



Warsztaty modelowania

02 – Analiza eksploracyjna danych (EDA)

opracowała Patrycja Naumczyk

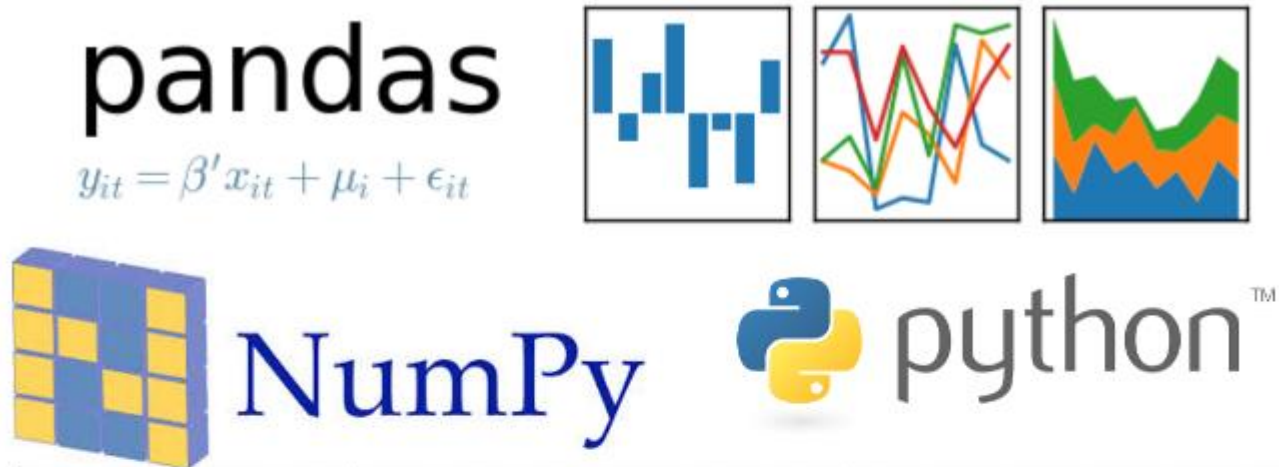
EDA - Exploratory Data Analysis

Absolutne minimum

1. Typy danych
2. Braki danych
3. Duplikaty
4. Statystyki opisowe (adekwatne do skali pomiarowej zmiennej)
5. Rozkłady zmiennych
6. Zależności między zmiennymi



Typy danych



Pandas dtype	Python type	NumPy type	Usage
object	str or mixed	string_, unicode_, mixed types	Text or mixed numeric and non-numeric values
int64	int	int_, int8, int16, int32, int64, uint8, uint16, uint32, uint64	Integer numbers
float64	float	float_, float16, float32, float64	Floating point numbers
bool	bool	bool_	True/False values
datetime64	NA	datetime64[ns]	Date and time values
timedelta[ns]	NA	NA	Differences between two datetimes
category	NA	NA	Finite list of text values

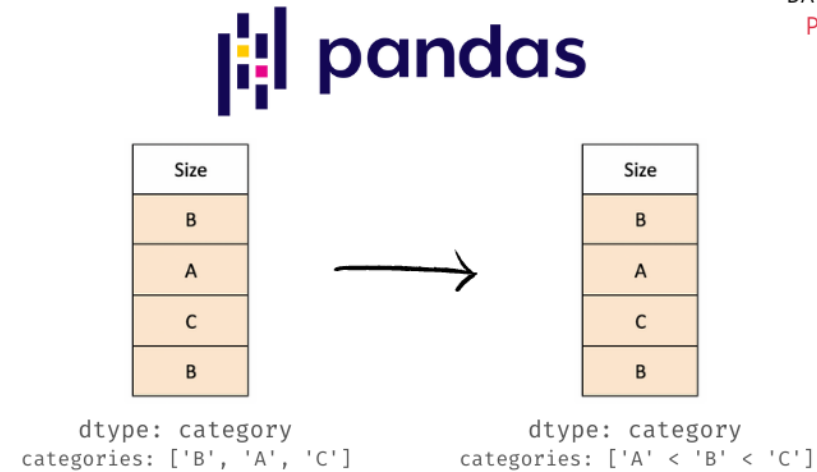
Kind of Data	pandas Data Type	Scalar	Array
TZ-aware datetime	DatetimeTZDtype	Timestamp	Datetimes
Timedeltas	(none)	Timedelta	Timedeltas
Period (time spans)	PeriodDtype	Period	Periods
Intervals	IntervalDtype	Interval	Intervals
Nullable Integer	Int64Dtype , ...	(none)	Nullable integer
Nullable Float	Float64Dtype , ...	(none)	Nullable float
Categorical	CategoricalDtype	(none)	Categoricals
Sparse	SparseDtype	(none)	Sparse
Strings	StringDtype	str	Strings
Nullable Boolean	BooleanDtype	bool	Nullable Boolean
PyArrow	ArrowDtype	Python Scalars or NA	PyArrow

