### **Module IV**

#### **Pandas**

- Pandas is a Python library for data analysis.
- Started by Wes McKinney in 2008 out of a need for a powerful and flexible quantitative analysis tool,
- Pandas is a high-level data manipulation tool
- It is built on the **Numpy package** and its key **data structure** is called the **DataFrame**.
- **DataFrames** allow you to store and manipulate tabular data in rows of observations and columns of variables.
- pandas has grown into one of the most popular Python libraries.
- It has an extremely active community of contributors.

There are several ways to create a DataFrame. One way way is to use a dictionary. For example:

#### 11. Create a very first Pandas dataframe (fromCSV)

#### Read Dataset:

Import pandas as pd pd.read\_csv("titanic.csv")

### Storing dataset in a Variable:

```
Import pandas as pd
titanic = pd.read_csv("titanic.csv")
titanic
```

### 12. Pandas display option and the methods head() & tail()

```
titanic
or
print(titanic)
<Output : first and last 5 rows>
```

### To see the max & min row setting:

```
pd.options.display.max_rows
<output : 60>
pd.options.display.min_rows
<output : 10>
```

### To change the min row setting:

```
pd.options.display.min_rows = 20

titanic

<output : display 20 rows if the number of rows are greater than max rows count>
```

### To change the max row setting:

```
pd.options.display.min_rows = 891
titanic
<output : display all 891 rows>
```

### Method head() & tail()

```
titanic.head()
<output : display first five rows of the dataset>
To see first 10 rows:
titanic.head(10)
<output : display first 10 rows of the dataset>
```

titanic.tail()

<output : display last five rows of the dataset>

To see last 10 rows:

titanic.tail(10)

<output : display last 10 rows of the dataset>

#### 13. First Data Inspection

#### Info Method:

### titanic.info()

```
titanic.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 9 columns):

# Column Non-Null Count Dtype
                           891 non-null
         survived
         pclass
sex
                           891 non-null
891 non-null
                                                          int64
                                                          object
         age
sibsp
                            714 non-null
891 non-null
                                                          float64
int64
         parch
fare
                           891 non-null
891 non-null
                                                          int64
float64
7 embarked 889 non-null object
8 deck 203 non-null object
dtypes: float64(2), int64(4), object(3)
memory usage: 62.8+ KB
```

#### **Describe Method:**

### titanic.describe()



#### 14. Built in Functions & Data Frame Attributes and Methods

### **Data Frame & Python Built-in Function**

```
In [4]: type(titanic)
Out[4]: pandas.core.frame.DataFrame
In [5]: len(titanic)
Out[5]: 891
In [6]:
         round(titanic, 0)
Out[6]:
               survived pclass
                                          sibsp parch fare embarked deck
                                 sex
                                      age
            0
                                male
                                      22.0
                                                         7.0
                                                                        NaN
                                                                          C
            1
                            1 female
                                      38.0
                                                     0 71.0
                                                                    C
                     1
            2
                     1
                                      26.0
                                               0
                                                     0
                                                         8.0
                            3 female
                                                                     S
                                                                        NaN
            3
                     1
                               female
                                      35.0
                                                        53.0
                                                                     S
                                                                          C
                     0
                            3
                                male
                                      35.0
                                               0
                                                     0
                                                         8.0
                                                                    S
                                                                        NaN
```

#### **Data Attributes**

### **DataFrame Methods**

13]:	sı	ırvived	pclass	sex	age	sibsp	parch	fare	embark	ed	decl
	0	0	3	male	22.0	1	0	7.2500		s	NaN
	1	1	1	female	38.0	1	0	71.2833		C	C
	2	1	3	female	26.0	0	0	7.9250	0	s	NaN
	3	1	1	female	35.0	1	0	53.1000	,	s	(
	4	0	3	male	35.0	0	0	8.0500		S	Nat
14]:	titan	ic.hea	d(n = 2	2)							
14]:	SII	rvived	nclass	sex	age	sibsp	narch	fare	embarked	l d	leck
			polass	301		зірэр	229232				
	0	0	3	male	22.0	1	0	7.2500	5	1	NaN
15]:	1 tit	1 anic.	1 info(	female	38.0	1	0	71.2833	(		C
15]:	tit <cl< td=""><td>anic.</td><td>info( panda ex: 8</td><td>female ) s.core 91 ent</td><td>38.0</td><td>1 ame.Da</td><td>o ataFra</td><td>71.2833 ame'&gt;</td><td></td><td></td><td></td></cl<>	anic.	info( panda ex: 8	female ) s.core 91 ent	38.0	1 ame.Da	o ataFra	71.2833 ame'>			
15]:	tit <cl< td=""><td>anic. ass ' geInc</td><td>1 info( panda</td><td>female s.core 91 ent (total</td><td>38.0 e.fra</td><td>1 ame.Da</td><td>o taFra to 890</td><td>71.2833 ame'&gt;</td><td></td><td></td><td></td></cl<>	anic. ass ' geInc	1 info( panda	female s.core 91 ent (total	38.0 e.fra	1 ame.Da	o taFra to 890	71.2833 ame'>			
15]:	tit <cl Ran Dat</cl 	anic. ass ' geInd a col Col	info( panda: ex: 89 umns umn	s.core 91 ent (total Non-	38.0 e.fracries l 9 c	ame.Da s, 0 t column	o taFra to 890 ns):	71.2833 ame'>			
15]:	tit <cl Ran Dat #</cl 	anic. ass ' geInc a col Col sur	info( panda: ex: 80 umns umn vived	s.core 91 ent (total Non-	38.0 e.fracries L 9 c Null	ame.Da s, 0 t column Cour	o taFra to 890 is): it Dt	71.2833 ame'> b type nt64			
15]:	tit <cl Ran Dat #</cl 	anic. ass ' geInc a col Col sur pcl	info( panda: lex: 89 umns umn  vived	s.core 91 ent (total Non- 891 891	38.0  e.fra cries 1 9 c -Null non-	ame.Da	o taFra to 890 ns): nt Dt ir	71.2833 ame'> b type nt64 nt64			
15]:	tit <cl Rand Dat # 1 2</cl 	anic. ass ' geInc a col Col sur pcl	info( panda: ex: 80 umns umn vived ass	s.core 91 ent (total Non- 891 891	38.0 e.fra cries l 9 c Null non- non-	ame.Da	o taFra to 890 is): it Dt ir ot	71.2833  ame'>  cype nt64 nt64 oject			
15]:	tit <cl Ran Dat # 1 2 3</cl 	anic. ass ' geInc a col Col sur pcl sex age	info( panda ex: 89 umns umn vived ass	s.core 91 ent (total Non- 891 891 891	38.0 e.fra cries l 9 c -Null non- non-	ame.Da	o taFra to 890 ns): nt Dt ir ob	71.2833  ame'>  type nt64 nt64 oject Loat64			
15]:	tit <cl Rand Dat # 1 2</cl 	anic. ass ' geInc a col Col sur pcl sex age	info( panda: ex: 89 umns umn vived ass	s.core 91 ent (total Non- 891 891 891 714	38.0 e.fra cries l 9 c Null non- non- non- non-	ame.Da	o taFra to 896 ns): nt Dt ir ob	71.2833  ame'>  cype nt64 nt64 oject			
15]:	tit <cl Ran Dat #  0 1 2 3 4 5</cl 	anic. ass ' geInc a col Col sur pcl sex age	info( panda: ex: 80 umns umn vived ass	s.core 91 ent (total Non- 891 891 891 714 891	38.0 e.fracries l 9 c Null non- non- non- non- non- non-	nme.Da	o taFra to 896 ns): nt Dt ir ob ir	71.2833  ame'>  cype  nt64 nt64 cject Loat64			
15]:	1 tit <cl Ran Dat # </cl 	anic. ass ' geInd a col col sur pcl sex age sib par far	info( panda: ex: 80 umns umn vived ass esp ech ee	s.core 91 ent (total Non- 891 891 891 891 891 891	38.0 e.fra cries l 9 c Null non- non- non- non- non- non- non- no	nme.Da	o intariant of int	71.2833  ame'>  type  164  164  164  164  164  164			

```
In [16]: titanic.min()

Out[16]: survived 0 pclass 1 sex female age 0.42 sibsp 0 parch 0 fare 0 dtype: object
```

