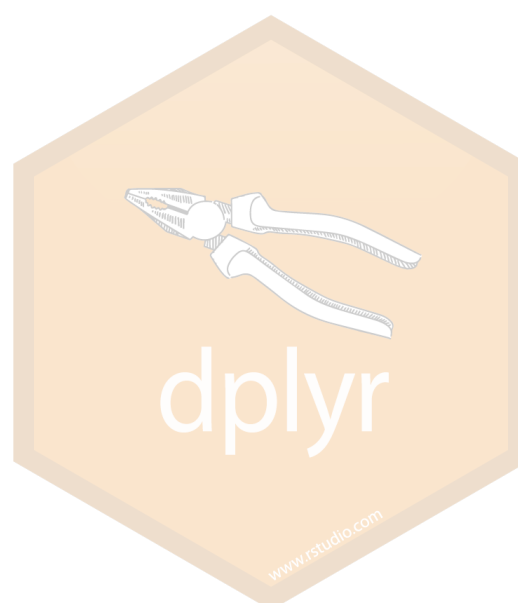
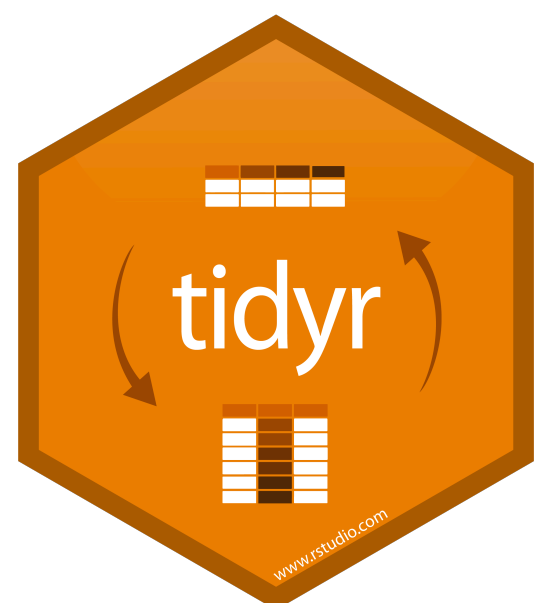
variables

country	year	cases	population
Afghanistan	1999	745	19997071
Afghanistan	2000	2666	20095360
Brazil	1999	37737	17206362
Brazil	2000	80488	17404898
China	1999	212258	1272015272
China	2000	210766	128048583

observations

1. Each variable has its own column

2. Each observation has its own row



Rules for tidy data

Data tidying

Data transformation

country	year	cases	population
Afghanistan	1999	745	19987071
Afghanistan	2000	2666	20595360
Brazil	1999	37737	172006362
Brazil	2000	80488	174504898
China	1999	212258	1272915272
China	2000	213766	1280425583

variables

country	year	cases	population
Afghanistan	1999	745	19987071
Afghanistan	2000	2666	20595360
Brazil	1999	37737	172006362
Brazil	2000	80488	174504898
China	1999	212258	1272915272
China	2000	213766	1280425583

observations

country	year	cases	population
Afghanistan	99	745	19987071
Afghanistan	00	2666	20595360
Brazil	99	37737	172006362
Brazil	00	80488	174504898
China	99	212258	1272915272
China	00	213766	1280425583

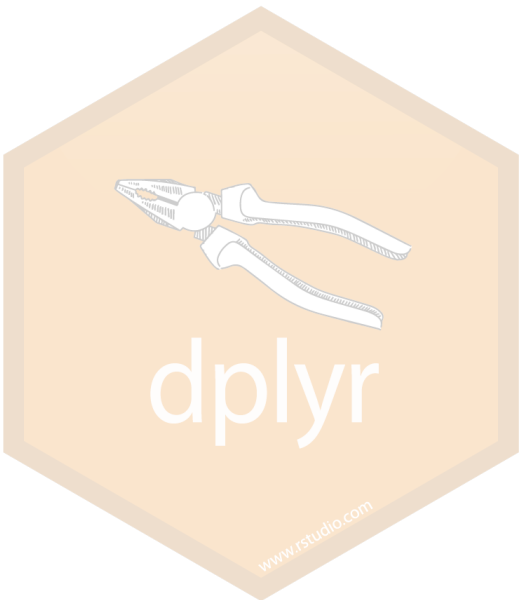
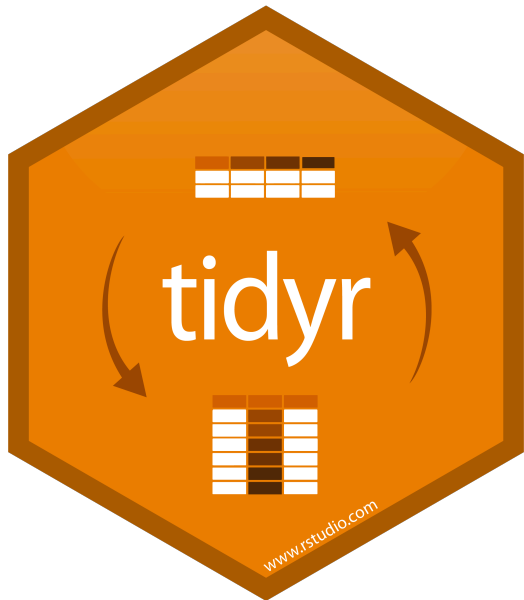
values

1. Each variable has its own column

2. Each observation has its own row

3. Each values has its own cell

Rules for tidy data



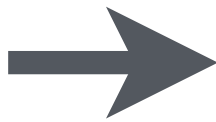
Data tidying

Data transformation

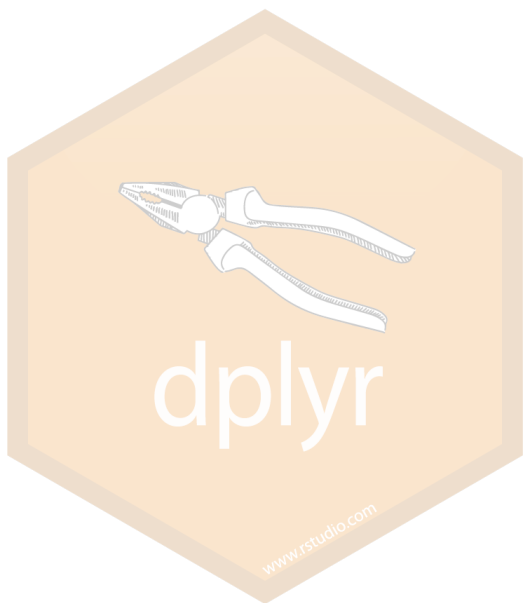
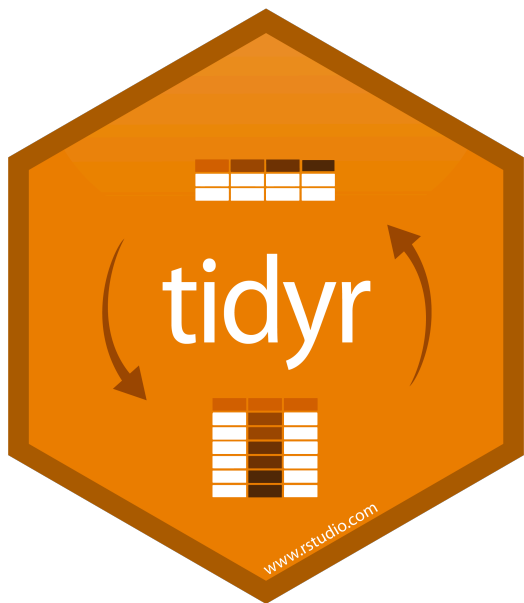
```
pivot_longer(data, cols, names_to = "name",  
values_to = "value", values_drop_na = FALSE)
```

table4a

country	1999	2000
A	0.7K	2K
B	37K	80K
C	212K	213K



country	year	cases
A	1999	0.7K
B	1999	37K
C	1999	212K
A	2000	2K
B	2000	80K
C	2000	213K



