

French given names per year per department

Edited by : **Oussama Oulkaid**

October, 2021

1. Introduction

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. given names data set of INSEE, we choose this dataset because it is sufficiently large, the analysis cannot be done by hand, the structure is simple.

We will use the *tidyverse* for this analysis. The file **dpt2019.csv** contains the data.

```
# The environment
library(tidyverse)
library(ggplot2)
```

2. Build the Dataframe from file

```
FirstNames <- read_delim("dpt2020.csv", delim =";")
```

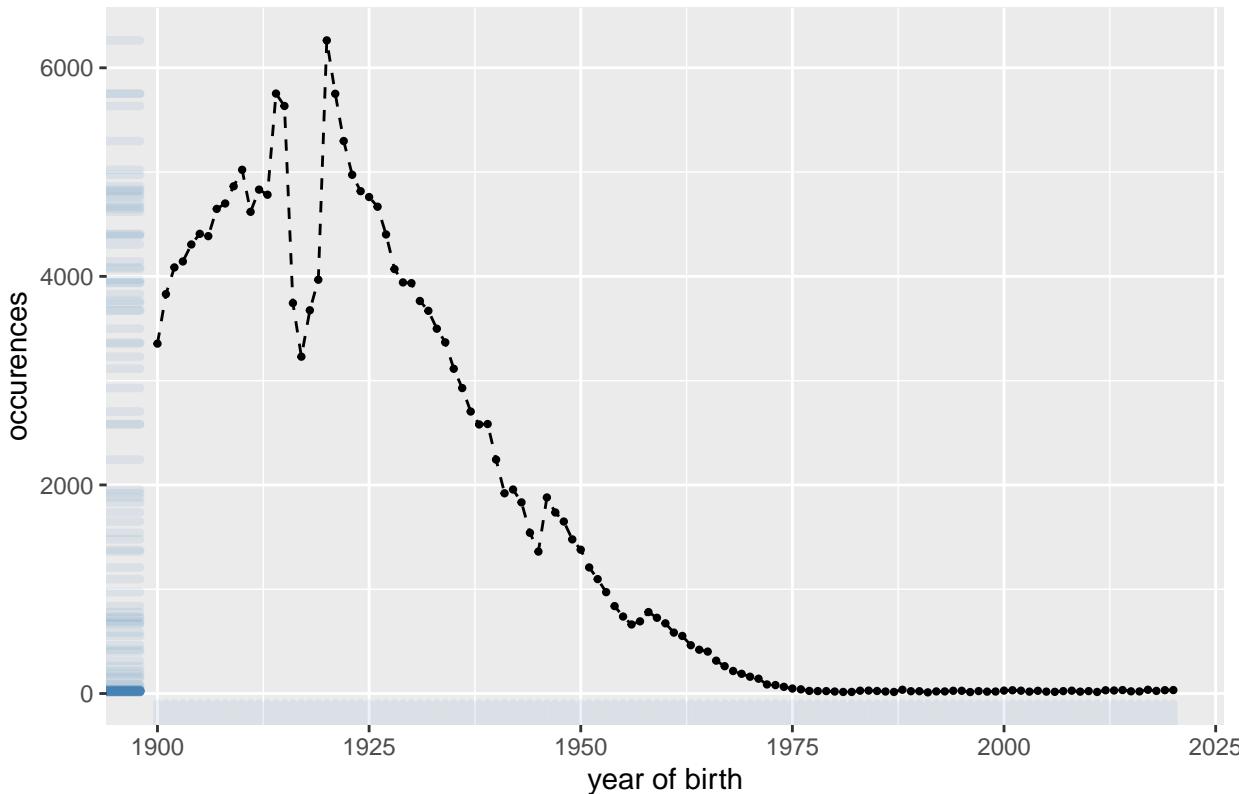
```
FirstNames
```

```
## # A tibble: 3,727,553 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
## # ... with 3,727,543 more rows
```

3. Analysing Firstnames' frequencies

We first choose an example Firstname (ALBERT), and we analyse its frequency.

occurrences of the name ALBERT = f (year of birth)

