# 20 Maart 2018

Herman Teirlinck,
01.17 - Clara Peeters

# What have I done?!?

```r
library(rvest)
library(dplyr)
library(magrittr)


waarnemingen <- read_html("https://waarnemingen.be/")
waarnemingen %>%
    html_nodes("table") %>%
    .[[3]] %>%
    html_table() %>%
    set_colnames(c("intro", "datum", "count", "soort", "photo", "gebied")) %>%
    select(datum, count, soort, gebied) %>%
    slice(-1)
```

Note: Rvest is a ' *tidyverse-adjacent*' package

DATA
MANIPULATION

PART 2

# Data Wrangling
## with dplyr and tidyr
### Cheat Sheet
**R** Studio

## Tidy Data - A foundation for wrangling in R

In a tidy data set:

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

&

Each **observation** is saved in its own **row**

Tidy data complements R's **vectorized operations**. R will automatically preserve observations as you manipulate variables. No other format works as intuitively with R.

M * A

## Syntax - Helpful conventions for wrangling

**dplyr::tbl_df(iris)**
Converts data to tbl class. tbl's are easier to examine than data frames. R displays only the data that fits onscreen:

```
Source: local data frame [150 x 5]

  Sepal.Length Sepal.Width Petal.Length
1          5.1         3.5          1.4
2          4.9         3.0          1.4
3          4.7         3.2          1.3
4          4.6         3.1          1.5
5          5.0         3.6          1.4
..         ...         ...          ...
Variables not shown: Petal.Width (dbl),
    Species (fctr)
```

**dplyr::glimpse(iris)**
Information dense summary of tbl data.

**utils::View(iris)**
View data set in spreadsheet-like display (note capital V).

**dplyr::%>%**
Passes object on left hand side as first argument (or . argument) of function on righthand side.

```
x %>% f(y)    is the same as f(x, y)
y %>% f(x, ., z)    is the same as f(x, y, z )
```

"Piping" with %>% makes code more readable, e.g.

```
iris %>%
  group_by(Species) %>%
  summarise(avg = mean(Sepal.Width)) %>%
  arrange(avg)
```
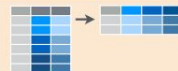
## Reshaping Data - Change the layout of a data set

