# **Data Cleaning Project**



**Robson Silva** 

## Info dos dados:

Case Number	Location	Fatal (Y/N)	href	
Date	Activity	Time	Case Number.1	
Year	Name	Species	Case Number.2	
Туре	Sex	Investigator or Source	Original Order	
Country	Age	pdf	Unnamed: 22	
Area	Injury	Href formula	Unnamed: 23	

### Iniciando com o básico:

```
attacks.drop duplicates(inplace=True)
In [5]:
In [6]:
            attacks.reset index(inplace=True)
In [8]:
            attacks.drop(columns='index', inplace=True)
In [7]:
            attacks.info()
        <class 'pandas.core.frame.DataFrame'>
        RangeIndex: 6312 entries, 0 to 6311
        Data columns (total 25 columns):
        index
                                  6312 non-null int64
        Case Number
                                  6310 non-null object
                                  6302 non-null object
        Date
                                  6300 non-null float64
        Year
                                  6298 non-null object
        Type
                                  6252 non-null object
        Country
                                  5847 non-null object
        Area
                                  5762 non-null object
        Location
                                  5758 non-null object
        Activity
                                  6092 non-null object
        Name
        Sex
                                  5737 non-null object
        Age
                                  3471 non-null object
                                  6274 non-null object
        Injury
```



### Colunas deletáveis:

```
href formula
                                 6301 non-null object
       href
                                 6302 non-null object
       Case Number.1
                                 6302 non-null object
       Case Number.2
                                 6302 non-null object
                                 6309 non-null float64
       original order
       Unnamed: 22
                                 1 non-null object
       Unnamed: 23
                                 2 non-null object
       dtypes: float64(2), int64(1), object(22)
       memory usage: 1.2+ MB
           attacks.drop(columns=['Unnamed: 22', 'Unnamed: 23'], axis=1, inplace=True)
n [9]:
```

# Ampliando a visualização.

```
In [10]: 1 pd.options.display.max_rows
Out[10]: 60
In [159]: 1 pd.options.display.max_rows = 200
```

## Organizando o trabalho.

Procurei por colunas mais organizadas e decidi voltar no final do projeto.

- Country
- Area
- Location

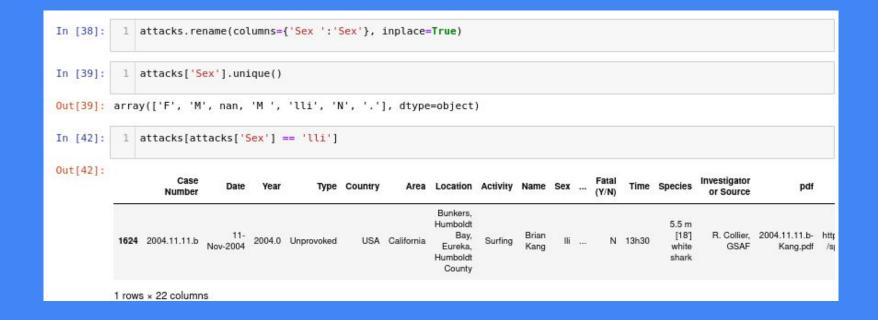
# Substituição manual de alguns valores da Coluna Fatal (Y/N).

```
In [24]: 1 attacks['Fatal (Y/N)'].replace({'M':'N', 'N':'N', 'N':'N', 'y':'Y','2017':np.nan}, inplace=True)
```

## Coluna 'Activity'

```
1 attacks['Activity'].describe()
In [26]:
Out[26]: count
                      5758
         unique
                      1532
                   Surfing
         top
         freq
                       971
         Name: Activity, dtype: object
In [31]:
             attacks['Activity'].unique()
Out[31]: array(['Paddling', 'Standing', 'Surfing', ...,
                'Crew swimming alongside their anchored ship',
                '4 men were bathing', 'Wreck of large double sailing canoe'],
               dtype=object)
In [30]:
          1 attacks['Activity'].isnull().sum()
Out[30]: 554
```

### Coluna 'Sex'



attacks['Sex'].replace({'M ':'M', 'N':'M',}, inplace=True)
Mais de 500 linhas nulas.