### **Method Chaining**

```
In [17]: titanic.mean()
Out[17]: survived
                       0.383838
          pclass
                        2.308642
                      29.699118
          age
          sibsp
                       0.523008
          parch
                       0.381594
          fare
                      32.204208
          dtype: float64
In [18]: titanic.mean().port_values()
Out[18]: parch
                       0.381594
         survived
                       0.383838
         sibsp
                       0.523008
```

In [19]: titanic.mean().sort\_values().head(2)

Out[19]: parch 0.381594

pclass

age

fare

survived 0.383838

dtype: float64

dtype: float64

### Resources:

https://pandas.pydata.org/pandas-

docs/stable/reference/api/pandas.DataFrame.html

2.308642

29.699118

32.204208

https://pandas.pydata.org/pandas-

docs/stable/reference/api/pandas.Series.html

### 15. Make it easy: TAB Completion and Tooltip

```
In [2]: import pandas as pd

In []: pd.read

pd.read_clipboard

In []: pd.read_csy
pd.read_excel
pd.read_feather
In []: pd.read_fwf
pd.read_gbq
pd.read_hdf
pd.read_html
pd.read_json
In []: pd.read_msgpack

In []:
```

```
In [2]: import pandas as pd

In []: pd.read_c|
pd.read_clipboard
In []: pd.read_csy
```

```
In [2]: import pandas as pd

In []: pd.read_csv("tit")

titanic.csv

titanic_clean.csv
titanic_complete.csv
titanic_dummy_csv
titanic_imp.csv
titanic_prep.csv
titanic_prep.csv
titanic_raw.csv
titanic_slice.xls
```

```
In [5]: titanic.sort_values()
Out[5]:
          Signature: titanic.sort_values(by, axis=0, ascending=True, inplace=False, kind='quickso
          rt', na_position='last')
          Docstring:
          Sort by the values along either axis
          Parameters
                  by : str or list of str
                     Name or list of names to sort by.
                                                  1 21.0750
                                                                    S NaN
                                   2.0
                                                  2 11.1333
                                                                    s
                                                                      NaN
                          3 female 27.0
                                                     30.0708
                                                                      NaN
          10
                                                     16.7000
                                   4.0
                                                     26.5500
                                                                         C
```

```
In [2]: import pandas as pd

In [4]: | titanic = pd.read_csv("titanic.csv")

In [ ]: tit

titanic

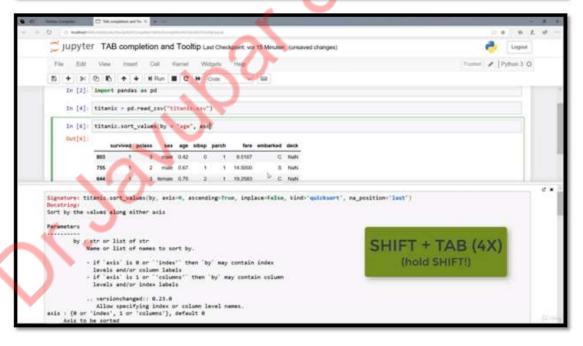
In [ ]: titanic_csv
titanic_clean.csv
titanic_complete.csv
In [ ]: titanic_dummy_csv
titanic_imp.csv
titanic_imp.csv
In [ ]: titanic_prep.csv
```

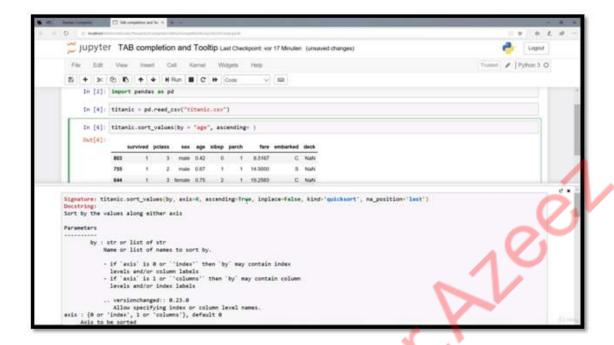
```
In [2]: import pandas as pd

In [4]: titanic = pd.read_csv("titanic.csv")

In [5]: titanic.sort|
Out[5]: titanic.sort_index
titanic.sort_values
sex age sibsp page
```