The `pivot_*` functions

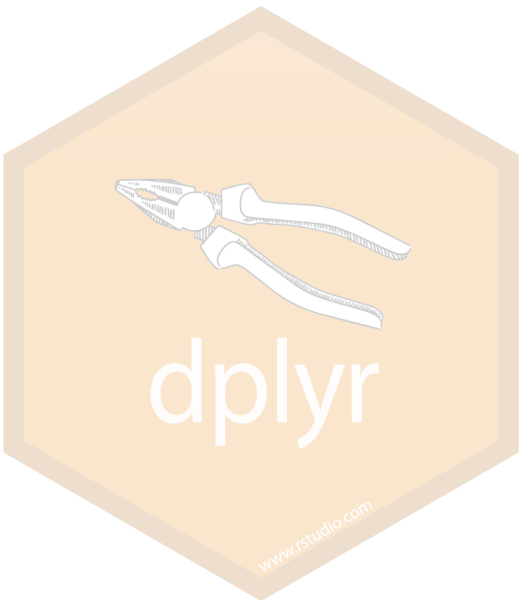
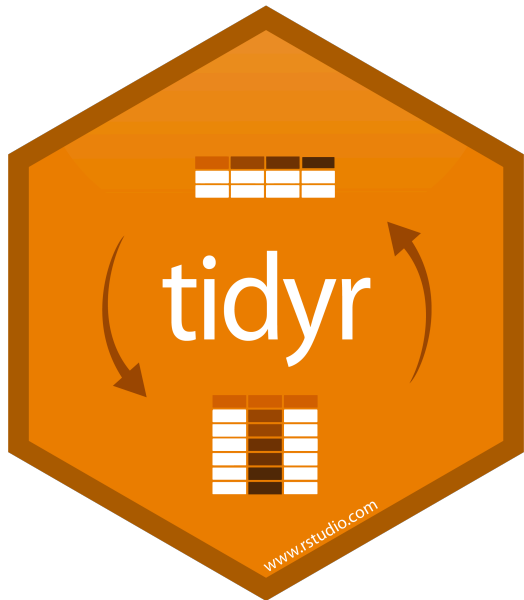
Data tidying

Data transformation

```
pivot_wider(data, names_from = "name",  
values_from = "value")
```

table2

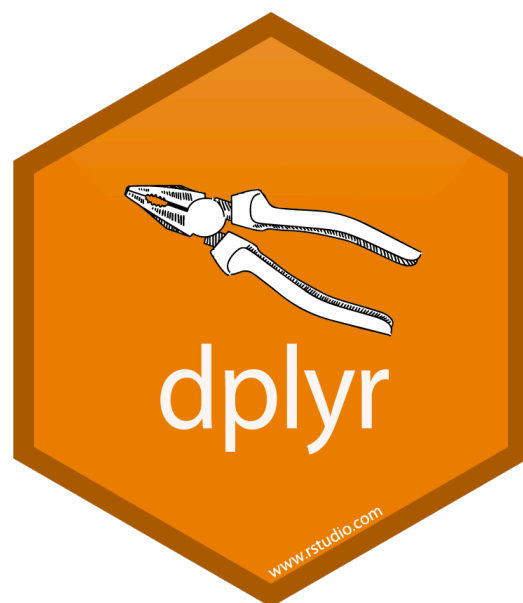
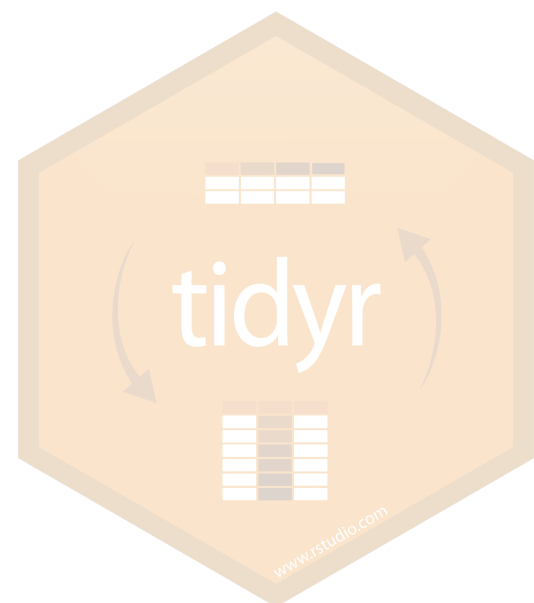
country	year	type	count		country	year	cases	pop
A	1999	cases	0.7K	➔	A	1999	0.7K	19M
A	1999	pop	19M		A	2000	2K	20M
A	2000	cases	2K		B	1999	37K	172M
A	2000	pop	20M		B	2000	80K	174M
B	1999	cases	37K		C	1999	212K	1T
B	1999	pop	172M		C	2000	213K	1T
B	2000	cases	80K					
B	2000	pop	174M					
C	1999	cases	212K					
C	1999	pop	1T					
C	2000	cases	213K					
C	2000	pop	1T					



