

# Tidyverse tutorial 2 - More advanced operations

Ariane Ducellier

10/05/2023

Load R packages.

```
library(tidyverse)
```

```
## -- Attaching core tidyverse packages ----- tidyverse 2.0.0 --
## v dplyr      1.1.3      v readr      2.1.4
## v forcats    1.0.0      v stringr   1.5.0
## v ggplot2    3.4.3      v tibble    3.2.1
## v lubridate  1.9.2      v tidyr     1.3.0
## v purrr      1.0.2
## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()     masks stats::lag()
## i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to become errors
```

## 1. Dealing with missing data

```
header <- c("age", "workclass", "fnlwgt", "education",
            "education_num", "marital_status", "occupation",
            "relationship", "race", "sex", "capital_gain",
            "capital_loss", "hours_per_week", "native_country", "target")
df <- read_csv("https://archive.ics.uci.edu/ml/machine-learning-databases/adult/adult.data",
              col_names=header, trim_ws=TRUE)
```

```
## Rows: 32561 Columns: 15
## -- Column specification -----
## Delimiter: ","
## chr (9): workclass, education, marital_status, occupation, relationship, rac...
## dbl (6): age, fnlwgt, education_num, capital_gain, capital_loss, hours_per_week
##
## i Use 'spec()' to retrieve the full column specification for this data.
## i Specify the column types or set 'show_col_types = FALSE' to quiet this message.
```

```
df <-df %>%
  mutate(workclass = na_if(workclass, "?"),
         occupation = na_if(occupation, "?"),
         native_country = na_if(native_country, "?"))
```

## 1.1 Filling values with previous value

```
df_fill11 <- df %>%  
  fill(workclass, occupation, native_country, .direction="down")
```

## 1.2 Filling values with most frequent value

For categorical variables.

```
m_freq_workcls <- names(table(df$workclass))[which.max(table(df$workclass))]  
m_freq_occup <- names(table(df$occupation))[which.max(table(df$occupation))]  
df_fill12 <- df %>%  
  replace_na(list(workclass = m_freq_workcls,  
                  occupation = m_freq_occup))
```

## 1.3 Dropping rows with missing values

```
df_no_na <- df %>%  
  drop_na()
```