**tidyr::gather(cases, "year", "n", 2:4)**
Gather columns into rows.

**tidyr::spread(pollution, size, amount)**
Spread rows into columns.

**tidyr::separate(storms, date, c("y", "m", "d"))**
Separate one column into several.

**tidyr::unite(data, col, ..., sep)**
Unite several columns into one.

**dplyr::data_frame(a = 1:3, b = 4:6)**
Combine vectors into data frame (optimized).

**dplyr::arrange(mtcars, mpg)**
Order rows by values of a column (low to high).

**dplyr::arrange(mtcars, desc(mpg))**
Order rows by values of a column (high to low).

**dplyr::rename(tb, y = year)**
Rename the columns of a data frame.

## Subset Observations (Rows)

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

**dplyr::distinct(iris)**
Remove duplicate rows.

**dplyr::sample_frac(iris, 0.5, replace = TRUE)**
Randomly select fraction of rows.

**dplyr::sample_n(iris, 10, replace = TRUE)**
Randomly select n rows.

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

**dplyr::top_n(storms, 2, date)**
Select and order top n entries (by group if grouped data).

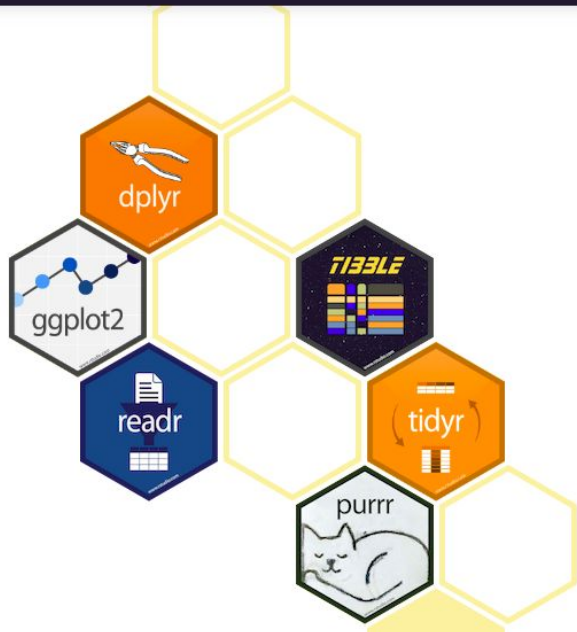| Logic in R - ?Comparison, ?base::Logic | | |
|---|---|---|
| < | Less than | != | Not equal to |
| > | Greater than | %in% | Group membership |
| == | Equal to | is.na | Is NA |
| <= | Less than or equal to | !is.na | Is not NA |
| >= | Greater than or equal to | &,|,!,xor,any,all | Boolean operators |

## Subset Variables (Columns)

**dplyr::select(iris, Sepal.Width, Petal.Length, Species)**
Select columns by name or helper function.

### Helper functions for select - ?select

**select(iris, contains("."))**
Select columns whose name contains a character string.

**select(iris, ends_with("Length"))**
Select columns whose name ends with a character string.

**select(iris, everything())**
Select every column.

**select(iris, matches(".t."))**
Select columns whose name matches a regular expression.

**select(iris, num_range("x", 1:5))**
Select columns named x1, x2, x3, x4, x5.

**select(iris, one_of(c("Species", "Genus")))**
Select columns whose names are in a group of names.

**select(iris, starts_with("Sepal"))**
Select columns whose name starts with a character string.

**select(iris, Sepal.Length:Petal.Width)**
Select all columns between Sepal.Length and Petal.Width (inclusive).

**select(iris, -Species)**
Select all columns except Species.

devtools::install_github("rstudio/EDAWR") for data sets          Learn more with browseVignettes(package = c("dplyr", "tidyr")) • dplyr 0.4.0• tidyr 0.2.0 • Updated: 1/15

## Summarise Data          ## Make New Variables          ## Combine Data Sets

# R packages for data science

The tidyverse is an opinionated **collection of R packages** designed for data science. All packages share an underlying design philosophy, grammar, and data structures.

Install the complete tidyverse with:

```
install.packages("tidyverse")
```

## Learn the tidyverse

See how the tidyverse makes data science faster, easier and more fun with "R for Data Science". Read it **online**,

Install the package suite:

```
install.packages("tidyverse")
```

Load the package suite:

```
library(tidyverse)
```

# TIDY?!?

See https://inbo.github.io/dwc-in-R/tidy.html#14

# Share your snippets during the coding session!

Go to [https://hackmd.io/7Yd3NsCFTwqHbRnHZbhlzg](https://hackmd.io/7Yd3NsCFTwqHbRnHZbhlzg) and post your code in between backticks:

*For example*:

```
library(dplyr)

my_data <- ...
```

Excel might contain column names with capital letters, spaces, etc., which can be annoying to select:

```
brandganzen <- read_excel("./data/20180123_brandganzen.xlsx")

brandganzen %>% select(`Locatie vangst`) # Ugh
```

With `janitor` your column names can be cleaned (lowercase, underscores instead of spaces). In addition, you can remove empty rows:

```
library(janitor) # Also tidyverse, but not loaded by default

brandganzen %>%
  remove_empty_rows() %>% # Additional step to remove empty rows
  clean_names() -> brandganzen

brandganzen %>% select(locatie_vangst)
```
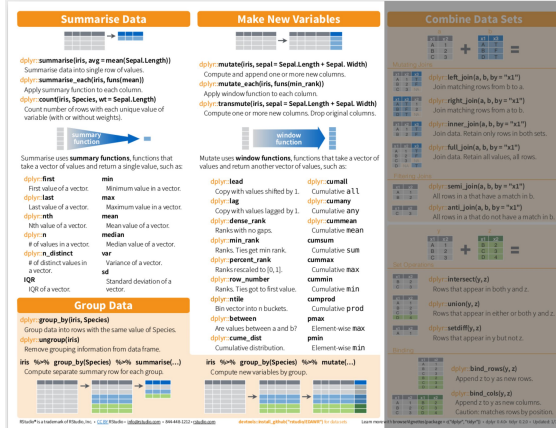
# recap/showcase



