

Regresja liniowa

Prognozowanie wartości danej zmiennej na podstawie innej.

Będę prognozował wartość 'Mieszkań oddanych do użytkowania' na podstawie 'roku'.

Spis treści

- Pobranie danych
- Wczytanie danych
- Wstępna obróbka danych
- Zmiana typu obiekt na float
- Współczynnik korelacji Pearsona
- Przygotowanie danych do modelu
- Model
- Korzystanie z modelu
- Wykres - regresja liniowa
- Znaczenie graficznej reprezentacji danych

Pobranie danych

<https://bdm.stat.gov.pl/> -> MIESZKANIA

ROczneKWARTALNEMIESIĘCZNE

Mieszkania

Przejdź do ustawień, aby wyświetlić pełny zakres dostępnych danych.											
	Wskaźnik	Jednostka	2014	2015	2016	2017	2018	2019	2020		
Mieszkania, na których budowę wydano pozwolenie na budowę zgłoszenie z projektem budowlanym	tyś.	3,8	196,9	188,8	211,5	250,9	257,6	268,8	216,1		
Mieszkania, których budowę rozpoczęto	tyś.	7,4	140,3	168,4	173,9	206,2	219,3	207,3	223,8		
Mieszkania w budowie	tyś.	3,7	69,9	70,8	73,2	78,9	79,7	82,6	88,0		
Mieszkania oddane do użytkownika	tyś.	3,7	102,2	102,7	102,3	118,3	104,1	107,4	120,4		

Arkusz programu Microsoft Excel (xlsx)

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X
1	Wskaźnik - wskaźnik/roczny																							
2	Wskaźnik - wskaźnik/roczny																							
3	Wskaźnik	Wskaźnik	2014	2015	2016	2017	2018	2019	2020	2021	2022	2023	2024	2025	2026	2027	2028	2029	2030	2031	2032	2033	2034	2035
4	Mieszkania, na których budowę wydano pozwolenie na budowę zgłoszenie z projektem budowlanym	tyś.	196,9	188,8	211,5	250,9	257,6	268,8	216,1	218,4	218,4	218,4	218,4	218,4	218,4	218,4	218,4	218,4	218,4	218,4	218,4	218,4	218,4	218,4
5	Mieszkania, których budowę rozpoczęto	tyś.	140,3	168,4	173,9	206,2	219,3	207,3	223,8	223,8	223,8	223,8	223,8	223,8	223,8	223,8	223,8	223,8	223,8	223,8	223,8	223,8	223,8	223,8
6	Mieszkania w budowie	tyś.	69,9	70,8	73,2	78,9	79,7	82,6	88,0	88,0	88,0	88,0	88,0	88,0	88,0	88,0	88,0	88,0	88,0	88,0	88,0	88,0	88,0	88,0
7	Mieszkania oddane do użytkownika	tyś.	102,2	102,7	102,3	118,3	104,1	107,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4
8	Mieszkania oddane do użytkownika	tyś.	102,2	102,7	102,3	118,3	104,1	107,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4
9	Mieszkania oddane do użytkownika	tyś.	102,2	102,7	102,3	118,3	104,1	107,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4
10	Mieszkania oddane do użytkownika	tyś.	102,2	102,7	102,3	118,3	104,1	107,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4
11	Mieszkania oddane do użytkownika	tyś.	102,2	102,7	102,3	118,3	104,1	107,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4
12	Mieszkania oddane do użytkownika	tyś.	102,2	102,7	102,3	118,3	104,1	107,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4	120,4

Wczytanie danych

```
In [1]: import pandas as pd
import matplotlib.pyplot as plt
import numpy as np
from sklearn import linear_model

In [2]: dane = pd.read_excel('C:\\Users\\V6\\Desktop\\Mieszkania.xlsx', header=3)
dane
Out[2]:
```