<https://pandas.pydata.org/docs/reference/arrays.html>

https://pandas.pydata.org/docs/user_guide/basics.html#basics-dtypes

Typy danych

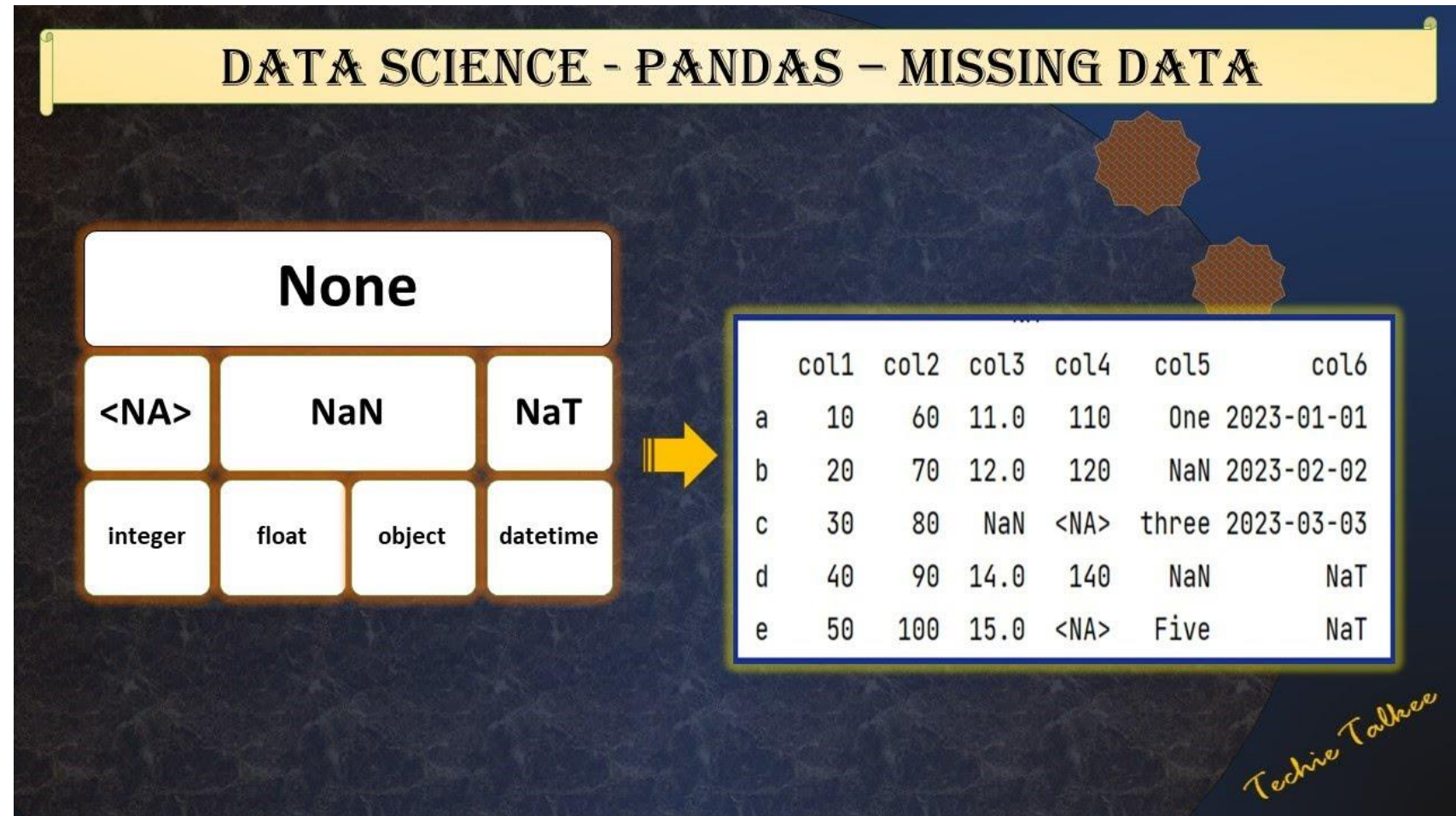
1. Atrybut `.dtypes` / `.dtype`
2. Metody
 - a) `astype()`
 - b) `to_numeric()`
 - c) `to_datetime()`
 - d) `select_dtypes()`