The `pivot_*` functions

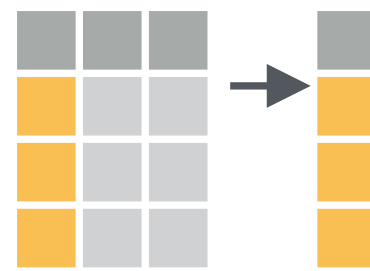
Data tidying

Data transformation

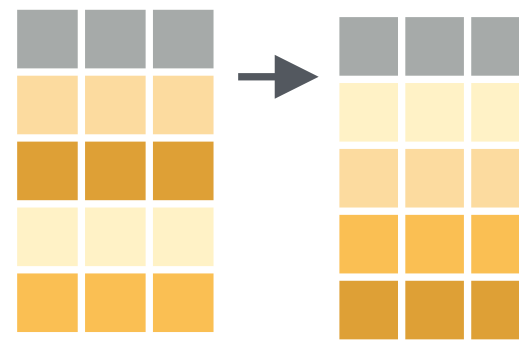


Data tidying

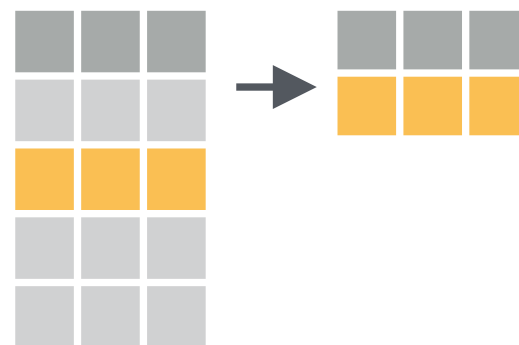
Data transformation



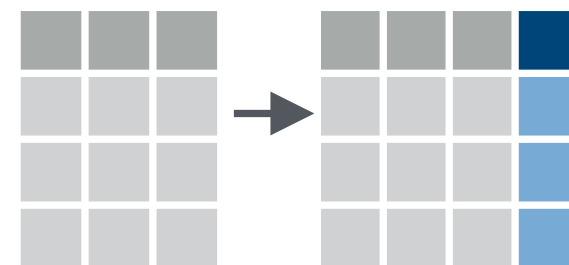
select(.data, ...) Extract columns as a table.
`select(mtcars, mpg, wt)`



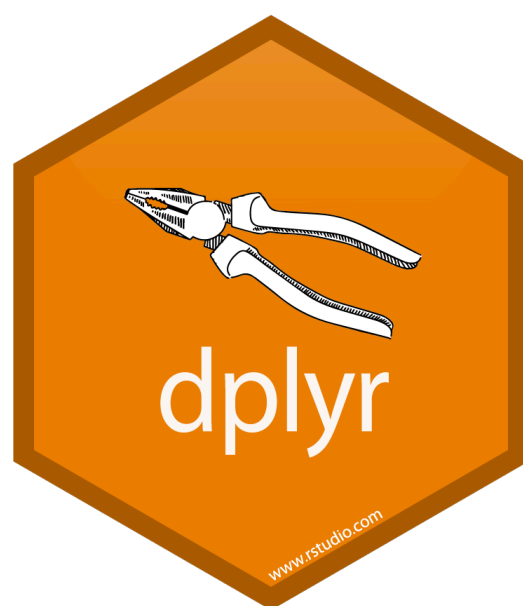
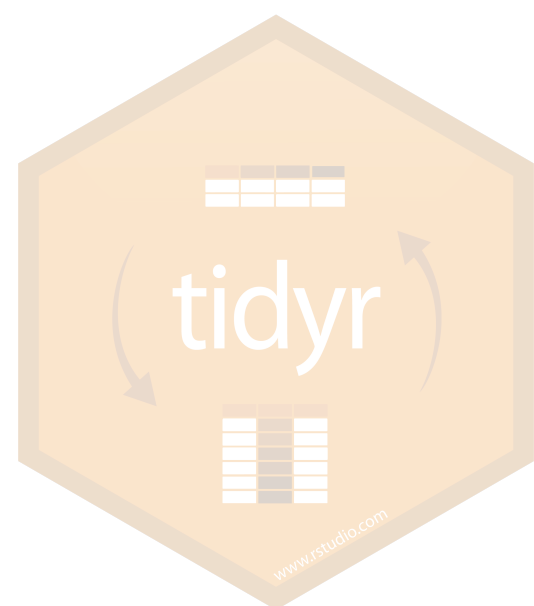
arrange(.data, ..., .by_group = FALSE) Order rows by values of a column or columns (low to high), use with **desc()** to order from high to low.
`arrange(mtcars, mpg)`
`arrange(mtcars, desc(mpg))`



filter(.data, ..., .preserve = FALSE) Extract rows that meet logical criteria.
`filter(mtcars, mpg > 20)`



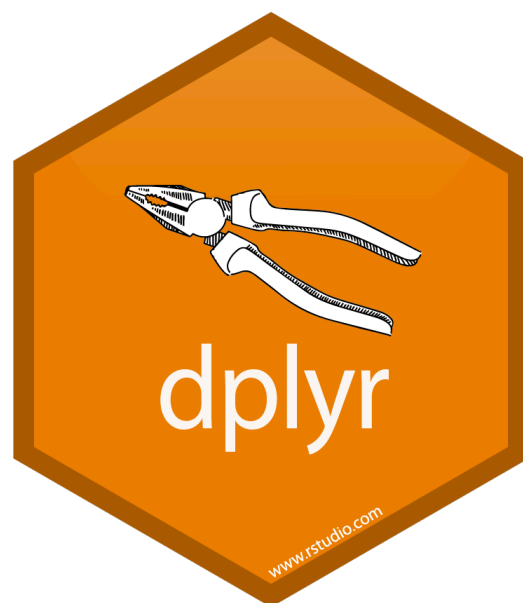
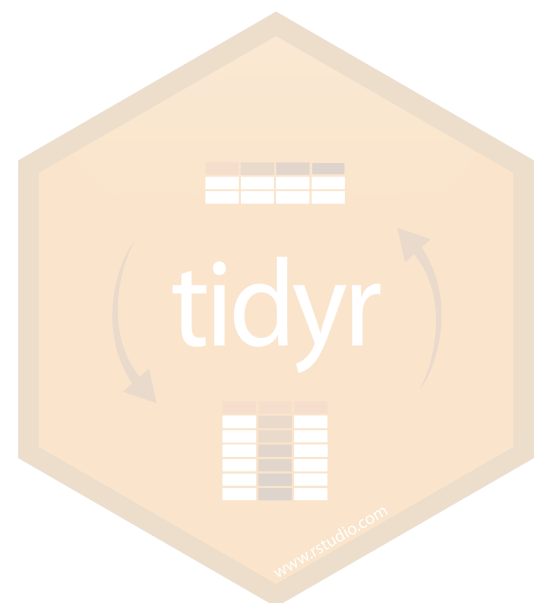
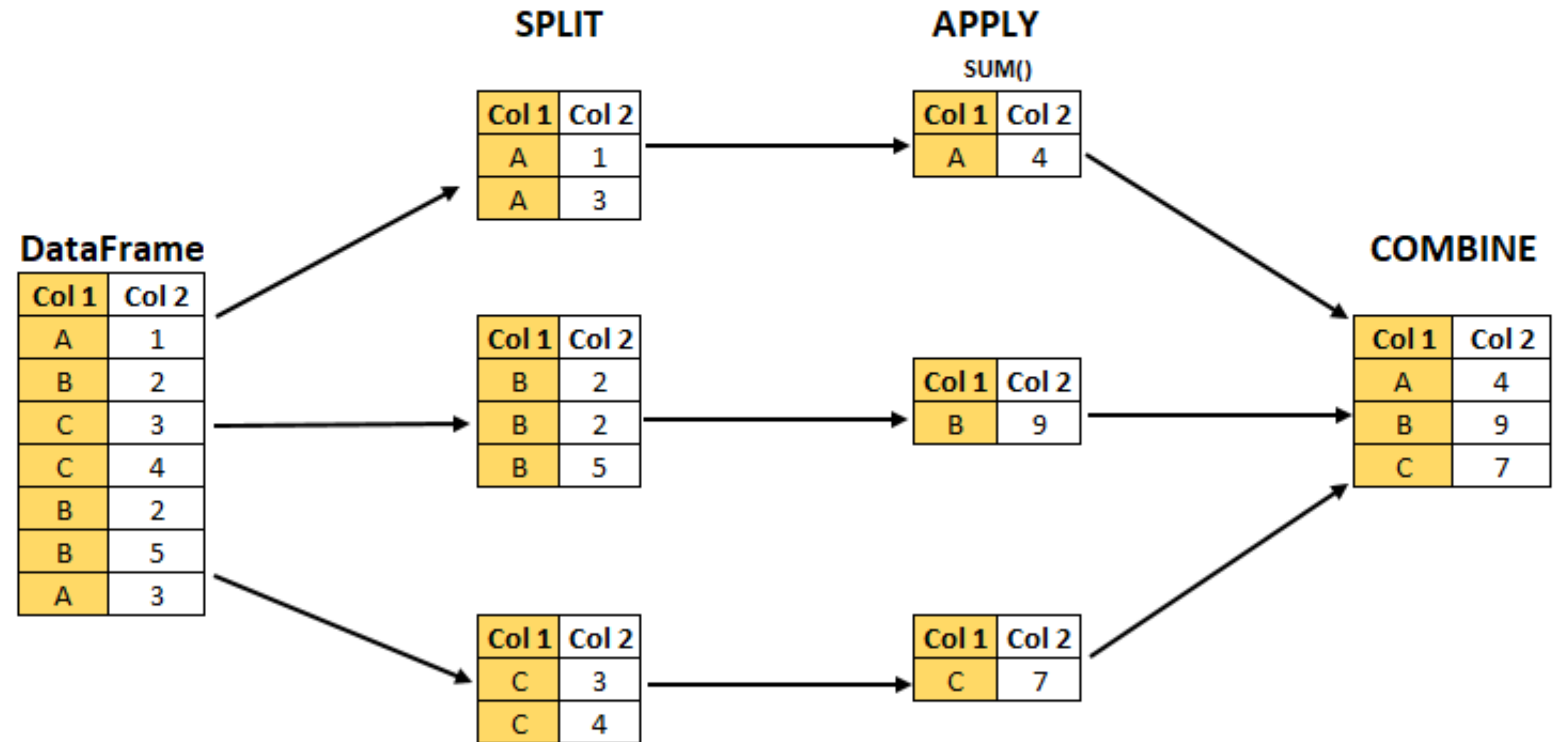
mutate(.data, ..., .keep = "all", .before = NULL, .after = NULL) Compute new column(s). Also **add_column()**, **add_count()**, and **add_tally()**.
`mutate(mtcars, gpm = 1 / mpg)`



