# Warsztaty modelowania

## 01 – biblioteka pandas

opracowała

Patrycja Naumczyk

# O czym będzie?

1. Numpy – czyli bebechy pandas
2. Pandas – struktury danych
   a) Serie vs tabele
   b) Indeksowanie
3. Tworzenie nowych kolumn i przypisywanie wartości
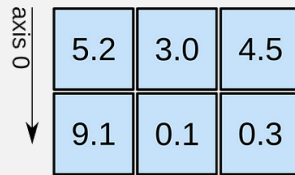4. Funkcje agregujące
5. Method chaining

# Series i dataframe

1. pd.Series()
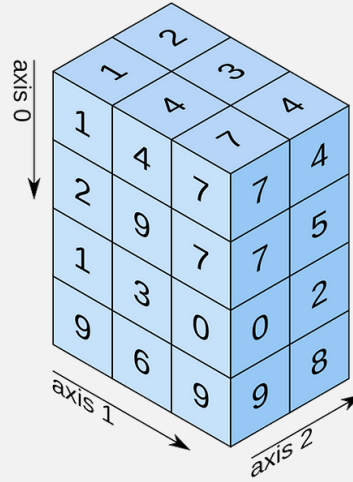2. pd.DataFrame()

# Indeksowanie



| axis=1 or columns | | |
|---|---|---|
| A | B | C |
| 0A | 0B | 0C |
| 1A | 1B | 1C |
| 2A | 2B | 2C |
| 3A | 3B | 3C |
| 4A | 4B | 4C |

pandas.Index

axis=0 or index

pandas.Index    panda.Series

| name | region | sales | expenses |
|---|---|---|---|
| William | East | 50000 | 42000 |
| Emma | North | 52000 | 43000 |
| Sofia | East | 90000 | 50000 |
| Markus | South | 34000 | 44000 |
| Edward | West | 42000 | 38000 |
| Thomas | West | 72000 | 39000 |
| Ethan | South | 49000 | 42000 |
| Olivia | West | 55000 | 60000 |
| Arun | West | 67000 | 39000 |
| Anika | East | 65000 | 44000 |
| Paulo | South | 67000 | 45000 |

An index is like a group of row labels

Column Index position

| | 0 | 1 | 2 | 3 | 4 | |
|---|---|---|---|---|---|---|
| | EmpID | Skill | Age | Pay | Name | Column Index values |
| 0 | 0 | 21 | Python | 35 | 15000.0 | Indhu |
| 1 | 1 | 12 | python | 27 | NaN | Karthi |
| 2 | 2 | 15 | JavaScript | 32 | 5000.0 | Palani |
| 3 | 3 | 7 | JavaScript | 27 | 12000.0 | Sarvesh |
| 4 | 4 | 10 | PYTHON | 25 | 1000.0 | Bindhu |

Row index position

Row Index values

# Indeksowanie cd.

1. Ustawianie dowolnego indeksu:
   a) set_index()
2. Czyszczenie indeksu:
   a) reset_index()
3. Odwołanie wprost:
   a) df[ ]

# Indeksowanie cd.

1. Ustawianie dowolnego indeksu:
   a) set_index()
2. Czyszczenie indeksu:
   a) reset_index()
3. Odwołanie wprost:
   a) df[ ]
4. Odwołanie pozycyjne (jak w numpy):
   a) iloc[ ]


df


df.iloc[0]


df1


df1.iloc[[0,2]]


df1

Step 2


df1.iloc[0:4:2,0:2]

# Indeksowanie cd.

1. Ustawianie dowolnego indeksu:
   a) set_index()
2. Czyszczenie indeksu:
   a) reset_index()
3. Odwołanie wprost:
   a) df[ ]
4. Odwołanie pozycyjne (jak w numpy):
   a) iloc[ ]
5. Odwołanie do etykiet:
   a) loc[ ]

## df

| | EmpID | Skill | Age | Pay | Name |
|---|---|---|---|---|---|
| 0 | 21 | Python | 35 | 15000.0 | Indhu |
| 1 | 12 | python | 27 | NaN | Karthi |
| 2 | 15 | JavaScript | 32 | 5000.0 | Palani |
| 3 | 7 | JavaScript | 27 | 12000.0 | Sarvesh |
| 4 | 10 | PYTHON | 25 | 1000.0 | Bindhu |

