| Cleaning data | <<CleaningData.csv>> |
|---|---|

| Empty Cells | Dropna() to remove empty cells |
|---|---|
| | Fillna() to fill something in empty cells |

| Dropna() | 16 60 '2020/12/16' 98 12<br>17 60 '2020/12/17' 100 12<br>18 45 '2020/12/18' 90 11<br>19 60 '2020/12/19' 103 12<br>20 45 '2020/12/20' 97 12<br>21 60 '2020/12/21' 108 13<br>22 45 NaN 100 11<br>23 60 '2020/12/23' 130 10<br>24 45 '2020/12/24' 105 13<br>25 60 '2020/12/25' 102 12<br>26 60 12/26/2020 100 12<br>27 60 '2020/12/27' 92 11<br>28 60 '2020/12/28' 103 13<br>29 60 '2020/12/29' 100 13 | import pandas as pd<br><br>df = pd.read_csv('CleaningData.csv')<br><br>new_df = df.dropna()<br><br>print(new_df.to_string()) | Duration Date Pulse Maxpuls<br>0 60 '2020/12/01' 110 13<br>1 60 '2020/12/02' 117 14<br>2 60 '2020/12/03' 103 13<br>3 45 '2020/12/04' 109 17<br>4 45 '2020/12/05' 117 14<br>5 60 '2020/12/06' 102 12<br>6 60 '2020/12/07' 110 13<br>7 450 '2020/12/08' 104 13<br>8 30 '2020/12/09' 109 13<br>9 60 '2020/12/10' 98 12<br>10 60 '2020/12/11' 103 14<br>11 60 '2020/12/12' 100 12<br>12 60 '2020/12/12' 100 12<br>13 60 '2020/12/13' 106 12<br>14 60 '2020/12/14' 104 13<br>15 60 '2020/12/15' 98 12<br>16 60 '2020/12/16' 98 12<br>17 60 '2020/12/17' 100 12<br>19 60 '2020/12/19' 103 12<br>20 45 '2020/12/20' 97 12<br>21 60 '2020/12/21' 108 13<br>23 60 '2020/12/23' 130 10<br>24 45 '2020/12/24' 105 11<br>25 60 '2020/12/25' 102 12<br>26 60 12/26/2020 100 12<br>27 60 '2020/12/27' 92 11<br>29 60 '2020/12/29' 100 13 |
|---|---|---|---|

Note= We can see its replacing new dataframe not the original one

| To replace null values in original dataframe use inplace=True | import pandas as pd<br><br>df = pd.read_csv('data.csv')<br><br>df.dropna(inplace = True)<br><br>print(df.to_string()) |
|---|---|

| Fillna() | import pandas as pd<br><br>df = pd.read_csv('data.csv')<br><br>df.fillna(130, inplace = True)<br>454rf |
|---|---|

Since we have less data its not good to remove entire row

| | | | | |
|---|---|---|---|---|
| For numerical values we will calculate mean value for coulmns | x = df["Calories"].mean() df["Calories"].fillna(x, inplace = True) | `print(x)` `304.68` | We got x value as 304.68 and it will filled with null values in df |  |

| | | | |
|---|---|---|---|
| Our dataset has wrong date format |  | df['Date'] = pd.to_datetime(df['Date']) print(df.to_string()) |  |

| | | |
|---|---|---|
| Its better to drop null values in date | df.dropna(subset=['Date'], inplace = True) | ```
20      45 2020-12-20      97      125    243.00
21      60 2020-12-21     108      131    364.20
23      60 2020-12-23     130      101    300.00
```  So 22 value is removed |

| | |
|---|---|
| Wrong data or unmatchable to respective column values | ```
    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
```  Here one value contains 450 m which may not be appropriate so we will reduce it two 45 using loc |

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| print(df.to_string()) | | | 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 | 45 | 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 |

| | | | | |
|---|---|---|---|---|
| Remove Duplicates | To check duplicate values | Df.duplicated() | 0 | False |
| | | | 1 | False |
| | | | 2 | False |
| | | | 3 | False |
| | | | 4 | False |
| | | | 5 | False |
| | | | 6 | False |
| | | | 7 | False |
| | | | 8 | False |
| | | | 9 | False |
| | | | 10 | False |
| | | | 11 | False |
| | | | 12 | True |
| | | | 13 | False |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| 12 row has dublicate values | df.drop_duplicates() | | 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 |
| | To drop dublicates | 13 | 60 | 2020-12-13 | 106 | 128 | 345.30 |

Data