

French given names per year per department

Lucas Mello Schnorr, Jean-Marc Vincent

October, 2022

The problem context

The aim of the activity is to develop a methodology to answer a specific question on a given dataset.

The dataset is the set of Firstname given in France on a large period of time. <https://www.insee.fr/fr/statistiques/2540004>, we choose this dataset because it is sufficiently large, you can't do the analysis by hand, the structure is simple

You need to use the *tidyverse* for this analysis. Unzip the file *dpt2020_txt.zip* (to get the **dpt2020.csv**). Read in R with this code. Note that you might need to install the **readr** package with the appropriate command.

Download Raw Data from the website

```
file = "dpt2021_csv.zip"
if(!file.exists(file)){
  download.file("https://www.insee.fr/fr/statistiques/fichier/2540004/dpt2021_csv.zip",
    destfile=file)
}
unzip(file)
```

Build the Dataframe from file

```
library(tidyverse)
```

```
## Warning: package 'tidyr' was built under R version 4.3.2
```

```
## Warning: package 'readr' was built under R version 4.3.2
```

```
## Warning: package 'stringr' was built under R version 4.3.2
```

```
## -- Attaching core tidyverse packages ----- tidyverse 2.0.0 --
## v dplyr      1.1.3      v readr      2.1.5
## v forcats    1.0.0      v stringr    1.5.1
## v ggplot2    3.4.3      v tibble     3.2.1
## v lubridate  1.9.3      v tidyr      1.3.1
```

```
## 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
```

```
FirstNames <- read_delim("dpt2021.csv",delim=";")
```

```
## Rows: 3784673 Columns: 5
## -- Column specification -----
## Delimiter: ";"
## chr (3): preusuel, annais, dpt
## dbl (2): sexe, nombre
##
## 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.
```

Let see summary about the data.

```
# Summary of the data
data_summary <- summary(FirstNames)
print(data_summary)
```

```
##      sexe      preusuel      annais      dpt
## Min.   :1.000   Length:3784673   Length:3784673   Length:3784673
## 1st Qu.:1.000   Class :character   Class :character   Class :character
## Median :2.000   Mode  :character   Mode  :character   Mode  :character
## Mean    :1.535
## 3rd Qu.:2.000
## Max.    :2.000
##      nombre
## Min.     : 3.0
## 1st Qu.: 4.0
## Median  : 7.0
## Mean    : 23.1
## 3rd Qu.: 18.0
## Max.    :6307.0
```

```
head_table <- head(FirstNames,10)
print(head_table)
```

```
## # A tibble: 10 x 5
##   sexe preusuel      annais dpt  nombre
##   <dbl> <chr>      <chr> <chr> <dbl>
## 1     1 _PRENOMS_RARES 1900   02     7
## 2     1 _PRENOMS_RARES 1900   04     9
## 3     1 _PRENOMS_RARES 1900   05     8
## 4     1 _PRENOMS_RARES 1900   06    23
## 5     1 _PRENOMS_RARES 1900   07     9
## 6     1 _PRENOMS_RARES 1900   08     4
```

```
## 7      1 _PRENOMS_RARES 1900  09      6
## 8      1 _PRENOMS_RARES 1900  10      3
## 9      1 _PRENOMS_RARES 1900  11     11
## 10     1 _PRENOMS_RARES 1900  12      7
```

```
library(conflicted)
```

```
## Warning: package 'conflicted' was built under R version 4.3.2
```

```
conflict_prefer("filter", "dplyr")
```

```
## [conflicted] Will prefer dplyr::filter over any other package.
```

```
conflict_prefer("lag", "dplyr")
```

```
## [conflicted] Will prefer dplyr::lag over any other package.
```

