# R WORKSHOP
# ON MILIEU INTERIEUR DATASETS

## DATA WRANGLING

Vincent Rouilly, Institut Pasteur Paris.

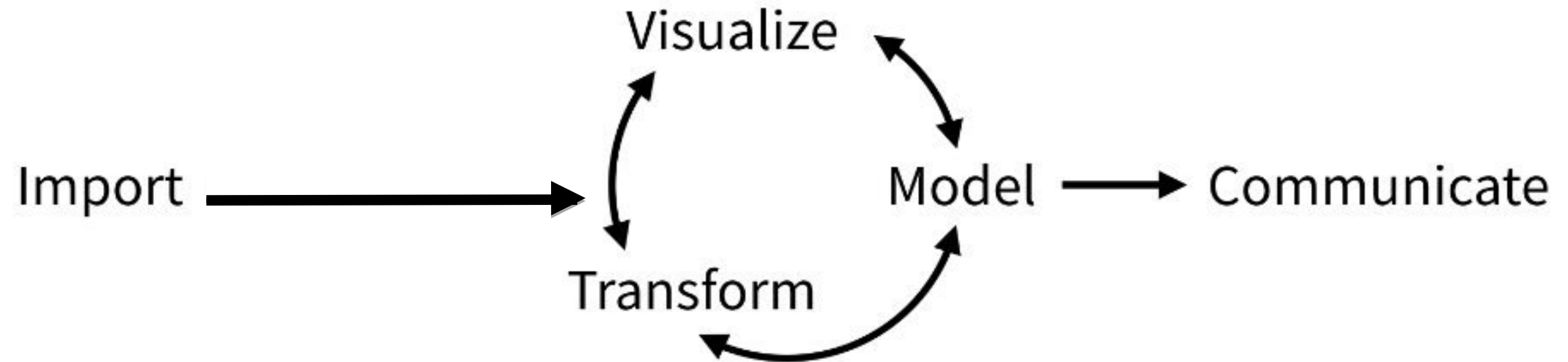# R WORKSHOP PROGRAM

| Times | Tuesday, June 17th. | Wednesday, June 18th. | Thursday, June 19th | Friday, June 20th |
|---|---|---|---|---|
| 9am | **Lecture**: Welcome & Workshop Introduction<br><br>Vincent Rouilly | **Lecture**: Milieu Intérieur Healthy Donor Studies<br><br>Darragh Duffy | **Lecture**: Comparing Immune Phenotypes Across Cohorts<br><br>Etienne Villain | **Lecture**: Transcriptional gene regulation & SES effects<br>Anthony Bertrand |
| 10am | Break | Break | Break | Break |
| 10.15am | RStudio set-up & Student introduction | Data Visualization | PCA and Clustering | Linear models |
| 12pm | Lunch | Lunch | Lunch | Lunch |
| 1pm | Data Wrangling | PCA and Clustering | Statistical tests (NHST) | Linear models |
| 2.30pm | Break | Break | Break | Break |
| 2.45pm | Data Wrangling / Data Visualization | PCA and Clustering | Statistical tests (NHST) | Linear models |
| 4pm | End of day | End of day | End of day | End of workshop |

# TYPICAL DATA ANALYSIS WORKFLOW

# TIDY DATA ANALYSIS

# Data import with the tidyverse : : CHEATSHEET

## Read Tabular Data with readr

**read_\***(file, col_names = TRUE, col_types = NULL, col_select = NULL, id = NULL, locale, n_max = Inf, skip = 0, na = c("", "NA"), guess_max = min(1000, n_max), show_col_types = TRUE) See **?read_delim**

**read_delim**("file.txt", delim = "|") Read files with any delimiter. If no delimiter is specified, it will automatically guess.
To make file.txt, run: write_file("A|B|C\n1|2|3\n4|5|NA", file = "file.txt")

**read_csv**("file.csv") Read a comma delimited file with period decimal marks.
write_file("A,B,C\n1,2,3\n4,5,NA", file = "file.csv")

**read_csv2**("file2.csv") Read semicolon delimited files with comma decimal marks.
write_file("A;B;C\n1,5;2;3\n4,5;5;NA", file = "file2.csv")