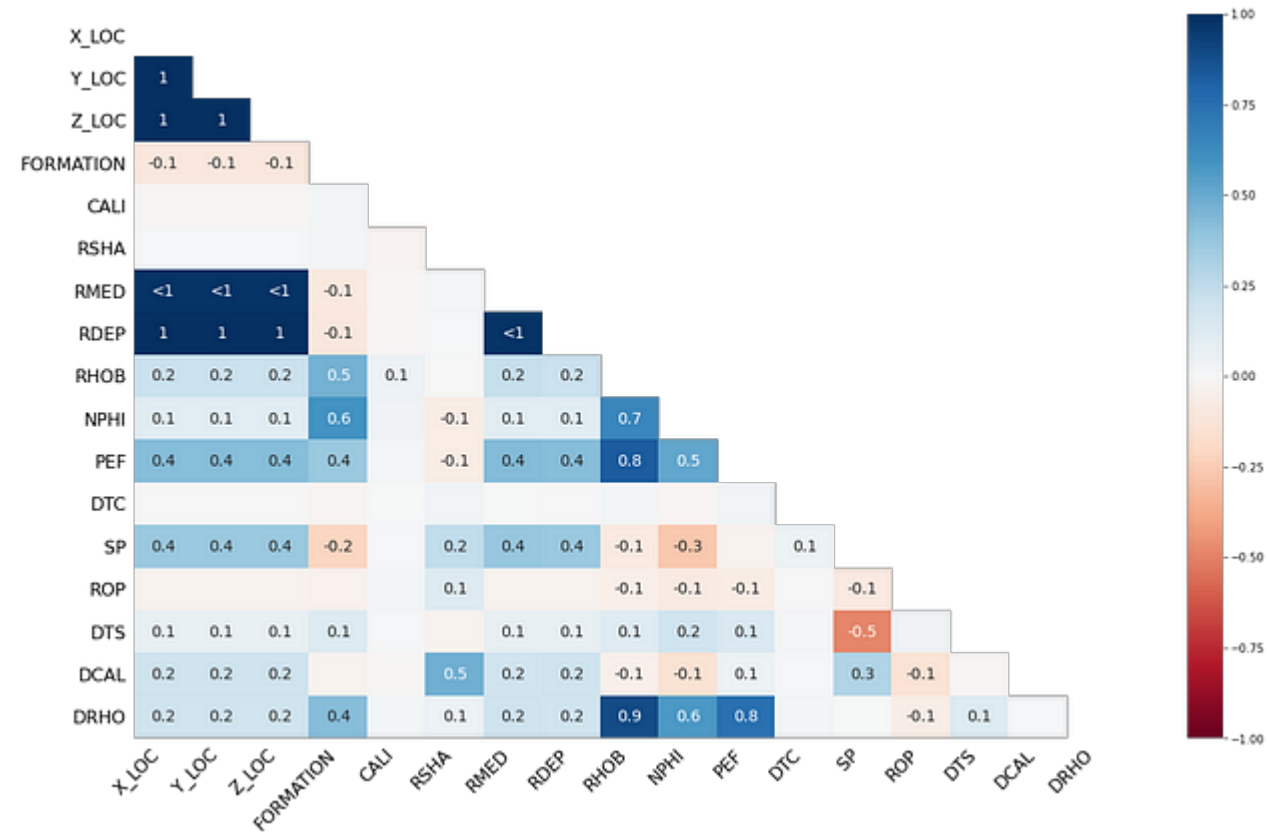