Simple transformations

Data tidying

Data transformation



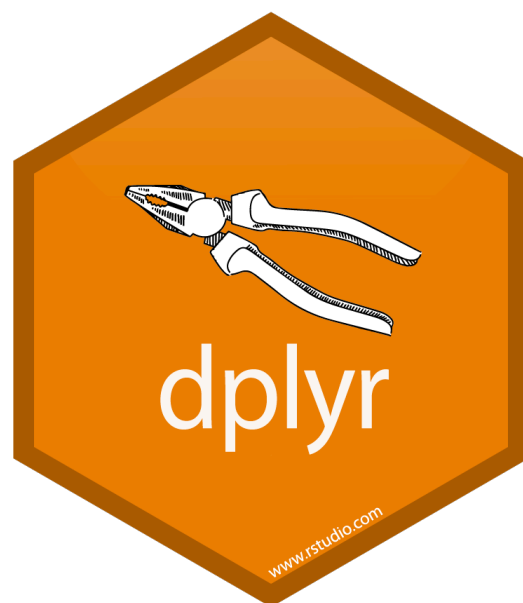
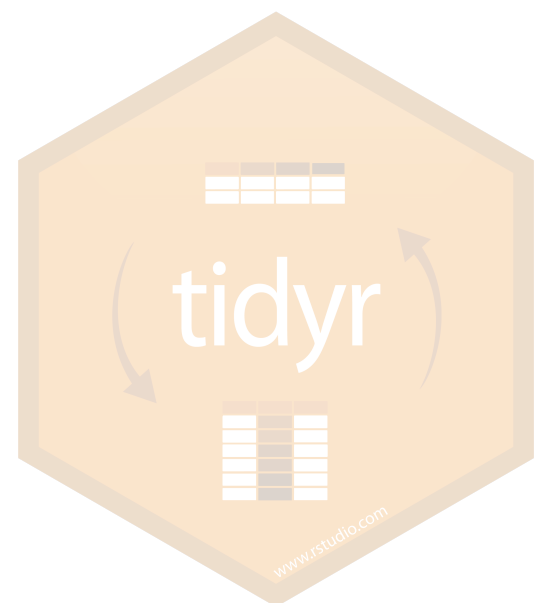
split-apply-combine