**read_tsv**("file.tsv") Read a tab delimited file. Also **read_table()**.
**read_fwf**("file.tsv", fwf_widths(c(2, 2, NA))) Read a fixed width file.
write_file("A\tB\tC\n1\t2\t3\n4\t5\tNA\n", file = "file.tsv")

### USEFUL READ ARGUMENTS

**No header**
read_csv("file.csv", col_names = FALSE)

**Provide header**
read_csv("file.csv",
    col_names = c("x", "y", "z"))

**Read multiple files into a single table**
read_csv(c("f1.csv", "f2.csv", "f3.csv"),
    id = "origin_file")

**Skip lines**
read_csv("file.csv", skip = 1)

**Read a subset of lines**
read_csv("file.csv", n_max = 1)

**Read values as missing**
read_csv("file.csv", na = c("1"))

**Specify decimal marks**
read_delim("file2.csv", locale =
    locale(decimal_mark = ","))

## Save Data with readr

**write_\***(x, file, na = "NA", append, col_names, quote, escape, eol, num_threads, progress)

**write_delim**(x, file, delim = " ") Write files with any delimiter.

**write_csv**(x, file) Write a comma delimited file.

**write_csv2**(x, file) Write a semicolon delimited file.

**write_tsv**(x, file) Write a tab delimited file.

---

One of the first steps of a project is to import outside data into R. Data is often stored in tabular formats, like csv files or spreadsheets.

The front page of this sheet shows how to import and save text files into R using **readr**.

The back page shows how to import spreadsheet data from Excel files using **readxl** or Google Sheets using **googlesheets4**.

### OTHER TYPES OF DATA
Try one of the following packages to import other types of files:

- **haven** - SPSS, Stata, and SAS files
- **DBI** - databases
- **jsonlite** - json
- **xml2** - XML
- **httr** - Web APIs
- **rvest** - HTML (Web Scraping)
- **readr::read_lines()** - text data

## Column Specification with readr

Column specifications define what data type each column of a file will be imported as. By default readr will generate a column spec when a file is read and output a summary.

**spec**(x) Extract the full column specification for the given imported data frame.

```
spec(x)
# cols(
#   age = col_integer(),
#   edu = col_character(),
#   earn = col_double()
# )
```
age is an integer
edu is a character
earn is a double (numeric)

### COLUMN TYPES
Each column type has a function and corresponding string abbreviation.

- **col_logical()** - "l"
- **col_integer()** - "i"
- **col_double()** - "d"
- **col_number()** - "n"
- **col_character()** - "c"
- **col_factor**(levels, ordered = FALSE) - "f"
- **col_datetime**(format = "") - "T"
- **col_date**(format = "") - "D"
- **col_time**(format = "") - "t"
- **col_skip()** - "-", "_"
- **col_guess()** - "?"

### USEFUL COLUMN ARGUMENTS

**Hide col spec message**
read_*(file, show_col_types = FALSE)

**Select columns to import**
Use names, position, or selection helpers.
read_*(file, col_select = c(age, earn))

**Guess column types**
To guess a column type, read_*() looks at the first 1000 rows of data. Increase with **guess_max**.
read_*(file, guess_max = Inf)

### DEFINE COLUMN SPECIFICATION

**Set a default type**
read_csv(
    file,
    col_type = list(.default = col_double())
)

**Use column type or string abbreviation**
read_csv(
    file,
    col_type = list(x = col_double(), y = "l", z = "_")
)

**Use a single string of abbreviations**
# col types: skip, guess, integer, logical, character
read_csv(
    file,
    col_type = "_?ilc"
)

# Import Spreadsheets

## with readxl

### READ EXCEL FILES



**read_excel(**path, sheet = NULL, range = NULL**)**
Read a .xls or .xlsx file based on the file extension.
See front page for more read arguments. Also
**read_xls()** and **read_xlsx()**.
read_excel("excel_file.xlsx")

### READ SHEETS

**read_excel(**path, **sheet = NULL)** Specify which sheet
to read by position or name.
read_excel(path, sheet = 1)
read_excel(path, sheet = "s1")

**excel_sheets(**path**)** Get a
vector of sheet names.
excel_sheets("excel_file.xlsx")

To **read multiple sheets:**
1. Get a vector of sheet
   names from the file path.
2. Set the vector names to
   be the sheet names.
3. Use purrr::map() and
   purrr::list_rbind() to read
   multiple files into one
   data frame.

```
path <- "your_file_path.xlsx"
path |>
  excel_sheets() |>
  set_names() |>
  map(read_excel, path = path) |>
  list_rbind()
```

### OTHER USEFUL EXCEL PACKAGES

For functions to write data to Excel files, see:
- **openxlsx**
- **writexl**

For working with non-tabular Excel data, see:
- **tidyxl**

---

### READXL COLUMN SPECIFICATION

Column specifications define what data type
each column of a file will be imported as.

