## import pandas

```python
import pandas as pd
```

## Read CSV Files

CSV files (comma separated files) are used to store big data sets.

```python
df = pd.read_csv('data.csv')
# By default, when you print a DataFrame, you will only get the first 5 rows, and the last 5 rows
df
```

|     | Duration | Pulse | Maxpulse | Calories |
|-----|----------|-------|----------|----------|
| 0   | 60       | 110   | 130      | 409.1    |
| 1   | 60       | 117   | 145      | 479.0    |
| 2   | 60       | 103   | 135      | 340.0    |
| 3   | 45       | 109   | 175      | 282.4    |
| 4   | 45       | 117   | 148      | 406.0    |
| ... | ...      | ...   | ...      | ...      |
| 164 | 60       | 105   | 140      | 290.8    |
| 165 | 60       | 110   | 145      | 300.0    |
| 166 | 60       | 115   | 145      | 310.2    |
| 167 | 75       | 120   | 150      | 320.4    |
| 168 | 75       | 125   | 150      | 330.4    |

169 rows × 4 columns

```
# to print the entire DataFrame
print(df.to_string())
```

```
    Duration  Pulse  Maxpulse  Calories
0         60    110       130     409.1
1         60    117       145     479.0
2         60    103       135     340.0
3         45    109       175     282.4
4         45    117       148     406.0
5         60    102       127     300.0
6         60    110       136     374.0
7         45    104       134     253.3
8         30    109       133     195.1
9         60     98       124     269.0
10        60    103       147     329.3
11        60    100       120     250.7
12        60    106       128     345.3
13        60    104       132     379.3
14        60     98       123     275.0
15        60     98       120     215.2
16        60    100       120     300.0
17        45     90       112       NaN
18        60    103       123     323.0
19        45     97       125     243.0
20        60    108       131     364.2
21        45    100       119     282.0
22        60    130       101     300.0
23        45    105       132     246.0
24        60    102       126     334.5
25        60    100       120     250.0
26        60     92       118     241.0
27        60    103       132       NaN
28        60    100       132     280.0
29        60    102       129     380.3
30        60     92       115     243.0
31        45     90       112     180.1
32        60    101       124     299.0
33        60     93       113     223.0
34        60    107       136     361.0
35        60    114       140     415.0
36        60    102       127     300.0
37        60    100       120     300.0
```

| 38 | 60 | 100 | 120 | 300.0 |
|---|---|---|---|---|
| 39 | 45 | 104 | 129 | 266.0 |
| 40 | 45 | 90 | 112 | 180.1 |
| 41 | 60 | 98 | 126 | 286.0 |
| 42 | 60 | 100 | 122 | 329.4 |
| 43 | 60 | 111 | 138 | 400.0 |
| 44 | 60 | 111 | 131 | 397.0 |
| 45 | 60 | 99 | 119 | 273.0 |
| 46 | 60 | 109 | 153 | 387.6 |
| 47 | 45 | 111 | 136 | 300.0 |
| 48 | 45 | 108 | 129 | 298.0 |
| 49 | 60 | 111 | 139 | 397.6 |
| 50 | 60 | 107 | 136 | 380.2 |
| 51 | 80 | 123 | 146 | 643.1 |
| 52 | 60 | 106 | 130 | 263.0 |
| 53 | 60 | 118 | 151 | 486.0 |
| 54 | 30 | 136 | 175 | 238.0 |
| 55 | 60 | 121 | 146 | 450.7 |
| 56 | 60 | 118 | 121 | 413.0 |
| 57 | 45 | 115 | 144 | 305.0 |
| 58 | 20 | 153 | 172 | 226.4 |
| 59 | 45 | 123 | 152 | 321.0 |
| 60 | 210 | 108 | 160 | 1376.0 |
| 61 | 160 | 110 | 137 | 1034.4 |
| 62 | 160 | 109 | 135 | 853.0 |
| 63 | 45 | 118 | 141 | 341.0 |
| 64 | 20 | 110 | 130 | 131.4 |
| 65 | 180 | 90 | 130 | 800.4 |
| 66 | 150 | 105 | 135 | 873.4 |
| 67 | 150 | 107 | 130 | 816.0 |
| 68 | 20 | 106 | 136 | 110.4 |
| 69 | 300 | 108 | 143 | 1500.2 |
| 70 | 150 | 97 | 129 | 1115.0 |
| 71 | 60 | 109 | 153 | 387.6 |
| 72 | 90 | 100 | 127 | 700.0 |
| 73 | 150 | 97 | 127 | 953.2 |
| 74 | 45 | 114 | 146 | 304.0 |
| 75 | 90 | 98 | 125 | 563.2 |
| 76 | 45 | 105 | 134 | 251.0 |
| 77 | 45 | 110 | 141 | 300.0 |
| 78 | 120 | 100 | 130 | 500.4 |
| 79 | 270 | 100 | 131 | 1729.0 |
| 80 | 30 | 159 | 182 | 319.2 |
| 81 | 45 | 149 | 169 | 344.0 |

| | | | | |
|---|---|---|---|---|
| 82 | 30 | 103 | 139 | 151.1 |
| 83 | 120 | 100 | 130 | 500.0 |
| 84 | 45 | 100 | 120 | 225.3 |
| 85 | 30 | 151 | 170 | 300.0 |
| 86 | 45 | 102 | 136 | 234.0 |
| 87 | 120 | 100 | 157 | 1000.1 |
| 88 | 45 | 129 | 103 | 242.0 |
| 89 | 20 | 83 | 107 | 50.3 |
| 90 | 180 | 101 | 127 | 600.1 |
| 91 | 45 | 107 | 137 | NaN |
| 92 | 30 | 90 | 107 | 105.3 |
| 93 | 15 | 80 | 100 | 50.5 |
| 94 | 20 | 150 | 171 | 127.4 |
| 95 | 20 | 151 | 168 | 229.4 |
| 96 | 30 | 95 | 128 | 128.2 |
| 97 | 25 | 152 | 168 | 244.2 |
| 98 | 30 | 109 | 131 | 188.2 |
| 99 | 90 | 93 | 124 | 604.1 |
| 100 | 20 | 95 | 112 | 77.7 |
| 101 | 90 | 90 | 110 | 500.0 |
| 102 | 90 | 90 | 100 | 500.0 |
| 103 | 90 | 90 | 100 | 500.4 |
| 104 | 30 | 92 | 108 | 92.7 |
| 105 | 30 | 93 | 128 | 124.0 |
| 106 | 180 | 90 | 120 | 800.3 |
| 107 | 30 | 90 | 120 | 86.2 |
| 108 | 90 | 90 | 120 | 500.3 |
| 109 | 210 | 137 | 184 | 1860.4 |
| 110 | 60 | 102 | 124 | 325.2 |
| 111 | 45 | 107 | 124 | 275.0 |
| 112 | 15 | 124 | 139 | 124.2 |
| 113 | 45 | 100 | 120 | 225.3 |
| 114 | 60 | 108 | 131 | 367.6 |
| 115 | 60 | 108 | 151 | 351.7 |
| 116 | 60 | 116 | 141 | 443.0 |
| 117 | 60 | 97 | 122 | 277.4 |
| 118 | 60 | 105 | 125 | NaN |
| 119 | 60 | 103 | 124 | 332.7 |
| 120 | 30 | 112 | 137 | 193.9 |
| 121 | 45 | 100 | 120 | 100.7 |
| 122 | 60 | 119 | 169 | 336.7 |
| 123 | 60 | 107 | 127 | 344.9 |
| 124 | 60 | 111 | 151 | 368.5 |
| 125 | 60 | 98 | 122 | 271.0 |