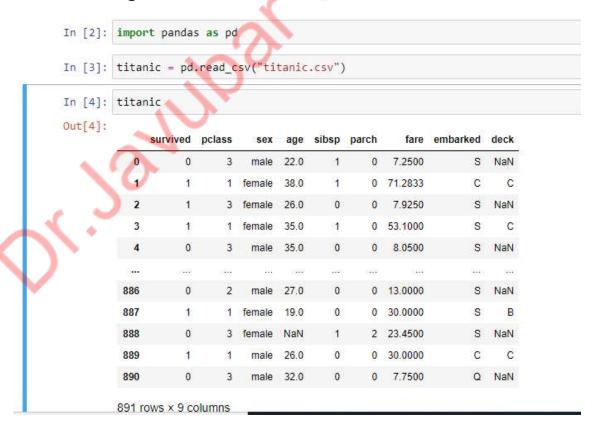
```
In [2]:
         import pandas as pd
In [4]:
         titanic = pd.read_csv("titanic.csv")
In [6]:
         titanic.sort_values(by = "age")
Out[6]:
              survived pclass
                                sex age sibsp parch
                                                           fare embarked deck
          803
                                male 0.42
                                                          8.5167
                                                                        C NaN
          755
                                male 0.67
                                              1
                                                         14.5000
                                                                           NaN
                            3 female 0.75
                                                         19.2583
                                                                           NaN
          469
                     1
                            3 female 0.75
                                                        19.2583
                                                                           NaN
           78
                                male 0.83
                                                         29,0000
                                                                           NaN
          831
                                male 0.83
                                              1
                                                        18.7500
                                                                           NaN
          305
                                male 0.92
                                                     2 151.5500
                                                                        S
                                                                             C
          827
                     1
                            2
                                male 1.00
                                                         37,0042
                                                                        C NaN
                            3 female 1.00
                                                        15.7417
                                                                           NaN
          164
                                male 1.00
                                               4
                                                         39.6875
                                                                           NaN
          183
                                male 1.00
                                                         39.0000
                                                                        S
                                                                              F
```





# 18. Selecting Column

### Selecting all the columns of a dataset



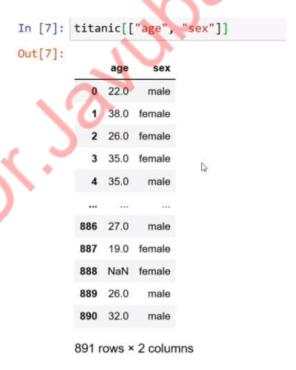
### Selection of a single column "age "

```
In [4]: titanic["age"]
Out[4]: 0
                22.0
                38.0
         2
                26.0
         3
                35.0
                35.0
         886
                27.0
         887
                19.0
         888
                 NaN
         889
                26.0
         890
                32.0
         Name: age, Length: 891, dtype: float64
```

# Finding the pandas series

```
In [5]: type(titanic["age"])
Out[5]: pandas.core.series.Series
```

# Selecting two or more column, it require to use list



### Checking the type of the dataframe

```
In [8]: type(titanic[["age", "sex"]])
Out[8]: pandas.core.frame.DataFrame
```

We can also change the column sequence

```
In [7]: titanic[["sex", "age"]]
Out[7]:
                age
                        sex
             0 22.0
                       male
                38.0
                     female
                26.0
                     female
                35.0
                     female
                35.0
                       male
           886
                27.0
                       male
           887
                19.0
                     female
           888
                NaN
                     female
           889
                26.0
                       male
           890 32.0
                       male
```

# 19. SELECTING ONE COLUMN WITH DOT NOTATION

### Selecting one Column with "dot notation"

# We can check two panda series by using equal function

# Selecting one Column with "dot notation"

In [6]: titanic.age.equals(titanic["age"])
Out[6]: True

# 20. Zero Based Indexing and Negative Indexing

# Zero-based Indexing

### column index positions

		0	1	2	3	4
row index positions		Nationality	Club V	World_Champion	Height	Goals_2018
Tow index positions	Player					
0	Lionel Messi	Argentina	FC Barcelona	False	1.70	45
1	Cristiano Ronaldo	Portugal	Juventus FC	False	1.87	44
2	Neymar Junior	Brasil	Paris SG	False	1.75	28
3	Kylian Mbappe	France	Paris SG	True	1.78	21
4	Manuel Neuer	Germany	FC Bayern	True	1.93	0

# Negative Indexing

### column index positions

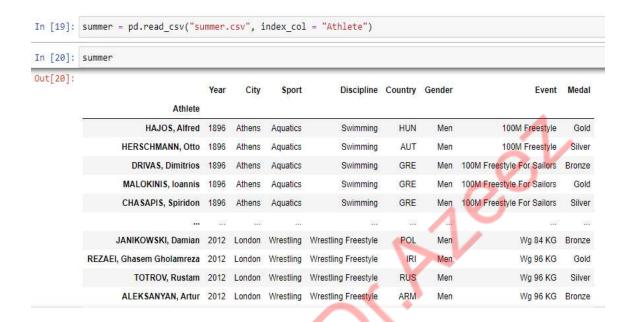
			4	-5	-4	-3	-2	-1
		Y		0	1	2	3	4
row index po	sition	S		Nationality	Club	World_Champion	Height	Goals_2018
			Player					
$\sim$ 1 $\sim$	-5	0	Lionel Messi	Argentina	FC Barcelona	False	1.70	45
	-4	1	Cristiano Ronaldo	Portugal	Juventus FC	False	1.87	44
	-3	2	Neymar Junior	Brasil	Paris SG	False	1.75	28
	-2	3	Kylian Mbappe	France	Paris SG	True	1.78	21
	-1	4	Manuel Neuer	Germany	FC Bayern	True	1.93	0

### 21. Selecting rows with iloc (position based indexing)

# Position-based Indexing and Slicing with iloc[]

```
import pandas as pd
  In [1]:
               summer = pd.read_csv("summer.csv")
               summer
  In [ ]:
In [14]: summer
Out[14]:
                                                                          Gender
               Year
                     City
                            Sport
                                      Discipline
                                                             Athlete Country
                                                                                               Event
                                                                                                    Medal
            0 1896
                                                         HAJOS, Alfred
                                                                      HUN
                   Athens
                          Aquatics
                                       Swimming
                                                                             Men
                                                                                         100M Freestyle
                                                                                                      Gold
            1 1896
                                       Swimming
                                                     HERSCHMANN, Otto
                                                                      AUT
                   Athens
                          Aquatics
                                                                             Men
                                                                                         100M Freestyle
                                                       DRIVAS, Dimitrios
                                                     MALOKINIS, Ioannis
            3 1896 Athens
                                       Swimming
                                                                      GRE
                                                                                 100M Freestyle For Sailors
                                                     CHASAPIS, Spiridon
            4 1896
                   Athens
                          Aquatics
                                       Swimming
                                                                      GRE
                                                                                 100M Freestyle For Sailors
                                                    JANIKOWSKI, Damian
                                                                      POL
                                                                                            Wg 84 KG Bronze
         31160 2012 London Wrestling Wrestling Freestyle
         31161 2012 London Wrestling Wrestling Freestyle REZAEI, Ghasem Gholamreza
                                                                       IRI
                                                                                            Wg 96 KG
         31162 2012 London Wrestling Wrestling Freestyle
                                                       TOTROV, Rustam
                                                                                            Wg 96 KG
                                                                      RUS
                                                                             Men
                                                                                                     Silver
         31163 2012 London Wrestling Wrestling Freestyle
                                                     ALEKSANYAN, Artur
                                                                      ARM
                                                                                            Wg 96 KG Bronze
                                                                             Men
         31164 2012 London Wrestling Wrestling Freestyle
                                                       LIDBERG, Jimmy
                                                                                            Wg 96 KG Bronze
        31165 rows × 9 columns
 In [15]: summer.info()
             <class 'pandas.core.frame.DataFrame'>
              RangeIndex: 31165 entries, 0 to 31164
             Data columns (total 9 columns):
                                    Non-Null Count Dtype
                                    31165 non-null
                                                         int64
                    Year
                    City
                                    31165 non-null object
               2
                    Sport
                                    31165 non-null object
                    Discipline 31165 non-null object
               4
                    Athlete
                                    31165 non-null object
               5
                    Country
                                    31161 non-null object
               6
                                    31165 non-null object
                    Gender
                    Event
                                    31165 non-null
                                                         object
               8
                    Medal
                                    31165 non-null object
             dtypes: int64(1), object(8)
             memory usage: 2.1+ MB
```