Use the **col_types** argument of **read_excel()** to
set the column specification.

**Guess column types**
To guess a column type, read_ excel() looks at
the first 1000 rows of data. Increase with the
**guess_max** argument.
read_excel(path, guess_max = Inf)

**Set all columns to same type, e.g. character**
read_excel(path, col_types = "text")

**Set each column individually**
```
read_excel(
  path,
  col_types = c("text", "guess", "guess","numeric")
)
```

### COLUMN TYPES

| logical | numeric | text | date | list |
|---------|---------|------|------|------|
| TRUE | 2 | hello | 1947-01-08 | hello |
| FALSE | 3.45 | world | 1956-10-21 | 1 |

- skip
- guess
- logical
- numeric
- text
- date
- list

Use **list** for columns that include multiple data
types. See **tidyr** and **purrr** for list-column data.

---

### CELL SPECIFICATION FOR READXL AND GOOGLESHEETS4



Use the **range** argument of **readxl::read_excel()** or
**googlesheets4::read_sheet()** to read a subset of cells from a
sheet.
read_excel(path, range = "Sheet1!B1:D2")
read_sheet(ss, range = "B1:D2")

Also use the range argument with cell specification functions
**cell_limits()**, **cell_rows()**, **cell_cols()**, and **anchored()**.

---

## with googlesheets4

### READ SHEETS



**read_sheet(**ss, sheet = NULL, range = NULL**)**
Read a sheet from a URL, a Sheet ID, or a dribble
from the googledrive package. See front page for
more read arguments. Same as **range_read()**.

### SHEETS METADATA

**URLs** are in the form:
https://docs.google.com/spreadsheets/d/
    **SPREADSHEET_ID**/edit#gid=**SHEET_ID**

**gs4_get(**ss**)** Get spreadsheet meta data.

**gs4_find(**...**)** Get data on all spreadsheet files.

**sheet_properties(**ss**)** Get a tibble of properties
for each worksheet. Also **sheet_names()**.

### WRITE SHEETS



**write_sheet(**data, ss =
NULL, sheet = NULL**)**
Write a data frame into a
new or existing Sheet.

**gs4_create(**name, ...,
sheets = NULL**)** Create a
new Sheet with a vector
of names, a data frame,
or a (named) list of data
frames.

**sheet_append(**ss, data,
sheet = 1**)** Add rows to
the end of a worksheet.

---

### GOOGLESHEETS4 COLUMN SPECIFICATION

Column specifications define what data type
each column of a file will be imported as.

Use the **col_types** argument of **read_sheet()/
range_read()** to set the column specification.

**Guess column types**
To guess a column type read_sheet()/
range_read() looks at the first 1000 rows of data.
Increase with **guess_max**.
read_sheet(path, guess_max = Inf)

**Set all columns to same type, e.g. character**
read_sheet(path, col_types = "c")

**Set each column individually**
# col types: skip, guess, integer, logical, character
read_sheets(ss, col_types = "_?ilc")

### COLUMN TYPES

| I | n | c | D | L |
|---|---|---|---|---|
| TRUE | 2 | hello | 1947-01-08 | hello |
| FALSE | 3.45 | world | 1956-10-21 | 1 |

- skip - "_" or "-"
- guess - "?"
- logical - "l"
- integer - "i"
- double - "d"
- numeric - "n"
- date - "D"
- datetime - "T"
- character - "c"
- list-column - "L"
- cell - "C" Returns
  list of raw cell data.

Use list for columns that include multiple data
types. See **tidyr** and **purrr** for list-column data.

### FILE LEVEL OPERATIONS

**googlesheets4** also offers ways to modify other
aspects of Sheets (e.g. freeze rows, set column
width, manage (work)sheets). Go to
**googlesheets4.tidyverse.org** to read more.