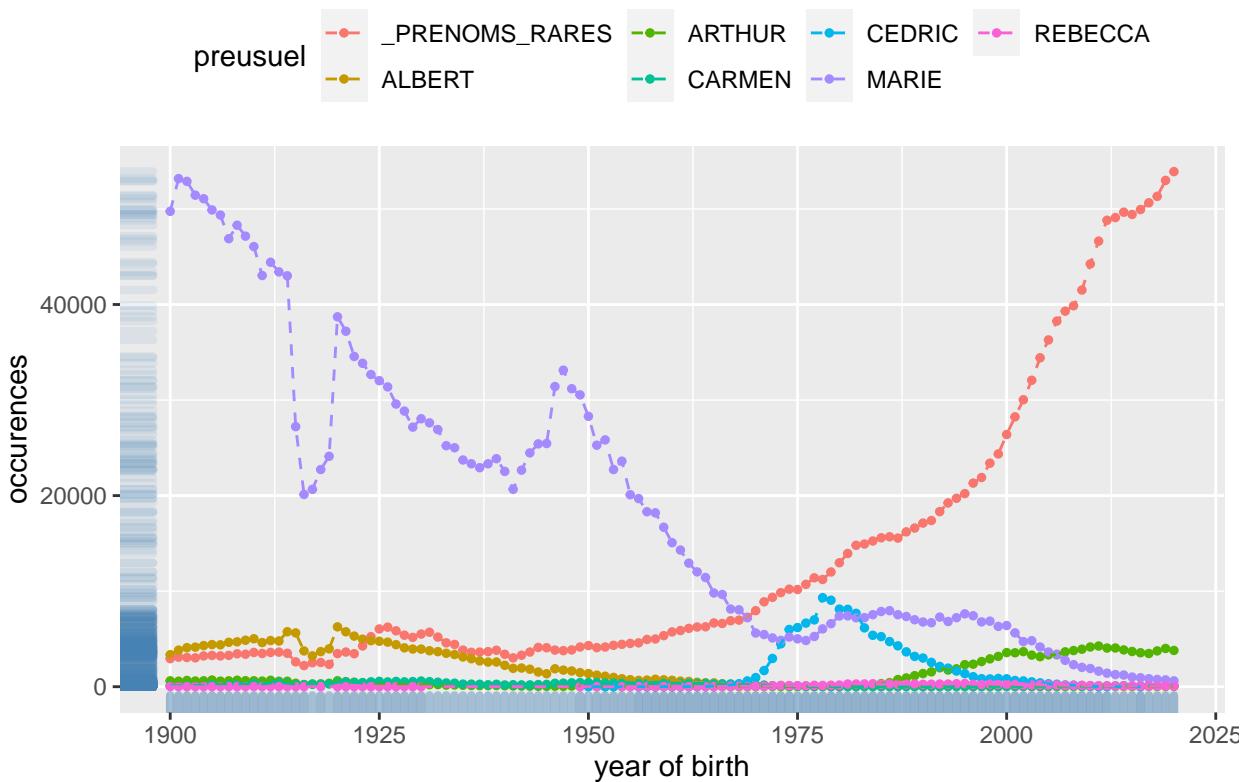


Comment :

After 1975, the number of new births with the Firstname ALBERT became very low. But the decline started since 1921.

