



THE UNIVERSITY
OF QUEENSLAND
AUSTRALIA

CREATE CHANGE

Introduction to data wrangling in R

Matthew Curry & Tomasz Zając

Institute for Social Science Research

The University of Queensland

ARDC Summer School 2024 Day 1

Social Science Stream

Acknowledgment of Country

The University of Queensland (UQ) acknowledges the Traditional Owners and their custodianship of the lands on which we meet.

We pay our respects to their Ancestors and their descendants, who continue cultural and spiritual connections to Country.

We recognise their valuable contributions to Australian and global society.



About the course

- Data manipulation methods in R
 - Subsetting data
 - Creating modifying variables
 - Merging datasets
 - Aggregating data

Why R?

- Powerful
- Designed for statistics and data science
- Free and open source
- Platform-independent
- Popular
- Great community support

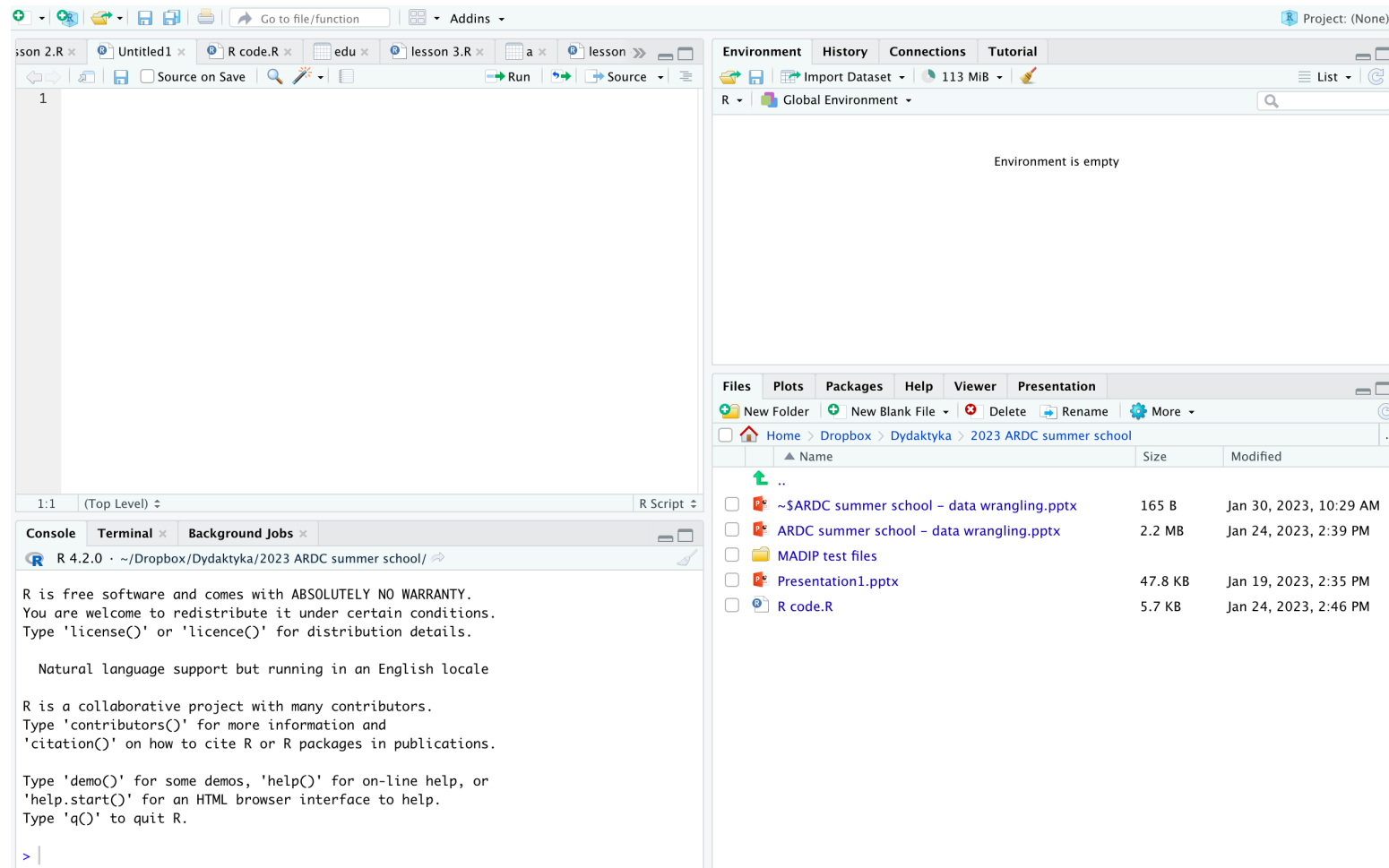


THE UNIVERSITY
OF QUEENSLAND
AUSTRALIA

CREATE CHANGE

R & RStudio

RStudio



Data types

- Booleans/logical (TRUE, FALSE)
- Numeric (1, 1.3, 5.9)
- Integer (1, 2, 10)
- Character – string (“a”, “abc”, “5”)
- Date
- Factor
- Polygons
- Missing: NA, NaN

Data objects

- Data.frame / tibble/ data.table
- Vector
- ...
- Matrix
- Array
- List

Some Useful Base R functions