For whole-file operations (e.g. renaming, sharing,
placing within a folder), see the tidyverse
package **googledrive** at
**googledrive.tidyverse.org**.

# Data transformation with dplyr : : **CHEATSHEET**

**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) |

## Summarize Cases

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**

**summarize**(.data, …) Compute table of summaries. mtcars |> summarize(avg = mean(mpg))

**count**(.data, …, wt = NULL, sort = FALSE, name = NULL) Count number of rows in each group defined by the variables in … Also **tally()**, **add_count()**, **add_tally()**. mtcars |> count(cyl)

## Group Cases

Use **group_by**(.data, …, .add = FALSE, .drop = TRUE) to create a "grouped" copy of a table grouped by columns in … dplyr functions will manipulate each "group" separately and combine the results.

mtcars |>
group_by(cyl) |>
summarize(avg = mean(mpg))

Use **rowwise**(.data, …) to group data into individual rows. dplyr functions will compute results for each row. Also apply functions to list-columns. See tidyr cheat sheet for list-column workflow.

starwars |>
rowwise() |>
mutate(film_count = length(films))

**ungroup**(x, …) Returns ungrouped copy of table.
g_mtcars <- mtcars |> group_by(cyl)
ungroup(g_mtcars)

## Manipulate Cases

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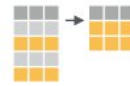
Row functions return a subset of rows as a new table.

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

**distinct**(.data, …, .keep_all = FALSE) Remove rows with duplicate values. mtcars |> distinct(gear)

**slice**(.data, …, .preserve = FALSE) Select rows by position. mtcars |> slice(10:15)

**slice_sample**(.data, …, n, prop, weight_by = NULL, replace = FALSE) Randomly select rows. Use n to select a number of rows and prop to select a fraction of rows. mtcars |> slice_sample(n = 5, replace = TRUE)

**slice_min**(.data, order_by, …, n, prop, with_ties = TRUE) and **slice_max()** Select rows with the lowest and highest values. mtcars |> slice_min(mpg, prop = 0.25)

**slice_head**(.data, …, n, prop) and **slice_tail()** Select the first or last rows. mtcars |> slice_head(n = 5)

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

### ADD CASES

**add_row**(.data, …, .before = NULL, .after = NULL) Add one or more rows to a table. cars |> add_row(speed = 1, dist = 1)

## Manipulate Variables

### EXTRACT VARIABLES

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

**pull**(.data, var = -1, name = NULL, …) Extract column values as a vector, by name or index. mtcars |> pull(wt)

**select**(.data, …) Extract columns as a table. mtcars |> select(mpg, wt)

**relocate**(.data, …, .before = NULL, .after = NULL) Move columns to new position. mtcars |> relocate(mpg, cyl, .after = last_col())

**Use these helpers with select() and across()**
e.g. mtcars |> select(mpg:cyl)

| contains(match) | num_range(prefix, range) | :, e.g., mpg:cyl |
| ends_with(match) | all_of(x)/any_of(x, …, vars) | !, e.g., !gear |
| starts_with(match) | matches(match) | everything() |

### MANIPULATE MULTIPLE VARIABLES AT ONCE

df <- tibble(x_1 = c(1, 2), x_2 = c(3, 4), y = c(4, 5))

**across**(.cols, .funs, …, .names = NULL) Summarize or mutate multiple columns in the same way. df |> summarize(across(everything(), mean))
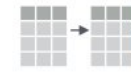
**c_across**(.cols) Compute across columns in row-wise data. df |> rowwise() |> mutate(x_total = sum(c_across(1:2)))

### MAKE NEW VARIABLES

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

**vectorized function**

**mutate**(.data, …, .keep = "all", .before = NULL, .after = NULL) Compute new column(s). Also **add_column()**. mtcars |> mutate(gpm = 1 / mpg) mtcars |> mutate(gpm = 1 / mpg, .keep = "none")

**rename**(.data, …) Rename columns. Use **rename_with()** to rename with a function. mtcars |> rename(miles_per_gallon = mpg)

# Vectorized Functions

## TO USE WITH MUTATE ()

**mutate()** applies vectorized functions to columns to create new columns. Vectorized functions take vectors as input and return vectors of the same length as output.