| | | | |
|---|---|---|---|---|
| 126 | 60 | 97 | 124 | 275.3 |
| 127 | 60 | 109 | 127 | 382.0 |
| 128 | 90 | 99 | 125 | 466.4 |
| 129 | 60 | 114 | 151 | 384.0 |
| 130 | 60 | 104 | 134 | 342.5 |
| 131 | 60 | 107 | 138 | 357.5 |
| 132 | 60 | 103 | 133 | 335.0 |
| 133 | 60 | 106 | 132 | 327.5 |
| 134 | 60 | 103 | 136 | 339.0 |
| 135 | 20 | 136 | 156 | 189.0 |
| 136 | 45 | 117 | 143 | 317.7 |
| 137 | 45 | 115 | 137 | 318.0 |
| 138 | 45 | 113 | 138 | 308.0 |
| 139 | 20 | 141 | 162 | 222.4 |
| 140 | 60 | 108 | 135 | 390.0 |
| 141 | 60 | 97 | 127 | NaN |
| 142 | 45 | 100 | 120 | 250.4 |
| 143 | 45 | 122 | 149 | 335.4 |
| 144 | 60 | 136 | 170 | 470.2 |
| 145 | 45 | 106 | 126 | 270.8 |
| 146 | 60 | 107 | 136 | 400.0 |
| 147 | 60 | 112 | 146 | 361.9 |
| 148 | 30 | 103 | 127 | 185.0 |
| 149 | 60 | 110 | 150 | 409.4 |
| 150 | 60 | 106 | 134 | 343.0 |
| 151 | 60 | 109 | 129 | 353.2 |
| 152 | 60 | 109 | 138 | 374.0 |
| 153 | 30 | 150 | 167 | 275.8 |
| 154 | 60 | 105 | 128 | 328.0 |
| 155 | 60 | 111 | 151 | 368.5 |
| 156 | 60 | 97 | 131 | 270.4 |
| 157 | 60 | 100 | 120 | 270.4 |
| 158 | 60 | 114 | 150 | 382.8 |
| 159 | 30 | 80 | 120 | 240.9 |
| 160 | 30 | 85 | 120 | 250.4 |
| 161 | 45 | 90 | 130 | 260.4 |
| 162 | 45 | 95 | 130 | 270.0 |
| 163 | 45 | 100 | 140 | 280.9 |
| 164 | 60 | 105 | 140 | 290.8 |
| 165 | 60 | 110 | 145 | 300.0 |
| 166 | 60 | 115 | 145 | 310.2 |
| 167 | 75 | 120 | 150 | 320.4 |
| 168 | 75 | 125 | 150 | 330.4 |

# Read JSON

- Big data sets are often stored, or extracted as JSON.

- JSON is plain text, but has the format of an object, and is well known in the world of programming, including Pandas.

- JSON objects have the same format as Python dictionaries.

In [4]:
```python
df_json = pd.read_json('data.json')
df_json
```

Out[4]:

|  | Duration | Pulse | Maxpulse | Calories |
|---|---|---|---|---|
| 0 | 60 | 110 | 130 | 409.1 |
| 1 | 60 | 117 | 145 | 479.0 |
| 2 | 60 | 103 | 135 | 340.0 |
| 3 | 45 | 109 | 175 | 282.4 |
| 4 | 45 | 117 | 148 | 406.0 |
| ... | ... | ... | ... | ... |
| 164 | 60 | 105 | 140 | 290.8 |
| 165 | 60 | 110 | 145 | 300.4 |
| 166 | 60 | 115 | 145 | 310.2 |
| 167 | 75 | 120 | 150 | 320.4 |
| 168 | 75 | 125 | 150 | 330.4 |

169 rows × 4 columns

# Viewing the Data

- `head()` : returns the headers and a specified number of rows, starting from the top.
- `tail()` : returns the headers and a specified number of rows, starting from the bottom.

In [5]:
```python
df.head() # by defaul, returns first rows
```

Out[5]:

| | Duration | Pulse | Maxpulse | Calories |
|---|---|---|---|---|
| 0 | 60 | 110 | 130 | 409.1 |
| 1 | 60 | 117 | 145 | 479.0 |
| 2 | 60 | 103 | 135 | 340.0 |
| 3 | 45 | 109 | 175 | 282.4 |
| 4 | 45 | 117 | 148 | 406.0 |

In [6]:
```python
df.head(10) # returns first 10 rows
```

Out[6]:

| | Duration | Pulse | Maxpulse | Calories |
|---|---|---|---|---|
| 0 | 60 | 110 | 130 | 409.1 |
| 1 | 60 | 117 | 145 | 479.0 |
| 2 | 60 | 103 | 135 | 340.0 |
| 3 | 45 | 109 | 175 | 282.4 |
| 4 | 45 | 117 | 148 | 406.0 |
| 5 | 60 | 102 | 127 | 300.0 |
| 6 | 60 | 110 | 136 | 374.0 |
| 7 | 45 | 104 | 134 | 253.3 |
| 8 | 30 | 109 | 133 | 195.1 |
| 9 | 60 | 98 | 124 | 269.0 |

```
In [7]:    df.tail() # by default, returns last 5 rows
```

Out[7]:

| | Duration | Pulse | Maxpulse | Calories |
|---|---|---|---|---|
| 164 | 60 | 105 | 140 | 290.8 |
| 165 | 60 | 110 | 145 | 300.0 |
| 166 | 60 | 115 | 145 | 310.2 |
| 167 | 75 | 120 | 150 | 320.4 |
| 168 | 75 | 125 | 150 | 330.4 |

```
In [8]:    df.tail(10) # returns last 10 rows
```

Out[8]:

| | Duration | Pulse | Maxpulse | Calories |
|---|---|---|---|---|
| 159 | 30 | 80 | 120 | 240.9 |
| 160 | 30 | 85 | 120 | 250.4 |
| 161 | 45 | 90 | 130 | 260.4 |
| 162 | 45 | 95 | 130 | 270.0 |
| 163 | 45 | 100 | 140 | 280.9 |
| 164 | 60 | 105 | 140 | 290.8 |
| 165 | 60 | 110 | 145 | 300.0 |
| 166 | 60 | 115 | 145 | 310.2 |
| 167 | 75 | 120 | 150 | 320.4 |
| 168 | 75 | 125 | 150 | 330.4 |

## Info About the Data

- `info()` : returns information about the data set.

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 169 entries, 0 to 168
Data columns (total 4 columns):
 #   Column    Non-Null Count  Dtype
---  ------    --------------  -----
 0   Duration  169 non-null    int64
 1   Pulse     169 non-null    int64
 2   Maxpulse  169 non-null    int64
 3   Calories  164 non-null    float64
dtypes: float64(1), int64(3)
memory usage: 5.4 KB
```

This information tells us that,

- there are 169 rows and 4 columns
- three columns have `int64` datatype and one column has `float64` data type.
- `no-null count` specifies how many data are no-null in each column. For example, `Calories` column contains 164 no-null data out of 169, that means, it has 5 rows with no value at all.

Note that, Empty values, or Null values, or Missing data should be handled carefully in data cleaning steps when analysing data.

## Cleaning Data

- Data cleaning means fixing bad data in data set.

- Bad data could be:

  - Empty cells
  - Data in wrong format
  - Wrong data
  - Duplicates

# Cleaning Empty Cells

- Empty cells can potentially give a wrong result when analyze data.