`20180222_survey_data_spreadsheet_tidy.csv`

1. Show min, max, mean weight per sex and species and save as a new object (df) `weight_per_species_sex`

2. Execute the following:
   a. Rename column `'weight_in_g'` to `'weight'`
   b. Replace `'weight'` values with values in kg
   c. Add column `'country'` with value `'US'`
   d. Save as new object `'data_kg_US'`

`20180123_brandganzen.xlsx`

1. How many adult geese per sex can you count (consider `'onbekend'` as a sex)?

```
Geslacht       n
  <chr>     <int>
1 Man          63
2 Onbekend     19
3 Vrouw       119
```

2. How many different catching methods were used in each location?

```
`Locatie vangst` n_methods
  <chr>               <int>
1 DEINZE                  1
2 DESTELBERGEN            1
...
```

# recap/showcase

`20180222_survey_data_spreadsheet_tidy.csv`

1. Show min, max, mean wei~~~~and
   species an~~~~object (df)
   `w~~~~_species_sex`

   *group_by & summarise*

2. Execute the following:
   a. Rename column `'weight_in_g'` to `'weight'`
   b. Replace ~~~~with values in kg

   *rename & mutate*

   c. ~~~~`country'` with value `'US'`
   d. Save as new object `'data_kg_US'`



`20180123_brandganzen.xlsx`

1. How many adult geese per~~~~you
   count (~~~~):

   *filter & group_by & count*

   ```
                    63
   2 Onbekend       19
   3 Vrouw         119
   ```

2. How many different catching ~~~~s
   were used in ~~~~

   *group_by & summarise*

   ```
                     1
   ~~BERGEN          1
   ...
   ```

# The ⬛ concept

We defined a number of challenges. If you were able to achieve a challenge, add a to ⬛ r laptop screen.

## The objective is that **everyone** achieves ⬛!

- Someone has more ⬛ than you? **Ask for help!**

- Someone has less ⬛ than you? **Provide help!**

- Download coding club material and work locally, not in sync with the Google drive

data

images

src

- Create new Rstudio project in the **/src** folder

New Project

Back    **Create Project from Existing Directory**

Project working directory:
~/projecten/2018_coding_club/src    Browse...

☐ Open in new session          Create Project    Cancel

- Download coding club material and work locally, not in sync with the Google drive
- Create new Rstudio project in the **src** folder…
- Use relative paths to data files:

```
> library(readr)
```

```
> read_csv2("../data/20180123_gent_groeiperwijk.csv")
```

Read in the data set

20180123_gent_groeiperwijk.csv

This is NOT a *tidy* data set!

Make this a tidy data set:

```
   wijk                            year growth
   <chr>                          <int>  <int>
 1 Binnenstad                      1999   - 36
 2 Bloemekenswijk                  1999     12
 3 Brugse Poort - Rooigem          1999     85
 4 Dampoort                        1999    107
 5 Drongen                         1999      3
 6 Elisabethbegijnhof - Papegaai   1999   -  4
 7 Gentbrugge                      1999      4
 8 Kanaaldorpen en -zone           1999      5
 9 Ledeberg                        1999   -  4
10 Macharius - Heirnis             1999     47
# ... with 265 more rows
```

Read in the [20180222_surveys.csv](20180222_surveys.csv) and the [20180222_species.csv](20180222_species.csv) data.

Join the species information columns (genus, species , taxa) to the survey data set, using the common identifier.

