

```
In [30]: print("Author: Louis-Antoine Etchian")
```

Author: Louis-Antoine Etchian

```
In [3]: import numpy as np
import pandas as pd
```

Load data from file

Most often data will come from somewhere, often csv files, and using `pd.read_csv()` will allow smooth creation of DataFrames.

Let's load that same heart-attack.csv that we used in Numpy before:

```
In [4]: data_names=np.array(["mpg", "cylinders", "displacement", "horsepower", "weight",
data=pd.read_csv('auto-mpg.data', delim_whitespace=True, header=None, names = data_names)
```

After loading data, it is good practice to check what we have. Usually, the sequences is:

1. Check dimension
2. Peek at the first rows
3. Get info on data types and missing values
4. Summarize columns

```
In [5]: # Check dimension (rows, columns)
data.shape
```

Out[5]: (398, 9)

```
In [6]: # Peek at the first rows
data.head()
```

Out[6]:

	mpg	cylinders	displacement	horsepower	weight	acceleration	model year	origin	car name
0	18.0	8	307.0	130.0	3504.0	12.0	70	1	chevrolet chevelle malibu
1	15.0	8	350.0	165.0	3693.0	11.5	70	1	buick skylark 320
2	18.0	8	318.0	150.0	3436.0	11.0	70	1	plymouth satellite
3	16.0	8	304.0	150.0	3433.0	12.0	70	1	amc rebel sst
4	17.0	8	302.0	140.0	3449.0	10.5	70	1	ford torino

```
In [8]: # Column names are
data.columns
```

```
Out[8]: Index(['mpg', 'cylinders', 'displacement', 'horsepower', 'weight',
              'acceleration', 'model year', 'origin', 'car name'],
              dtype='object')
```

```
In [9]: # Get info on data types and missing values
data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 398 entries, 0 to 397
Data columns (total 9 columns):
 #   Column          Non-Null Count  Dtype
---  -
 0   mpg             398 non-null   float64
 1   cylinders        398 non-null   int64
 2   displacement     398 non-null   float64
 3   horsepower       398 non-null   object
 4   weight           398 non-null   float64
 5   acceleration     398 non-null   float64
 6   model year       398 non-null   int64
 7   origin           398 non-null   int64
 8   car name         398 non-null   object
dtypes: float64(4), int64(3), object(2)
memory usage: 28.1+ KB
```

Summarize values

What is the mean, std, min, max in each column?

```
In [10]: data.mean()
```

C:\Users\LA\AppData\Local\Temp\ipykernel_8952\531903386.py:1: FutureWarning: Dropping of nuisance columns in DataFrame reductions (with 'numeric_only=None') is deprecated; in a future version this will raise TypeError. Select only valid columns before calling the reduction.

```
data.mean()
```

```
Out[10]: mpg             23.514573
cylinders             5.454774
displacement         193.425879
weight              2970.424623
acceleration         15.568090
model year           76.010050
origin                1.572864
dtype: float64
```

```
In [11]: # where are the other columns? Check data types
data.dtypes
```

```
Out[11]: mpg          float64
cylinders         int64
displacement      float64
horsepower        object
weight            float64
acceleration      float64
model year        int64
origin            int64
car name          object
dtype: object
```

Notice that many columns are of type object, which is not a number. Maybe this has to do with missing values? We know from peeking at the first rows that there are '?' values in there. Let's replace these with the string NaN for not-a-number.

```
In [12]: # replace '?' with 'NaN'
data = data.replace({'?': 'NaN'})
data.head()
```

```
Out[12]:
```

	mpg	cylinders	displacement	horsepower	weight	acceleration	model year	origin	car name
0	18.0	8	307.0	130.0	3504.0	12.0	70	1	chevrolet chevelle malibu
1	15.0	8	350.0	165.0	3693.0	11.5	70	1	buick skylark 320
2	18.0	8	318.0	150.0	3436.0	11.0	70	1	plymouth satellite
3	16.0	8	304.0	150.0	3433.0	12.0	70	1	amc rebel sst
4	17.0	8	302.0	140.0	3449.0	10.5	70	1	ford torino

Pandas knows that 'NaN' probably means that numbers are missing. Now we can convert the data type from object to float

```
In [13]: # convert dtypes (Only horsepower seems appropriate)
data = data.astype({'horsepower': float})
data.dtypes
```

```
Out[13]: mpg                float64
cylinders                int64
displacement            float64
horsepower              float64
weight                  float64
acceleration            float64
model year              int64
origin                  int64
car name                 object
dtype: object
```

We could have loaded the data with the `na_values` argument to indicate that '?' means missing number:

```
In [14]: data=pd.read_csv('auto-mpg.data', delim_whitespace=True, header=None, names = data.dtypes)
```

```
Out[14]: mpg                float64
cylinders                int64
displacement            float64
horsepower              float64
weight                  float64
acceleration            float64
model year              int64
origin                  int64
car name                 object
dtype: object
```