```
In [10]:    # Let's take a new data set which contains bad data
            ddf = pd.read_csv('dirtydata.csv')
            ddf
```

Out[10]:

|    | Duration | Date | Pulse | Maxpulse | Calories |
|----|----------|------|-------|----------|----------|
| 0  | 60  | '2020/12/01' | 110 | 130 | 409.1 |
| 1  | 60  | '2020/12/02' | 117 | 145 | 479.0 |
| 2  | 60  | '2020/12/03' | 103 | 135 | 340.0 |
| 3  | 45  | '2020/12/04' | 109 | 175 | 282.4 |
| 4  | 45  | '2020/12/05' | 117 | 148 | 406.0 |
| 5  | 60  | '2020/12/06' | 102 | 127 | 300.0 |
| 6  | 60  | '2020/12/07' | 110 | 136 | 374.0 |
| 7  | 450 | '2020/12/08' | 104 | 134 | 253.3 |
| 8  | 30  | '2020/12/09' | 109 | 133 | 195.1 |
| 9  | 60  | '2020/12/10' | 98  | 124 | 269.0 |
| 10 | 60  | '2020/12/11' | 103 | 147 | 329.3 |
| 11 | 60  | '2020/12/12' | 100 | 120 | 250.7 |
| 12 | 60  | '2020/12/12' | 100 | 120 | 250.7 |
| 13 | 60  | '2020/12/13' | 106 | 128 | 345.3 |
| 14 | 60  | '2020/12/14' | 104 | 132 | 379.3 |
| 15 | 60  | '2020/12/15' | 98  | 123 | 275.0 |
| 16 | 60  | '2020/12/16' | 98  | 120 | 215.2 |

| | | | | | |
|---|---|---|---|---|---|
| 17 | 60 | '2020/12/17' | 100 | 120 | 300.0 |
| 18 | 45 | '2020/12/18' | 90 | 112 | NaN |
| 19 | 60 | '2020/12/19' | 103 | 123 | 323.0 |
| 20 | 45 | '2020/12/20' | 97 | 125 | 243.0 |
| 21 | 60 | '2020/12/21' | 108 | 131 | 364.2 |
| 22 | 45 | NaN | 100 | 119 | 282.0 |
| 23 | 60 | '2020/12/23' | 130 | 101 | 300.0 |
| 24 | 45 | '2020/12/24' | 105 | 132 | 246.0 |
| 25 | 60 | '2020/12/25' | 102 | 126 | 334.5 |
| 26 | 60 | 20201226 | 100 | 120 | 250.0 |
| 27 | 60 | '2020/12/27' | 92 | 118 | 241.0 |
| 28 | 60 | '2020/12/28' | 103 | 132 | NaN |
| 29 | 60 | '2020/12/29' | 100 | 132 | 280.0 |
| 30 | 60 | '2020/12/30' | 102 | 129 | 380.3 |
| 31 | 60 | '2020/12/31' | 92 | 115 | 243.0 |

In [11]: 

```
ddf.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 32 entries, 0 to 31
Data columns (total 5 columns):
 #   Column    Non-Null Count  Dtype
---  ------    --------------  -----
 0   Duration  32 non-null     int64
 1   Date      31 non-null     object
 2   Pulse     32 non-null     int64
 3   Maxpulse  32 non-null     int64
 4   Calories  30 non-null     float64
dtypes: float64(1), int64(3), object(1)
memory usage: 1.4+ KB
```

## Remove Empty cells

- One way to deal with empty cells is to remove rows that contain empty cells (removing few rows are OK if the data set is big enough, because it does not have a big impact on the result).
- `dropna()` method can be used to remove empty rows in the data set.

In [12]:
```python
# By default, the dropna() method returns a new DataFrame, and will not change the original
new_df = ddf.dropna()
new_df
```

Out[12]:

| | Duration | Date | Pulse | Maxpulse | Calories |
|---|---|---|---|---|---|
| 0 | 60 | '2020/12/01' | 110 | 130 | 409.1 |
| 1 | 60 | '2020/12/02' | 117 | 145 | 479.0 |
| 2 | 60 | '2020/12/03' | 103 | 135 | 340.0 |
| 3 | 45 | '2020/12/04' | 109 | 175 | 282.4 |
| 4 | 45 | '2020/12/05' | 117 | 148 | 406.0 |
| 5 | 60 | '2020/12/06' | 102 | 127 | 300.0 |
| 6 | 60 | '2020/12/07' | 110 | 136 | 374.0 |
| 7 | 450 | '2020/12/08' | 104 | 134 | 253.3 |
| 8 | 30 | '2020/12/09' | 109 | 133 | 195.1 |
| 9 | 60 | '2020/12/10' | 98 | 124 | 269.0 |
| 10 | 60 | '2020/12/11' | 103 | 147 | 329.3 |
| 11 | 60 | '2020/12/12' | 100 | 120 | 250.7 |
| 12 | 60 | '2020/12/12' | 100 | 120 | 250.7 |
| 13 | 60 | '2020/12/13' | 106 | 128 | 345.3 |
| 14 | 60 | '2020/12/14' | 104 | 132 | 379.3 |
| 15 | 60 | '2020/12/15' | 98 | 123 | 275.0 |

| | | | | | |
|---|---|---|---|---|---|
| 16 | 60 | '2020/12/16' | 98 | 120 | 215.2 |
| 17 | 60 | '2020/12/17' | 100 | 120 | 300.0 |
| 19 | 60 | '2020/12/19' | 103 | 123 | 323.0 |
| 20 | 45 | '2020/12/20' | 97 | 125 | 243.0 |
| 21 | 60 | '2020/12/21' | 108 | 131 | 364.2 |
| 23 | 60 | '2020/12/23' | 130 | 101 | 300.0 |
| 24 | 45 | '2020/12/24' | 105 | 132 | 246.0 |
| 25 | 60 | '2020/12/25' | 102 | 126 | 334.5 |
| 26 | 60 | 20201226 | 100 | 120 | 250.0 |
| 27 | 60 | '2020/12/27' | 92 | 118 | 241.0 |
| 29 | 60 | '2020/12/29' | 100 | 132 | 280.0 |
| 30 | 60 | '2020/12/30' | 102 | 129 | 380.3 |
| 31 | 60 | '2020/12/31' | 92 | 115 | 243.0 |

In [13]:
```python
new_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 29 entries, 0 to 31
Data columns (total 5 columns):
 #   Column    Non-Null Count  Dtype
---  ------    --------------  -----
 0   Duration  29 non-null     int64
 1   Date      29 non-null     object
 2   Pulse     29 non-null     int64
 3   Maxpulse  29 non-null     int64
 4   Calories  29 non-null     float64
dtypes: float64(1), int64(3), object(1)
memory usage: 1.4+ KB
```

## Replace Empty Values

- Another way of dealing with empty cells is to insert a new value instead.
- This way you do not have to delete entire rows just because of some empty cells.
- The `fillna()` method allows us to replace empty cells with a value.

In [14]:
```python
# Replace null values with 999
new_df = ddf.fillna(999)
new_df
```

Out[14]:

| | Duration | Date | Pulse | Maxpulse | Calories |
|---|---|---|---|---|---|
| 0 | 60 | '2020/12/01' | 110 | 130 | 409.1 |
| 1 | 60 | '2020/12/02' | 117 | 145 | 479.0 |
| 2 | 60 | '2020/12/03' | 103 | 135 | 340.0 |
| 3 | 45 | '2020/12/04' | 109 | 175 | 282.4 |
| 4 | 45 | '2020/12/05' | 117 | 148 | 406.0 |
| 5 | 60 | '2020/12/06' | 102 | 127 | 300.0 |
| 6 | 60 | '2020/12/07' | 110 | 136 | 374.0 |
| 7 | 450 | '2020/12/08' | 104 | 134 | 253.3 |
| 8 | 30 | '2020/12/09' | 109 | 133 | 195.1 |
| 9 | 60 | '2020/12/10' | 98 | 124 | 269.0 |
| 10 | 60 | '2020/12/11' | 103 | 147 | 329.3 |
| 11 | 60 | '2020/12/12' | 100 | 120 | 250.7 |
| 12 | 60 | '2020/12/12' | 100 | 120 | 250.7 |
| 13 | 60 | '2020/12/13' | 106 | 128 | 345.3 |
| 14 | 60 | '2020/12/14' | 104 | 132 | 379.3 |
| 15 | 60 | '2020/12/15' | 98 | 123 | 275.0 |
| 16 | 60 | '2020/12/16' | 98 | 120 | 215.2 |

|    | Duration | Date | Pulse | Maxpulse | Calories |
|----|----------|------|-------|----------|----------|
| 17 | 60 | '2020/12/17' | 100 | 120 | 300.0 |
| 18 | 45 | '2020/12/18' | 90 | 112 | 999.0 |
| 19 | 60 | '2020/12/19' | 103 | 123 | 323.0 |
| 20 | 45 | '2020/12/20' | 97 | 125 | 243.0 |
| 21 | 60 | '2020/12/21' | 108 | 131 | 364.2 |
| 22 | 45 | 999 | 100 | 119 | 282.0 |
| 23 | 60 | '2020/12/23' | 130 | 101 | 300.0 |
| 24 | 45 | '2020/12/24' | 105 | 132 | 246.0 |
| 25 | 60 | '2020/12/25' | 102 | 126 | 334.5 |
| 26 | 60 | 20201226 | 100 | 120 | 250.0 |
| 27 | 60 | '2020/12/27' | 92 | 118 | 241.0 |
| 28 | 60 | '2020/12/28' | 103 | 132 | 999.0 |
| 29 | 60 | '2020/12/29' | 100 | 132 | 280.0 |
| 30 | 60 | '2020/12/30' | 102 | 129 | 380.3 |
| 31 | 60 | '2020/12/31' | 92 | 115 | 243.0 |

In [15]:
```python
new_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 32 entries, 0 to 31
Data columns (total 5 columns):
 #   Column    Non-Null Count  Dtype
---  ------    --------------  -----
 0   Duration  32 non-null     int64
 1   Date      32 non-null     object
 2   Pulse     32 non-null     int64
 3   Maxpulse  32 non-null     int64
 4   Calories  32 non-null     float64
dtypes: float64(1), int64(3), object(1)
memory usage: 1.4+ KB
```