### 'Href' e 'href formula' colunas.

```
href = attacks[attacks['href formula'] != attacks['href']]
In [49]:
In [50]:
               href.reset index()
Out[50]:
                                                                                                                                  Investigator
               index
                                                              Country
                                                                                 Location Activity
                                                                                                                                    or Source
                        Number
                                                                                                                            White
                                                                       New South
                                                                                    Martin
                                                                                                                                     B. Myatt,
                                                                                                                                                    20
                  50 2018.01.13
                                         2018.0 Unprovoked AUSTRALIA
                                                                                                                            shark,
                                                                                                                                       GSAF
                                                                                                                            3.5 m
```

70 rows x 23 columns

### Eliminando 'href'

```
for i in range(len(href)):
    print(i)
    print(href['href'].iloc[i])

thtp://sharkattackfile.net/spreadsheets/pdf_directory/http://sharkattackfile.net/spreadsheets/pdf_directory/2018.0

1.13-Stewart.pdf
http://sharkattackfile.net/spreadsheets/pdf_directory/2018.01.13-Stewart.pdf

1
http://sharkattackfile.net/spreadsheets/pdf_directory/2018.01.13-Stewart.pdf

2
http://sharkattackfile.net/spreadsheets/pdf_directory/http://sharkattackfile.net/spreadsheets/pdf_directory/2017.08.27-Brundler.pdf
http://sharkattackfile.net/spreadsheets/pdf_directory/2017.08.27-Brundler.pdf

2
http://sharkattackfile.net/spreadsheets/pdf_directory/http://sharkattackfile.net/spreadsheets/pdf_directory/2017.06.05-FrenchPolynesia.pdf
http://sharkattackfile.net/spreadsheets/pdf_directory/2017.06.05-FrenchPolynesia.pdf

3
http://sharkattackfile.net/spreadsheets/pdf_directory/2017.06.05-FrenchPolynesia.pdf
```

```
In [ ]: 1 attacks.drop(columns='href', axis=1, inplace=True)
In [56]: 1 attacks.rename(columns={'href formula':'Link'}, inplace=True)
In [57]: 1 attacks.dropna(thresh=9, axis=0, inplace=True)
In [58]: 1 attacks.reset_index(inplace=True)
```

## Trabalhando com 'Type'.

Type possui apenas valores 4 nulos.

### Colunas: Case Numbers

```
In [76]:
                attacks[attacks['Case Number.1'] != attacks['Case Number']]
Out[76]:
                 index
                                                                                    Location
                                                                                                Activity
                                                                                                                                     Time Species
                                                                 Country
                                                                                                          Name ...
                                                                             Area
                            Number
                                                                          Eastern
                                                                                                                     Lacerations
                                                                 SOUTH
                                                                                   St. Francis
                                                                                                                                             White
                         2018.04.03
                                            2018.0 Unprovoked
                                                                                                 Surfing
                                                                                                                                  N 15h00
                                                                                                                     to left knee
                                                                                                                                             shark Tra
                                                                                                                     & lower leg
                                                            In [78]:
                                                                            attacks['Case Number'].iloc[5488] = '1905.09.06'
24 rows x 22 columns
                                                            In [79]:
                                                                            attacks[attacks['Case Number'].isnull()]
                                                            Out[79]:
                                                                                        Date Year Type Country Area Location Activity Name ... Injury
                                                                       0 rows x 22 columns
```

```
In [80]: 1 attacks.drop(columns=['Case Number.1','Case Number.2'], inplace=True)
```

### Coluna 'Year'

[94]:

[94]:

```
In [95]: 1 attacks['Year'] = attacks['Year'].apply(int)
```

```
        index
        Case Number
        Date
        Year
        Type
        Country
        Area

        6177
        6177
        0000.0214
        Ca. 214 B.C.
        0.0
        Unprovoked
        NaN
        Ionian Sea

        6178
        6178
        0000.0336
        Ca. 336.B.C.
        0.0
        Unprovoked
        GREECE
        Piraeus
```

6297	6.0	ND.0005	Before 1903	0	Ung
6298	5.0	ND.0004	Before 1903	0	Ung

attacks[attacks['Year'] == 0]

```
In [87]:
               attacks[attacks['Year'].isnull()]
Out[87]:
                index
                                     Date Year
                                                    Type
                                                           Country
                                                                         Area
                          Number
                                  Reported
                  187 2017.01.08.R
                                          NaN
                                                   Invalid AUSTRALIA Queensland
                                 Jan-2017
                                  Reported
                 6079 1836.08.19.R
                                          NaN Unprovoked ENGLAND Cumberland
                                 Aug-1836
               attacks['Year'].iloc[187] = 2017
In [90]:
In [93]:
              attacks['Year'].iloc[6079] = 1836
```

