# Data Massage - R

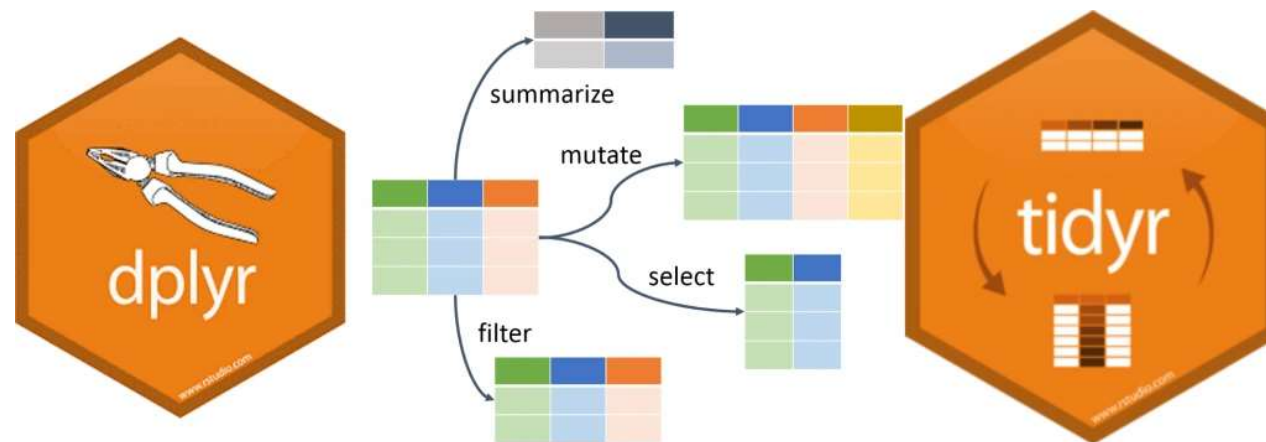**Introduction Data Massage with R**

**IT Academy**

# Index

- Data massaging with R

- What is tidyr?

- What is dplyr?

# Data massaging with R

Data Massaging is the process of extracting data to remove unneeded information or clean up a dataset to get into a usable format.

- To do it in R, we can use **dplyr** and **tidyr** packages

# What is tidyr?

**tidyr is a package that makes it easy to "tidy" your data.**

Tidy data is data that's easy to work with: it's easy to munge (with **dplyr**), visualise (with **ggplot2** or ggvis) and model (with R's hundreds of modelling packages). The two most important properties of tidy data are:

- Each variable forms a column
- Each observation forms a row



Tidy data graphic from R for Data Science

# What is dplyr?

**dplyr is a package which provides a set of tools for efficiently manipulating datasets in R**

You can track developement progress at https://github.com/tidyverse/dplyr

## dplyr

`CRAN 1.0.7`  `R-CMD-check passing`  `codecov 86%`  `R-CMD-check passing`

### Overview

dplyr is a grammar of data manipulation, providing a consistent set of verbs that help you solve the most common data manipulation challenges:

- `mutate()` adds new variables that are functions of existing variables
- `select()` picks variables based on their names.
- `filter()` picks cases based on their values.
- `summarise()` reduces multiple values down to a single summary.
- `arrange()` changes the ordering of the rows.

These all combine naturally with `group_by()` which allows you to perform any operation "by group". You can learn more about them in `vignette("dplyr")`. As well as these single-table verbs, dplyr also provides a variety of two-table verbs, which you can learn about in `vignette("two-table")`.

# Basics from tidyr and dplyr

| Package | Function | Use |
|---|---|---|
| dplyr | select | select variables/columns |
| dplyr | filter | select observations/rows |
| dplyr | Mutate | transform or recode variables |
| dplyr | summarize / summarise | summarize data |
| dplyr | group_by | identify subgroups for further processing |
| tidyr | gather | convert wide format dataset to long format |
| tidyr | spread | convert long format dataset to wide format |

# You are going to use:

# Data Transformation with dplyr : : **CHEAT SHEET**

**dplyr** functions work with pipes and expect **tidy data**. In tidy data:

Each **variable** is in its own **column**

&

Each **observation**, or **case**, is in its own **row**

**pipes**
x %>% f(y)
becomes f(x, y)

## Summarise Cases

These apply **summary functions** to columns to create a new table of summary statistics. Summary functions take vectors as input and return one value (see back).

**summary function**

**summarise**(.data, ...)
Compute table of summaries.
*summarise(mtcars, avg = mean(mpg))*

**count**(x, ..., wt = NULL, sort = FALSE)
Count number of rows in each group defined by the variables in ... Also **tally()**.
*count(iris, Species)*

**VARIATIONS**

**summarise_all()** - Apply funs to every column.
**summarise_at()** - Apply funs to specific columns.
**summarise_if()** - Apply funs to all cols of one type.