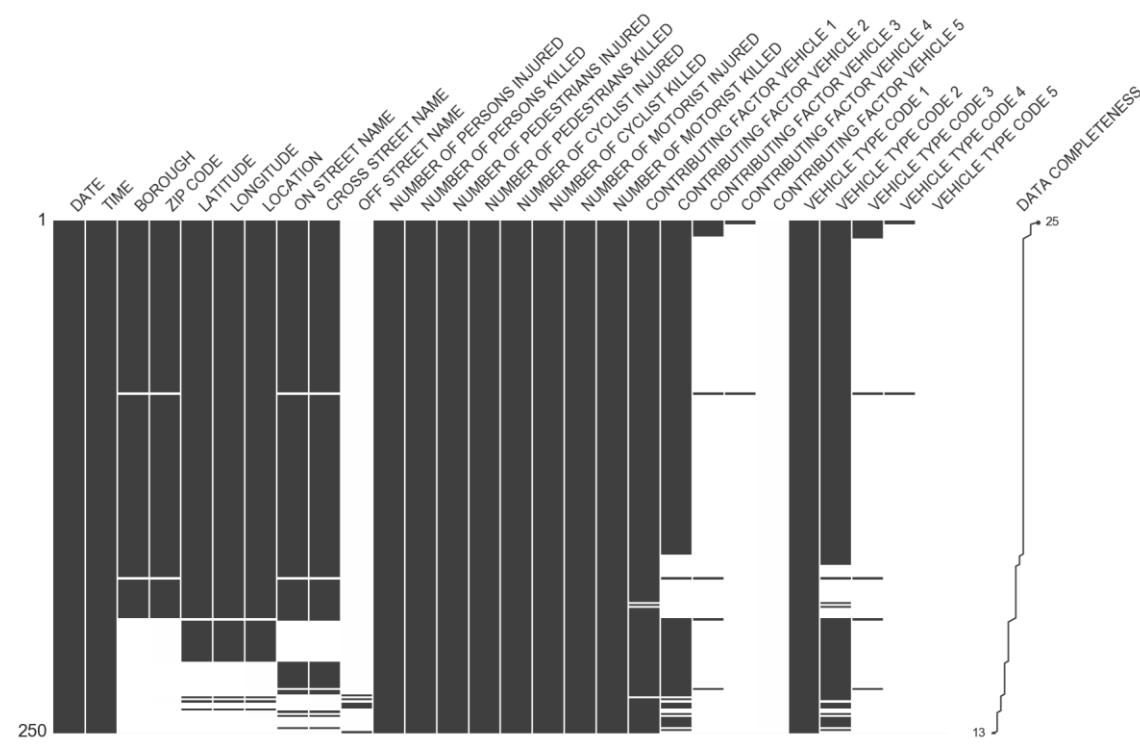
Set Category Order