### Coluna Year, mais limpa, sem nenhum elemento nulo.

```
In [459]:
              attacks['Year'].unique()
Out[459]: array([2018, 2017, 2016, 2015, 2014, 2013, 2012, 2011, 2010, 2009, 2008,
                 2007, 2006, 2005, 2004, 2003, 2002, 2001, 2000, 1999, 1998, 1997,
                 1996, 1995, 1984, 1994, 1993, 1992, 1991, 1990, 1989, 1969, 1988,
                 1987, 1986, 1985, 1983, 1982, 1981, 1980, 1979, 1978, 1977, 1976,
                 1975, 1974, 1973, 1972, 1971, 1970, 1968, 1967, 1966, 1965, 1964,
                 1963, 1962, 1961, 1960, 1959, 1958, 1957, 1956, 1955, 1954, 1953,
                 1952, 1951, 1950, 1949, 1948, 1848, 1947, 1946, 1945, 1944, 1943,
                 1942, 1941, 1940, 1939, 1938, 1937, 1936, 1935, 1934, 1933, 1932,
                 1931, 1930, 1929, 1928, 1927, 1926, 1925, 1924, 1923, 1922, 1921,
                 1920, 1919, 1918, 1917, 1916, 1915, 1914, 1913, 1912, 1911, 1910,
                 1909, 1908, 1907, 1906, 1905, 1904, 1903, 1902, 1901, 1900, 1899,
                 1898, 1897, 1896, 1895, 1894, 1893, 1892, 1891, 1890, 1889, 1888,
                 1887, 1886, 1885, 1884, 1883, 1882, 1881, 1880, 1879, 1878, 1877,
                 1876, 1875, 1874, 1873, 1872, 1871, 1870, 1869, 1868, 1867, 1866,
                 1865, 1864, 1863, 1862, 1861, 1860, 1859, 1858, 1857, 1856, 1855,
                 1853, 1852, 1851, 1850, 1849, 1847, 1846, 1845, 1844, 1842, 1841,
                 1840, 1839, 1837, 1836, 1835, 1834, 1832, 1831, 1830, 1829, 1828,
                 1827, 1826, 1825, 1823, 1822, 1819, 1818, 1817, 1816, 1815, 1812,
                 1811, 1810, 1808, 1807, 1805, 1804, 1803, 1802, 1801, 1800, 1797,
                 1792, 1791, 1788, 1787, 1786, 1785, 1784, 1783, 1780, 1779, 1776,
                 1771, 1767, 1764, 1758, 1753, 1751, 1749, 1755, 1748, 1742, 1738,
                 1733, 1723, 1721, 1703, 1700, 1642, 1638, 1637, 1617, 1595, 1580,
                 1555, 1554, 1543, 500, 77,
                                                  5, -214, -336, -493, -725,
```

'0' São 37 valores que serão feitos manualmente.

## Limpando 'Year'

```
In [371]:
              def rob string(x):
                  if re.findall(r'.*B\.C*.', x):
                      return int(''.join(re.findall(f'\d+', x)))*-1
                  elif re.findall(r'.*[bB]efore.*', x):
                      return ''.join(re.findall(f'\d+', x))
                  else:
                      return 0
In [368]:
              texto = 'Ca, 214 B.C.'
In [369]:
              rob string(texto)
Out[369]: -214
In [370]:
              tex = 'Before Mar-1956'
In [362]:
              rob string(tex)
Out[362]: 1956
```

5 attacks.loc[ attacks['Year']==0, 'Year' ] = 'Teste'

## Limpando 'Name'

```
In [472]: 1 len(attacks['Name'].apply(lambda x:re.findall(r'.*male.*',x)).sum())
Out[472]: 809
```

```
In [494]: 1 attacks['Sex'].isnull().sum()
Out[494]: 559
```

```
In [496]: 1 attacks.loc[(attacks['Name']==r'.*^male.*') & (attacks['Sex'].isnull())]

Out[496]: 

Index Case Number Date Year Type Country Area Location Activity Name ... Age Injury Fatal (Y/N) Time Spent Orows × 21 columns
```

# Conclusões e melhorias