```
In [16]:   # Replace only for a specific column
           # Replace NULL values in the "Calories" columns with the number 888:
           ddf2 = pd.read_csv('dirtydata.csv')
           ddf2['Calories'].fillna(888, inplace = True)
           ddf2
```

Out[16]:

| | Duration | Date | Pulse | Maxpulse | Calories |
|---|---|---|---|---|---|
| 0 | 60 | '2020/12/01' | 110 | 130 | 409.1 |
| 1 | 60 | '2020/12/02' | 117 | 145 | 479.0 |
| 2 | 60 | '2020/12/03' | 103 | 135 | 340.0 |
| 3 | 45 | '2020/12/04' | 109 | 175 | 282.4 |
| 4 | 45 | '2020/12/05' | 117 | 148 | 406.0 |
| 5 | 60 | '2020/12/06' | 102 | 127 | 300.0 |
| 6 | 60 | '2020/12/07' | 110 | 136 | 374.0 |
| 7 | 450 | '2020/12/08' | 104 | 134 | 253.3 |
| 8 | 30 | '2020/12/09' | 109 | 133 | 195.1 |
| 9 | 60 | '2020/12/10' | 98 | 124 | 269.0 |
| 10 | 60 | '2020/12/11' | 103 | 147 | 329.3 |
| 11 | 60 | '2020/12/12' | 100 | 120 | 250.7 |
| 12 | 60 | '2020/12/12' | 100 | 120 | 250.7 |
| 13 | 60 | '2020/12/13' | 106 | 128 | 345.3 |
| 14 | 60 | '2020/12/14' | 104 | 132 | 379.3 |
| 15 | 60 | '2020/12/15' | 98 | 123 | 275.0 |
| 16 | 60 | '2020/12/16' | 98 | 120 | 215.2 |
| 17 | 60 | '2020/12/17' | 100 | 120 | 300.0 |
| 18 | 45 | '2020/12/18' | 90 | 112 | 888.0 |

| | | | | | |
|---|---|---|---|---|---|
| 19 | 60 | '2020/12/19' | 103 | 123 | 323.0 |
| 20 | 45 | '2020/12/20' | 97 | 125 | 243.0 |
| 21 | 60 | '2020/12/21' | 108 | 131 | 364.2 |
| 22 | 45 | NaN | 100 | 119 | 282.0 |
| 23 | 60 | '2020/12/23' | 130 | 101 | 300.0 |
| 24 | 45 | '2020/12/24' | 105 | 132 | 246.0 |
| 25 | 60 | '2020/12/25' | 102 | 126 | 334.5 |
| 26 | 60 | 20201226 | 100 | 120 | 250.0 |
| 27 | 60 | '2020/12/27' | 92 | 118 | 241.0 |
| 28 | 60 | '2020/12/28' | 103 | 132 | 888.0 |
| 29 | 60 | '2020/12/29' | 100 | 132 | 280.0 |
| 30 | 60 | '2020/12/30' | 102 | 129 | 380.3 |
| 31 | 60 | '2020/12/31' | 92 | 115 | 243.0 |

## Replace Using Mean, Median, or Mode

In [17]:
```python
# using mean() method
# Mean = the average value (the sum of all values divided by number of values).

ddf3 = pd.read_csv('dirtydata.csv')
x = ddf3["Calories"].mean()
x
```

Out[17]: 304.68

In [18]:
```python
ddf3["Calories"].fillna(x, inplace = True)
ddf3
```

Out[18]:

| | Duration | Date | Pulse | Maxpulse | Calories |
|---|---|---|---|---|---|

| | | | | | |
|---|---|---|---|---|---|
| 0 | 60 | '2020/12/01' | 110 | 130 | 409.10 |
| 1 | 60 | '2020/12/02' | 117 | 145 | 479.00 |
| 2 | 60 | '2020/12/03' | 103 | 135 | 340.00 |
| 3 | 45 | '2020/12/04' | 109 | 175 | 282.40 |
| 4 | 45 | '2020/12/05' | 117 | 148 | 406.00 |
| 5 | 60 | '2020/12/06' | 102 | 127 | 300.00 |
| 6 | 60 | '2020/12/07' | 110 | 136 | 374.00 |
| 7 | 450 | '2020/12/08' | 104 | 134 | 253.30 |
| 8 | 30 | '2020/12/09' | 109 | 133 | 195.10 |
| 9 | 60 | '2020/12/10' | 98 | 124 | 269.00 |
| 10 | 60 | '2020/12/11' | 103 | 147 | 329.30 |
| 11 | 60 | '2020/12/12' | 100 | 120 | 250.70 |
| 12 | 60 | '2020/12/12' | 100 | 120 | 250.70 |
| 13 | 60 | '2020/12/13' | 106 | 128 | 345.30 |
| 14 | 60 | '2020/12/14' | 104 | 132 | 379.30 |
| 15 | 60 | '2020/12/15' | 98 | 123 | 275.00 |
| 16 | 60 | '2020/12/16' | 98 | 120 | 215.20 |
| 17 | 60 | '2020/12/17' | 100 | 120 | 300.00 |
| 18 | 45 | '2020/12/18' | 90 | 112 | 304.68 |
| 19 | 60 | '2020/12/19' | 103 | 123 | 323.00 |
| 20 | 45 | '2020/12/20' | 97 | 125 | 243.00 |
| 21 | 60 | '2020/12/21' | 108 | 131 | 364.20 |
| 22 | 45 | NaN | 100 | 119 | 282.00 |
| 23 | 60 | '2020/12/23' | 130 | 101 | 300.00 |
| 24 | 45 | '2020/12/24' | 105 | 132 | 246.00 |

| | Duration | Date | Pulse | Maxpulse | Calories |
|---|---|---|---|---|---|
| **25** | 60 | '2020/12/25' | 102 | 126 | 334.50 |
| **26** | 60 | 20201226 | 100 | 120 | 250.00 |
| **27** | 60 | '2020/12/27' | 92 | 118 | 241.00 |
| **28** | 60 | '2020/12/28' | 103 | 132 | 304.68 |
| **29** | 60 | '2020/12/29' | 100 | 132 | 280.00 |
| **30** | 60 | '2020/12/30' | 102 | 129 | 380.30 |
| **31** | 60 | '2020/12/31' | 92 | 115 | 243.00 |

In [19]:
```python
# using median() method
# Median = the value in the middle, after you have sorted all values ascending.
ddf3 = pd.read_csv('dirtydata.csv')