### Changing the Index to "Athlete" Column.



### Selecting row using iloc

```
In [30]:
         summer.iloc[0]
Out[30]: Year
                                   1896
          City
                                 Athens
          Sport
                               Aquatics
          Discipline
                               Swimming
          Country
                                    HUN
          Gender
                                    Men
          Event
                        100M Freestyle
         Name: HAJOS, Alfred, dtype: object
```

### Selecting last row using negative indexing in iloc

```
In [31]:
         summer.iloc[-1]
Out[31]:
         Year
                                        2012
          City
                                      London
          Sport
                                   Wrestling
         Discipline
                        Wrestling Freestyle
         Country
         Gender
                                         Men
         Event
                                    Wg 96 KG
         Medal
                                      Bronze
         Name: LIDBERG, Jimmy, dtype: object
```

# Selecting first three rows using iloc by using list

[33]:	<pre>summer.iloc[[1,2,</pre>	3]]							
[33]:		Year	City	Sport	Discipline	Country	Gender	Event	Medal
	Athlete								
	HERSCHMANN, Otto	1896	Athens	Aquatics	Swimming	AUT	Men	100M Freestyle	Silver
	DRIVAS, Dimitrios	1896	Athens	Aquatics	Swimming	GRE	Men	100M Freestyle For Sailors	Bronze
	MALOKINIS, Ioannis	1896	Athens	Aquatics	Swimming	GRE	Men	100M Freestyle For Sailors	Gold

# Selecting first three rows using slicing in iloc

In [36]:	summer.iloc[1 : 4	]						Y	
Out[36]:		Year	City	Sport	Discipline	Country	Gender	Event	Medal
	Athlete						1,		
	HERSCHMANN, Otto	1896	Athens	Aquatics	Swimming	AUT	Men	100M Freestyle	Silver
	DRIVAS, Dimitrios	1896	Athens	Aquatics	Swimming	GRE	Men	100M Freestyle For Sailors	Bronze
	MALOKINIS, Ioannis	1896	Athens	Aquatics	Swimming	GRE	Men	100M Freestyle For Sailors	Gold

# Using slicing operation, selecting first three rows

In [37]:	<pre>summer.iloc[:3]</pre>			7	•				
Out[37]:		Year	City	Sport	Discipline	Country	Gender	Event	Medal
	Athlete		1						
	HAJOS, Alfred	1896	Athens	Aquatics	Swimming	HUN	Men	100M Freestyle	Gold
	HERSCHMANN, Otto	1896	Athens	Aquatics	Swimming	AUT	Men	100M Freestyle	Silver
	DRIVAS, Dimitrios	1896	Athens	Aquatics	Swimming	GRE	Men	100M Freestyle For Sailors	Bronze

# Using slicing operation, selecting last three rows

In [38]:	summer.iloc[-3:]								
Out[38]:		Year	City	Sport	Discipline	Country	Gender	Event	Medal
<b>Y</b>	Athlete								
	TOTROV, Rustam	2012	London	Wrestling	Wrestling Freestyle	RUS	Men	Wg 96 KG	Silver
	ALEKSANYAN, Artur	2012	London	Wrestling	Wrestling Freestyle	ARM	Men	Wg 96 KG	Bronze
	LIDBERG, Jimmy	2012	London	Wrestling	Wrestling Freestyle	SWE	Men	Wg 96 KG	Bronze

# 22. Slicing rows and columns with iloc(position-based indexing)

# Using head & tail to select first and last set of rows

[41]:	summer.head(10)								
[41]:		Year	City	Sport	Discipline	Country	Gender	Event	Medal
	Athlete								
	HAJOS, Alfred	1896	Athens	Aquatics	Swimming	HUN	Men	100M Freestyle	Gold
	HERSCHMANN, Otto	1896	Athens	Aquatics	Swimming	AUT	Men	100M Freestyle	Silver
	DRIVAS, Dimitrios	1896	Athens	Aquatics	Swimming	GRE	Men	100M Freestyle For Sailors	Bronze
	MALOKINIS, Ioannis	1896	Athens	Aquatics	Swimming	GRE	Men	100M Freestyle For Sailors	Gold
	CHASAPIS, Spiridon	1896	Athens	Aquatics	Swimming	GRE	Men	100M Freestyle For Sailors	Silver
	CHOROPHAS, Efstathios	1896	Athens	Aquatics	Swimming	GRE	Men	1200M Freestyle	Bronze
	HAJOS, Alfred	1896	Athens	Aquatics	Swimming	HUN	Men	1200M Freestyle	Gold
	ANDREOU, Joannis	1896	Athens	Aquatics	Swimming	GRE	Men	1200M Freestyle	Silver
	CHOROPHAS, Efstathios	1896	Athens	Aquatics	Swimming	GRE	Men	400M Freestyle	Bronze
	NEUMANN, Paul	1896	Athens	Aquatics	Swimming	AUT	Men	400M Freestyle	Gold