- `<-` or `=`
- `view()`
- `head()`
- `tail()`
- `table()`
- `prop.table()`
- `summary()`



Open RStudio

“Melbournian hipster coding in R”



R packages



<https://tidyverse.tidyverse.org>

Loading data

Data format	Package	Command
csv	base R	<code>read.csv()</code> – special case of <code>read.table()</code>
	readr	<code>read_csv()</code> - special case of <code>read_delim()</code>
	data.table	<code>fread()</code>
	arrow	<code>read_csv_arrow()</code>
xls, xlsx	readxl	<code>read_excel()</code>
SPSS, Stata	foreign	<code>read.spss()</code> , <code>read.dta()</code>
	haven	<code>read_spss()</code> , <code>read_dta()</code>
rds	base R	<code>readRDS()</code>

Example data

Example datasets with random data:

File 1 – Data on higher education completions

File 2 – Data on income tax

Goal:

Compute the gender pay gap (GPG) for each broad field of education for graduates from 2011:

- Extract data on the 2011 cohort of graduates
- Create required variables (e.g. broad field of education), clean the data
- Extract tax data
- Transform tax data to the right format
- Merge both datasets
- Aggregate data and compute the GPG



Working with R packages and example data

“R packages being produced”

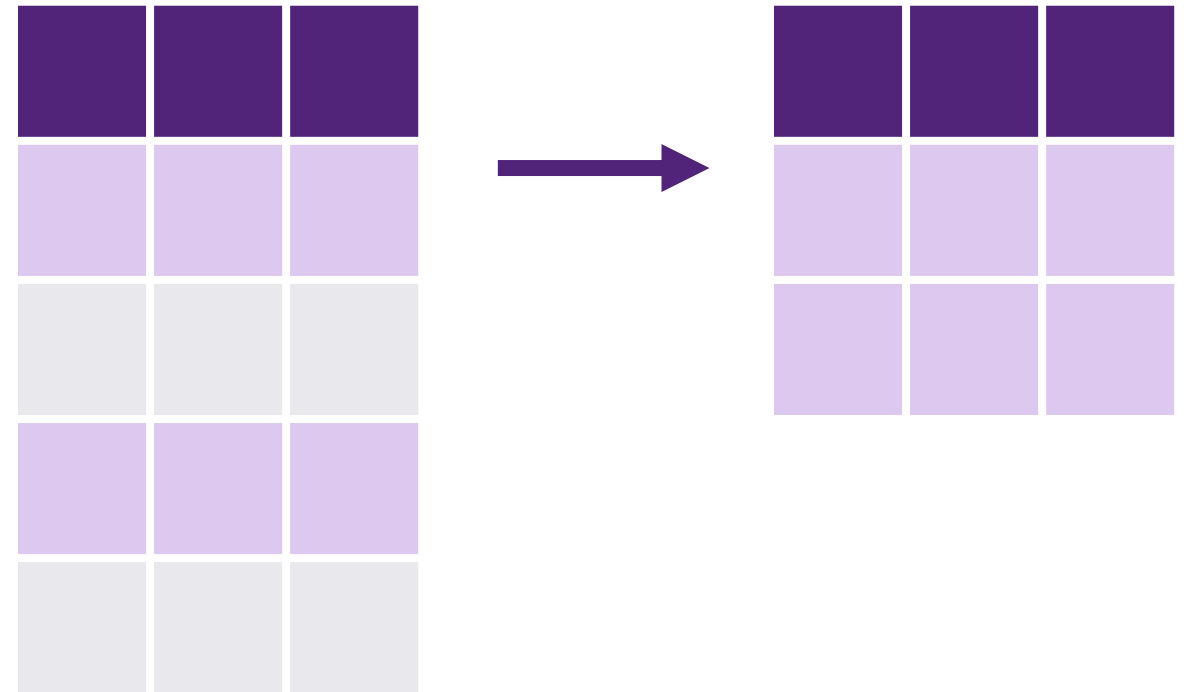


Filter rows/ extract cases

`filter(data, condition)`

Conditions:

- `x == y` equal to
- `x != y` is not equal to
- `x > y` (`>=`, `<=`, `<`) greater/smaller than
- `x %in% y` within
- `!is.na(x)` is missing
- `&` and
- `|` or



visualisations inspired by: <https://github.com/rstudio/cheatsheets/blob/main/data-transformation.pdf>

Select columns

`select(data, columns)`

tidy selection:

`col1, col2, col3`

`col2:col5`

`-col3`

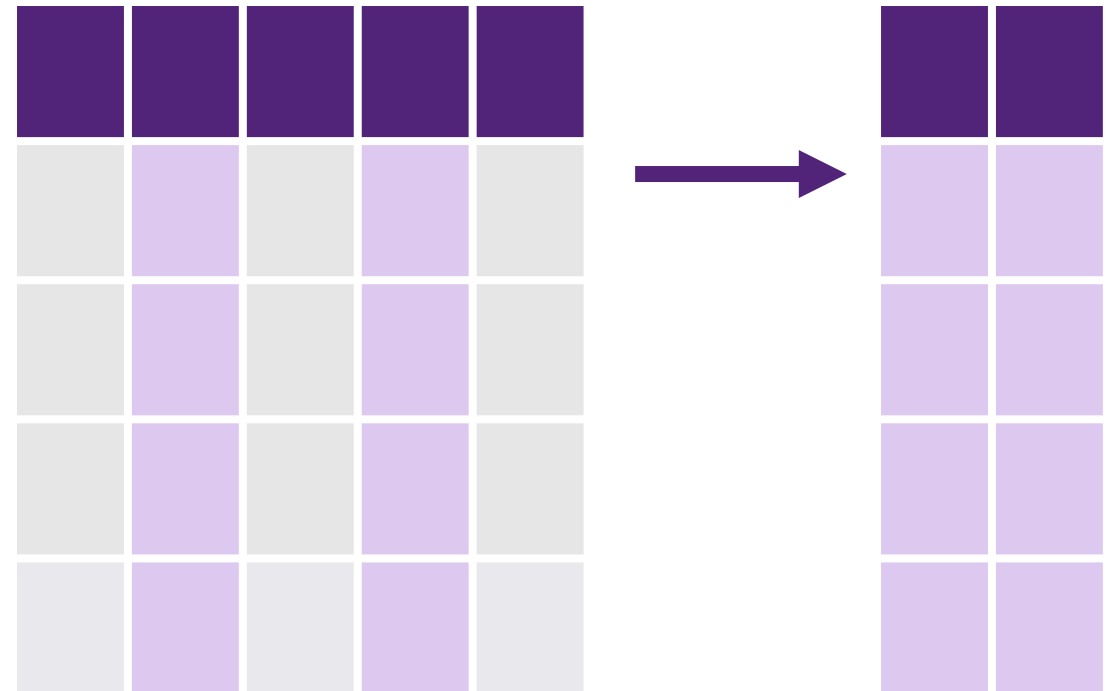
`col1 | col4:col8`

`starts_with()`

`ends_with()`

`contains()`

`where()`



Chaining – Pipe operator

Requires `magrittr` or `dplyr`

```
function4(function3(function2(function1(x, add_arg1), add_arg2), add_arg3), add_arg4)
```

```
x %>%
```

```
  function1(add_arg1) %>%
```

```
  function2(add_arg2) %>%
```

```
  function3(add_arg3) %>%
```

```
  function4(add_arg4)
```