Braki danych

1. Reprezentacja braku
 - a) np.nan
 - b) NaT
 - c) NA – string i nullable
 - d) None
2. Metody
 - a) isna() / isnull()
 - b) notna()
 - c) dropna()
 - d) fillna() / interpolate()

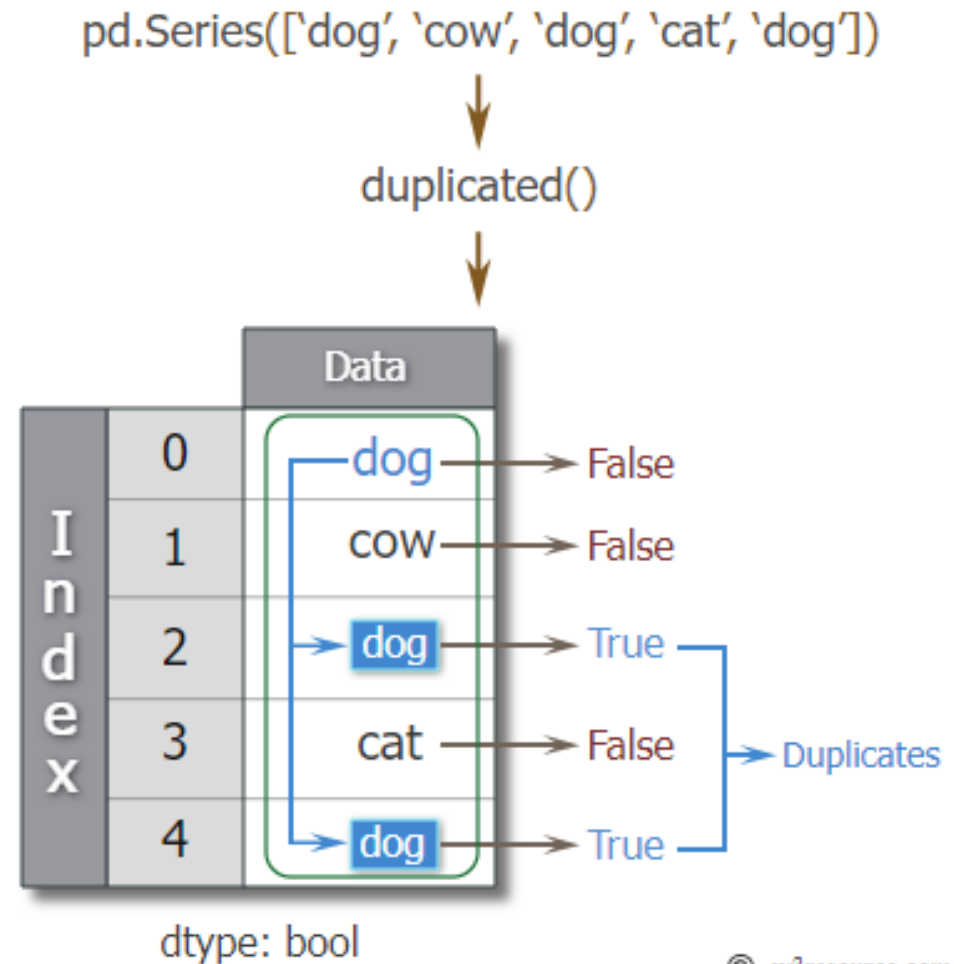


Braki danych - missingno



Duplikaty danych

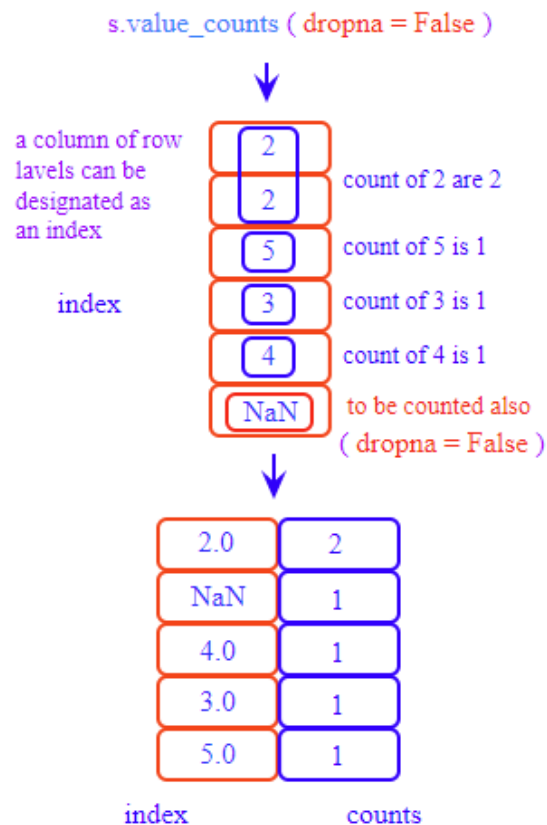
1. Flagi metod konstruowania obiektów
2. Metody badania indexu
 - a) `.index.is_unique`
 - b) `.index.duplicated()`
3. Metody badania treści
 - a) `duplicated()`
 - b) `drop_duplicates()`



Statystyki opisowe

1. Metody:

- describe()
- value_counts()
- unique() / nunique()