# To select a single value from a dataset using iloc.

ıt[41]:		Year	City	Sport	Discipline	Country	Gender	Event	Medal
	Athlete								
	HAJOS, Alfred	1896	Athens	Aquatics	Swimming	HUN	Men	100M Freestyle	Gold
	HERSCHMANN, Otto	1896	Athens	Aquatics	Swimming	AUT	Men	100M Freestyle	Silver
	DRIVAS, Dimitrios	1896	Athens	Aquatics	Swimming	GRE	Men	100M Freestyle For Sailors	Bronze
	MALOKINIS, Ioannis	1896	Athens	Aquatics	Swimming	GRE	Men	100M Freestyle For Sailors	Gold
	CHASAPIS, Spiridon	1896	Athens	Aquatics	Swimming	GRE	Men	100M Freestyle For Sailors	Silver
	CHOROPHAS, Efstathios	1896	Athens	Aquatics	Swimming	GRE	Men	1200M Freestyle	Bronze
	HAJOS, Alfred	1896	Athens	Aquatics	Swimming	HUN	Men	1200M Freestyle	Gold
	ANDREOU, Joannis	1896	Athens	Aquatics	Swimming	GRE	Men	1200M Freestyle	Silver
	CHOROPHAS, Efstathios	1896	Athens	Aquatics	Swimming	GRE	Men	400M Freestyle	Bronze
1	NEUMANN, Paul	1896	Athens	Aquatics	Swimming	AUT	Men	400M Freestyle	Gold

In [42]: summer.iloc[0, 4]

Out[42]: 'HUN'

In [43]: summer.iloc[0,:3]

Out[43]: Year 1896

City Athens Sport Aquatics

Name: HAJOS, Alfred, dtype: object

```
In [44]: summer.iloc[0,[0,2,5,7]]
Out[44]: Year
                         1896
          Sport
                     Aquatics
          Gender
                          Men
          Medal
                         Gold
          Name: HAJOS, Alfred, dtype: object
In [45]: summer.iloc[34:39,[0,2,5,7]]
Out[45]:
                                Year
                                       Sport Gender
                                                     Medal
                        Athlete
                GARRETT, Robert 1896 Athletics
                                                Men
                                                      Silver
                KELLNER, Gyula 1896 Athletics
                                                Men
                                                     Bronze
                 LOUIS, Spyridon 1896 Athletics
                                                Men
                                                       Gold
            VASILAKOS, Kharilaos 1896 Athletics
                                                Men
                                                      Silver
           DAMASKOS, Evangelos 1896 Athletics
                                                Men
                                                     Bronze
In [47]: summer.iloc[:, 4]
Out[47]: Athlete
                                         HUN
          HAJOS, Alfred
          HERSCHMANN, Otto
                                         AUT
          DRIVAS, Dimitrios
                                         GRE
          MALOKINIS, Ioannis
                                         GRE
          CHASAPIS, Spiridon
                                         GRE
          JANIKOWSKI, Damian
                                         POL
          REZAEI, Ghasem Gholamreza
                                         IRI
          TOTROV, Rustam
                                         RUS
          ALEKSANYAN, Artur
                                         ARM
          LIDBERG, Jimmy
                                         SWE
          Name: Country, Length: 31165, dtype: object
```

### Compare two selecting row outputs by using equals function

```
In [48]: summer.iloc[:, 4].equals(summer.Country)
Out[48]: True
```

# 23. Position based Indexing Cheat Sheets

### Position-based Indexing (iloc)



Zero-based indexing applies!

### Position-based Indexing (iloc) - Example 1



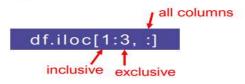
Output is an element ("Paris SG").

# Position-based Indexing (iloc) - Example 2



Output is a DataFrame.

# Position-based Indexing (iloc) - Example 3



		0	1	2	3	4
		Nationality	Club	World_Champion	Height	Goals_2018
	Player					
1	Cristiano Ronaldo	Portugal	Juventus FC	False	1.87	44
2	Neymar Junior	Brasil	Paris SG	False	1.75	28

Output is a DataFrame.

Position-based Indexing (iloc) - Example 4



Output is a Pandas Series.

# 24. LABEL BASED INDEXING AND SLICING WITH LOC[]

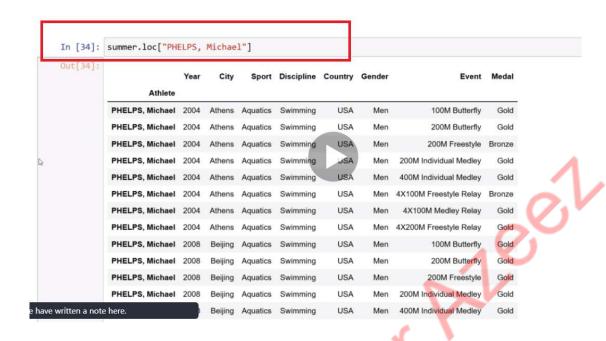
Out[31]: Discipline Country Gender Event Medal Athlete HAJOS, Alfred 1896 Athens Aquatics Swimming HUN Men 100M Freestyle Gold HERSCHMANN, Otto AUT GRE MALOKINIS, Ioannis Athens Aquatics Swimming GRE Men 100M Freestyle For Sailors Gold CHASAPIS, Spiridon 1896 Athens GRE 100M Freestyle For Sailors Silver JANIKOWSKI, Damian 2012 London Wrestling Wrestling Freestyle POL Men Wg 84 KG Bronze REZAEI, Ghasem Gholamreza 2012 London Wrestling Wrestling Freestyle IRI Wa 96 KG Men RUS Wg 96 KG TOTROV, Rustam 2012 London Wrestling Wrestling Freestyle ALEKSANYAN, Artur 2012 London Wrestling Wrestling Freestyle ARM Wg 96 KG Bronze LIDBERG, Jimmy 2012 London Wrestling Wrestling Freestyle SWE Wg 96 KG Bronze

31165 rows × 8 columns

# Selecting Rows with loc[ In [32]: summer.iloc[2]

Out[32]: Year 1896
City Athens
Sport Aquatics
Discipline Swimming
Country GRE
Gender Men
Event 100M Freestyle For Sailors
Medal Name: DRIVAS, Dimitrios, dtype: object

In [ ]: summer.loc[T'DRIVAS, Dimitrios"]



# 26. Indexing and Slicing with reindex()



```
In [6]: summer.loc[[0,5,30000],["Athlete","Medal"]]

Out[6]:

Athlete Medal

0 HAJOS, Alfred Gold

5 CHOROPHAS, Efstathios Bronze

30000 PAUTARAN, Maryna Bronze
```

### Reindex

# The summer table have rows upto 31165



\_\_\_\_\_

# Importing from CSV and first inspection

```
import pandas as pd
summer = pd.read_csv("summer.csv")
summer
summer.info()
summer.describe()
```

### Select one column

summer.Medal summer["Medal"]

# Select multiple Columns

summer["Year","Medal"]
summer.loc[:,["Year","Medal"]]

### Select positional rows

summer.iloc[10:21]

# **Select positional rows**

# Section 4 - Pandas Series and Index Objects

Welcome to the Section on Pandas Series and Pandas Index objects!

In this Section, you will learn how to work with and analyze single columns of a DataFrame (Pandas Series).