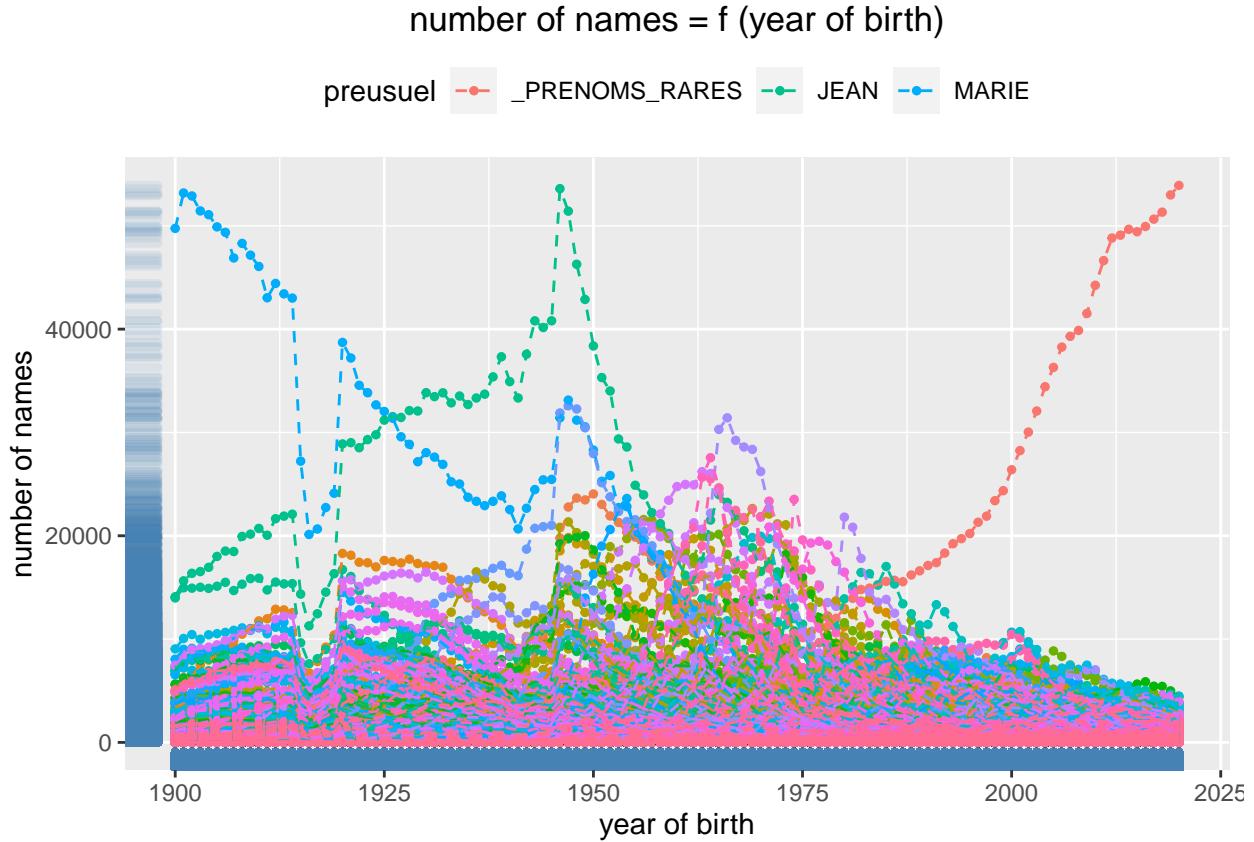
Now, let's compare several Firstnames' frequencies.

occurrences of sample names = f (year of birth)



Comment : It appears that the set of rare Firstnames is getting diversified, thus the increasing curve of the total number of this set. We note here that there are many samples for which the year of birth is not known (labeled XXXX), but we've ignored them in this analysis.

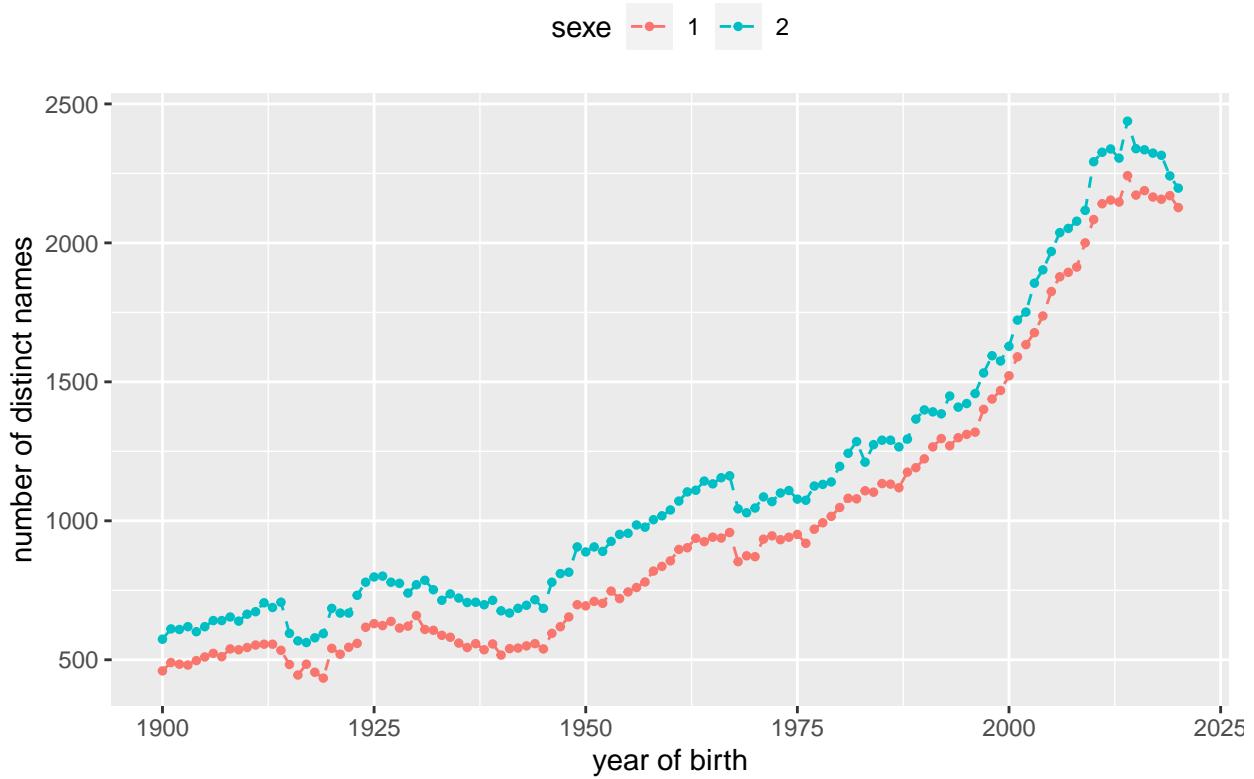
We can also plot all the names in the data set, in order to have an exhaustive overview over the distribution (event though not readable). We'll identify the three most used Firstnames since 1900 :