Unnamed: 0	Unnamed: 1	Unnamed: 2	Unnamed: 3	Unnamed: 4	Unnamed: 5	Unnamed: 6	Unnamed: 7	Unnamed: 8	Unnamed: 9	Unnamed: 10	Unnamed: 11	Unnamed: 12	Unnamed: 13	Unnamed: 14	Unnamed: 15
0	Wskaźnik	NaN	Jednostka	2000	2001	2002	2003	2004	2005	2006	...	2011	2012	2013	2014
1	Mieszkania, na których budowę wydano pozwolenie na budowę zgłoszenie z projektem budowlanym	NaN	tyś.	152,0	149,0	86,0	98,9	114,9	123,9	169,9	...	184,1	165,3	138,0	138,0
2	Mieszkania, których budowę rozpoczęto	NaN	tyś.	125,6	114,4	77,0	83,7	101,1	105,8	138,0	...	162,2	141,8	127,0	127,0
3	Mieszkania w budowie	NaN	tyś.	710,4	718,8	698,2	619,2	612,1	603,9	626,5	...	723,8	712,7	694,0	694,0
4	Mieszkania oddane do użytkownika	NaN	tyś.	87,8	106,0	97,6	162,7	106,1	114,1	115,4	...	131,0	152,9	145,0	145,0
5	w tym w budownictwie indywidualnym	tyś.	31,6	35,4	48,0	113,2	99,9	98,7	93,5	...	65,4	70,3	72,0	72,0	72,0
6	NaN	NaN	tyś.	24,4	25,8	15,4	12,0	9,4	8,2	9,0	...	3,8	4,2	3,0	3,0
7	w budownictwie społecznym	tyś.	24,6	34,6	26,4	28,7	29,1	37,6	42,0	...	56,9	74,4	65,0	65,0	65,0
8	NaN	NaN	tyś.	24,6	34,6	26,4	28,7	29,1	37,6	42,0	...	56,9	74,4	65,0	65,0

8 rows x 24 columns

Wstępna obróbka danych

```
In [3]: # usuniecie kolumn
dane = dane.drop('Unnamed: 0', axis=1)
dane = dane.drop('Unnamed: 1', axis=1)
dane = dane.drop('Unnamed: 2', axis=1)

In [4]: # usuniecie wierszy
dane = dane.drop([1,2,3,5,6,7])

In [5]: # zmiiana nazw indeksow
dane = dane.rename(index={0:"rok", 4:"mieszkania"})
dane.head(3)
Out[5]:
```

Unnamed: 3	Unnamed: 4	Unnamed: 5	Unnamed: 6	Unnamed: 7	Unnamed: 8	Unnamed: 9	Unnamed: 10	Unnamed: 11	Unnamed: 12	...	Unnamed: 14	Unnamed: 15	
rok	2000	2001	2002	2003	2004	2005	2006	2007	2008	2009	...	2011	2012
mieszkania	87,8	106,0	97,6	162,7	106,1	114,1	115,4	133,7	165,2	160,0	...	131,0	152,9

2 rows x 21 columns

```
In [6]: # transponowanie ranki danych
dane = dane.transpose()
dane
Out[6]:
```

	rok	mieszkania
Unnamed: 3	2000	87,8
Unnamed: 4	2001	106,0
Unnamed: 5	2002	97,6
Unnamed: 6	2003	162,7
Unnamed: 7	2004	106,1
Unnamed: 8	2005	114,1
Unnamed: 9	2006	115,4
Unnamed: 10	2007	133,7
Unnamed: 11	2008	165,2
Unnamed: 12	2009	160,0
Unnamed: 13	2010	135,8
Unnamed: 14	2011	131,0
Unnamed: 15	2012	152,9
Unnamed: 16	2013	145,1
Unnamed: 17	2014	143,2
Unnamed: 18	2015	147,7
Unnamed: 19	2016	163,3
Unnamed: 20	2017	178,3
Unnamed: 21	2018	185,1
Unnamed: 22	2019	207,4
Unnamed: 23	2020	221,4

Zmiana typu obiekt na float

```
In [7]: dane.info()
<class 'pandas.core.frame.DataFrame'>
Int64Index: 21 entries, Unnamed: 3 to Unnamed: 23
Data columns (total 24 columns):
# Column Non-Null Count Dtype
---
0 rok 21 non-null object
1 mieszkania 21 non-null object
dtypes: object(2)
memory usage: 1.1+ KB

Zmiana przecinka na kropke
In [8]: # zmiany 'komórki' w których występowałyby sam przecinek
# dane = dane.replace(",",".")

In [9]: # zmiany każdy przecinek w kolumnie 'mieszkania' na kropkę
dane["mieszkania"] = [x.replace(",",".") for x in dane["mieszkania"]]
dane.head(3)
Out[9]:
```

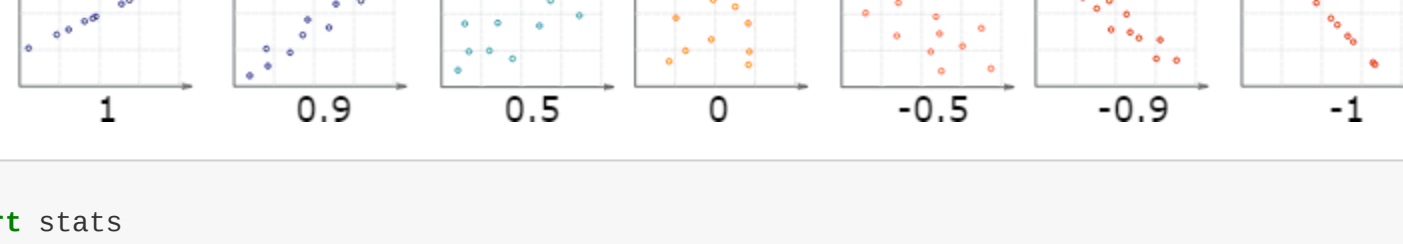
Unnamed: 3	rok	mieszkania
Unnamed: 3	2000	87.8
Unnamed: 4	2001	106.0
Unnamed: 5	2002	97.6