1. Choose a firstname and analyse its frequency along time. Compare several firstnames frequency.

Answer 1:

```
# Load required libraries
library(tidyverse)

# Read data with semicolon delimiter
FirstNames <- read_delim("dpt2021.csv", delim = ";", show_col_types = FALSE)

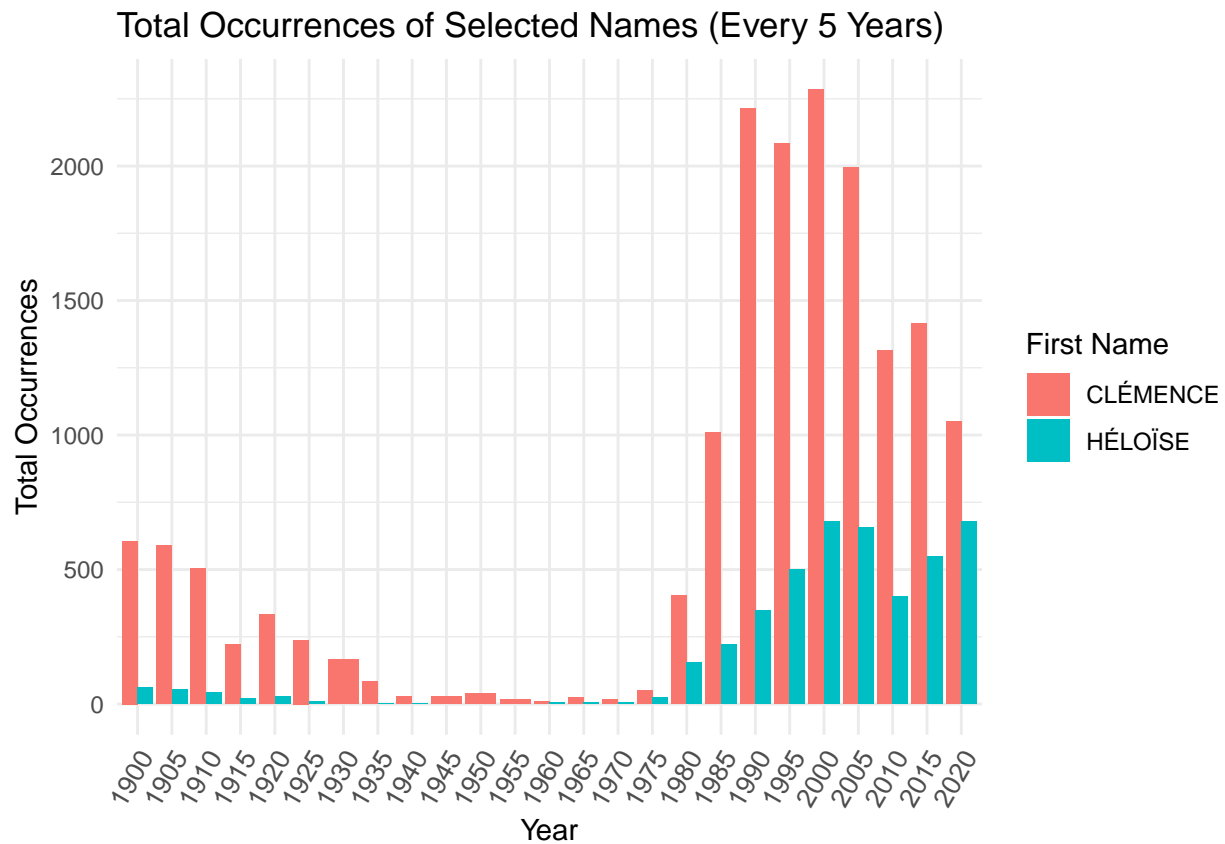
# Selected names
selected_names <- c("CLÉMENCE", "HÉLOÏSE")

# Filter data for the selected names and valid years
filtered_data <- FirstNames %>%
  filter(preusuel %in% selected_names) %>%
  filter(str_detect(annais, "\\d{4}$")) %>%
  mutate(annais_numeric = as.numeric(annais)) %>%
  filter(annais_numeric %% 5 == 0)

# Group and summarize the filtered data
grouped_data <- filtered_data %>%
  group_by(annais, preusuel) %>%
  summarise(total_occurrences = sum(nombre, na.rm = TRUE), .groups = "drop")

# Create a bar chart with bars next to each other for each year
ggplot(grouped_data, aes(x = annais, y = total_occurrences, fill = preusuel)) +
  geom_bar(stat = "identity", position = "dodge") +
  labs(title = "Total Occurrences of Selected Names (Every 5 Years)",
```

```
x = "Year",
y = "Total Occurrences",
fill = "First Name") +
theme_minimal() +
theme(axis.text.x = element_text(size = 10, angle = 60, hjust = 1))
```



```
# Adjust the size parameter and angle
```

We can see from this that CLEMENCE is the more frequent compared to HELOISE.

