# Data wrangling

Importing and data cleaning

**Jordan Creed**
Moffitt Cancer Center

June 16, 2021

# Elements of "tidy" data

Data comes in all shapes and sizes, but not all data organization is made equal!

3 elements of tidy data:

1. Each variable has its own column
2. Each observation has its own row
3. Each value has it own cell

👉 For more information please see chapter 11 of R for Data Science

# Elements of "tidy" data

When working with tidy data, we can use the same tools in similar ways for different datasets...

...but working with untidy data often means reinventing the wheel with one-time approaches that are hard to iterate or reuse.

# TCGA Data

# The Cancer Genome Atlas (TCGA)

TCGA is the large cancer genomics program spanning 33 cancer types and multiple institutions.

- HNSC - Head and neck squamous cell carcinoma
- KIRC - Kidney renal clear cell carcinoma

Data is saved as two `.csv`s: `data/tcga-clinical.csv` and `data/tcga-gene-exp.csv`

# Importing files

# Import method depends on file type!

- Most text files : `readr::read_delim()` or `readr::read_csv()`
- Microsoft excel files : `readxl::read_excel()`
- R data files : `load()`
- SAS data files: `haven::read_sas()`

🧑‍💻 remember to copy code from the import gui to your scripts!

# Data Import : : CHEAT SHEET

R's **tidyverse** is built around **tidy data** stored in **tibbles**, which are enhanced data frames.

The front side of this sheet shows how to read text files into R with **readr**.

The reverse side shows how to create tibbles with **tibble** and to layout tidy data with **tidyr**.

### OTHER TYPES OF DATA
Try one of the following packages to import other types of files
- **haven** - SPSS, Stata, and SAS files
- **readxl** - excel files (.xls and .xlsx)
- **DBI** - databases
- **jsonlite** - json
- **xml2** - XML
- **httr** - Web APIs
- **rvest** - HTML (Web Scraping)

## Save Data

Save **x**, an R object, to **path**, a file path, as:

**Comma delimited file**
write_csv(x, path, na = "NA", append = FALSE, col_names = !append)

**File with arbitrary delimiter**
write_delim(x, path, delim = " ", na = "NA", append = FALSE, col_names = !append)

**CSV for excel**
write_excel_csv(x, path, na = "NA", append = FALSE, col_names = !append)

**String to file**
write_file(x, path, append = FALSE)

**String vector to file, one element per line**
write_lines(x, path, na = "NA", append = FALSE)

**Object to RDS file**
write_rds(x, path, compress = c("none", "gz", "bz2", "xz"), ...)

**Tab delimited files**
write_tsv(x, path, na = "NA", append = FALSE, col_names = !append)

## Read Tabular Data - These functions share the common arguments:

read_*(file, col_names = TRUE, col_types = NULL, locale = default_locale(), na = c("", "NA"), quoted_na = TRUE, comment = "", trim_ws = TRUE, skip = 0, n_max = Inf, guess_max = min(1000, n_max), progress = interactive())

**Comma Delimited Files**
read_csv("file.csv")
To make file.csv run:
write_file(x = "a,b,c\n1,2,3\n4,5,NA", path = "file.csv")

**Semi-colon Delimited Files**
read_csv2("file2.csv")
write_file(x = "a;b;c\n1;2;3\n4;5;NA", path = "file2.csv")

**Files with Any Delimiter**
read_delim("file.txt", delim = "|")
write_file(x = "a|b|c\n1|2|3\n4|5|NA", path = "file.txt")

**Fixed Width Files**
read_fwf("file.fwf", col_positions = c(1, 3, 5))
write_file(x = "a b c\n1 2 3\n4 5 NA", path = "file.fwf")

**Tab Delimited Files**
read_tsv("file.tsv") Also read_table().
write_file(x = "a\tb\tc\n1\t2\t3\n4\t5\tNA", path = "file.tsv")

### USEFUL ARGUMENTS

**Example file**
write_file("a,b,c\n1,2,3\n4,5,NA","file.csv")
f <- "file.csv"

**No header**
read_csv(f, col_names = FALSE)

**Provide header**
read_csv(f, col_names = c("x", "y", "z"))

**Skip lines**
read_csv(f, skip = 1)

**Read in a subset**
read_csv(f, n_max = 1)

**Missing Values**
read_csv(f, na = c("1", "."))

## Read Non-Tabular Data

**Read a file into a single string**
read_file(file, locale = default_locale())

**Read each line into its own string**
read_lines(file, skip = 0, n_max = -1L, na = character(), locale = default_locale(), progress = interactive())

**Read Apache style log files**
read_log(file, col_names = FALSE, col_types = NULL, skip = 0, n_max = -1, progress = interactive())

**Read a file into a raw vector**
read_file_raw(file)

**Read each line into a raw vector**
read_lines_raw(file, skip = 0, n_max = -1L, progress = interactive())

## Data types

readr functions guess the types of each column and convert types when appropriate (but will NOT convert strings to factors automatically).