| vectorized function |

### OFFSET
dplyr::**lag()** - offset elements by 1
dplyr::**lead()** - offset elements by -1

### CUMULATIVE AGGREGATE
dplyr::**cumall()** - cumulative all()
dplyr::**cumany()** - cumulative any()
    **cummax()** - cumulative max()
dplyr::**cummean()** - cumulative mean()
    **cummin()** - cumulative min()
    **cumprod()** - cumulative prod()
    **cumsum()** - cumulative sum()

### RANKING
dplyr::**cume_dist()** - proportion of all values <=
dplyr::**dense_rank()** - rank w ties = min, no gaps
dplyr::**min_rank()** - rank with ties = min
dplyr::**ntile()** - bins into n bins
dplyr::**percent_rank()** - min_rank scaled to [0,1]
dplyr::**row_number()** - rank with ties = "first"

### MATH
    **+, -, \*, /, ^, %/%, %%** - arithmetic ops
    **log(), log2(), log10()** - logs
    **<, <=, >, >=, !=, ==** - logical comparisons
dplyr::**between()** - x >= left & x <= right
dplyr::**near()** - safe == for floating point numbers

### MISCELLANEOUS
dplyr::**case_when()** - multi-case if_else()
```
starwars |>
  mutate(type = case_when(
    height > 200 | mass > 200 ~ "large",
    species == "Droid"        ~ "robot",
    TRUE                      ~ "other")
  )
```
dplyr::**coalesce()** - first non-NA values by element across a set of vectors
dplyr::**if_else()** - element-wise if() + else()
dplyr::**na_if()** - replace specific values with NA
    **pmax()** - element-wise max()
    **pmin()** - element-wise min()

# Summary Functions

## TO USE WITH SUMMARIZE ()

**summarize()** applies summary functions to columns to create a new table. Summary functions take vectors as input and return single values as output.

| summary function |

### COUNT
dplyr::**n()** - number of values/rows
dplyr::**n_distinct()** - # of uniques
    **sum(!is.na())** - # of non-NAs

### POSITION
    **mean()** - mean, also **mean(!is.na())**
    **median()** - median

### LOGICAL
    **mean()** - proportion of TRUEs
    **sum()** - # of TRUEs

### ORDER
dplyr::**first()** - first value
dplyr::**last()** - last value
dplyr::**nth()** - value in nth location of vector

### RANK
    **quantile()** - nth quantile
    **min()** - minimum value
    **max()** - maximum value

### SPREAD
    **IQR()** - Inter-Quartile Range
    **mad()** - median absolute deviation
    **sd()** - standard deviation
    **var()** - variance

# Row Names

Tidy data does not use rownames, which store a variable outside of the columns. To work with the rownames, first move them into a column.


tibble::**rownames_to_column()**
Move row names into col.
```
a <- mtcars |>
  rownames_to_column(var = "C")
```


tibble::**column_to_rownames()**
Move col into row names.
```
a |> column_to_rownames(var = "C")
```

Also tibble::**has_rownames()** and tibble::**remove_rownames()**.

# Combine Tables

## COMBINE VARIABLES



**bind_cols(**…, .name_repair**)** Returns tables placed side by side as a single table. Column lengths must be equal. Columns will NOT be matched by id (to do that look at Relational Data below), so be sure to check that both tables are ordered the way you want before binding.

## RELATIONAL DATA

Use a "**Mutating Join**" to join one table to columns from another, matching values with the rows that they correspond to. Each join retains a different combination of values from the tables.

**left_join(**x, y, by = NULL, copy = FALSE, suffix = c(".x", ".y"), …, keep = FALSE, na_matches = "na"**)** Join matching values from y to x.

**right_join(**x, y, by = NULL, copy = FALSE, suffix = c(".x", ".y"), …, keep = FALSE, na_matches = "na"**)** Join matching values from x to y.

**inner_join(**x, y, by = NULL, copy = FALSE, suffix = c(".x", ".y"), …, keep = FALSE, na_matches = "na"**)** Join data. Retain only rows with matches.

**full_join(**x, y, by = NULL, copy = FALSE, suffix = c(".x", ".y"), …, keep = FALSE, na_matches = "na"**)** Join data. Retain all values, all rows.