Mutate – add/ modify variables

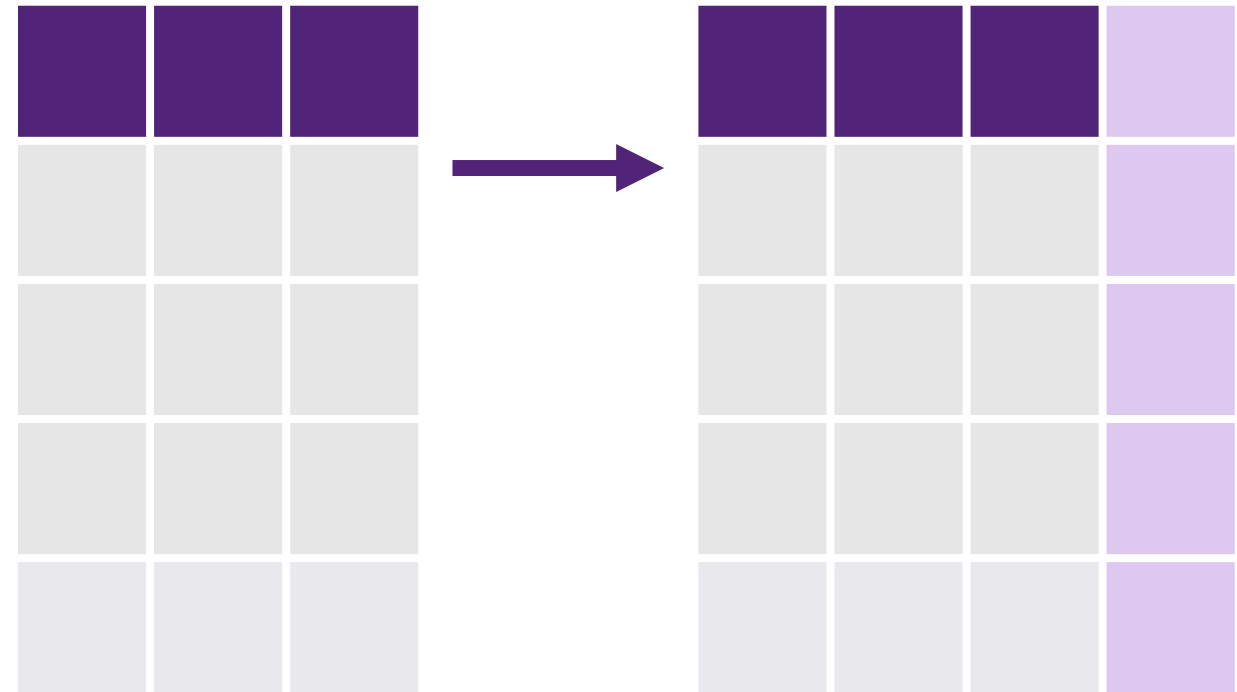
```
mutate(data, new_var = function(old_var))
```

`across()` – applies the same transformation to multiple columns

```
mutate(across(col1:col3, ~ .x^2))
```

```
mutate(across(where(is.numeric),  
              ~ .x^2))
```

```
mutate(across(where(is.numeric),  
              list(sq = ~.x^2)))
```



Reshaping Data

`pivot_longer()`

wide → long
transformations

wide

id	x	y	z
1	a	c	e
2	b	d	f

long

id	key	val
1	x	a
2	x	b
1	y	c
2	y	d
1	z	e
2	z	f

`pivot_wider()`

long → wide
transformations

Image source: <https://www.garrickadenbuie.com/project/tidyexplain/#spread-and-gather>

Merging data – mutating joins

x - left dataset

ID	V1	V2
1		
2		
3		

y - right dataset

ID	V3	V4
1		
2		
4		

left_join(x,y)

ID	V1	V2	V3	V4
1				
2				
3			NA	NA

right_join(x,y)

ID	V1	V2	V3	V4
1				
2				
4	NA	NA		

inner_join(x,y)

ID	V1	V2	V3	V4
1				
2				

full_join(x,y)

ID	V1	V2	V3	V4
1				
2				
3			NA	NA
4	NA	NA		

Collapsing data

original data

ID	V1	V2

grouped data
`group_by()`

ID	V1	V2

collapsed data
`summarise()`

ID	V1	m V2

aggregated value as new
column - `mutate()`

ID	V1	V2	m v2

Additional resources

https://github.com/AURIN-OFFICE/HASS_Summer_School

<https://cran.r-project.org>

<https://posit.co/resources/cheatsheets/>

<https://stackoverflow.com>

Google/ Bing/ DuckDuckGo/...

ChatGPT