© w3resource.com

```
df = pd.DataFrame( { 'categorical': pd.Categorical(['s', 't', 'u']),
                    'numeric': [2, 3, 4],
                    'object': ['p', 'q', 'r'] } )
df.describe()
```

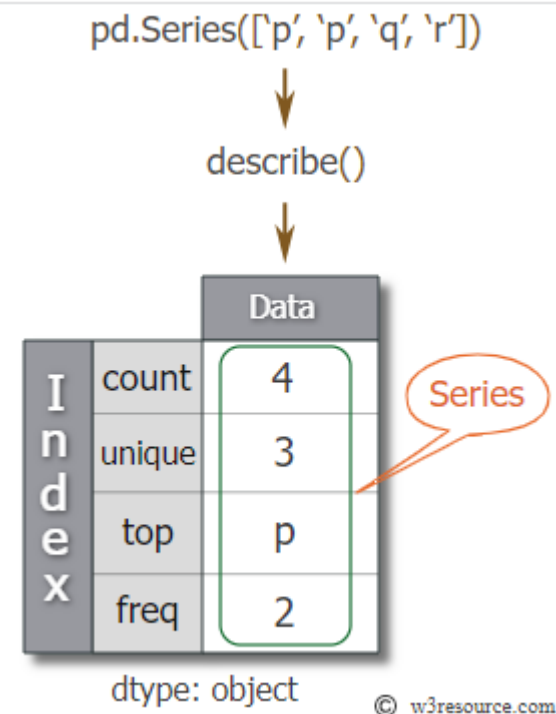
↓

`'numeric': [2, 3, 4]`

DataFrame		
df	numeric	
count	3.0	3 numbers
mean	3.0	mean or average
std	1.0	Standard Deviation
min	2.0	minimum value
25%	2.5	25th percentiles
50%	3.0	50th percentiles
75%	3.5	75th percentiles
max	4.0	maximum value

N.B. : by default describe returns only numeric field

© w3resource.com



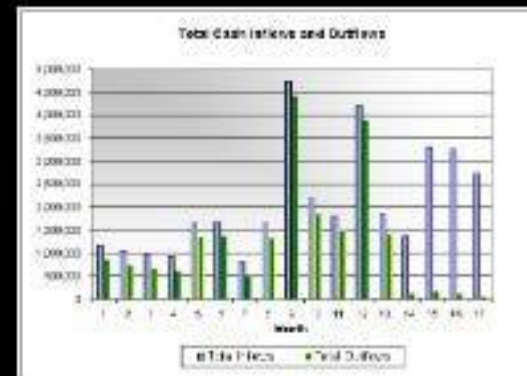


WHEN DATA IS IN TABLE FORM

ID	NAME	CLASS	MARK	SEX
1	John Doe	Four	75	female
2	Max Ruhn	Three	85	male
3	Arnold	Three	95	male
4	Kevin Star	Four	80	female
5	John Mike	Four	60	female
6	Alex John	Four	55	male
7	My John Rob	Fifth	78	male
8	Arnold	Five	85	male
9	Tex Cry	Six	76	male
10	Big John	Four	55	female