x = ddf3["Calories"].median()
x
```

Out[19]: 291.2

In [20]:
```python
ddf3["Calories"].fillna(x, inplace = True)
ddf3
```

Out[20]:

| | Duration | Date | Pulse | Maxpulse | Calories |
|---|---|---|---|---|---|
| **0** | 60 | '2020/12/01' | 110 | 130 | 409.1 |
| **1** | 60 | '2020/12/02' | 117 | 145 | 479.0 |
| **2** | 60 | '2020/12/03' | 103 | 135 | 340.0 |
| **3** | 45 | '2020/12/04' | 109 | 175 | 282.4 |
| **4** | 45 | '2020/12/05' | 117 | 148 | 406.0 |
| **5** | 60 | '2020/12/06' | 102 | 127 | 300.0 |
| **6** | 60 | '2020/12/07' | 110 | 136 | 374.0 |

| | | | | | |
|---|---|---|---|---|---|
| 7 | 450 | '2020/12/08' | 104 | 134 | 253.3 |
| 8 | 30 | '2020/12/09' | 109 | 133 | 195.1 |
| 9 | 60 | '2020/12/10' | 98 | 124 | 269.0 |
| 10 | 60 | '2020/12/11' | 103 | 147 | 329.3 |
| 11 | 60 | '2020/12/12' | 100 | 120 | 250.7 |
| 12 | 60 | '2020/12/12' | 100 | 120 | 250.7 |
| 13 | 60 | '2020/12/13' | 106 | 128 | 345.3 |
| 14 | 60 | '2020/12/14' | 104 | 132 | 379.3 |
| 15 | 60 | '2020/12/15' | 98 | 123 | 275.0 |
| 16 | 60 | '2020/12/16' | 98 | 120 | 215.2 |
| 17 | 60 | '2020/12/17' | 100 | 120 | 300.0 |
| 18 | 45 | '2020/12/18' | 90 | 112 | 291.2 |
| 19 | 60 | '2020/12/19' | 103 | 123 | 323.0 |
| 20 | 45 | '2020/12/20' | 97 | 125 | 243.0 |
| 21 | 60 | '2020/12/21' | 108 | 131 | 364.2 |
| 22 | 45 | NaN | 100 | 119 | 282.0 |
| 23 | 60 | '2020/12/23' | 130 | 101 | 300.0 |
| 24 | 45 | '2020/12/24' | 105 | 132 | 246.0 |
| 25 | 60 | '2020/12/25' | 102 | 126 | 334.5 |
| 26 | 60 | 20201226 | 100 | 120 | 250.0 |
| 27 | 60 | '2020/12/27' | 92 | 118 | 241.0 |
| 28 | 60 | '2020/12/28' | 103 | 132 | 291.2 |
| 29 | 60 | '2020/12/29' | 100 | 132 | 280.0 |
| 30 | 60 | '2020/12/30' | 102 | 129 | 380.3 |
| 31 | 60 | '2020/12/31' | 92 | 115 | 243.0 |

```
In [21]:   # using mode() method
           # Mode = the value that appears most frequently

           ddf3 = pd.read_csv('dirtydata.csv')

           x = ddf3["Calories"].mode()[0]
           x
```

Out[21]:  300.0

```
In [22]:   ddf3["Calories"].fillna(x, inplace = True)
           ddf3
```

Out[22]:

|    | Duration | Date         | Pulse | Maxpulse | Calories |
|----|----------|--------------|-------|----------|----------|
| 0  | 60       | '2020/12/01' | 110   | 130      | 409.1    |
| 1  | 60       | '2020/12/02' | 117   | 145      | 479.0    |
| 2  | 60       | '2020/12/03' | 103   | 135      | 340.0    |
| 3  | 45       | '2020/12/04' | 109   | 175      | 282.4    |
| 4  | 45       | '2020/12/05' | 117   | 148      | 406.0    |
| 5  | 60       | '2020/12/06' | 102   | 127      | 300.0    |
| 6  | 60       | '2020/12/07' | 110   | 136      | 374.0    |
| 7  | 450      | '2020/12/08' | 104   | 134      | 253.3    |
| 8  | 30       | '2020/12/09' | 109   | 133      | 195.1    |
| 9  | 60       | '2020/12/10' | 98    | 124      | 269.0    |
| 10 | 60       | '2020/12/11' | 103   | 147      | 329.3    |
| 11 | 60       | '2020/12/12' | 100   | 120      | 250.7    |
| 12 | 60       | '2020/12/12' | 100   | 120      | 250.7    |
```

| | | | | | |
|---|---|---|---|---|---|
| 13 | 60 | '2020/12/13' | 106 | 128 | 345.3 |
| 14 | 60 | '2020/12/14' | 104 | 132 | 379.3 |
| 15 | 60 | '2020/12/15' | 98 | 123 | 275.0 |
| 16 | 60 | '2020/12/16' | 98 | 120 | 215.2 |
| 17 | 60 | '2020/12/17' | 100 | 120 | 300.0 |
| 18 | 45 | '2020/12/18' | 90 | 112 | 300.0 |
| 19 | 60 | '2020/12/19' | 103 | 123 | 323.0 |
| 20 | 45 | '2020/12/20' | 97 | 125 | 243.0 |
| 21 | 60 | '2020/12/21' | 108 | 131 | 364.2 |
| 22 | 45 | NaN | 100 | 119 | 282.0 |
| 23 | 60 | '2020/12/23' | 130 | 101 | 300.0 |
| 24 | 45 | '2020/12/24' | 105 | 132 | 246.0 |
| 25 | 60 | '2020/12/25' | 102 | 126 | 334.5 |
| 26 | 60 | 20201226 | 100 | 120 | 250.0 |
| 27 | 60 | '2020/12/27' | 92 | 118 | 241.0 |
| 28 | 60 | '2020/12/28' | 103 | 132 | 300.0 |
| 29 | 60 | '2020/12/29' | 100 | 132 | 280.0 |
| 30 | 60 | '2020/12/30' | 102 | 129 | 380.3 |
| 31 | 60 | '2020/12/31' | 92 | 115 | 243.0 |

## Cleaning Wrong Format

- Cells with data of wrong format can make it difficult, or even impossible, to analyze data.

- Two options to fix it: remove the rows, or convert all cells in the columns into the same format.

```
In [23]:    # Data set
            ddf4 = pd.read_csv('dirtydata.csv')
            ddf4
```

Out[23]:

|    | Duration | Date | Pulse | Maxpulse | Calories |
|----|----------|------|-------|----------|----------|
| 0  | 60 | '2020/12/01' | 110 | 130 | 409.1 |
| 1  | 60 | '2020/12/02' | 117 | 145 | 479.0 |
| 2  | 60 | '2020/12/03' | 103 | 135 | 340.0 |
| 3  | 45 | '2020/12/04' | 109 | 175 | 282.4 |
| 4  | 45 | '2020/12/05' | 117 | 148 | 406.0 |
| 5  | 60 | '2020/12/06' | 102 | 127 | 300.0 |
| 6  | 60 | '2020/12/07' | 110 | 136 | 374.0 |
| 7  | 450 | '2020/12/08' | 104 | 134 | 253.3 |
| 8  | 30 | '2020/12/09' | 109 | 133 | 195.1 |
| 9  | 60 | '2020/12/10' | 98 | 124 | 269.0 |
| 10 | 60 | '2020/12/11' | 103 | 147 | 329.3 |
| 11 | 60 | '2020/12/12' | 100 | 120 | 250.7 |
| 12 | 60 | '2020/12/12' | 100 | 120 | 250.7 |
| 13 | 60 | '2020/12/13' | 106 | 128 | 345.3 |
| 14 | 60 | '2020/12/14' | 104 | 132 | 379.3 |
| 15 | 60 | '2020/12/15' | 98 | 123 | 275.0 |
| 16 | 60 | '2020/12/16' | 98 | 120 | 215.2 |
| 17 | 60 | '2020/12/17' | 100 | 120 | 300.0 |
| 18 | 45 | '2020/12/18' | 90 | 112 | NaN |
| 19 | 60 | '2020/12/19' | 103 | 123 | 323.0 |

| | | | | | |
|---|---|---|---|---|---|
| **20** | 45 | '2020/12/20' | 97 | 125 | 243.0 |
| **21** | 60 | '2020/12/21' | 108 | 131 | 364.2 |
| **22** | 45 | NaN | 100 | 119 | 282.0 |
| **23** | 60 | '2020/12/23' | 130 | 101 | 300.0 |
| **24** | 45 | '2020/12/24' | 105 | 132 | 246.0 |
| **25** | 60 | '2020/12/25' | 102 | 126 | 334.5 |
| **26** | 60 | 20201226 | 100 | 120 | 250.0 |
| **27** | 60 | '2020/12/27' | 92 | 118 | 241.0 |
| **28** | 60 | '2020/12/28' | 103 | 132 | NaN |
| **29** | 60 | '2020/12/29' | 100 | 132 | 280.0 |
| **30** | 60 | '2020/12/30' | 102 | 129 | 380.3 |
| **31** | 60 | '2020/12/31' | 92 | 115 | 243.0 |