## df.loc[0,['EmpID','Skill']]

```
EmpID         21
Skill     Python
Name: 0, dtype: object
```

## df.loc[[0],['EmpID','Skill']]

| | EmpID | Skill |
|---|---|---|
| 0 | 21 | Python |

## df1

| Name | EmpID | Skill | Age | Pay |
|---|---|---|---|---|
| Indhu | 21 | Python | 35 | 15000.0 |
| Karthi | 12 | python | 27 | NaN |
| Palani | 15 | JavaScript | 32 | 5000.0 |
| Sarvesh | 7 | JavaScript | 27 | 12000.0 |
| Bindhu | 10 | PYTHON | 25 | 1000.0 |

Step 2

## df1.loc[::2,"EmpID":"Age"]

| Name | EmpID | Skill | Age |
|---|---|---|---|
| Indhu | 21 | Python | 35 |
| Palani | 15 | JavaScript | 32 |
| Bindhu | 10 | PYTHON | 25 |

# Porównanie iloc[ ] i loc[ ]

| Input given in iloc | Return Type |
|---|---|
| 1.Both row_index and column_index given as single integer | Single value |
| 2. One input is given as single integer and other input is given as list of integer/integers | Series |
| 3. Both row_index and column_index given as list of integer/integers. | DataFrame |

| Input given in loc | Return Type |
|---|---|
| 1.Both row_index and column_index given as single label | Single value |
| 2. One input is given as single label and other input is given as list of label/labels | Series |
| 3. Both row_index and column_index given as list of label/labels. | DataFrame |

df.loc[0,"EmpID"] → 21    Return type -> single value

df.loc[0,["EmpID"]] →
```
EmpID    21
Name: 0, dtype: object
```
Return Type -Series

df.loc[[0],"EmpID"] →
```
0     21
Name: EmpID, dtype: int64
```
Return Type -Series

df.loc[[0],["EmpID"]] →
```
      EmpID
 0      21
```
Return type -DataFrame

# Indeksowanie cd.

1. Ustawianie dowolnego indeksu:
   a) set_index()
2. Czyszczenie indeksu:
   a) reset_index()
3. Odwołanie wprost:
   a) df[ ]
4. Odwołanie pozycyjne (jak w numpy):
   a) iloc[ ]
5. Odwołanie do etykiet:
   a) loc[ ]
6. Maski logiczne



The row index points to the corresponding row in each Series (axis = 0)

| | year | make | model | body | condition | odometer | color | interior | sellingprice |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 2015 | Kia | Sorento | SUV | 5.0 | 16639.0 | white | black | 21500 |
| 1 | 2015 | Kia | Sorento | SUV | 5.0 | 9393.0 | white | beige | 21500 |
| 2 | 2014 | BMW | 3 Series | Sedan | 4.5 | 1331.0 | gray | black | 30000 |
| 3 | 2015 | Volvo | S60 | Sedan | 4.1 | 14282.0 | white | black | 27750 |
| 4 | 2014 | BMW | 6 Series Gran Coupe | Sedan | 4.3 | 2641.0 | gray | black | 67000 |

The column index points to each individual Series (axis = 1)

Each column is a Pandas Series

```
cars["make"] == "BMW"
```

```
0    False
1    False
2    True
3    False
4    True
Name: make, dtype: bool
```

```
conditions = (cars["make"] == "BMW") & (cars["model"] == "3 Series")

cars[conditions]
```

| | year | make | model | body | condition | odometer | color | interior | sellingprice |
|---|---|---|---|---|---|---|---|---|---|
| 2 | 2014 | BMW | 3 Series | Sedan | 4.5 | 1331.0 | gray | black | 30000 |

# Tworzenie kolumn

1. Przypisanie wartości

```python
df['Course'] = 'Computer science'
df
```

|   | Name  | Age | City      | Country   | Course           |
|---|-------|-----|-----------|-----------|------------------|
| a | Jack  | 34  | Sydeny    | Australia | Computer science |
| b | Riti  | 30  | Delhi     | India     | Computer science |
| c | Tom   | 31  | Mumbai    | India     | Computer science |
| d | Neelu | 32  | Bangalore | India     | Computer science |
| e | John  | 16  | New York  | US        | Computer science |
| f | Mike  | 17  | las vegas | US        | Computer science |

```python
df.loc[:,'Grade'] = 'A'
df
```

|   | Name  | Age | City      | Country   | Course           | Grade |
|---|-------|-----|-----------|-----------|------------------|-------|
| a | Jack  | 34  | Sydeny    | Australia | Computer science | A     |
| b | Riti  | 30  | Delhi     | India     | Computer science | A     |
| c | Tom   | 31  | Mumbai    | India     | Computer science | A     |
| d | Neelu | 32  | Bangalore | India     | Computer science | A     |
| e | John  | 16  | New York  | US        | Computer science | A     |
| f | Mike  | 17  | las vegas | US        | Computer science | A     |