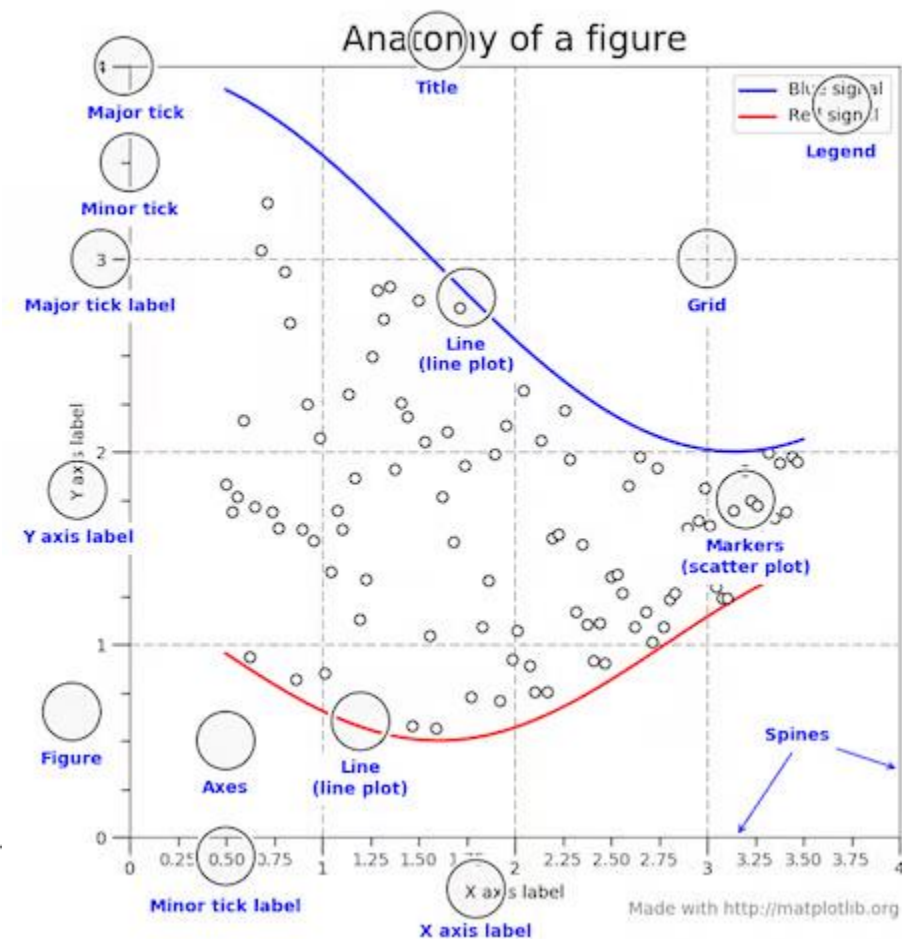
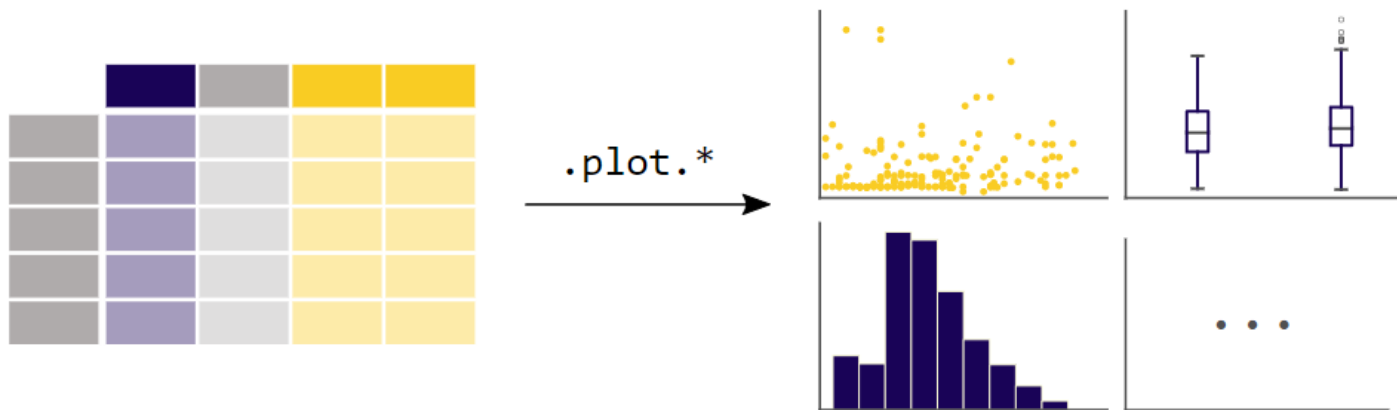
WHEN DATA IS IN PLOT



imgflip.com

Podstawowe wykresy pandas

- `'bar'` or `'barh'` for bar plots
- `'hist'` for histogram
- `'box'` for boxplot
- `'kde'` or `'density'` for density plots
- `'area'` for area plots
- `'scatter'` for scatter plots
- `'hexbin'` for hexagonal bin plots
- `'pie'` for pie plots



https://pandas.pydata.org/docs/user_guide/visualization.html