In [24]:

```python
# Convert all the cells in the "Date" column into dates
ddf4['Date'] = pd.to_datetime(ddf4['Date'])
ddf4
```

Out[24]:

| | Duration | Date | Pulse | Maxpulse | Calories |
|---|---|---|---|---|---|
| **0** | 60 | 2020-12-01 | 110 | 130 | 409.1 |
| **1** | 60 | 2020-12-02 | 117 | 145 | 479.0 |
| **2** | 60 | 2020-12-03 | 103 | 135 | 340.0 |
| **3** | 45 | 2020-12-04 | 109 | 175 | 282.4 |
| **4** | 45 | 2020-12-05 | 117 | 148 | 406.0 |
| **5** | 60 | 2020-12-06 | 102 | 127 | 300.0 |
| **6** | 60 | 2020-12-07 | 110 | 136 | 374.0 |
| **7** | 450 | 2020-12-08 | 104 | 134 | 253.3 |

| | | | | | |
|---|---|---|---|---|---|
| **8** | 30 | 2020-12-09 | 109 | 133 | 195.1 |
| **9** | 60 | 2020-12-10 | 98 | 124 | 269.0 |
| **10** | 60 | 2020-12-11 | 103 | 147 | 329.3 |
| **11** | 60 | 2020-12-12 | 100 | 120 | 250.7 |
| **12** | 60 | 2020-12-12 | 100 | 120 | 250.7 |
| **13** | 60 | 2020-12-13 | 106 | 128 | 345.3 |
| **14** | 60 | 2020-12-14 | 104 | 132 | 379.3 |
| **15** | 60 | 2020-12-15 | 98 | 123 | 275.0 |
| **16** | 60 | 2020-12-16 | 98 | 120 | 215.2 |
| **17** | 60 | 2020-12-17 | 100 | 120 | 300.0 |
| **18** | 45 | 2020-12-18 | 90 | 112 | NaN |
| **19** | 60 | 2020-12-19 | 103 | 123 | 323.0 |
| **20** | 45 | 2020-12-20 | 97 | 125 | 243.0 |
| **21** | 60 | 2020-12-21 | 108 | 131 | 364.2 |
| **22** | 45 | NaT | 100 | 119 | 282.0 |
| **23** | 60 | 2020-12-23 | 130 | 101 | 300.0 |
| **24** | 45 | 2020-12-24 | 105 | 132 | 246.0 |
| **25** | 60 | 2020-12-25 | 102 | 126 | 334.5 |
| **26** | 60 | 2020-12-26 | 100 | 120 | 250.0 |
| **27** | 60 | 2020-12-27 | 92 | 118 | 241.0 |
| **28** | 60 | 2020-12-28 | 103 | 132 | NaN |
| **29** | 60 | 2020-12-29 | 100 | 132 | 280.0 |
| **30** | 60 | 2020-12-30 | 102 | 129 | 380.3 |
| **31** | 60 | 2020-12-31 | 92 | 115 | 243.0 |

```
In [25]:    # removes rows with NaT: Not a Time

            ddf4.dropna(subset=['Date'], inplace = True)
            ddf4
```

Out[25]:

| | Duration | Date | Pulse | Maxpulse | Calories |
|---|---|---|---|---|---|
| **0** | 60 | 2020-12-01 | 110 | 130 | 409.1 |
| **1** | 60 | 2020-12-02 | 117 | 145 | 479.0 |
| **2** | 60 | 2020-12-03 | 103 | 135 | 340.0 |
| **3** | 45 | 2020-12-04 | 109 | 175 | 282.4 |
| **4** | 45 | 2020-12-05 | 117 | 148 | 406.0 |
| **5** | 60 | 2020-12-06 | 102 | 127 | 300.0 |
| **6** | 60 | 2020-12-07 | 110 | 136 | 374.0 |
| **7** | 450 | 2020-12-08 | 104 | 134 | 253.3 |
| **8** | 30 | 2020-12-09 | 109 | 133 | 195.1 |
| **9** | 60 | 2020-12-10 | 98 | 124 | 269.0 |
| **10** | 60 | 2020-12-11 | 103 | 147 | 329.3 |
| **11** | 60 | 2020-12-12 | 100 | 120 | 250.7 |
| **12** | 60 | 2020-12-12 | 100 | 120 | 250.7 |
| **13** | 60 | 2020-12-13 | 106 | 128 | 345.3 |
| **14** | 60 | 2020-12-14 | 104 | 132 | 379.3 |
| **15** | 60 | 2020-12-15 | 98 | 123 | 275.0 |
| **16** | 60 | 2020-12-16 | 98 | 120 | 215.2 |
| **17** | 60 | 2020-12-17 | 100 | 120 | 300.0 |
| **18** | 45 | 2020-12-18 | 90 | 112 | NaN |
| **19** | 60 | 2020-12-19 | 103 | 123 | 323.0 |

| | | Date | | | |
|---|---|---|---|---|---|
| 20 | 45 | 2020-12-20 | 97 | 125 | 243.0 |
| 21 | 60 | 2020-12-21 | 108 | 131 | 364.2 |
| 23 | 60 | 2020-12-23 | 130 | 101 | 300.0 |
| 24 | 45 | 2020-12-24 | 105 | 132 | 246.0 |
| 25 | 60 | 2020-12-25 | 102 | 126 | 334.5 |
| 26 | 60 | 2020-12-26 | 100 | 120 | 250.0 |
| 27 | 60 | 2020-12-27 | 92 | 118 | 241.0 |
| 28 | 60 | 2020-12-28 | 103 | 132 | NaN |
| 29 | 60 | 2020-12-29 | 100 | 132 | 280.0 |
| 30 | 60 | 2020-12-30 | 102 | 129 | 380.3 |
| 31 | 60 | 2020-12-31 | 92 | 115 | 243.0 |

## Cleaning Wrong Data

- Wrong data means the data value is wrong. For example, "Duration" column should have values 30 to 60, but someone written 450, which is wrong.

- Sometimes you can spot wrong data by looking at the data set, because you have an expectation of what it should be.

In [26]:
```python
# Data set
ddf5 = pd.read_csv('dirtydata.csv')
ddf5
```

Out[26]:

| | Duration | Date | Pulse | Maxpulse | Calories |
|---|---|---|---|---|---|
| 0 | 60 | '2020/12/01' | 110 | 130 | 409.1 |
| 1 | 60 | '2020/12/02' | 117 | 145 | 479.0 |
| 2 | 60 | '2020/12/03' | 103 | 135 | 340.0 |

| | | | | | |
|---|---|---|---|---|---|
| 3 | 45 | '2020/12/04' | 109 | 175 | 282.4 |
| 4 | 45 | '2020/12/05' | 117 | 148 | 406.0 |
| 5 | 60 | '2020/12/06' | 102 | 127 | 300.0 |
| 6 | 60 | '2020/12/07' | 110 | 136 | 374.0 |
| 7 | 450 | '2020/12/08' | 104 | 134 | 253.3 |
| 8 | 30 | '2020/12/09' | 109 | 133 | 195.1 |
| 9 | 60 | '2020/12/10' | 98 | 124 | 269.0 |
| 10 | 60 | '2020/12/11' | 103 | 147 | 329.3 |
| 11 | 60 | '2020/12/12' | 100 | 120 | 250.7 |
| 12 | 60 | '2020/12/12' | 100 | 120 | 250.7 |
| 13 | 60 | '2020/12/13' | 106 | 128 | 345.3 |
| 14 | 60 | '2020/12/14' | 104 | 132 | 379.3 |
| 15 | 60 | '2020/12/15' | 98 | 123 | 275.0 |
| 16 | 60 | '2020/12/16' | 98 | 120 | 215.2 |
| 17 | 60 | '2020/12/17' | 100 | 120 | 300.0 |
| 18 | 45 | '2020/12/18' | 90 | 112 | NaN |
| 19 | 60 | '2020/12/19' | 103 | 123 | 323.0 |
| 20 | 45 | '2020/12/20' | 97 | 125 | 243.0 |
| 21 | 60 | '2020/12/21' | 108 | 131 | 364.2 |
| 22 | 45 | NaN | 100 | 119 | 282.0 |
| 23 | 60 | '2020/12/23' | 130 | 101 | 300.0 |
| 24 | 45 | '2020/12/24' | 105 | 132 | 246.0 |
| 25 | 60 | '2020/12/25' | 102 | 126 | 334.5 |
| 26 | 60 | 20201226 | 100 | 120 | 250.0 |