## COLUMN MATCHING FOR JOINS

Use **by = c("col1", "col2", …)** to specify one or more common columns to match on.
left_join(x, y, by = "A")

Use a named vector, **by = c("col1" = "col2")**, to match on columns that have different names in each table.
left_join(x, y, by = c("C" = "D"))

Use **suffix** to specify the suffix to give to unmatched columns that have the same name in both tables.
left_join(x, y, by = c("C" = "D"),
suffix = c("1", "2"))

## COMBINE CASES



**bind_rows(**…, .id = NULL**)** Returns tables one on top of the other as a single table. Set .id to a column name to add a column of the original table names (as pictured).

Use a "**Filtering Join**" to filter one table against the rows of another.



**semi_join(**x, y, by = NULL, copy = FALSE, …, na_matches = "na"**)** Return rows of x that have a match in y. Use to see what will be included in a join.

**anti_join(**x, y, by = NULL, copy = FALSE, …, na_matches = "na"**)** Return rows of x that do not have a match in y. Use to see what will not be included in a join.

Use a "**Nest Join**" to inner join one table to another into a nested data frame.

**nest_join(**x, y, by = NULL, copy = FALSE, keep = FALSE, name = NULL, …**)** Join data, nesting matches from y in a single new data frame column.

## SET OPERATIONS

**intersect(x, y, …)**
Rows that appear in both x and y.

**setdiff(x, y, …)**
Rows that appear in x but not y.

**union(x, y, …)**
Rows that appear in x or y, duplicates removed). **union_all()** retains duplicates.

Use **setequal()** to test whether two data sets contain the exact same rows (in any order).

# Data tidying with tidyr : : **CHEATSHEET**

**Tidy data** is a way to organize tabular data in a consistent data structure across packages. A table is tidy if:
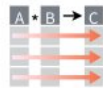
Each **variable** is in its own **column**   &   Each **observation**, or **case**, is in its own row

Access **variables** as **vectors**

Preserve **cases** in vectorized operations

## Tibbles

**AN ENHANCED DATA FRAME**

Tibbles are a table format provided by the **tibble** package. They inherit the data frame class, but have improved behaviors:

- **Subset** a new tibble with ], a vector with [[ and $.
- **No partial matching** when subsetting columns.
- **Display** concise views of the data on one screen.

**options**(tibble.print_max = n, tibble.print_min = m, tibble.width = Inf) Control default display settings.

**View()** or **glimpse()** View the entire data set.

### CONSTRUCT A TIBBLE

**tibble**(…) Construct by columns.
tibble(x = 1:3, y = c("a", "b", "c"))

**tribble**(…) Construct by rows.
tribble(~x,  ~y,
        1, "a",
        2, "b",
        3, "c")

Both make this tibble

```
A tibble: 3 × 2
      x     y
  <int> <chr>
1     1     a
2     2     b
3     3     c
```

**as_tibble**(x, …) Convert a data frame to a tibble.

**enframe**(x, name = "name", value = "value") Convert a named vector to a tibble. Also **deframe()**.

**is_tibble**(x) Test whether x is a tibble.

## Reshape Data - Pivot data to reorganize values into a new layout.

table4a

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

"Lengthen" data by collapsing several columns into two. Column names move to a new names_to column and values to a new values_to column.

pivot_longer(table4a, cols = 2:3, names_to = "year", values_to = "cases")

table2

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

The inverse of pivot_longer(). "Widen" data by expanding two columns into several. One column provides the new column names, the other the values.

pivot_wider(table2, names_from = type, values_from = count)

## Split Cells - Use these functions to split or combine cells into individual, isolated values.

table5

**unite**(data, col, …, sep = "_", remove = TRUE, na.rm = FALSE) Collapse cells across several columns into a single column.

unite(table5, century, year, col = "year", sep = "")

table3

**separate_wider_delim**(data, cols, delim, …, names = NULL, names_sep = NULL, names_repair = "check unique", too_few, too_many, cols_remove = TRUE) Separate each cell in a column into several columns. Also **separate_wider_regex()** and **separate_wider_position()**.

