# Exploring data #1

# Tidyverse and cheatsheets

# The "tidyverse"

So far, we have used a number of packages that are part of the *tidyverse*. The tidyverse is a collection of recent and developing packages for R, many written by Hadley Wickham.

"A giant among data nerds"

{ https://priceonomics.com/hadley-wickham-the-man-who-revolutionized-r/}

## Cheatsheets

RStudio has several very helpful **cheatsheets**. These are one-page sheets (front and back) that cover many of the main functions for a certain topic or task in R. These cheatsheets cover a lot of the main "tidyverse" functions.

You can access these directly from RStudio. Go to "Help" -> "Cheatsheets" and select the cheatsheet on the topic of interest.

You can find even more of these cheatsheets at https://www.rstudio.com/resources/cheatsheets/.

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

dplyr

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

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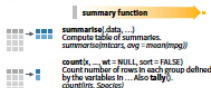
**Each observation, or case,** is in its own row

**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).
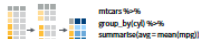
summary function

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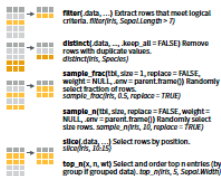
mtcars %>%
group_by(cyl) %>%
summarise(avg = mean(mpg))

**group_by(**_data, ..., add = FALSE_**)**
Returns copy of table grouped by ...
`g_iris <- group_by(iris, Species)`

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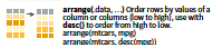Returns ungrouped copy of table.
`ungroup(g_iris)`

## Manipulate Cases

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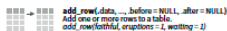Row functions return a subset of rows as a new table.

**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, now = 1:32)`

**rename(**_data, ..._**)** Rename columns.
`rename(iris, Length = Sepal.Length)`

**More reading / practice**

If you would like more reading and practice on what we've covered so far on transforming data, see chapter 5 of the "R for Data Science" book suggested at the start of the course.

As a reminder, that is available at:

http://r4ds.had.co.nz