This worked nicely. Now we can describe all columns, meaning printing basic statistics. Note that by default Pandas ignores NaN, whereas Numpy does not.

```
In [15]: data.describe() # ignores NaN
```

```
Out[15]:
```

	mpg	cylinders	displacement	horsepower	weight	acceleration	model year
count	398.000000	398.000000	398.000000	392.000000	398.000000	398.000000	398.000000
mean	23.514573	5.454774	193.425879	104.469388	2970.424623	15.568090	76.010050
std	7.815984	1.701004	104.269838	38.491160	846.841774	2.757689	3.697627
min	9.000000	3.000000	68.000000	46.000000	1613.000000	8.000000	70.000000
25%	17.500000	4.000000	104.250000	75.000000	2223.750000	13.825000	73.000000
50%	23.000000	4.000000	148.500000	93.500000	2803.500000	15.500000	76.000000
75%	29.000000	8.000000	262.000000	126.000000	3608.000000	17.175000	79.000000
max	46.600000	8.000000	455.000000	230.000000	5140.000000	24.800000	82.000000

We could be interested by these statistics in each of the genders. To get these, we first group values by gender, then ask for the description. We will only look at age for clarity

```
In [21]: #Let's try with the mpg for each origin
data.groupby(by='origin').describe().mpg
```

Out[21]:

	count	mean	std	min	25%	50%	75%	max
origin								
1	249.0	20.083534	6.402892	9.0	15.0	18.5	24.00	39.0
2	70.0	27.891429	6.723930	16.2	24.0	26.5	30.65	44.3
3	79.0	30.450633	6.090048	18.0	25.7	31.6	34.05	46.6

Find NaNs

How many NaNs in each column?

We can ask which entries are null, which produces a boolean array

```
In [23]: data.isnull()
```

Out[23]:

	mpg	cylinders	displacement	horsepower	weight	acceleration	model year	origin	car name
0	False	False	False	False	False	False	False	False	False
1	False	False	False	False	False	False	False	False	False
2	False	False	False	False	False	False	False	False	False
3	False	False	False	False	False	False	False	False	False
4	False	False	False	False	False	False	False	False	False
...
393	False	False	False	False	False	False	False	False	False
394	False	False	False	False	False	False	False	False	False
395	False	False	False	False	False	False	False	False	False
396	False	False	False	False	False	False	False	False	False
397	False	False	False	False	False	False	False	False	False

398 rows × 9 columns

Applying `sum()` to this boolean array will count the number of `True` values in each column

```
In [24]: data.isnull().sum()
```

```
Out[24]: mpg                0
cylinders                0
displacement             0
horsepower               6
weight                  0
acceleration             0
model year              0
origin                  0
car name                 0
dtype: int64
```

We get complementary information from `info()`

```
In [25]: data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 398 entries, 0 to 397
Data columns (total 9 columns):
#   Column          Non-Null Count  Dtype
---  -
0   mpg              398 non-null    float64
1   cylinders        398 non-null    int64
2   displacement     398 non-null    float64
3   horsepower       392 non-null    float64
4   weight           398 non-null    float64
5   acceleration     398 non-null    float64
6   model year      398 non-null    int64
7   origin           398 non-null    int64
8   car name         398 non-null    object
dtypes: float64(5), int64(3), object(1)
memory usage: 28.1+ KB
```

We can fill (replace) these missing values, for example with the minimum value in each column

```
In [26]: data.fillna(data.min()).describe()
```

```
Out[26]:
```

	mpg	cylinders	displacement	horsepower	weight	acceleration	model year
count	398.000000	398.000000	398.000000	398.000000	398.000000	398.000000	398.000000
mean	23.514573	5.454774	193.425879	103.587940	2970.424623	15.568090	76.010050
std	7.815984	1.701004	104.269838	38.859575	846.841774	2.757689	3.697627
min	9.000000	3.000000	68.000000	46.000000	1613.000000	8.000000	70.000000
25%	17.500000	4.000000	104.250000	75.000000	2223.750000	13.825000	73.000000
50%	23.000000	4.000000	148.500000	92.000000	2803.500000	15.500000	76.000000
75%	29.000000	8.000000	262.000000	125.000000	3608.000000	17.175000	79.000000
max	46.600000	8.000000	455.000000	230.000000	5140.000000	24.800000	82.000000

Count unique values (a histogram)

We finish off, with our good friend the histogram

```
In [27]: data['mpg'].value_counts()
```

```
Out[27]: 13.0    20
          14.0    19
          18.0    17
          15.0    16
          26.0    14
          ..
          31.9     1
          16.9     1
          18.2     1
          22.3     1
          44.0     1
          Name: mpg, Length: 129, dtype: int64
```

```
In [ ]:
```