

In [197]:

```
import pandas as pd
import os
import numpy as np
import matplotlib as plt
%matplotlib inline
```

Скачал данные с <https://github.com/wesm/pydata-book/tree/2nd-edition/datasets/babynames>
(<https://github.com/wesm/pydata-book/tree/2nd-edition/datasets/babynames>)

План такой:

- 1) Получить названия директорий откуда пандас будет считывать файлы,
- 2) Путем танцев с бубнами, костылями и пандасами создать датафрейм,
- 3) Вытащить из названий файлов Год анализа,
- 4) Понять, что это даже не начало задания на визуализацию

In [170]:

```
yourpath = r'C:\Users\a.dikov\OTUS\5\datasets\babynames'

directory = []

for root, dirs, files in os.walk(yourpath, topdown=False):
    for name in files:
        directory.append(os.path.join(root, name))
        print(os.path.join(root, name))

#1 часть готова
```

[illegible]

[illegible]

C:\Users\a.dikov\OTUS\5\datasets\babynames\yob2001.txt
C:\Users\a.dikov\OTUS\5\datasets\babynames\yob2002.txt
C:\Users\a.dikov\OTUS\5\datasets\babynames\yob2003.txt
C:\Users\a.dikov\OTUS\5\datasets\babynames\yob2004.txt
C:\Users\a.dikov\OTUS\5\datasets\babynames\yob2005.txt
C:\Users\a.dikov\OTUS\5\datasets\babynames\yob2006.txt
C:\Users\a.dikov\OTUS\5\datasets\babynames\yob2007.txt
C:\Users\a.dikov\OTUS\5\datasets\babynames\yob2008.txt
C:\Users\a.dikov\OTUS\5\datasets\babynames\yob2009.txt
C:\Users\a.dikov\OTUS\5\datasets\babynames\yob2010.txt

In [172]:

```
directory = directory[1:]  
directory  
#обрезаем для удобства
```

Out[172]:

[illegible]

[illegible]


```
'C:\\Users\\a.dikov\\OTUS\\5\\datasets\\babynames\\yob2000.txt',  
'C:\\Users\\a.dikov\\OTUS\\5\\datasets\\babynames\\yob2001.txt',  
'C:\\Users\\a.dikov\\OTUS\\5\\datasets\\babynames\\yob2002.txt',  
'C:\\Users\\a.dikov\\OTUS\\5\\datasets\\babynames\\yob2003.txt',  
'C:\\Users\\a.dikov\\OTUS\\5\\datasets\\babynames\\yob2004.txt',  
'C:\\Users\\a.dikov\\OTUS\\5\\datasets\\babynames\\yob2005.txt',  
'C:\\Users\\a.dikov\\OTUS\\5\\datasets\\babynames\\yob2006.txt',  
'C:\\Users\\a.dikov\\OTUS\\5\\datasets\\babynames\\yob2007.txt',  
'C:\\Users\\a.dikov\\OTUS\\5\\datasets\\babynames\\yob2008.txt',  
'C:\\Users\\a.dikov\\OTUS\\5\\datasets\\babynames\\yob2009.txt',  
'C:\\Users\\a.dikov\\OTUS\\5\\datasets\\babynames\\yob2010.txt']
```

In [210]:

```
dfc = pd.DataFrame(columns=['Name', 'Sex', 'Amount', 'Year'])  
  
for i in directory:  
    df = pd.read_csv(i, names=['Name', 'Sex', 'Amount'], delimiter=',' )  
    df['Year'] = i  
    df['Year'] = df['Year'].map(lambda x: x.lstrip(r'C:\\Users\\a.dikov\\OTUS\\5\\datasets\\b  
abynames\\yob').rstrip('.txt'))  
    dfc = dfc.append(df)  
  
# 2 и 3 часть готовы
```

In [211]:

```
dfc
```

Out[211]:

	Name	Sex	Amount	Year
0	Mary	F	7065	1880
1	Anna	F	2604	1880
2	Emma	F	2003	1880
3	Elizabeth	F	1939	1880
4	Minnie	F	1746	1880
5	Margaret	F	1578	1880
6	Ida	F	1472	1880
7	Alice	F	1414	1880
8	Bertha	F	1320	1880
9	Sarah	F	1288	1880
10	Annie	F	1258	1880
11	Clara	F	1226	1880
12	Ella	F	1156	1880
13	Florence	F	1063	1880
14	Cora	F	1045	1880
15	Martha	F	1040	1880
16	Laura	F	1012	1880
17	Nellie	F	995	1880
18	Grace	F	982	1880
19	Carrie	F	949	1880
20	Maude	F	858	1880
21	Mabel	F	808	1880
22	Bessie	F	794	1880
23	Jennie	F	793	1880
24	Gertrude	F	787	1880
25	Julia	F	783	1880
26	Hattie	F	769	1880
27	Edith	F	768	1880
28	Mattie	F	704	1880
29	Rose	F	700	1880
...
33808	Zaviyon	M	5	2010
33809	Zaybrien	M	5	2010
33810	Zayshawn	M	5	2010
33811	Zayyan	M	5	2010
33812	Zeal	M	5	2010
33813	Zealan	M	5	2010

	Name	Sex	Amount	Year
33814	Zecharia	M	5	2010
33815	Zeferino	M	5	2010
33816	Zekariah	M	5	2010
33817	Zeki	M	5	2010
33818	Zeriah	M	5	2010
33819	Zeshan	M	5	2010
33820	Zhyier	M	5	2010
33821	Zildjian	M	5	2010
33822	Zinn	M	5	2010
33823	Zishan	M	5	2010
33824	Ziven	M	5	2010
33825	Zmari	M	5	2010
33826	Zoren	M	5	2010
33827	Zuhaib	M	5	2010
33828	Zyeire	M	5	2010
33829	Zygmunt	M	5	2010
33830	Zykerion	M	5	2010
33831	Zylar	M	5	2010
33832	Zylin	M	5	2010
33833	Zymaire	M	5	2010
33834	Zyonne	M	5	2010
33835	Zyquarius	M	5	2010
33836	Zyran	M	5	2010
33837	Zzyzx	M	5	2010

1690784 rows × 4 columns

ЗАДАНИЕ 1: Сгруппируйте данные по полу и году и визуализируйте общую динамику рождаемости обоих полов

In []:

```
dfc['tech_col'] = '1' # добавляем 1 чтобы по ней суммировать, скорее всего неправильно так делать, но что поделать `\_(\ツ)_/`
```