We will differentiate between numerical columns (integers or floats) and non-numerical columns (strings).

### Analyzing Columns / Pandas Series

non-n	umerical c	olumn	nume	rical co	lumn
	Nationality	Club	World_Champion	Height	Goals_2018
Player					
Lionel Messi	Argentina	FC Barcelona	False	1.70	45
Christiano Ronaldo	Portugal	Juventus FC	False	1.87	44
Neymar Junior	Brasil	Paris SG	False	1.75	28
Kylian Mbappe	France	Paris SG	True	1.78	21
Manuel Neuer	Germany	FC Bayern	True	1.93	0

Next, you will learn how to work with and analyze Pandas Index objects and how to change single or multiple labels.

We will differentiate between (Row) Index and Column Index.

# Analyzing and modifying Pandas Index objects



### PANDAS SERIES

### FIRST STEP WITH PANDAS SERIES

```
In [7]: import pandas as pd
In [8]: titanic = pd.read_csv("titanic.csv")
In [ ]: titanic
In [ ]: titanic.info()
In [ ]: titanic["age"]
In [ ]: type(titanic["age"])
In [ ]: titanic["age"].equals (titanic.age)
In [13]: age = titanic["age"]
In [ ]: age.head()
In [ ]: age.tail()
In [ ]: age.dtype
In [ ]: age.shape
In [ ]: len(age)
In [ ]: age.index
In [ ]: age.info()
In [ ]: age.to_frame().info()
```

# Analyzing Numerical Series with unique(), nunique() and value\_counts()

```
In [ ]: import pandas as pd
In [ ]: titanic = pd.read_csv("titanic.csv")
In [ ]: titanic.csv
In [ ]: age = titanic["age"]
In [ ]: age
In [ ]: age.describe()
In [ ]: age.count()
In [ ]: age.size
In [ ]: len(age)
In [ ]: age.sum()
In [ ]: sum(age)
In [ ]: age.mean()
In [ ]: age.median()
In [ ]: age.std()
In [ ]: age.min()
In [ ]: age.max()
In [ ]: age.unique()
In [ ]: len(age.unique())
In [ ]: age.nunique()
In [ ]: age.value_counts()
```

# Analyzing Non-Numerical Series with unique(), nunique() and value\_counts()

```
In [ ]: import pandas as pd
In [ ]: summer = pd.read_csv("summer.csv")
In [ ]: summer.head()
In [ ]: summer.tail()
In [ ]: summer.info()
In [ ]: athlete = summer("athlete")
In [ ]: athlete.head()
In [ ]: athlete.tail(5)
In [ ]: tyep(athlete)
In [ ]: athlete.dtype
 In [ ]: athlete.shape
 In [ ]: athlete.describe()
 In [ ]: athlete.size
 In [ ]: athlete.count()
 In [ ]: athlete.min()
 In [ ]: athlete.unique()
 In [ ]: len(athlete.unique())
 In [ ]: athlete.nunique()
 In [ ]: athlete.value_counts()
```

#### **PANDAS SERIES**

# **Creating Pandas Series Part-1**

```
In [1]: import pandas as pd
 In [2]: summer = pd.read csv("summer.csv")
 In [3]: summer.head()
 Out[3]:
            Year
                                               Athlete Country Gender
                                                                                       Medal
                  City
                         Sport Discipline
                                                                                 Event
          0 1896 Athens Aquatics Swimming
                                           HAJOS, Alfred
                                                        HUN
                                                                           100M Freestyle
          1 1896 Athens Aquatics Swimming HERSCHMANN, Otto
                                                        AUT
                                                                           100M Freestyle
                                                                                        Silver
                                                               Men
          2 1896 Athens Aquatics Swimming
                                         DRIVAS, Dimitrios
                                                        GRE
                                                               Men 100M Freestyle For Sailors Bronze
          3 1896 Athens Aquatics Swimming
                                       MALOKINIS, Ioannis
                                                        GRE
                                                                    100M Freestyle For Sailors
                                       CHASAPIS, Spiridon
                                                        GRE
                                                                    100M Freestyle For Sailors
          4 1896 Athens Aquatics Swimming
                                                                                       Silver
In [6]: summer["Athlete"]
Out[6]: 0
                                   HAJOS, Alfred
                                HERSCHMANN, Otto
          1
          2
                               DRIVAS, Dimitrios
          3
                              MALOKINIS, Ioannis
          4
                              CHASAPIS, Spiridon
                              JANIKOWSKI, Damian
          31160
          31161
                     REZAEI, Ghasem Gholamreza
          31162
                                  TOTROV, Rustam
                               ALEKSANYAN, Artur
          31163
          31164
                                  LIDBERG, Jimmy
          Name: Athlete, Length: 31165, dtype: object
In [7]:
          summer.Athlete
Out[7]:
                                   HAJOS, Alfred
                                HERSCHMANN, Otto
                               DRIVAS, Dimitrios
          2
          3
                              MALOKINIS, Ioannis
          4
                              CHASAPIS, Spiridon
          31160
                              JANIKOWSKI, Damian
          31161
                     REZAEI, Ghasem Gholamreza
                                  TOTROV, Rustam
          31162
          31163
                               ALEKSANYAN, Artur
          31164
                                  LIDBERG, Jimmy
          Name: Athlete, Length: 31165, dtype: object
```

```
In [8]: summer.iloc[0]
Out[8]: Year
                                 1896
         City
                               Athens
         Sport
                             Aquatics
        Discipline
                             Swimming
        Athlete
                        HAJOS, Alfred
         Country
         Gender
         Event
                       100M Freestyle
        Medal
                                 Gold
        Name: 0, dtype: object
```

# Importing from CSV

```
In [9]: pd.read_csv("summer.csv", usecols=["Athlete"], squeeze = True)
Out[9]: 0
                              HAJOS, Alfred
                          HERSCHMANN, Otto
        1
        2
                         DRIVAS, Dimitrios
        3
                         MALOKINIS, Ioannis
        4
                         CHASAPIS, Spiridon
        31160
                         JANIKOWSKI, Damian
        31161
                 REZAEI, Ghasem Gholamreza
        31162
                             TOTROV, Rustam
        31163
                          ALEKSANYAN, Artur
        31164
                             LIDBERG, Jimmy
        Name: Athlete, Length: 31165, dtype: object
```

# Creating Scratch from pd.Series()

```
In [12]: pd.Series([10,25,6,36,2])
Out[12]: 0
              10
         1
              25
         2
               6
         3
              36
         4
         dtype: int64
In [14]: pd.Series([10,25,6,36,2], index = ["Mon","Tue", "Wed", "Thu","Fri"])
Out[14]: Mon
                 25
         Tue
                 6
         Wed
         Thu
                 36
         Fri
                 2
         dtype: int64
In [15]: pd.Series([10,25,6,36,2], index = ["Mon","Tue", "Wed", "Thu", "Fri"], name = "sales")
Out[15]: Mon
                10
         Tue
                 25
         Wed
                 6
         Thu
                36
         Fri
                 2
         Name: sales, dtype: int64
```