Data tidying

Data transformation



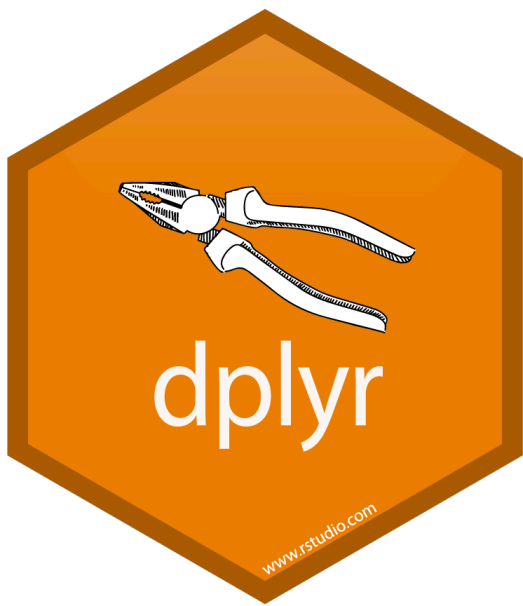
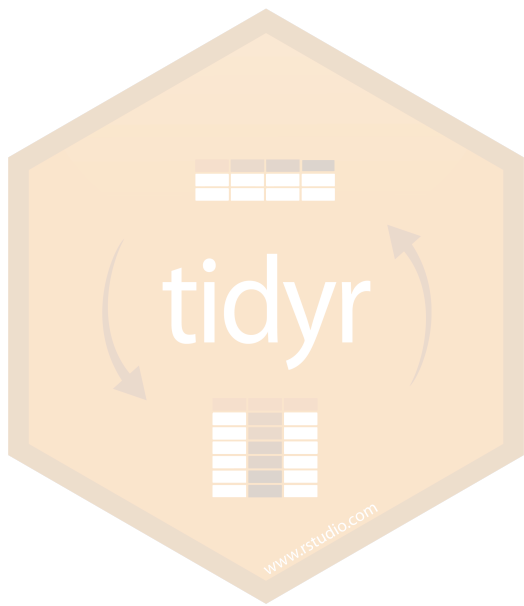
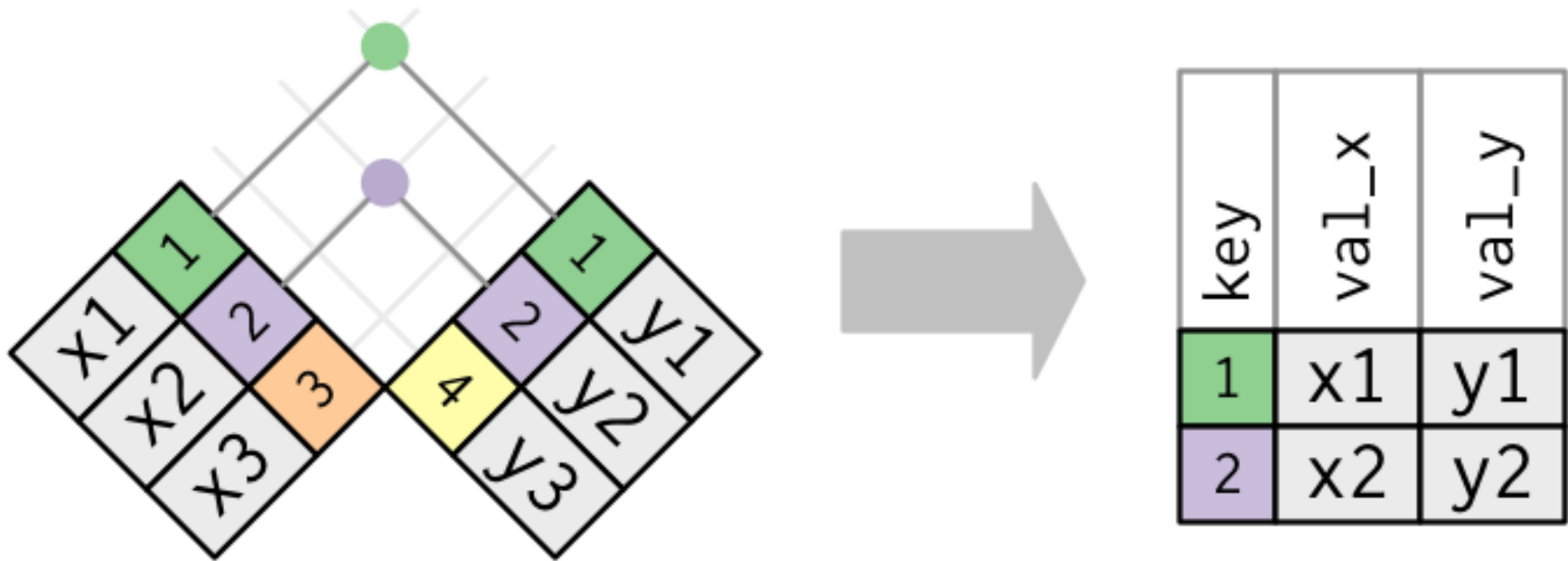
```
mtcars %>%  
  group_by(cyl) %>%  
  summarise(avg = mean(mpg))
```



