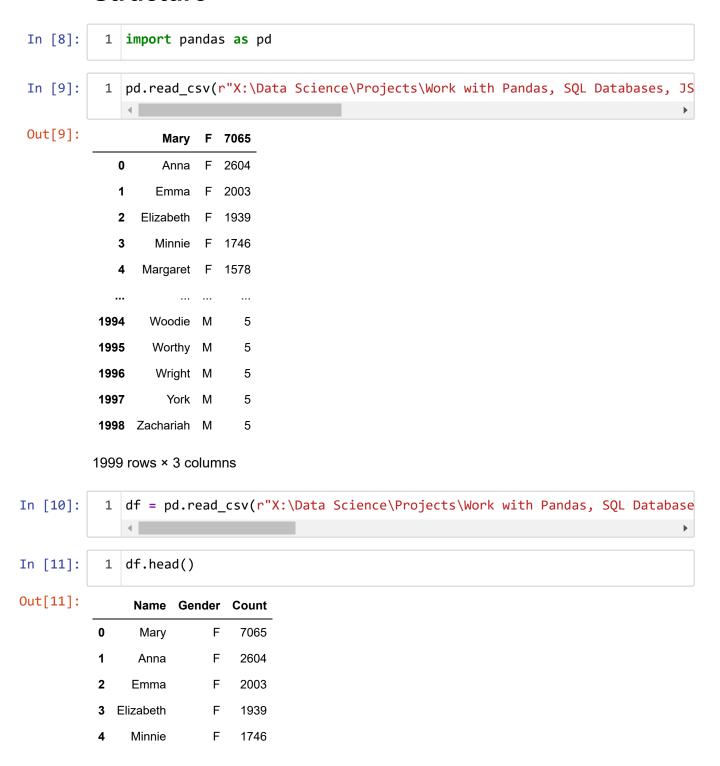
## Importing one file and unserstanding Data Structure



```
In [12]:
           1 df.info()
         <class 'pandas.core.frame.DataFrame'>
         RangeIndex: 2000 entries, 0 to 1999
         Data columns (total 3 columns):
              Column Non-Null Count Dtype
                      -----
          0
                      2000 non-null
                                      object
              Name
                                      object
          1
              Gender 2000 non-null
          2
              Count
                      2000 non-null
                                      int64
         dtypes: int64(1), object(2)
         memory usage: 47.0+ KB
```

## **Importing & Merging many files**

```
df_1880 = pd.read_csv('yob1880.txt', header = None, names = ['Names', "Gen
In [13]:
            2 df 1880.head()
Out[13]:
                Names Gender Count
           0
                                7065
                 Mary
           1
                 Anna
                                2604
           2
                Emma
                                2003
                                 1939
           3
              Elizabeth
                            F
                                1746
                Minnie
In [14]:
               df_1881 = pd.read_csv('yob1881.txt', header = None, names = ['Names', "Gen
               df_1881
Out[14]:
                   Names
                          Gender
                                 Count
              0
                               F
                                   6919
                    Mary
              1
                               F
                    Anna
                                    2698
              2
                   Emma
                               F
                                   2034
                               F
              3
                 Elizabeth
                                    1852
                 Margaret
                               F
                                    1658
           1930
                   Wiliam
                                      5
           1931
                   Wilton
                               Μ
                                      5
                                      5
           1932
                    Wing
           1933
                    Wood
                                      5
           1934
                   Wright
                                      5
                               Μ
```

1935 rows × 3 columns

In [17]: 1 pd.concat(objs = [df\_1880, df\_1881], axis = 0)

Out[17]:

	Names	Gender	Count
0	Mary	F	7065
1	Anna	F	2604
2	Emma	F	2003
3	Elizabeth	F	1939
4	Minnie	F	1746
1930	Wiliam	М	5
1931	Wilton	М	5
1932	Wing	М	5
1933	Wood	М	5
1934	Wright	М	5

3935 rows × 3 columns

Out[21]:

	Year	Names	Gender	Count
0	1880	Mary	F	7065
1	1880	Anna	F	2604
2	1880	Emma	F	2003
3	1880	Elizabeth	F	1939
4	1880	Minnie	F	1746
3930	1881	Wiliam	М	5
3931	1881	Wilton	М	5
3932	1881	Wing	М	5
3933	1881	Wood	М	5
3934	1881	Wright	М	5

3935 rows × 4 columns

```
pd.read_csv("yob{}.txt".format(1880), header = None, names = ["Name", "Gen
In [25]:
Out[25]:
                    Name Gender Count
                                    7065
              0
                     Mary
                               F
                                   2604
              1
                               F
                     Anna
              2
                    Emma
                                    2003
              3
                 Elizabeth
                               F
                                    1939
              4
                   Minnie
                                    1746
           1995
                                      5
                   Woodie
                               Μ
           1996
                   Worthy
                                      5
           1997
                    Wright
                               Μ
                                      5
           1998
                     York
                                      5
                               Μ
           1999 Zachariah
                               Μ
                                      5
          2000 rows × 3 columns
In [28]:
               years = list(range(1880,2019))
In [29]:
            1 years
Out[29]: [1880,
           1881,
           1882,
           1883,
           1884,
           1885,
           1886,
           1887,
           1888,
           1889,
           1890,
           1891,
           1892,
           1893,
           1894,
           1895,
           1896,
           1897,
           1898,
            1000
```

```
In [40]:
           1
              dataframes = []
           2
              for year in years:
           3
                   data = pd.read_csv("yob{}.txt".format(year), header = None, names = ["
                   dataframes.append(data)
In [41]:
           1 dataframes
           2294
                   winston
                                         5
           2295
                      York
                                 Μ
                                         5
           2296
                Zachariah
                                 Μ
           [2297 rows x \ 3 \ columns],
                      Name Gender
                                    Count
           0
                      Mary
                                 F
                                     9128
           1
                                     3994
                      Anna
           2
                      Emma
                                 F
                                     2728
           3
                 Elizabeth
                                     2582
           4
                  Margaret
                                 F
                                     2204
           2289
                    Wallie
                                 Μ
                                         5
                   Willian
           2290
           2291
                      Wirt
                                         5
           2292
                       Yee
                                 Μ
                                         5
           2293
                       Zeb
           [2294 rows x \ 3 \ columns],
                      Name Gender Count
In [42]:
           1 len(dataframes)
Out[42]: 139
In [44]:
           1 | df = pd.concat(dataframes, axis = 0, keys = years, names = ["Year"]).dropl
```

```
In [45]:
           1 df
```

Out	[45]	١
ouc	T 7	

	Year	Name	Gender	Count
0	1880	Mary	F	7065
1	1880	Anna	F	2604
2	1880	Emma	F	2003
3	1880	Elizabeth	F	1939
4	1880	Minnie	F	1746
1957041	2018	Zylas	М	5
1957042	2018	Zyran	М	5
1957043	2018	Zyrie	М	5
1957044	2018	Zyron	М	5
1957045	2018	Zzyzx	М	5

1957046 rows × 4 columns

```
In [47]:
           1 df.info()
```

<class 'pandas.core.frame.DataFrame'>

RangeIndex: 1957046 entries, 0 to 1957045

Data columns (total 4 columns):

Column Dtype

-----

0 int64 Year 1 Name object

2

Gender object

int64 3 Count

dtypes: int64(2), object(2) memory usage: 59.7+ MB

```
In [48]:
           1 df.to_csv("us_baby_names.csv", index = False)
```

In [49]: 1 pd.read\_csv("us\_baby\_names.csv")

Out[49]:

	Year	Name	Gender	Count
0	1880	Mary	F	7065
1	1880	Anna	F	2604
2	1880	Emma	F	2003
3	1880	Elizabeth	F	1939
4	1880	Minnie	F	1746
1957041	2018	Zylas	М	5
1957042	2018	Zyran	М	5
1957043	2018	Zyrie	М	5
1957044	2018	Zyron	М	5
1957045	2018	Zzyzx	М	5

1957046 rows × 4 columns

In [ ]:

1