# **Creating Pandas Series Part-2**

# from Numpy Array

dtype: int64

```
In [19]: import pandas as pd
         import numpy as np
In [22]: sales = np.array([10,25,6,36,2])
         sales
Out[22]: array([10, 25, 6, 36, 2])
In [23]: pd.Series(sales)
Out[23]: 0
              25
               6
              36
         dtype: int32
          from List
 In [24]: sales = [10,25,6,36,2]
 In [25]: pd.Series(sales)
 Out[25]: 0
               10
               25
                6
               36
          dtype: int64
          from Dictionary
 In [27]: dic = {"Mon":10, "Tue":25, "Wed":6, "Thu":36, "Fri":2}
 In [29]: sales = pd.Series(dic)
 In [30]: sales
 Out[30]: Mon
                 10
                 25
          Tue
          Wed
                 6
          Thu
                 36
          Fri
                  2
```

```
In [31]: pd.Series(dic, index = ["Fri", "Sat", "Sun", "Mon", "Tue", "Wed"])
Out[31]: Fri
                2.0
        Sat
                NaN
                NaN
        Sun
               10.0
        Mon
               25.0
        Tue
        Wed
                6.0
        dtype: float64
In [32]: pd.Series(dic, index = [1,2,3,4,5])
Out[32]: 1
            NaN
            NaN
            NaN
            NaN
            NaN
        dtype: float64
        Indexing and Slicing Pandas Series
In [ ]: import pandas as pd
In [ ]: titanic = pd.read_csv("titanic.csv")
In [ ]: titanic.head()
In [ ]: titanic.tail()
In [ ]: age = titanic.age
In [ ]: age.head()
In [ ]: age.tail()
```

In [ ]: age.index

In [ ]: age[0]

In [ ]: age[2]

In [ ]: age.iloc[-1]

In [ ]: age.iloc[:3]

In [ ]: age.loc[:3]

In [ ]: age[890]

# # Sorting and Introduction to the inplace-parameter

```
In [ ]: import pandas as pd
In [ ]: dic = {1:10, 3:25, 2:6, 4:36, 5:2, 6:0, 7:None}
dic
In [ ]: sales = pd.Series(dic)
sales
In [ ]: sales.sort_index()
In [ ]: sales
In [ ]: sales
In [ ]: sales.sort_index(ascending = True, inplace = True)
In [ ]: sales
In [ ]: sales
In [ ]: sales.sort_values()
In [ ]: sales
```

# nlargest(), nsmallest(), idxmax() and idxmin()

```
In [ ]: import pandas as pd
In [ ]: titanic = pd.read_csv("titanic.csv")
In [ ]: titanic.head()
In [ ]: age = titanic.age
In [ ]: age
In [ ]: age.head()
In [ ]: age.sort_values(ascending=False).head(5)
In [ ]: age.nlargest()
In [ ]: age.nsmallest()
In [ ]: age.idxmax()
In [ ]: age.idxmin()
```

#### 5.1 Section DataFrame Basics II

### In this Section, you will learn how to

- 1. filter DataFrames by one and by many conditions.
- 2. perform advanced filtering techniques
- 3. remove rows and columns from a DataFrame
- 4. add new columns to a DataFrame

### DATAFRAME BASICS II

### FILTERING DATAFRAMES WITH ONE CONDITION

```
In [ ]: import pandas as pd
In [ ]: titanic = pd.read_csv("titanic.csv")
In [ ]: titanic
In [ ]: titanic.loc[titanic.sex =="male"]
In [ ]: t1 = titanic.loc[titanic.sex =="male"]
In [ ]: t1
```

# DATAFRAME BASICS II

### ## FILTERING DATAFRAMES WITH MANY CONDITIONS (AND)

```
In [ ]: import pandas as pd
In [ ]: titanic = pd.read_csv("titanic.csv")
In [ ]: titanic
In [ ]: titanic.loc[titanic.sex =="male"]
In [ ]: mask1 = titanic.sex =="male"
In [ ]: mask2 = titanic.age > 14
In [ ]: mask1
In [ ]: mask2
In [ ]: mask2
```

```
In [ ]: output = titanic.loc[mask1 & mask2]
In [ ]: output
In [ ]: output.info()
In [ ]: output.describe()
```

### ## FILTERING DATAFRAMES WITH MANY CONDITIONS (OR)

```
In [ ]: import pandas as pd
In [ ]: titanic = pd.read_csv("titanic.csv")
In [ ]: titanic
In [ ]: titanic.loc[titanic.sex =="male"]
In [ ]: mask1 = titanic.sex =="male"
In [ ]: mask2 = titanic.age > 14
In [ ]: mask1
In [ ]: mask2
In [ ]: (mask1 & mask2).head()
In [ ]: output = titanic.loc[mask1 | mask2]
In [ ]: output
In [ ]: output.info()
In [ ]: output.describe()
```

### ## ADVANCED FILTERING WITH BETWEEN(), ISIN() AND ~

```
In [ ]: import pandas as pd
In [ ]: summer = pd.read_csv("summer.csv")
In [ ]: summer
In [ ]: og_1988 = summer.loc[summer.Year == 1988]
In [ ]: og_1988.head()
In [ ]: og_1988.tail()
In [ ]: og_1988.info()
In [ ]: og_since1992 = summer.loc[summer.Year >= 1992]
In [ ]: og_since1992
In [ ]: og_since1992.head()
In [ ]: og_since1992.tail()
In [ ]: summer.Year.between(1960,1969).head()
In [ ]: og_60s = summer.loc[summer.Year.between(1960,1969, inclusive = True)]
In [ ]: og_60s
In [ ]: my_favoriate_games = [1972,1996]
In [ ]: summer.Year.isin(my_favoriate_games).head()
In [ ]: og_72_96 = summer.loc[summer.Year.isin(my_favoriate_games)]
```

```
In [ ]: og_72_96
In [ ]: og_not_72_96 = summer.loc[~summer.Year.isin(my_favoriate_games)]
In [ ]: og_not_72_96
In [ ]: og_not_72_96.Year.unique()
In [ ]: mask1 = summer.Year == 1960
In [ ]: mask2 = summer.Year == 1996
In [ ]: og_or_72_96 = summer.loc[mask1 | mask2]
In [ ]: og_or_72_96
In [ ]: test1 = summer.loc[summer.Year == 1960]
In [ ]: test1
In [ ]: summer
```