| | Duration | Date | Pulse | Maxpulse | Calories |
|---|---|---|---|---|---|
| 27 | 60 | '2020/12/27' | 92 | 118 | 241.0 |
| 28 | 60 | '2020/12/28' | 103 | 132 | NaN |
| 29 | 60 | '2020/12/29' | 100 | 132 | 280.0 |
| 30 | 60 | '2020/12/30' | 102 | 129 | 380.3 |
| 31 | 60 | '2020/12/31' | 92 | 115 | 243.0 |

In the above dataset, in row 7 of `Duration` column, the value should be "45" instead of "450".

In [27]:
```python
# Set "Duration" = 45 in row 7

ddf5.loc[7, 'Duration'] = 45
ddf5
```

Out[27]:

| | Duration | Date | Pulse | Maxpulse | Calories |
|---|---|---|---|---|---|
| 0 | 60 | '2020/12/01' | 110 | 130 | 409.1 |
| 1 | 60 | '2020/12/02' | 117 | 145 | 479.0 |
| 2 | 60 | '2020/12/03' | 103 | 135 | 340.0 |
| 3 | 45 | '2020/12/04' | 109 | 175 | 282.4 |
| 4 | 45 | '2020/12/05' | 117 | 148 | 406.0 |
| 5 | 60 | '2020/12/06' | 102 | 127 | 300.0 |
| 6 | 60 | '2020/12/07' | 110 | 136 | 374.0 |
| 7 | 45 | '2020/12/08' | 104 | 134 | 253.3 |
| 8 | 30 | '2020/12/09' | 109 | 133 | 195.1 |
| 9 | 60 | '2020/12/10' | 98 | 124 | 269.0 |
| 10 | 60 | '2020/12/11' | 103 | 147 | 329.3 |
| 11 | 60 | '2020/12/12' | 100 | 120 | 250.7 |
| 12 | 60 | '2020/12/12' | 100 | 120 | 250.7 |

| | | | | | |
|---|---|---|---|---|---|
| **13** | 60 | '2020/12/13' | 106 | 128 | 345.3 |
| **14** | 60 | '2020/12/14' | 104 | 132 | 379.3 |
| **15** | 60 | '2020/12/15' | 98 | 123 | 275.0 |
| **16** | 60 | '2020/12/16' | 98 | 120 | 215.2 |
| **17** | 60 | '2020/12/17' | 100 | 120 | 300.0 |
| **18** | 45 | '2020/12/18' | 90 | 112 | NaN |
| **19** | 60 | '2020/12/19' | 103 | 123 | 323.0 |
| **20** | 45 | '2020/12/20' | 97 | 125 | 243.0 |
| **21** | 60 | '2020/12/21' | 108 | 131 | 364.2 |
| **22** | 45 | NaN | 100 | 119 | 282.0 |
| **23** | 60 | '2020/12/23' | 130 | 101 | 300.0 |
| **24** | 45 | '2020/12/24' | 105 | 132 | 246.0 |
| **25** | 60 | '2020/12/25' | 102 | 126 | 334.5 |
| **26** | 60 | 20201226 | 100 | 120 | 250.0 |
| **27** | 60 | '2020/12/27' | 92 | 118 | 241.0 |
| **28** | 60 | '2020/12/28' | 103 | 132 | NaN |
| **29** | 60 | '2020/12/29' | 100 | 132 | 280.0 |
| **30** | 60 | '2020/12/30' | 102 | 129 | 380.3 |
| **31** | 60 | '2020/12/31' | 92 | 115 | 243.0 |

```
In [28]:    """
            Loop through all values in the "Duration" column.

            If the value is higher than 120, set it to 120:
            """
            for x in ddf5.index:
              if ddf5.loc[x, "Duration"] > 120:
                ddf5.loc[x, "Duration"] = 120

            ddf5
```

Out[28]:

| | Duration | Date | Pulse | Maxpulse | Calories |
|---|---|---|---|---|---|
| 0 | 60 | '2020/12/01' | 110 | 130 | 409.1 |
| 1 | 60 | '2020/12/02' | 117 | 145 | 479.0 |
| 2 | 60 | '2020/12/03' | 103 | 135 | 340.0 |
| 3 | 45 | '2020/12/04' | 109 | 175 | 282.4 |
| 4 | 45 | '2020/12/05' | 117 | 148 | 406.0 |
| 5 | 60 | '2020/12/06' | 102 | 127 | 300.0 |
| 6 | 60 | '2020/12/07' | 110 | 136 | 374.0 |
| 7 | 45 | '2020/12/08' | 104 | 134 | 253.3 |
| 8 | 30 | '2020/12/09' | 109 | 133 | 195.1 |
| 9 | 60 | '2020/12/10' | 98 | 124 | 269.0 |
| 10 | 60 | '2020/12/11' | 103 | 147 | 329.3 |
| 11 | 60 | '2020/12/12' | 100 | 120 | 250.7 |
| 12 | 60 | '2020/12/12' | 100 | 120 | 250.7 |
| 13 | 60 | '2020/12/13' | 106 | 128 | 345.3 |
| 14 | 60 | '2020/12/14' | 104 | 132 | 379.3 |
| 15 | 60 | '2020/12/15' | 98 | 123 | 275.0 |

| | | | | | |
|---|---|---|---|---|---|
| **16** | 60 | '2020/12/16' | 98 | 120 | 215.2 |
| **17** | 60 | '2020/12/17' | 100 | 120 | 300.0 |
| **18** | 45 | '2020/12/18' | 90 | 112 | NaN |
| **19** | 60 | '2020/12/19' | 103 | 123 | 323.0 |
| **20** | 45 | '2020/12/20' | 97 | 125 | 243.0 |
| **21** | 60 | '2020/12/21' | 108 | 131 | 364.2 |
| **22** | 45 | NaN | 100 | 119 | 282.0 |
| **23** | 60 | '2020/12/23' | 130 | 101 | 300.0 |
| **24** | 45 | '2020/12/24' | 105 | 132 | 246.0 |
| **25** | 60 | '2020/12/25' | 102 | 126 | 334.5 |
| **26** | 60 | 20201226 | 100 | 120 | 250.0 |
| **27** | 60 | '2020/12/27' | 92 | 118 | 241.0 |
| **28** | 60 | '2020/12/28' | 103 | 132 | NaN |
| **29** | 60 | '2020/12/29' | 100 | 132 | 280.0 |
| **30** | 60 | '2020/12/30' | 102 | 129 | 380.3 |
| **31** | 60 | '2020/12/31' | 92 | 115 | 243.0 |

In [29]:
```python
"""
OR, Remove all the rows which have wrong data.

Delete rows where "Duration" is higher than 120:
"""

for x in ddf5.index:
  if ddf5.loc[x, "Duration"] > 120:
    ddf5.drop(x, inplace = True)
ddf5
```