po zmianie przecinka na kropke dopiero mogą zmienić typ kolumn z obiekt na int/float

```
In [10]: #zmiana object(string) na int
dane["rok"] = pd.to_numeric(dane["rok"])
#dane["rok"] = pd.to_numeric(dane["rok"], downcast='float') - gdybym chciał zamienić na float
dane["mieszkania"] = pd.to_numeric(dane["mieszkania"])

In [11]: dane.info()
<class 'pandas.core.frame.DataFrame'>
Int64Index: 21 entries, Unnamed: 3 to Unnamed: 23
Data columns (total 24 columns):
# Column Non-Null Count Dtype
---
0 rok 21 non-null int64
1 mieszkania 21 non-null float64
dtypes: float64(1), int64(1)
memory usage: 1.1+ KB
```

Współczynnik korelacji Pearsona



```
In [12]: # przykład
from scipy import stats
a = np.array([5,3,6])
b = np.array([2,3,4])
stats.pearsonr(a, b)
Out[12]: (0.8666254937844388, 0.3333333333333333)

In [13]: from scipy import stats
a = dane["rok"].values
b = dane["mieszkania"].values
stats.pearsonr(a, b)
# Odp. r = 0,84
Out[13]: (0.8413042761590577, 1.7891785643153179e-06)

In [14]: c = stats.pearsonr(a, b)
print(f"r czyli współczynnik korelacji Pearsona wynosi: ",round(c[0], 2))
r czyli współczynnik korelacji Pearsona wynosi: 0.84
```

Przygotowanie danych do modelu

```
In [15]: x = dane["rok"].values
y = dane["mieszkania"].values

In [16]: X
Out[16]: array([2000, 2001, 2002, 2003, 2004, 2005, 2006, 2007, 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015, 2016, 2017, 2018, 2019, 2020],
      dtype=int64)

In [17]: y
Out[17]: array([ 87.8, 106. ,  97.6, 162.7, 106.1, 114.1, 115.4, 133.7, 165.2,
        160. , 135.8, 131. , 152.9, 145.1, 143.2, 147.7, 163.3, 178.3,
        185.1, 207.4, 221.4])

In [18]: x = x.reshape((21, 1))
#y = y.reshape((21, 1))

In [19]: X
Out[19]: array([[2000],
      [2001],
      [2002],
      [2003],
      [2004],
      [2005],
      [2006],
      [2007],
      [2008],
      [2009],
      [2010],
      [2011],
      [2012],
      [2013],
      [2014],
      [2015],
      [2016],
      [2017],
      [2018],
      [2019],
      [2020]], dtype=int64)
```

Model

```
In [20]: # stworzenie modelu
model = linear_model.LinearRegression()

In [21]: # dopasowanie modelu do danych
model.fit(X, y)
Out[21]: LinearRegression()
```

Korzystanie z modelu

Ile zostanie oddanych mieszkań do użytkowania w 1930 roku?

```
In [22]: model.predict([[1930]])
Out[22]: array([-233.29610391])

Odp. W 1930 roku zostało oddanych -233 tysięcy mieszkań do użytkowania. 🤖
```

WAŻNE



Posiadanie linii regresji nie oznacza, że można podstawić w miejsce X dowolną wartość i przeprowadzić trafną prognozę Y. Tworzenie prognoz na podstawie wartości x, które wykraczają poza zakres posiadanych danych, jest niedopuszczalne. Statystycy nazywają to **ekstrapolacją**, uważają na badaczy, którzy próbują stawiać tezy wykraczające poza zakres danych.

Dlatego w miejsce X będę wstawiał tylko 2000-2020

Wykres - regresja liniowa

```
In [23]: yPred=model.predict(x)
plt.figure(figsize=(10,10))
plt.scatter(x,y)
plt.plot(x,yPred,color='green',lw=3)
plt.xticks(range(2000,2021), rotation = 45)
for i, txt in enumerate(y):
    plt.annotate(txt, (x[i]-0.5, y[i]+1.5), color='red')
Out[23]:
```