# ## ADVANCED FILTERING WITH any() AND all()

```
In [ ]: import pandas as pd
In [ ]: titanic = pd.read_csv("titanic.csv")
In [ ]: titanic.head()
In [ ]: titanic.tail()
In [ ]: titanic.sex == "male"
In [ ]: (titanic.sex == "male").any()
In [ ]: (titanic.sex == "male").all()
In [ ]: (titanic.age == 80.0).any()
In [ ]: pd.Series([-1,0.5,1,-0.1,0]).all()
In [ ]: titanic.fare.all()
```

# Removing Columns ¶

```
In [ ]: import pandas as pd
In [ ]: summer = pd.read_csv("summer.csv")
In [ ]: summer.head()
In [ ]: summer.drop(columns = "Sport")
In [ ]: summer.head()
In [ ]: summer.drop(columns = ["Sport", "Discipline"])
In [ ]: summer.head()
In [ ]: summer.head()
In [ ]: summer.head()
In [ ]: summer.loc(columns = ["Sport", "Discipline"], inplace = True)
In [ ]: summer.head()
```

# **Removing Rows**

```
In [ ]: import pandas as pd
In [ ]: summer = pd.read_csv("summer.csv", index_col = "Athlete")
In [ ]: summer.head()
In [ ]: summer.drop(index = "HAJOS, Alfred")
In [ ]: summer.head()
In [ ]: summer.drop(index = "HAJOS, Alfred", inplace = True)
In [ ]: summer.head()
In [ ]: summer.head()
In [ ]: summer.head()
In [ ]: summer = pd.read_csv("summer.csv", index_col = "Athlete")
In [ ]: summer.head()
In [ ]: summer.head()
In [ ]: summer.head()
In [ ]: summer.head()
```

# **Removing Rows**

```
In [ ]: import pandas as pd
In [ ]: titanic = pd.read_csv("titanic.csv")
In [ ]: titanic.head()
In [ ]: titanic["Zeros"] = 0
In [ ]: titanic.head()
In [ ]: titanic = pd.read_csv("titanic.csv")
In [ ]: titanic.head()
```

### DATAFRAME BASICS II

### Creating columns based on other columns

```
In [ ]: import pandas as pd
In [ ]: titanic = pd.read_csv("titanic.csv")
In [ ]: titanic.head()
In [ ]: 1912-titanic.age
In [ ]: titanic["YOB"] = 1912-titanic.age
In [ ]: titanic.head()
In [ ]: titanic.sibsp + titanic.parch
In [ ]: titanic["relatives"] = titanic.sibsp + titanic.parch
```

```
In [ ]: titanic.head()
In [ ]: titanic.drop(columns = ["sibsp", "parch"], inplace = True)
In [ ]: titanic.head()
In [ ]: inflation_factor = 10
In [ ]: titanic.fare*10
In [ ]: titanic.fare = titanic.fare*10
In [ ]: titanic.head()
```

# Adding columns with insert()

```
In [ ]: import pandas as pd
In [ ]: titanic = pd.read_csv("titanic.csv")
In [ ]: titanic.head()
In [ ]: titanic["Test"] = "Test"
In [ ]: titanic.head()
In [ ]: relatives = titanic.sibsp + titanic.parch relatives.head()
In [ ]: titanic.insert(loc = 6, column = "relatives", value = relatives)
In [ ]: titanic.head()
```

#### DataFrame Basics II ¶

#### Exercise 5: Filtering DataFrames & Adding/Removing Rows and Columns

```
In []: #run the cell!
import pandas as pd

In []: #run the cell!
cars = pd.read_csv("cars.csv")

In []: #run the cell!
cars

40. Check for all elements in the column "origin" whether they are equal to "europe". Save the boolean Series in the variable mask1! Fill in the gaps!
```

In [ ]: mask1 = cars.origin == "europe"
mask1

```
41. Check for all elements in the column "mpg" whether they are smaller than 20. Save the boolean Series in the variable mask2!
 In [ ]: mask2 = cars.mpg < 20
           42. Filter the cars DataFrame for all cars from europe (use mask1) and save the subset in the variable europe! Fill in the gaps!
 In [ ]: europe = cars.loc[mask1].copy()
           Inspect! What is the name of the second car?
 In [ ]: # run the cell!
           europe
           The second car is a ... peugeot 504!
           43. Get some meta information on the DataFrame europe. How many cars are from europe?
 In [ ]: europe.info()
            .. 70 cars are from europe!
             44. Filter the DataFrame cars for all cars from europe with low fuel efficiency (mpg lower than 20). Use mask1 and mask2!
                 Save the subset in the variable europe_le! Fill in the gaps!
   In [ ]: europe_le = cars.loc[mask1 & mask2].copy()
             Inspect! What is the name of the least efficient car from europe?
   In [ ]: #run the cell!
             europe_le
             The least efficient european car is ... peugeot 604sl
             45. Filter the DataFrame cars for all cars with an mpg between 10 and 15 (both ends inclusive!).
                 Save the subset in the variable mpg_10_15! Fill in the gaps!
   In [ ]: mpg_10_15 = cars.loc[cars.mpg.between(10,15)].copy()
            Inspect! The first car is ...?
   In [ ]: # run the cell!
            mpg_10_15
         The first car is... buick skylark 320!
         Inspect! How many cars are in the subset?
In [ ]: #run the cell!
         mpg_10_15.info()
         There are ... 68 cars!
          46. Filter the Dataframe cars for all cars that are not built in the years 73 and 74, and only the columns "mpg" and "name"!
              Save the subset in the variable not_73_74. Fill in the gaps!
In [ ]: not_73_74 = cars.loc[~cars.model_year.isin([73, 74]), ["mpg", "name"]].copy()
In [ ]: #run the cell!
         not 73 74
         Inspect! How many cars are in this subset?
In [ ]: # run the cell!
         not_73_74.info()
```

Well Done!

The "I\_per\_100km" value for the audi 100 ls is ... 9.80!

### Section 6: Manipulating Elements in a DataFrame / Slice +++Important...

If you want to avoid the most commonly made mistakes in Pandas, this is probably the most important Section of the course!

There are two types of mistakes that even experienced Pandas coders are faced with:

- 1. Intention: Manipulating elements in a DataFrame. But: The code only manipulates these elements in a slice/subset. The original DataFrame remains unchanged.
- 2. Intention: Manipulating elements in a slice/copy of a DataFrame. But: The code also manipulates these elements in the original DataFrame.

Both types are dangerous and in some cases very hard to detect. In this Section, you will learn how Pandas work under the hood. And you will learn some simple rules to avoid these pitfalls.