split-apply-combine

Data tidying

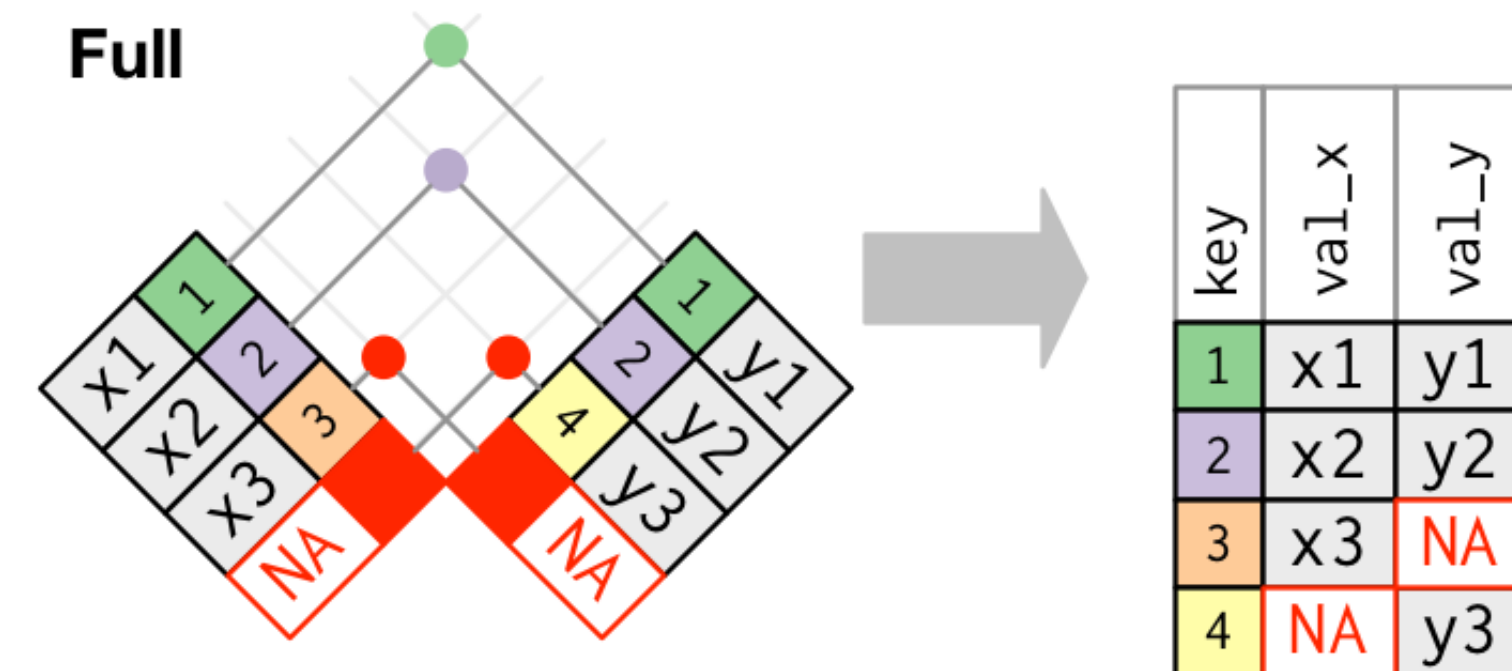
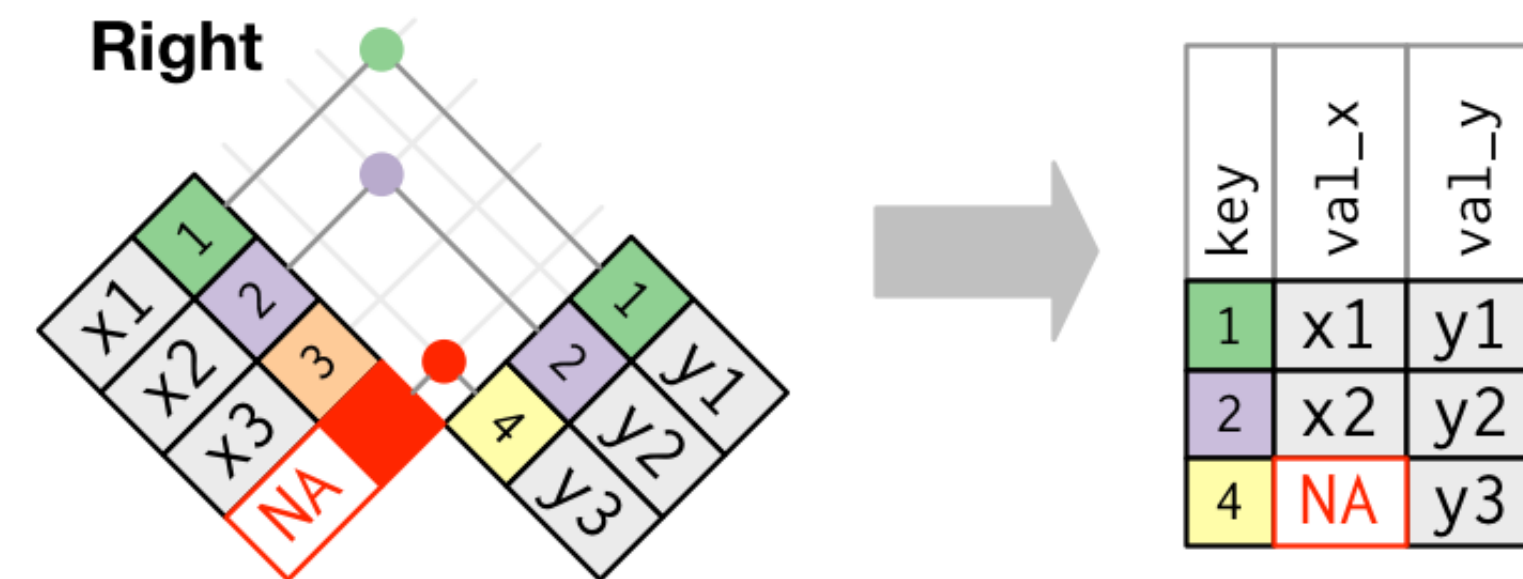
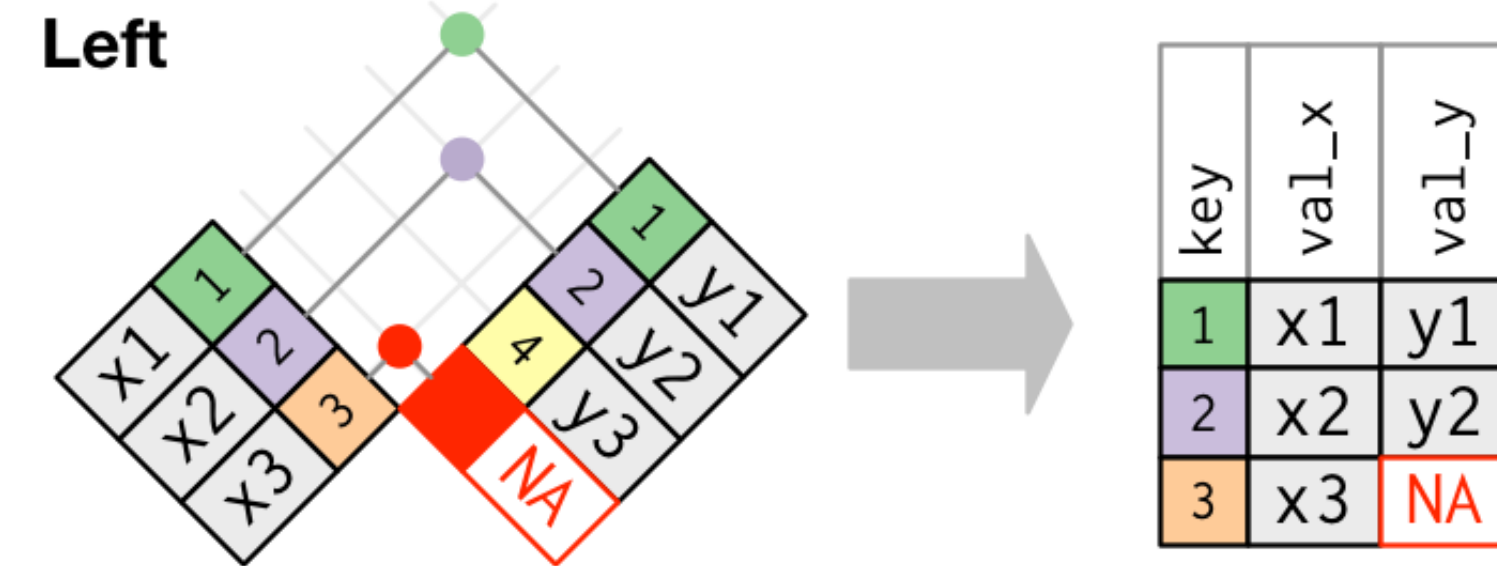
Data transformation



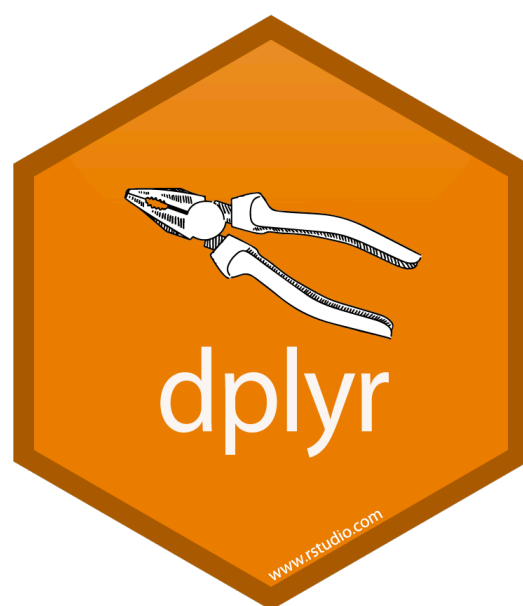
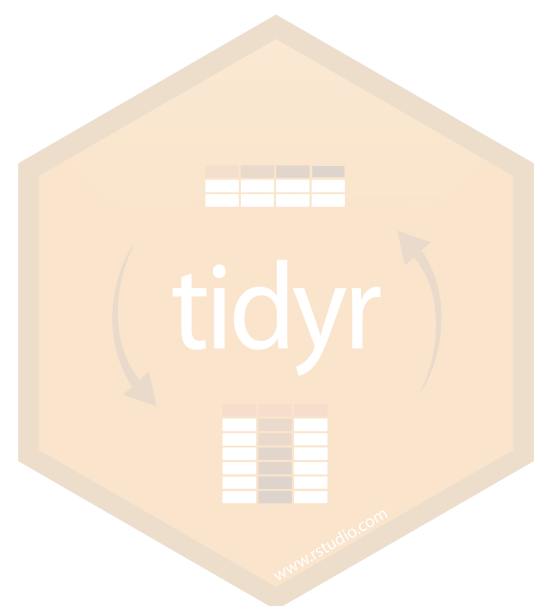
Mutating joins - Inner join