Out[29]:

| Duration | Date | Pulse | Maxpulse | Calories |
|---|---|---|---|---|

|    |    |              |     |     |       |
|----|----|--------------|-----|-----|-------|
| 0  | 60 | '2020/12/01' | 110 | 130 | 409.1 |
| 1  | 60 | '2020/12/02' | 117 | 145 | 479.0 |
| 2  | 60 | '2020/12/03' | 103 | 135 | 340.0 |
| 3  | 45 | '2020/12/04' | 109 | 175 | 282.4 |
| 4  | 45 | '2020/12/05' | 117 | 148 | 406.0 |
| 5  | 60 | '2020/12/06' | 102 | 127 | 300.0 |
| 6  | 60 | '2020/12/07' | 110 | 136 | 374.0 |
| 7  | 45 | '2020/12/08' | 104 | 134 | 253.3 |
| 8  | 30 | '2020/12/09' | 109 | 133 | 195.1 |
| 9  | 60 | '2020/12/10' | 98  | 124 | 269.0 |
| 10 | 60 | '2020/12/11' | 103 | 147 | 329.3 |
| 11 | 60 | '2020/12/12' | 100 | 120 | 250.7 |
| 12 | 60 | '2020/12/12' | 100 | 120 | 250.7 |
| 13 | 60 | '2020/12/13' | 106 | 128 | 345.3 |
| 14 | 60 | '2020/12/14' | 104 | 132 | 379.3 |
| 15 | 60 | '2020/12/15' | 98  | 123 | 275.0 |
| 16 | 60 | '2020/12/16' | 98  | 120 | 215.2 |
| 17 | 60 | '2020/12/17' | 100 | 120 | 300.0 |
| 18 | 45 | '2020/12/18' | 90  | 112 | NaN   |
| 19 | 60 | '2020/12/19' | 103 | 123 | 323.0 |
| 20 | 45 | '2020/12/20' | 97  | 125 | 243.0 |
| 21 | 60 | '2020/12/21' | 108 | 131 | 364.2 |
| 22 | 45 | NaN          | 100 | 119 | 282.0 |
| 23 | 60 | '2020/12/23' | 130 | 101 | 300.0 |

| | Duration | Date | Pulse | Maxpulse | Calories |
|---|---|---|---|---|---|
| 24 | 45 | '2020/12/24' | 105 | 132 | 246.0 |
| 25 | 60 | '2020/12/25' | 102 | 126 | 334.5 |
| 26 | 60 | 20201226 | 100 | 120 | 250.0 |
| 27 | 60 | '2020/12/27' | 92 | 118 | 241.0 |
| 28 | 60 | '2020/12/28' | 103 | 132 | NaN |
| 29 | 60 | '2020/12/29' | 100 | 132 | 280.0 |
| 30 | 60 | '2020/12/30' | 102 | 129 | 380.3 |
| 31 | 60 | '2020/12/31' | 92 | 115 | 243.0 |

## Cleanning Duplicates

- Duplicate rows are rows that have been registered more than one time.

In [30]:

```
ddf5
```

Out[30]:

| | Duration | Date | Pulse | Maxpulse | Calories |
|---|---|---|---|---|---|
| 0 | 60 | '2020/12/01' | 110 | 130 | 409.1 |
| 1 | 60 | '2020/12/02' | 117 | 145 | 479.0 |
| 2 | 60 | '2020/12/03' | 103 | 135 | 340.0 |
| 3 | 45 | '2020/12/04' | 109 | 175 | 282.4 |
| 4 | 45 | '2020/12/05' | 117 | 148 | 406.0 |
| 5 | 60 | '2020/12/06' | 102 | 127 | 300.0 |
| 6 | 60 | '2020/12/07' | 110 | 136 | 374.0 |
| 7 | 45 | '2020/12/08' | 104 | 134 | 253.3 |
| 8 | 30 | '2020/12/09' | 109 | 133 | 195.1 |
| 9 | 60 | '2020/12/10' | 98 | 124 | 269.0 |

| | | | | | |
|---|---|---|---|---|---|
| 10 | 60 | '2020/12/11' | 103 | 147 | 329.3 |
| 11 | 60 | '2020/12/12' | 100 | 120 | 250.7 |
| 12 | 60 | '2020/12/12' | 100 | 120 | 250.7 |
| 13 | 60 | '2020/12/13' | 106 | 128 | 345.3 |
| 14 | 60 | '2020/12/14' | 104 | 132 | 379.3 |
| 15 | 60 | '2020/12/15' | 98 | 123 | 275.0 |
| 16 | 60 | '2020/12/16' | 98 | 120 | 215.2 |
| 17 | 60 | '2020/12/17' | 100 | 120 | 300.0 |
| 18 | 45 | '2020/12/18' | 90 | 112 | NaN |
| 19 | 60 | '2020/12/19' | 103 | 123 | 323.0 |
| 20 | 45 | '2020/12/20' | 97 | 125 | 243.0 |
| 21 | 60 | '2020/12/21' | 108 | 131 | 364.2 |
| 22 | 45 | NaN | 100 | 119 | 282.0 |
| 23 | 60 | '2020/12/23' | 130 | 101 | 300.0 |
| 24 | 45 | '2020/12/24' | 105 | 132 | 246.0 |
| 25 | 60 | '2020/12/25' | 102 | 126 | 334.5 |
| 26 | 60 | 20201226 | 100 | 120 | 250.0 |
| 27 | 60 | '2020/12/27' | 92 | 118 | 241.0 |
| 28 | 60 | '2020/12/28' | 103 | 132 | NaN |
| 29 | 60 | '2020/12/29' | 100 | 132 | 280.0 |
| 30 | 60 | '2020/12/30' | 102 | 129 | 380.3 |
| 31 | 60 | '2020/12/31' | 92 | 115 | 243.0 |

In the above dataset, rows 11 and 12 are duplicates.

```python
# To remove duplicates, use the drop_duplicates() method

ddf5.drop_duplicates(inplace = True)
ddf5
```

| | Duration | Date | Pulse | Maxpulse | Calories |
|---|---|---|---|---|---|
| 0 | 60 | '2020/12/01' | 110 | 130 | 409.1 |
| 1 | 60 | '2020/12/02' | 117 | 145 | 479.0 |
| 2 | 60 | '2020/12/03' | 103 | 135 | 340.0 |
| 3 | 45 | '2020/12/04' | 109 | 175 | 282.4 |
| 4 | 45 | '2020/12/05' | 117 | 148 | 406.0 |
| 5 | 60 | '2020/12/06' | 102 | 127 | 300.0 |
| 6 | 60 | '2020/12/07' | 110 | 136 | 374.0 |
| 7 | 45 | '2020/12/08' | 104 | 134 | 253.3 |
| 8 | 30 | '2020/12/09' | 109 | 133 | 195.1 |
| 9 | 60 | '2020/12/10' | 98 | 124 | 269.0 |
| 10 | 60 | '2020/12/11' | 103 | 147 | 329.3 |
| 11 | 60 | '2020/12/12' | 100 | 120 | 250.7 |
| 13 | 60 | '2020/12/13' | 106 | 128 | 345.3 |
| 14 | 60 | '2020/12/14' | 104 | 132 | 379.3 |
| 15 | 60 | '2020/12/15' | 98 | 123 | 275.0 |
| 16 | 60 | '2020/12/16' | 98 | 120 | 215.2 |
| 17 | 60 | '2020/12/17' | 100 | 120 | 300.0 |
| 18 | 45 | '2020/12/18' | 90 | 112 | NaN |
| 19 | 60 | '2020/12/19' | 103 | 123 | 323.0 |
| 20 | 45 | '2020/12/20' | 97 | 125 | 243.0 |

| | | | | | |
|---|---|---|---|---|---|
| **21** | 60 | '2020/12/21' | 108 | 131 | 364.2 |
| **22** | 45 | NaN | 100 | 119 | 282.0 |
| **23** | 60 | '2020/12/23' | 130 | 101 | 300.0 |
| **24** | 45 | '2020/12/24' | 105 | 132 | 246.0 |
| **25** | 60 | '2020/12/25' | 102 | 126 | 334.5 |
| **26** | 60 | 20201226 | 100 | 120 | 250.0 |
| **27** | 60 | '2020/12/27' | 92 | 118 | 241.0 |
| **28** | 60 | '2020/12/28' | 103 | 132 | NaN |
| **29** | 60 | '2020/12/29' | 100 | 132 | 280.0 |
| **30** | 60 | '2020/12/30' | 102 | 129 | 380.3 |
| **31** | 60 | '2020/12/31' | 92 | 115 | 243.0 |