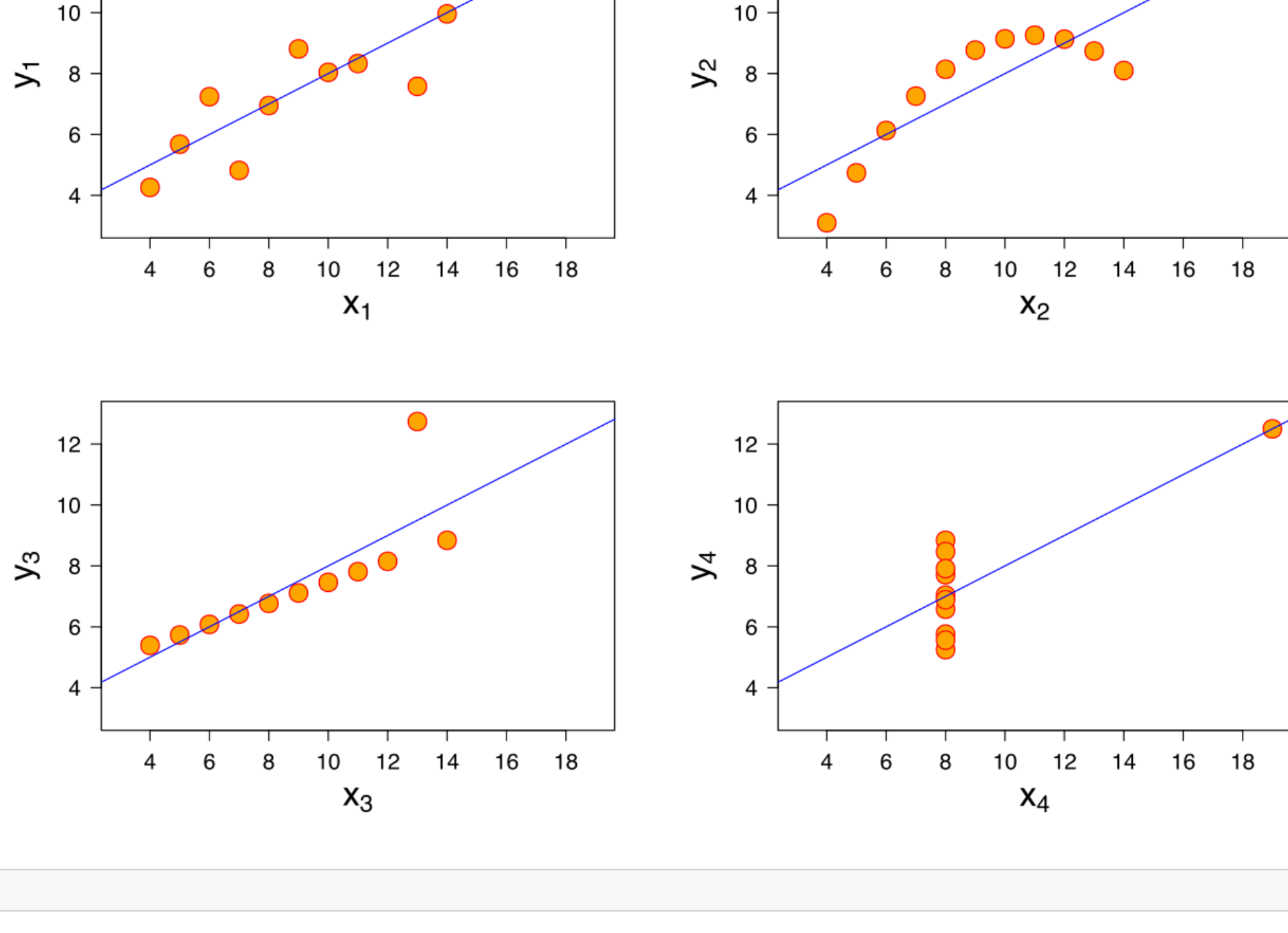
```
In [24]: import sklearn.metrics as sm
print(f"średni błąd kwadratowy =", round(sm.mean_absolute_error(y,yPred),2))
print(f"średni kwadrat błędów =", round(sm.mean_squared_error(y,yPred),2))
print(f"Medianny błąd absolutny =", round(sm.median_absolute_error(y,yPred),2))
print(f"R2 błędów =", round(sm.r2_score(y,yPred),2))
średni błąd kwadratowy = 14.42
średni kwadrat błędów = 339.92
Medianny błąd absolutny = 10.93
R2 błędów = 0.71
```

Znaczenie graficznej reprezentacji danych

https://pl.wikipedia.org/wiki/Kwantet_Anscombe%E2%80%9999

Kwartet Anscombe'a

- zestaw czterech zestawów danych o identycznych cechach statystycznych, takich jak średnia arytmetyczna, wariancja, współczynnik korelacji czy równanie regresji liniowej. Układ tych danych został stworzony w 1973 roku aby ukazać znaczenie graficznej reprezentacji danych przy okazji ich analizy statystycznej.



```
In [ ]:
```