Data tidying

Data transformation



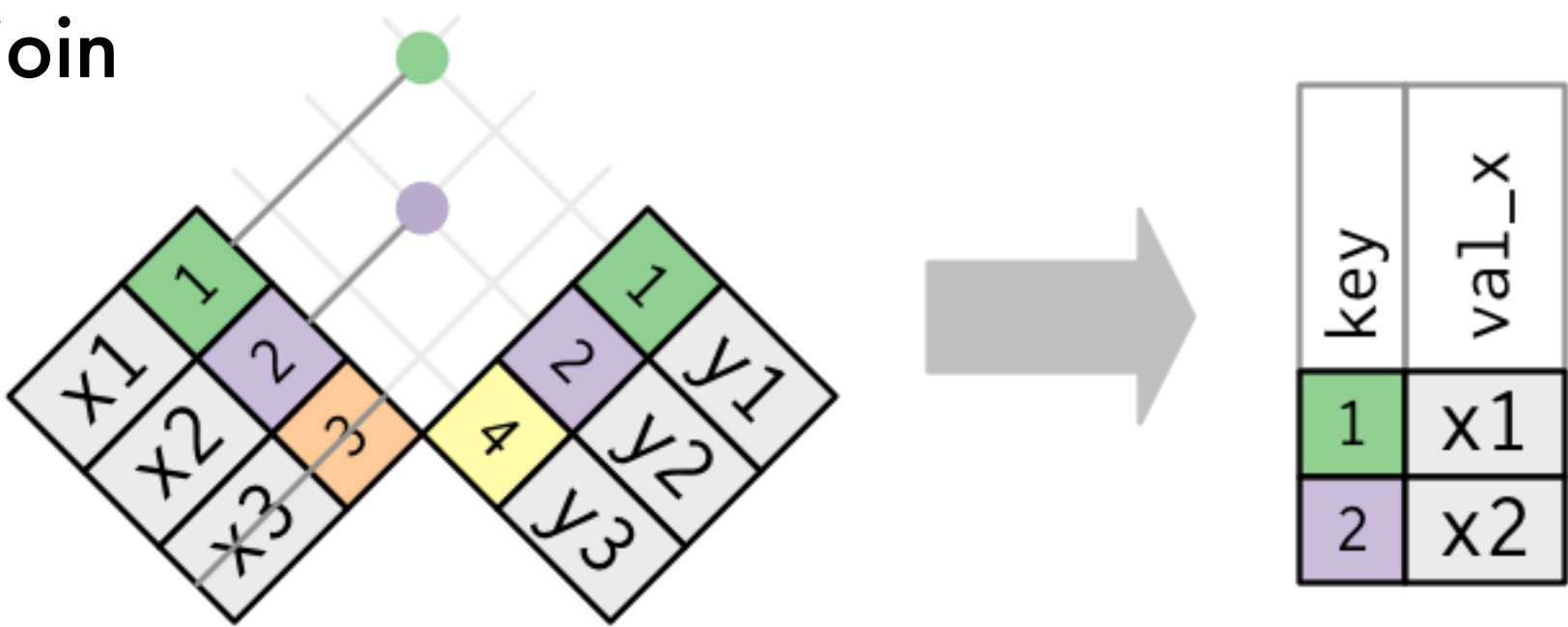
Mutating joins - Outer joins



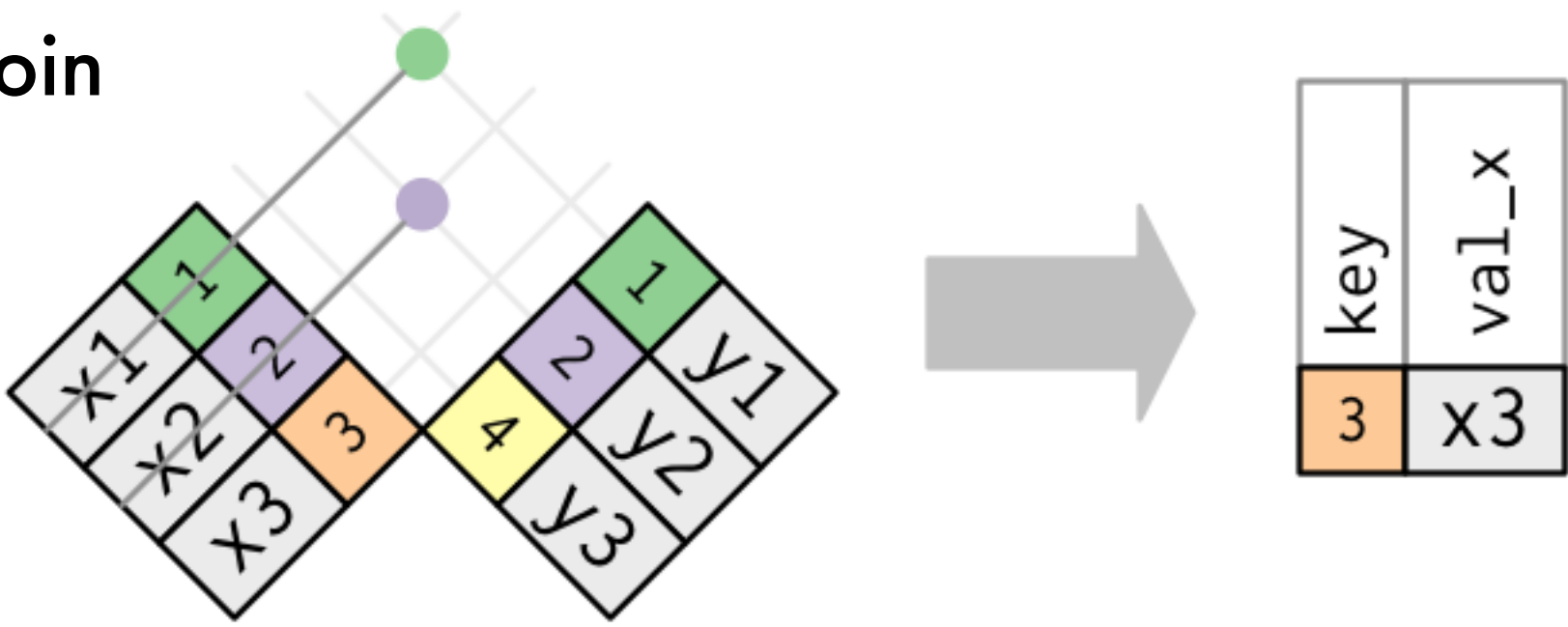
Data tidying

Data transformation

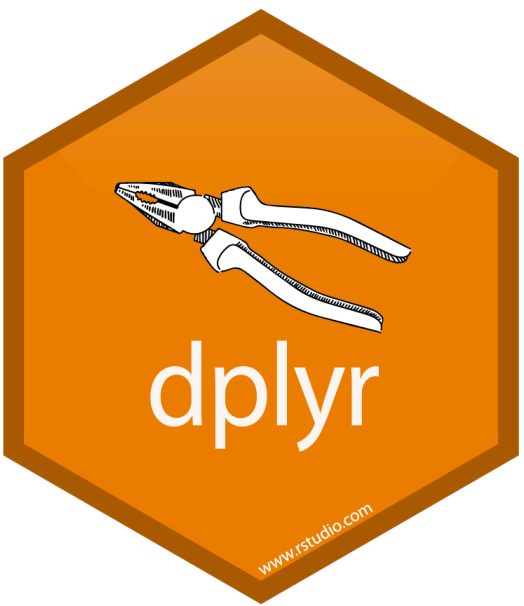
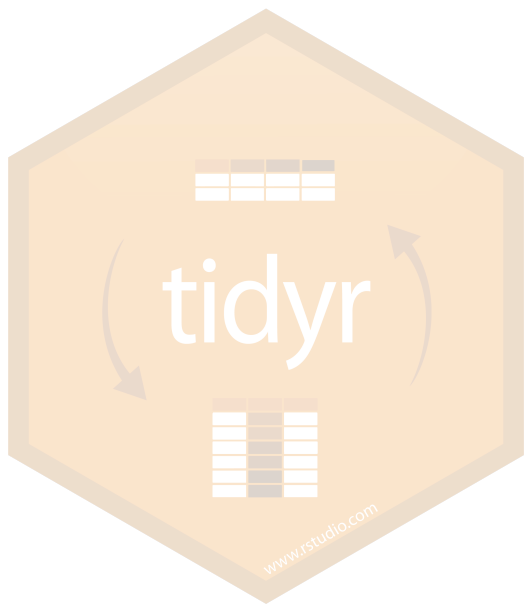
Semi join



Anti join