Comment : JEAN and MARIE are the most used Firstnames since 1900. The sum of all sets of rare names end up to form the biggest set of birth names. What we conclude from this graph is that the global set of names is getting diversified.

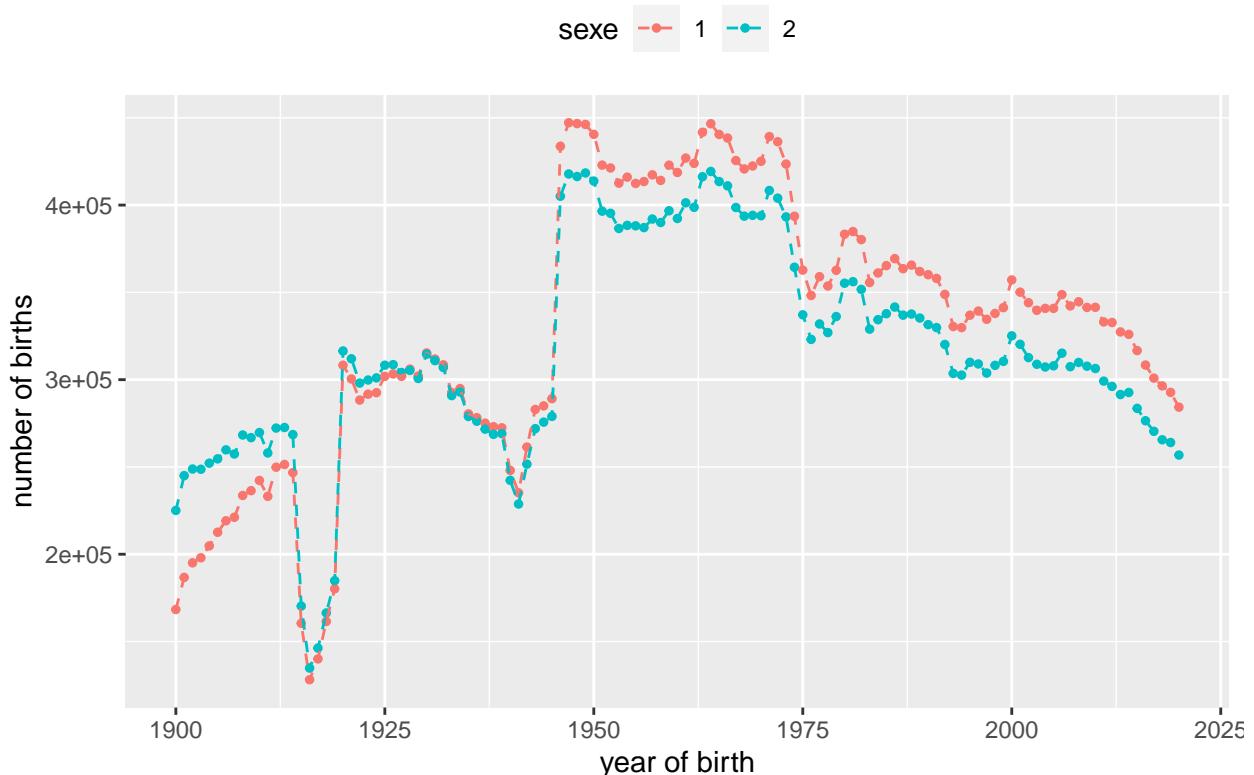
This might lead to think about the evolution of the diversity of birth names (how many distinct names). As shown on the following graph :

