# Przetwarzanie i wizualizacja danych

prof. UAM dr hab. Tomasz Górecki

tomasz.gorecki@amu.edu.pl

Zakład Statystyki Matematycznej i Analizy Danych
Wydział Matematyki i Informatyki
Uniwersytet im. Adama Mickiewicza w Poznaniu

A data set is **tidy** iff:

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

2. Each **case** is in its own **row**

3. Each **value** is in its own **cell**



## tame data

| baker | cinnamon_1 | cardamom_2 | nutmeg_3 |
|---|---|---|---|
| Emma | 1 | 0 | 1 |
| Harry | 1 | 1 | 1 |
| Ruby | 1 | 0 | 1 |
| Zainab | 0 | NA | 0 |

| trial | Emma | Harry | Ruby | Zainab |
|---|---|---|---|---|
| cinnamon_1 | 1 | 1 | 1 | 0 |
| cardamom_2 | 0 | 1 | 0 | NA |
| nutmeg_3 | 1 | 1 | 1 | 0 |

## tidy data

| baker | spice | order | correct |
|---|---|---|---|
| Emma | Cinnamon | 1 | 1 |
| Harry | Cinnamon | 1 | 1 |
| Ruby | Cinnamon | 1 | 1 |
| Zainab | Cinnamon | 1 | 0 |
| Emma | Cardamom | 2 | 0 |
| Harry | Cardamom | 2 | 1 |
| Ruby | Cardamom | 2 | 0 |
| Zainab | Cardamom | 2 | NA |
| Emma | Nutmeg | 3 | 1 |
| Harry | Nutmeg | 3 | 1 |
| Ruby | Nutmeg | 3 | 1 |
| Zainab | Nutmeg | 3 | 0 |

# Data tidying with tidyr : : **CHEAT SHEET**

**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.

**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")

**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.

**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 = "")

**separate**(data, col, into, sep = "[^[:alnum:]]+", remove = TRUE, convert = FALSE, extra = "warn", fill = "warn", …) Separate each cell in a column into several columns. Also **extract**().

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

**separate_rows**(data, …, sep = "[^[:alnum:].]+", convert = FALSE) Separate each cell in a column into several rows.

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

## Expand Tables

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

**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)

**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).

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

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

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

R Studio

dplyr jest pakietem zaprojektowanym do wydajnego (napisany w C++) i efektywnego manipulowania danymi. Dodatkowo wszystko co robimy na ramce danych możemy też zrobić na innych obiektach np. tabelach w bazie danych. Jest on obecnie częścią większej grupy pakietów zwanej tidyverse.

# 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) |
| **ends_with**(match) | **one_of**(…) |
| **matches**(match) | **starts_with**(match) |

:, e.g. mpg:cyl
-, e.g. -Species

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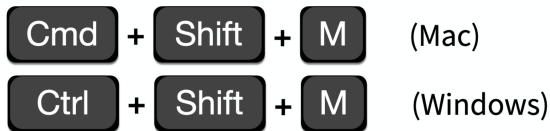
**rename**(.data, …) Rename columns.
*rename(iris, Length = Sepal.Length)*

R Studio

## Vector Functions

**TO USE WITH MUTATE ()**

**mutate()** and **transmute()** apply 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 |

**OFFSETS**

dplyr::**lag()** - Offset elements by 1
dplyr::**lead()** - Offset elements by -1

**CUMULATIVE AGGREGATES**

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

**RANKINGS**

dplyr::**cume_dist()** - Proportion of all values <=
dplyr::**dense_rank()** - rank with 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

**MISC**

dplyr::**case_when()** - multi-case if_else()
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()
dplyr::**recode()** - Vectorized switch()
dplyr::**recode_factor()** - Vectorized switch() for factors

## Summary Functions

**TO USE WITH SUMMARISE ()**

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

> | summary function |

**COUNTS**

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

**LOCATION**

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

**LOGICALS**

    **mean()** - Proportion of TRUE's
    **sum()** - # of TRUE's

**POSITION/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.

**rownames_to_column()**
Move row names into col.
a <- rownames_to_column(iris, var = "C")

**column_to_rownames()**
Move col in row names.
column_to_rownames(a, var = "C")

Also **has_rownames()**, **remove_rownames()**

## Combine Tables

**COMBINE VARIABLES**

Use **bind_cols()** to paste tables beside each other as they are.

**bind_cols(...)** Returns tables placed side by side as a single table.
BE SURE THAT ROWS ALIGN.

Use a "**Mutating Join**" to join one table to columns from another, matching values for each row. Each join retains a different combination of values from the tables.

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

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

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

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

Use **by = c("col1", "col2")** to specify the column(s) to match on.
left_join(x, y, by = "A")

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

Use **suffix** to specify suffix to give to duplicate column names.
left_join(x, y, by = c("C" = "D"), suffix = c("1", "2"))

**COMBINE CASES**

Use **bind_rows()** to paste tables below each other as they are.

**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)

**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).

**EXTRACT ROWS**

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

**semi_join**(x, y, by = NULL, ...)
Return rows of x that have a match in y. USEFUL TO SEE WHAT WILL BE JOINED.

**anti_join**(x, y, by = NULL, ...)
Return rows of x that do not have a match in y. USEFUL TO SEE WHAT WILL NOT BE JOINED.

- filter,
- arrange,
- % > % – strumienie.

## Shortcut to type %>%

| Cmd | + | Shift | + | M |  (Mac)

| Ctrl | + | Shift | + | M |  (Windows)

- select,
- mutate,
- summarise,
- group_by.