A message shows the type of each column in the result.

```
## Parsed with column specification:
## cols(
##     age = col_integer(),
##     sex = col_character(),
##     earn = col_double()
## )
```

age is an integer
sex is a character
earn is a double (numeric)

1. Use **problems()** to diagnose problems.
   x <- read_csv("file.csv"); problems(x)

2. Use a col_ function to guide parsing.
   - **col_guess()** - the default
   - **col_character()**
   - **col_double()**, **col_euro_double()**
   - **col_datetime**(format = "") Also
     **col_date**(format = ""), **col_time**(format = "")
   - **col_factor**(levels, ordered = FALSE)
   - **col_integer()**
   - **col_logical()**
   - **col_number()**, **col_numeric()**
   - **col_skip()**
   x <- read_csv("file.csv", col_types = cols(
     A = col_double(),
     B = col_logical(),
     C = col_factor()))

3. Else, read in as character vectors then parse with a parse_ function.
   - **parse_guess()**
   - **parse_character()**
   - **parse_datetime()** Also **parse_date()** and **parse_time()**
   - **parse_double()**
   - **parse_factor()**
   - **parse_integer()**
   - **parse_logical()**
   - **parse_number()**
   x$A <- parse_number(x$A)

# Data cleaning with dplyr

# dplyr verbs and syntax

select() subset columns

- select(tibble_name, variable1, variable2)

mutate() create new variables/columns

- mutate(tibble_name, new_variable = variable1*2)

# common `select()` options

`:` selects a range of columns

`-` selects every column but those specified

`starts_with()` and `ends_with()` selects columns whose names start/end with the specified string

`contains()` selects columns whose names contain the specified string

# dplyr verbs and syntax

`filter()` subset rows

- `filter(tibble_name, variable1 == "specific value")`

`summarise()` aggregates rows

- `summarise(tibble_name, avg_value = mean(variable1))`

`arrange()` orders rows

- `arrange(tibble_name, variable1)`

Operations can be chained together with the pipe operator (`%>%`)

# logical tests for `filter()`

- `==` : equal
- `!=` : not equal
- `<` & `<=` : less than & less than or equal to
- `>` & `>=` : greater than & greater than or equal to
- `|` : or
- `&` : and
- `!` : not
- `%in%` : in the set

# dplyr helpers

`n()` the number of rows

`n_distinct()` the number of unique values for a variable

`group_by()` collects observations by a common value
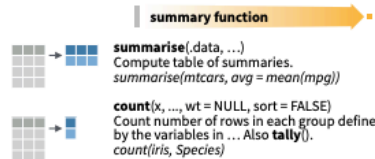
# Data Transformation with dplyr : : CHEAT SHEET

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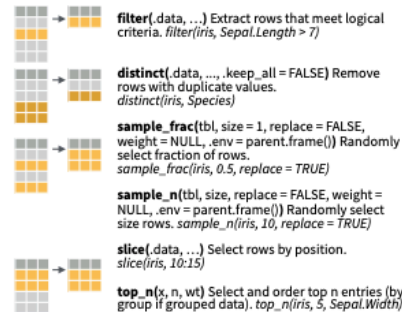
mtcars %>%
group_by(cyl) %>%
summarise(avg = mean(mpg))

**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

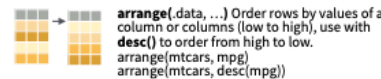Row functions return a subset of rows as a new table.

**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() | ! | & | |

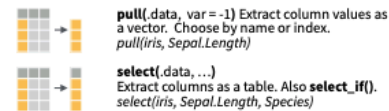See **?base::Logic** and **?Comparison** for help.

### 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)*

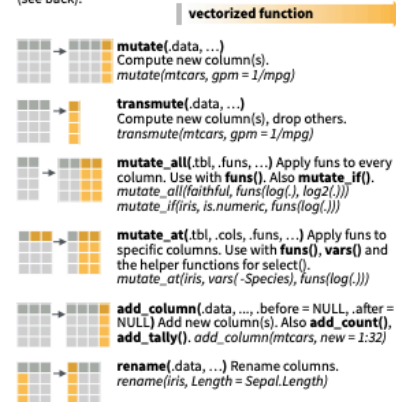**select**(.data, ...) Extract columns as a table. Also **select_if()**. *select(iris, Sepal.Length, Species)*

**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)*