Compare the result when applying the different commands to join the data...



inner_join(x, y)

left_join(x, y)

full_join(x, y)

right_join(x, y)

For the [20180123_observations_NPHK_cameratrapping.csv](20180123_observations_NPHK_cameratrapping.csv) data:

- count the observed humans for each month:

```
  sequenceMonth humans_observed
          <int>           <int>
1             5               1
2             6               1
3             7              38
4             8             153
5             9              38
6            10              25
```

- add an additional column with the counts for each animal-deploymentID combination

```
   sequenceDay sequenceMonth sequenceYear deploymentSamplingPoint animalVernacularName animalCount  point_animal_counts
         <int>         <int>        <int> <chr>                   <chr>                      <int>                <int>
1            7             7         2017 JW_0090                 Ass                            1                   12
2            6             7         2017 JW_0090                 Ass                            3                   12
3            7             7         2017 JW_0090                 Ass                            1                   12
4           15             7         2017 JW_0090                 Ass                            1                   12
5           16             7         2017 JW_0090                 Ass                            2                   12
6           27             7         2017 JW_0090                 Ass                            2                   12
7           27             7         2017 JW_0090                 Ass                            1                   12
8           27             7         2017 JW_0090                 Ass                            1                   12
9           13             5         2017 JW_0020                 Beech Marten                   1                    3
10          13             5         2017 JW_0020                 Beech Marten                   1                    3
# ... with 3,750 more rows
```

More challenges!

For the [stierkikker](#) formulieren data, derive all the columns concerning
`blankvoorn` and remove those rows for which all values are `NA`:
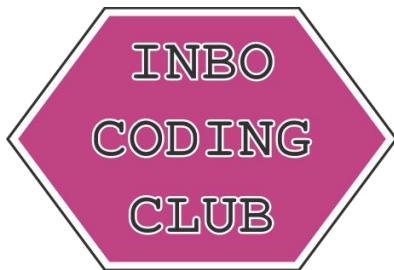
```
# A tibble: 95 x 12
   `Fuik 1 - Bijvangs…  `Fuik 2 - Bijvangs…  `Fuik 3 - Bijvangs…  `Fuik 4 - Bijvangs…  `Fuik 5 - Bijvang…  `Fuik 6 - Bijvang…
   <chr>                <chr>                <chr>                <chr>                <chr>               <chr>
 1 NA                   5-10cm               NA                   NA                   NA                  NA
 2 ??                   ??                   NA                   NA                   NA                  NA
 3 ??                   NA                   NA                   NA                   NA                  NA
 4 ??                   NA                   NA                   NA                   NA                  NA
 5 ??                   NA                   NA                   NA                   NA                  NA
 6 ??                   NA                   NA                   NA                   NA                  NA
 7 5-10cm               NA                   NA                   NA                   NA                  NA
 8 5-10cm               NA                   NA                   NA                   NA                  NA
 9 5-10cm               5-10cm               5-10cm               NA                   NA                  NA
10 NA                   5-10cm               5-10cm               5-10cm               5-10cm              5-10cm
# ... with 85 more rows, and 6 more variables: `Fuik 7 - Bijvangst [Blankvoorn]`   <chr>, `Fuik 8 - Bijvangst
#   [Blankvoorn]`  <chr>, `Fuik 9 - Bijvangst [Blankvoorn]`   <chr>, `Fuik 10 - Bijvangst [Blankvoorn]`   <chr>, `Fuik 11 -
#   Bijvangst [Blankvoorn]`  <chr>, `Fuik 12 - Bijvangst [Blankvoorn]`   <chr>
```

For the `20180123_rainfall_klemskerke_clean.csv` data, calculate the yearly rainfall sum from 2012 till 2016:

```
 year                  value
   <dttm>               <dbl>
1 2012-01-01 00:00:00    934
2 2013-01-01 00:00:00    701
3 2014-01-01 00:00:00    727
4 2015-01-01 00:00:00    789
5 2016-01-01 00:00:00    775
```

**INBO CODING CLUB**

Zaal: Herman Teirlinck - 01.21 - Jeanne Brabants
Datum: 26/04/2018, van 10:00 tot 12:00
*(registratie aangekondigd via DG_useR@inbo.be)*