

LAB: Bivariate analysis

2024-01-31

- M1 MIDS & MFA
- [Université Paris Cité](#)
- Année 2023-2024
- [Course Homepage](#)



- Moodle

```
to_be_loaded <- c("tidyverse",
                  "glue",
                  "magrittr",
                  "lobstr",
                  "arrow",
                  "ggforce",
                  "vcd",
                  "ggmosaic",
                  "htr",
                  "cowplot",
                  "patchwork"
)

for (pck in to_be_loaded) {
  if (!require(pck, character.only = T)) {
    install.packages(pck, repos="http://cran.rstudio.com/")
    stopifnot(require(pck, character.only = T))
  }
}
```

Objectives

In Exploratory analysis of tabular data, bivariate analysis is the second step. It consists in exploring, summarizing, visualizing pairs of columns of a dataset.

Bivariate techniques depend on the types of columns we are facing.

For *numerical/numerical* samples

- Scatter plots

- Smoothed lineplots (for example linear regression)
- 2-dimensional density plots

For *categorical/categorical* samples : mosaicplots and variants

For *numerical/categorical* samples

- Boxplots per group
- Histograms per group
- Density plots per group
- Quantile-Quantile plots


Dataset

Once again we rely on the Census dataset.

Since 1948, the US Census Bureau carries out a monthly Current Population Survey, collecting data concerning residents aged above 15 from 150000 households. This survey is one of the most important sources of information concerning the american workforce. Data reported in file `Recensement.txt` originate from the 2012 census.

Load the data into the session environment and call it `df`. Take advantage of the fact that we saved the result of our data wrangling job in a self-documented file format. Download a `parquet` file from the following URL:

<https://stephane-v-boucheron.fr/data/Recensement.parquet>

 Use `httr::GET()` and `WriteBin()`.

Categorical/Categorical pairs

```
df |>
  select(where(is.factor)) |>
  head()
```

```
# A tibble: 6 x 9
  SEXE REGION STAT_MARI SYNDICAT CATEGORIE NIV_ETUDES NB_PERS NB_ENF REV_FOYER
<fct> <fct> <fct> <fct> <fct> <fct> <fct> <fct> <fct>
1 F NE C non "Administ~ Bachelor 2 0 [35000-4~
2 M W M non "Building~ 12 years ~ 2 0 [17500-2~
3 M S C non "Administ~ Associate~ 2 0 [75000-1~
4 M NE D oui "Services" 12 years ~ 4 1 [17500-2~
5 M W M non "Services" 9 years s~ 8 1 [75000-1~
6 M NW C non "Services" 12 years ~ 6 0 [1e+05-1~
```

Explore the connection between `CATEGORIE` and `SEX`. Compute the 2-ways contingency table using `table()`, and `count()` from `dplyr`.

Use `tibble::as_tibble()` to transform the output of `table()` into a dataframe/tibble.

Use `tidyr::pivot_wider()` so as to obtain a wide (but messy) tibble with the same the same shape as the output of `table()`. Can you spot a difference?

Use `mosaicplot()` from base R to visualize the contingency table.

Use `geom_mosaic` from `ggmosaic` to visualize the contingency table

- Make the plot as readable as possible
- Reorder `CATEGORIE` according to counts
- Collapse rare levels of `CATEGORIE` (consider that a level is rare if it has less than 40 occurrences). Use tools from `forcats`.

Testing association

Chi-square independence/association test

Categorical/Numerical pairs

Grouped boxplots

Plot boxplots of `AGE` according to `NIV_ETUDES`

Draw density plots of `AGE`, facet by `NIV_ETUDES` and `SEXE`

Collapse rare levels of `NIV_ETUDES` and replay.

Numerical/Numerical pairs

Make a scatterplot of `SAL_HOR` with respect to `AGE`

pairs from base R

`ggpairs()`

Useful links

- [rmarkdown](#)
- [dplyr](#)
- [ggplot2](#)
- *R Graphic Cookbook*. Winston Chang. O' Reilly.
- [A blog on ggplot object](#)
- `skimr`
- `vcd`
- `ggmosaic`
- `ggforce`
- `arrow`
- `httr`