## Group Cases

Use **group_by()** to create a "grouped" copy of a table. dplyr functions will manipulate each "group" separately and then combine the results.
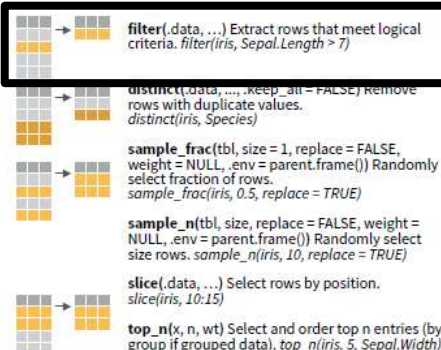
mtcars %>%
group_by(cyl) %>%
summarise(avg = mean(mpg))

**group_by**(.data, ..., add = FALSE)
Returns copy of table grouped by ...
*g_iris <- group_by(iris, Species)*

**ungroup**(x, ...)
Returns ungrouped copy of table.
*ungroup(g_iris)*

## Manipulate Cases

### EXTRACT CASES
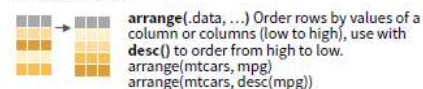
Row functions return a subset of rows as a new table.

**filter**(.data, ...) Extract rows that meet logical criteria. *filter(iris, Sepal.Length > 7)*

**distinct**(.data, ..., .keep_all = FALSE) Remove rows with duplicate values.
*distinct(iris, Species)*

**sample_frac**(tbl, size = 1, replace = FALSE, weight = NULL, .env = parent.frame()) Randomly select fraction of rows.
*sample_frac(iris, 0.5, replace = TRUE)*

**sample_n**(tbl, size, replace = FALSE, weight = NULL, .env = parent.frame()) Randomly select size rows. *sample_n(iris, 10, replace = TRUE)*

**slice**(.data, ...) Select rows by position.
*slice(iris, 10:15)*

**top_n**(x, n, wt) Select and order top n entries (by group if grouped data). *top_n(iris, 5, Sepal.Width)*

### Logical and boolean operators to use with filter()

| < | <= | is.na() | %in% | \| | xor() |
|---|----|---------|------|-----|-------|
| > | >= | !is.na() | ! | & | |

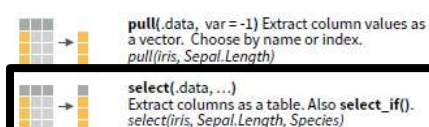See **?base::Logic** and **?Comparison** for help.

### ARRANGE CASES

**arrange**(.data, ...) 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))*

### ADD CASES

**add_row**(.data, ..., .before = NULL, .after = NULL) Add one or more rows to a table.
*add_row(faithful, eruptions = 1, waiting = 1)*

## Manipulate Variables

### EXTRACT VARIABLES

Column functions return a set of columns as a new vector or table.

**pull**(.data, var = -1) Extract column values as a vector. Choose by name or index.
*pull(iris, Sepal.Length)*

**select**(.data, ...)
Extract columns as a table. Also **select_if()**.
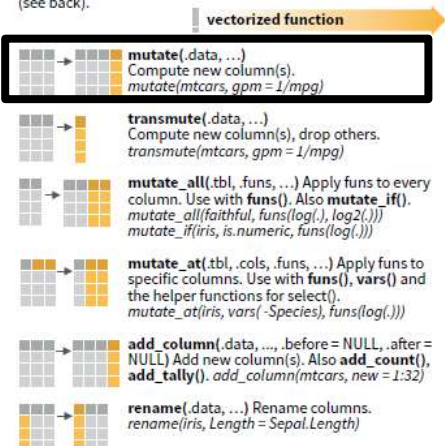*select(iris, Sepal.Length, Species)*

Use these helpers with select (),
e.g. select(iris, starts_with("Sepal"))

| **contains**(match) | **num_range**(prefix, range) | :, e.g. mpg:cyl |
| **ends_with**(match) | **one_of**(...) | -, e.g. -Species |
| **matches**(match) | **starts_with**(match) | |

### MAKE NEW VARIABLES

These apply **vectorized functions** to columns. Vectorized funs take vectors as input and return vectors of the same length as output (see back).

**vectorized function**

**mutate**(.data, ...)
Compute new column(s).
*mutate(mtcars, gpm = 1/mpg)*

**transmute**(.data, ...)
Compute new column(s), drop others.
*transmute(mtcars, gpm = 1/mpg)*

**mutate_all**(.tbl, .funs, ...) Apply funs to every column. Use with **funs()**. Also **mutate_if()**.
*mutate_all(faithful, funs(log(.), log2(.)))*
*mutate_if(iris, is.numeric, funs(log(.)))*

**mutate_at**(.tbl, .cols, .funs, ...) Apply funs to specific columns. Use with **funs()**, **vars()** and the helper functions for select().
*mutate_at(iris, vars( -Species), funs(log(.)))*

**add_column**(.data, ..., .before = NULL, .after = NULL) Add new column(s). Also **add_count()**, **add_tally()**. *add_column(mtcars, new = 1:32)*

**rename**(.data, ...) Rename columns.
*rename(iris, Length = Sepal.Length)*

**barcelona.cat/barcelonactiva**