diversity of names = f (year of birth)



Analyzing the Birth Rate :

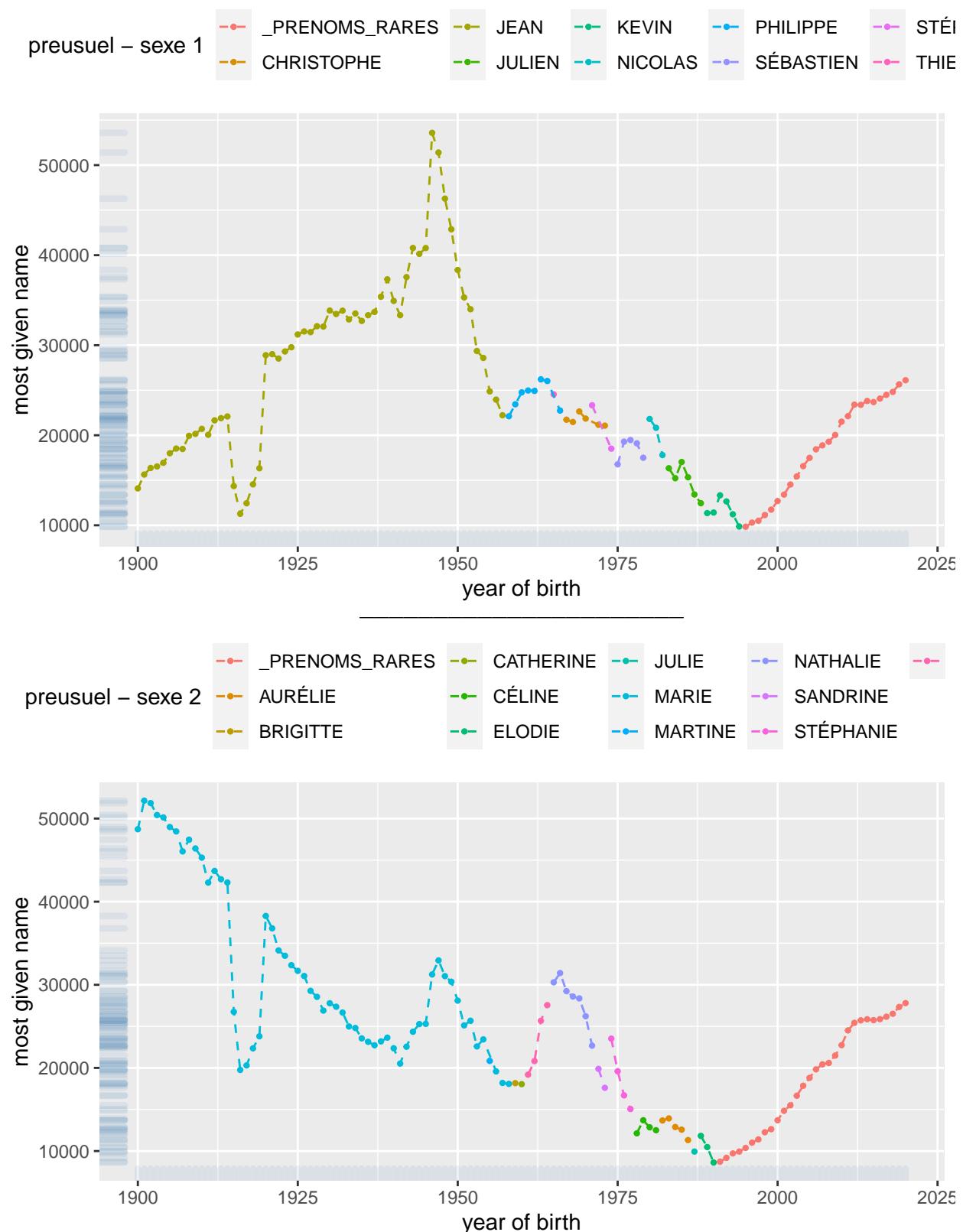
Birth Rate = f (year of birth)