separate(table3, rate, sep = "/", into = c("cases", "pop"))

table3

**separate_longer_delim**(data, cols, delim, .., width, keep_eampty) Separate each cell in a column into several rows.

separate_longer_delim(table3, rate, sep = "/")

## Expand Tables

Create new combinations of variables or identify implicit missing values (combinations of variables not present in the data).

x

**expand**(data, …) Create a new tibble with all possible combinations of the values of the variables listed in … Drop other variables.
expand(mtcars, cyl, gear, carb)

x

**complete**(data, …, fill = list()) Add missing possible combinations of values of variables listed in … Fill remaining variables with NA.
complete(mtcars, cyl, gear, carb)

## Handle Missing Values

Drop or replace explicit missing values (NA).

x

**drop_na**(data, …) Drop rows containing NA's in … columns.
drop_na(x, x2)

x

**fill**(data, …, .direction = "down") Fill in NA's in … columns using the next or previous value.
fill(x, x2)

x

**replace_na**(data, replace) Specify a value to replace NA in selected columns.
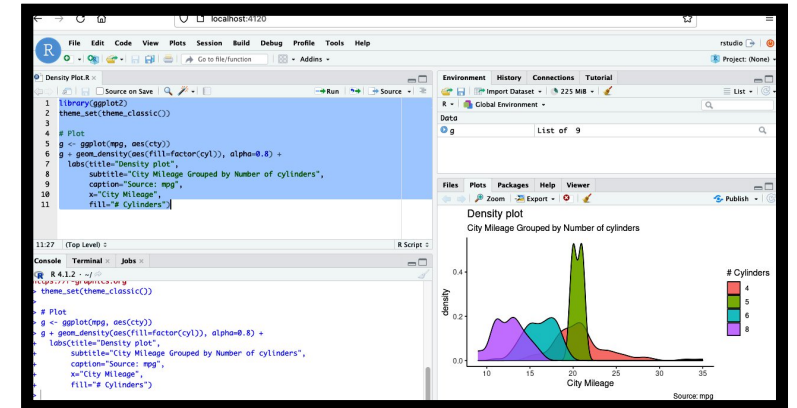replace_na(x, list(x2 = 2))

posit®

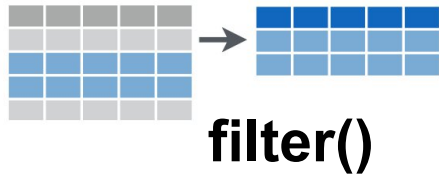# GENERAL WORKSHOP PRINCIPLES



Topic Introduction
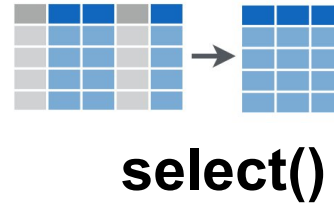Overview

Cookbook
R Recipes
Synthetic Data
Use case

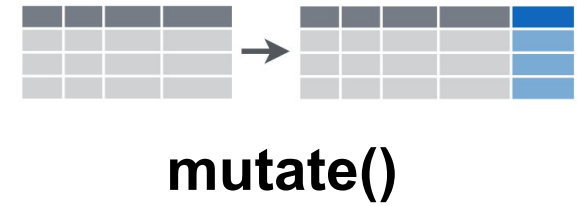Application
to MI Datasets

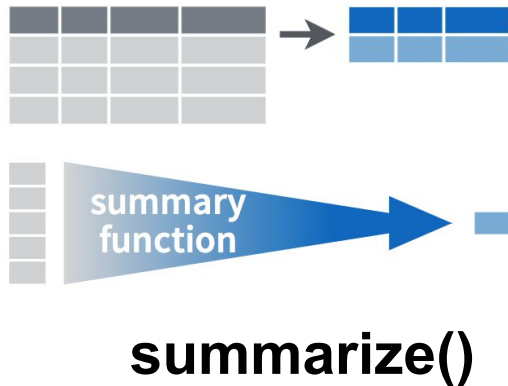# DATA WRANGLING OBJECTIVES

**Subset Observations** (Rows)

**filter()**

**Subset Variables** (Columns)

**select()**

**Make New Variables**

**mutate()**

**Summarise Data**

summary function

**summarize()**

**Group by rows**

**group_by()**

# TIDY DATA

# THANK YOU


Questions? Comments?