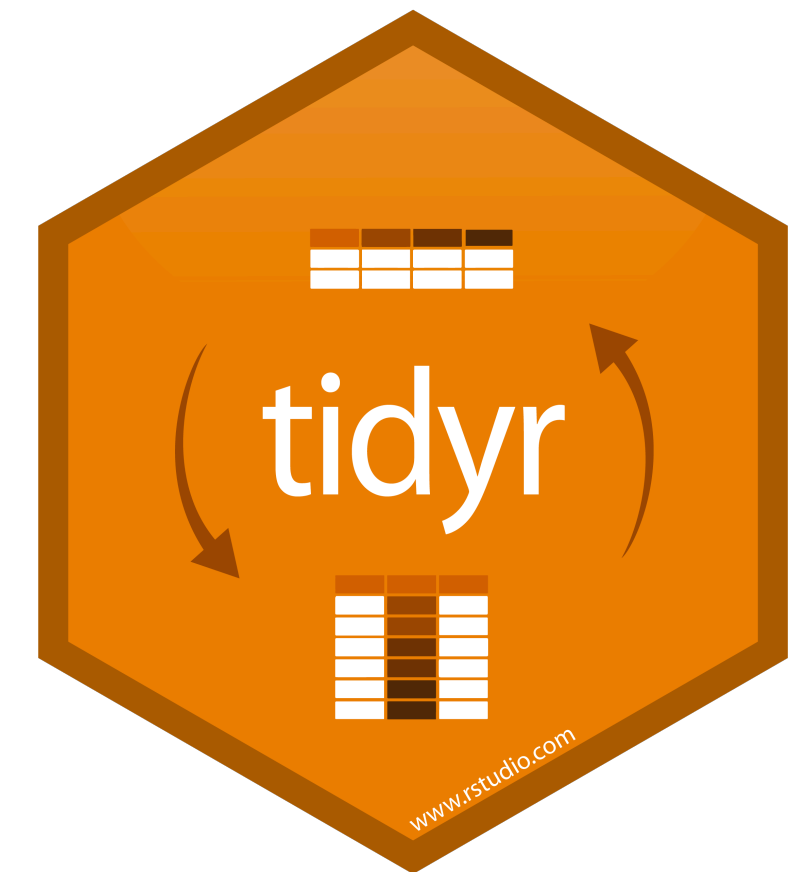
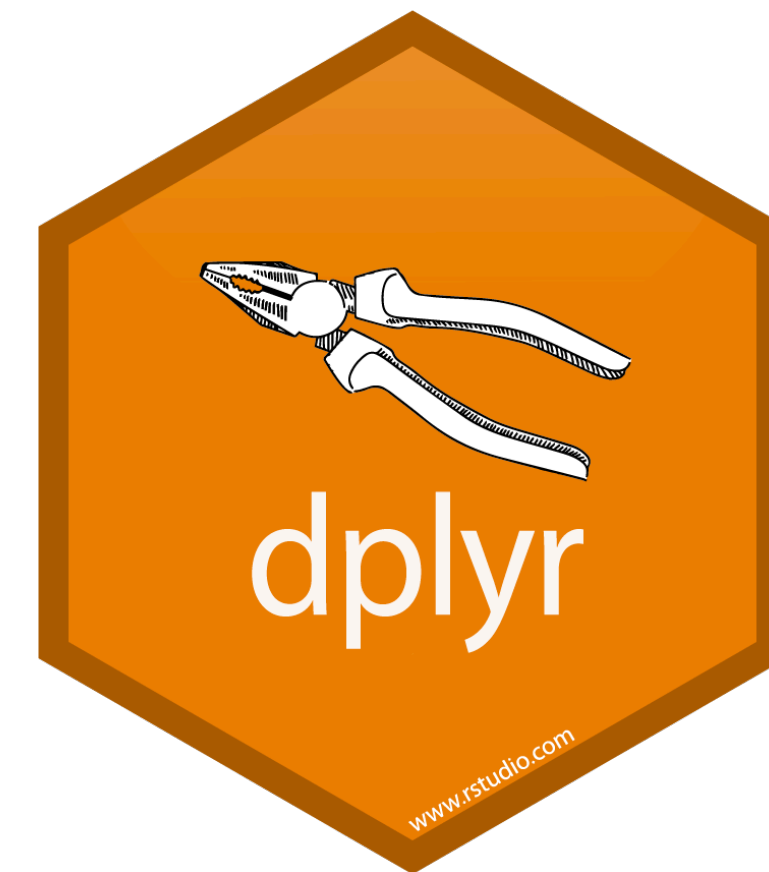


Filtering joins



Programme for today

- Introductory lecture
- Hands-on tutorial on data wrangling
- Individual exercises



At the end of the day you should be familiar with transforming and tidying data in R using the dplyr and tidy packages