Comment : We observe a sudden drop in the birth rate, both around 1915 and 1940. Which might be a clear manifestation of respectively the world war 1 (1914-1918) and the world war 2 (1939-1945).

4. Compute the most given firstname per year

most given name = f (year of birth)

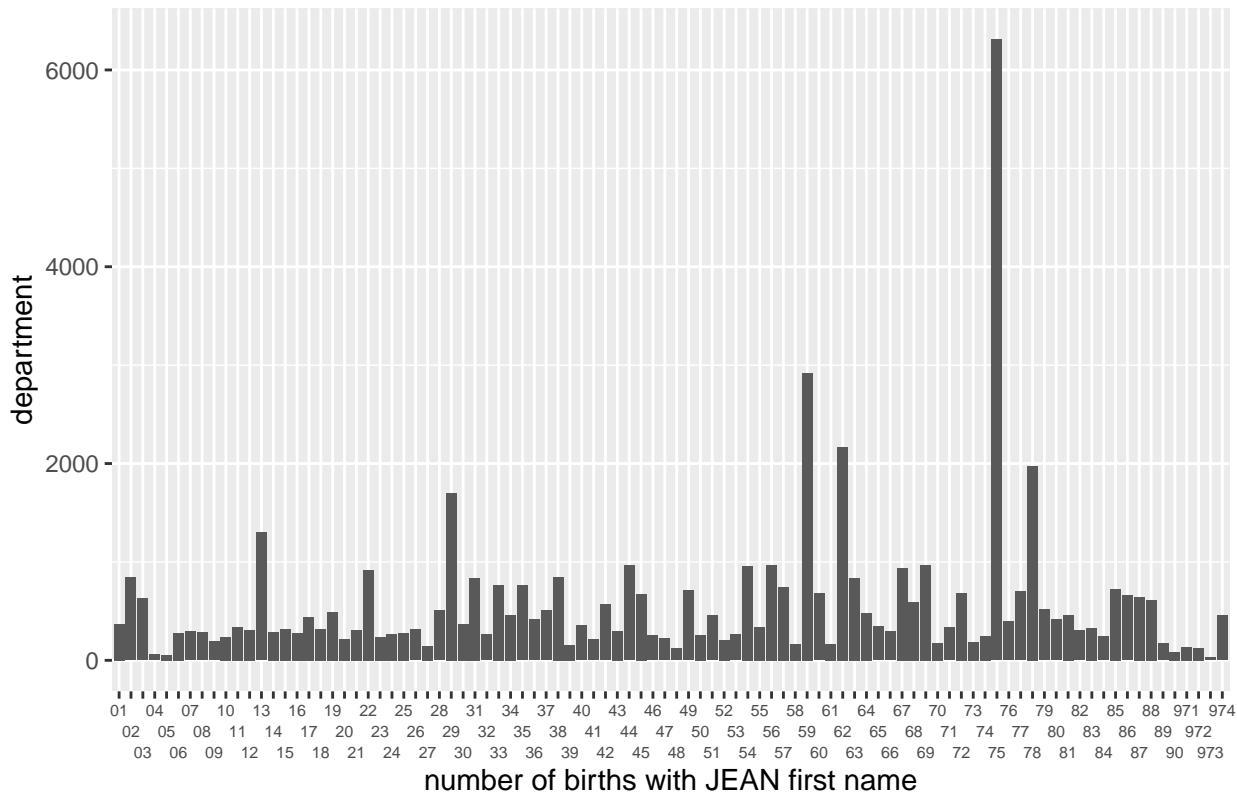


Synthesis : TODO

5. Any correlation between the first name and the localization (department) ?

Let's take for instance the most used name on 1946 (JEAN) ; 53584 (the peak), and see its distribution over the departments.

number of births with JEAN first name = f (department)



Comment : TODO