# Tworzenie kolumn

1. Przypisanie wartości
2. Operacje arytmetyczne
   a) Operatory matematyczne (+ - * /)
   b) Metody (add() sub() mul() div() )

# Tworzenie kolumn

1. Przypisanie wartości
2. Operacje arytmetyczne
   a) Operatory matematyczne (+ - * /)
   b) Metody (add() sub() mul() div() )
3. Metoda assign()

```
df = df.assign(Year='3')
df
```

|   | Name | Age | City | Country | Course | Grade | Year |
|---|------|-----|------|---------|--------|-------|------|
| a | Jack | 34 | Sydeny | Australia | Computer science | A | 3 |
| b | Riti | 30 | Delhi | India | Computer science | A | 3 |
| c | Tom | 31 | Mumbai | India | Computer science | A | 3 |
| d | Neelu | 32 | Bangalore | India | Computer science | A | 3 |
| e | John | 16 | New York | US | Computer science | A | 3 |
| f | Mike | 17 | las vegas | US | Computer science | A | 3 |

# Tworzenie kolumn

1. Przypisanie wartości
2. Operacje arytmetyczne
3. Metoda assign()
4. Przypisanie warunkowe:
    a) Funkcja np.where()
    b) Funkcja np.select()

The
function
name

A True/False
condition that
determines the
output

```
np.where(condition, output-if-true, output-if-false)
```

The output if
condition is True

The output if
condition is False

## numpy.select

```
numpy.select(condlist, choicelist, default=0)
```

Return an array drawn from elements in choicelist, depending on conditions.

Parameters: **condlist** : *list of bool ndarrays*

The list of conditions which determine from which array in *choicelist* the output elements are taken. When multiple conditions are satisfied, the first one encountered in *condlist* is used.

**choicelist** : *list of ndarrays*

The list of arrays from which the output elements are taken. It has to be of the same length as *condlist*.

**default** : *scalar, optional*

The element inserted in *output* when all conditions evaluate to False.

# Funkcje agregujące

1. Operacje matematyczne i statystyczne

THE `.sum()` METHOD SUMS THE
VALUES OF A VARIABLE OR DATAFRAME

The name of the series you
want to operate on

Optional
parameters

```
your_series.sum(…)
```

The method name

| name | sales |
|------|-------|
| Arun | 67000 |
| Edward | 42000 |
| William | 50000 |
| Emma | 52000 |
| Sofia | 90000 |

`dataframe.sales.sum()`

301,000

The name of your
Pandas dataframe

The method name

```
your_dataframe.column.sum(…)
```

The name of the column
in the dataframe that
you want to operate on

Optional
parameters

| name | sales |
|------|-------|
| Arun | 67000 |
| Edward | 42000 |
| William | 50000 |
| Emma | 52000 |
| Sofia | 90000 |

`dataframe.sales.mean()`

| mean |
|------|
| 60200 |

# Funkcje agregujące

1. Operacje matematyczne i statystyczne
2. Metoda agg()

# Funkcje agregujące

1. Operacje matematyczne i statystyczne
2. Metoda agg()
3. Metoda groupby()

# Funkcje agregujące

# Funkcje agregujące

1. Operacje matematyczne i statystyczne
2. Metoda agg()
3. Metoda groupby()
4. Funkcja pd.merge()

# Funkcje agregujące

1. Operacje matematyczne i statystyczne
2. Metoda agg()
3. Metoda groupby()
4. Funkcja pd.merge()
5. Funkcja pd.concat()

# Method chaining

Open
parenthesis

The name the
DataFrame you
want to operate on

```
(dataframe
    .pandas_method()
    .pandas_method()
)
```

Multiple Pandas
methods on
separate lines

Close
parenthesis

## Method Chaining

Most pandas methods return a DataFrame so that another pandas method can be applied to the result. This improves readability of code.

```
df = (pd.melt(df)
        .rename(columns={
            'variable' : 'var',
            'value' : 'val'})
        .query('val >= 200')
    )
```