2. Establish by gender the most given firstname by year. Analyse the evolution of the most frequent firstname.

Answer 2:

```
# Remove rows with NA values in the 'annais' column
FirstNames <- FirstNames[complete.cases(FirstNames$annais), ]

# Find the most frequent name for each gender in each year
most_frequent_names <- FirstNames %>%
  group_by(sexe, annais, preusuel) %>%
  summarise(total_occurrences = sum(nombre), .groups = "drop") %>%
```

```

group_by(sexe, annais) %>%
  arrange(desc(total_occurrences)) %>%
  slice(1) %>%
  ungroup()

# Convert 'annais' to numeric
most_frequent_names$annais <- as.numeric(as.character(most_frequent_names$annais))

## Warning: NAs introduced by coercion

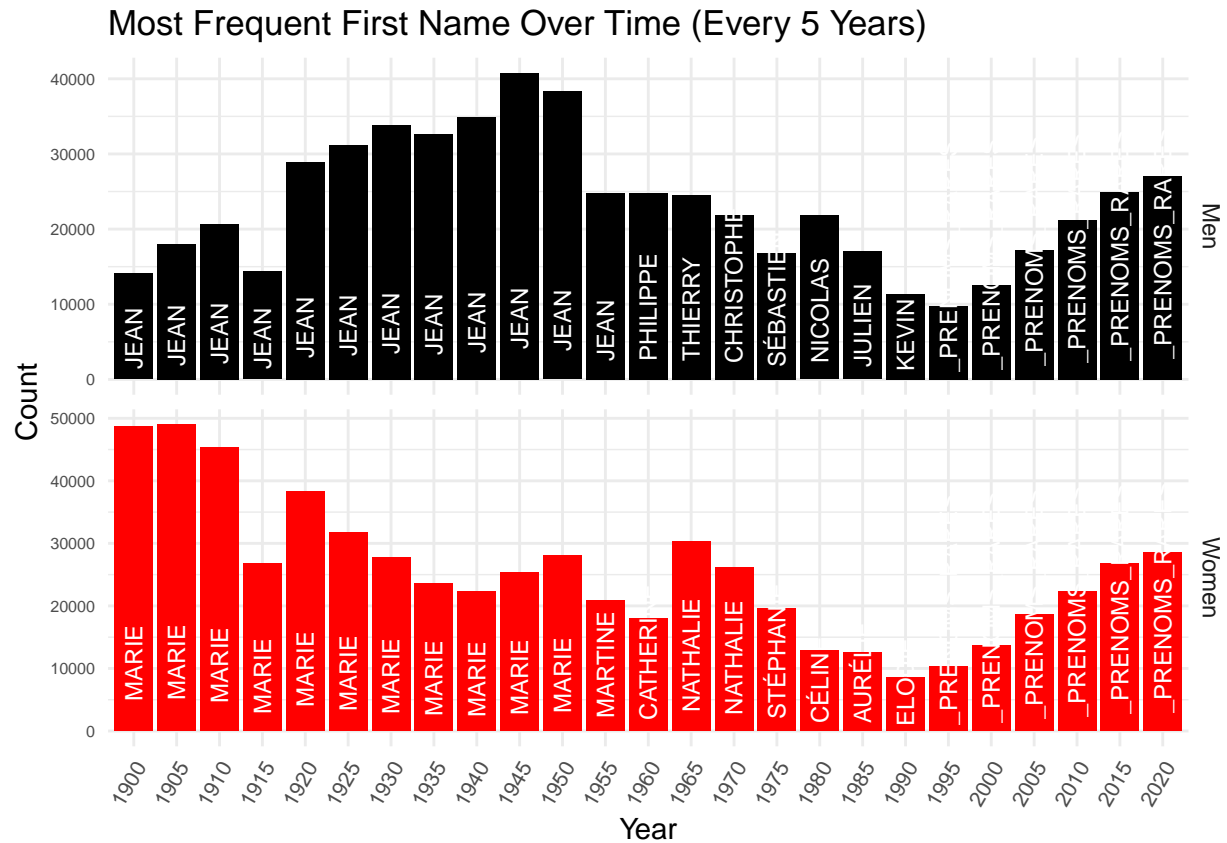
most_frequent_names <- most_frequent_names %>%
  filter(annais %% 5 == 0)

# Determine the color for each year based on gender
most_frequent_names <- most_frequent_names %>%
  mutate(color = ifelse(sexe == 1, "black", "red"))

# Create a bar chart with facets for each gender
p<-ggplot(most_frequent_names, aes(x = as.factor(annais),
                                   y = total_occurrences, fill = color, label = preusuel)) +
  geom_bar(stat = "identity") +
  geom_text(position = position_stack(vjust = 0.1), size = 3,
            color = "white", angle = 90, hjust = 0) + # Add labels vertically with white text
  labs(title = "Most Frequent First Name Over Time (Every 5 Years)",
        x = "Year",
        y = "Count") +
  scale_fill_identity() +
  theme_minimal() +
  theme(axis.text.x = element_text(size = 8, angle = 60, hjust = 1),
        axis.text.y = element_text(size = 6)) + # Adjust text size for y-axis
  facet_grid(sexe ~ ., scales = "free_y",
             labeller = labeller(sexe = c("1" = "Men", "2" = "Women")))) +
  theme(legend.position = "none") # Remove default legend

print(p)

```



We can see how Marie is the most frequent name in women , and Jean in men.

3. Optional : Which department has a larger variety of names along time ? Is there some sort of geographical correlation with the data?