

SPRAWOZDANIE

Zajęcia: Nauka o danych I

Prowadzący: prof. dr hab. Vasyl Martsenyuk

Laboratorium Nr 3 Data 19.10.2024 Temat: Wykorzystanie pakietu Pandas do manipulacji i przetwarzania danych w Pythonie	Bartosz Bieniek Informatyka II stopień, stacjonarne, 1 semestr, gr.A
--	---

1. Polecenie

Zadanie wymagało wykonania następujących kroków:

1. Wczytywanie danych i wyświetlanie podstawowych informacji
2. Obliczanie podstawowych statystyk
3. Identyfikacja i obsługa brakujących danych
4. Wykrywanie wartości odstających
5. Analiza zależności między kolumnami
6. Przekształcanie danych

2. Opis programu opracowanego:

<https://github.com/mindgoner/Studia/tree/master/Nauka%20o%20Danych/Laboratorium%203>

```
Laboratorium 3 > | | Laboratorium3.ipynb > # 4. Wykrywanie wartości odstających
Generate + Code + Markdown Run All Restart Clear All Outputs Jupyter Variables Outline Python 3.11.9

# 1. Wczytywanie danych i wyświetlanie podstawowych informacji
import pandas as pd

file_path = 'E2-opiekunowie.csv'
df = pd.read_csv(file_path)
print(df.head())
print(df.info())
print(df.describe())

[0] ✓ OK Python

...
0 Izabela Wojcik MAZOWIECKIE Warszawa Warszawa - Praga-Polnoc Gmina \
1 Paulina Pomorska MAZOWIECKIE Warszawa Warszawa - Polnoc Otwock
2 Anna Spiechow MAZOWIECKIE Otwock Otwock
3 Angelika Rogaczewska POMORSKIE Gdansk Gdansk
4 Paulina Adamus LODZKIE tdd tdd-Grdmiecie

Kod pocztowy Liczba miejsc Liczba dzieci zapisanych \
0 03-484 8 5.0
1 02-621 8 8.0
2 05-400 8 8.0
3 88-464 5 5.0
4 90-260 5 8.0

Liczba miejsc ze sredkow FERS Liczba dzieci zapisanych na miejsca z FERS \
0 0.0 0.0
1 0.0 0.0
2 NaN NaN
3 NaN NaN
4 0.0 0.0

Opłata miesięczna za pobyt Opłata miesięczna na miejscach z KPO/FERS \
0 NaN 0.0 zł 0.0 zł
1 1999.0 zł 0.0 zł
2 1999.0 zł 0.0 zł
...
25% 0.000000
50% 0.000000
75% 0.000000
max 2.000000
Output is truncated. View as a scrollable element or open in a text editor. Adjust cell output settings.
```

Rys. 1. Fragment kodu i wynik jego wykonania

Podano dane w formacie CSV, które zostały wczytane do DataFrame. Wyświetlono pierwsze kilka wierszy, informacje o strukturze danych oraz podstawowe statystyki opisowe.

```
# 2. Obliczanie podstawowych statystyk
print("2.1. Średnia liczba miejsc: "+str(df['Liczba miejsc'].mean()))
print("2.2. Średnia liczba zapisanych dzieci: "+str(df['Liczba dzieci zapisanych'].mean()))
print("2.3. Suma miejsc ze środków FERS: "+str(df['Liczba miejsc ze środków FERS'].sum()))
df['Opłata miesięczna za pobyt'] = pd.to_numeric(df['Opłata miesięczna za pobyt'], errors='coerce')
print("2.4. Maksymalna opłata miesięczna za pobyt: "+str(df['Opłata miesięczna za pobyt'].max()))
df['Opłata za wyżywienie - dzienna'] = pd.to_numeric(df['Opłata za wyżywienie - dzienna'], errors='coerce')
print("2.5. Minimalna dzienna opłata za wyżywienie: "+str(df['Opłata za wyżywienie - dzienna'].min()))
```

```
✓ 0/0
1. Średnia liczba miejsc: 6.23887775511022
2. Średnia liczba zapisanych dzieci: 4.5449183308053
3. Suma miejsc ze środków FERS: 515.0
4. Maksymalna opłata miesięczna za pobyt: 3700.0
5. Minimalna dzienna opłata za wyżywienie: 0.0
```

Rys. 2. Fragment kodu i wynik jego wykonania

Na podstawie wybranych kolumn obliczono średnie, sumy i minimalne/maksymalne wartości. Statystyki te pomogły zrozumieć rozkład danych w zbiorze.

```
# 3. Identyfikacja i obsługa brakujących danych
print("3.1. Liczba brakujących wartości w każdej kolumnie: "+str(df.isnull().sum()))
missing_percentage = (df.isnull().sum() / len(df)) * 100
print("3.2. Procent brakujących danych w każdej kolumnie:")
print(missing_percentage.head(5))

# Usunięcie danych gdzie brakuje ponad 50% kolumn
threshold = 50
columns_to_drop = missing_percentage[missing_percentage > threshold].index
df_cleaned = df.drop(columns=columns_to_drop)
print("3.3. Kolumny usunięte z powodu dużej liczby braków (> 50%):")
print(", ".join(columns_to_drop))

# Uzupełnianie brakujących w kolumnie "Opłata miesięczna za pobyt" danych średnią wartością
df_cleaned['Opłata miesięczna za pobyt'] = df_cleaned['Opłata miesięczna za pobyt'].fillna(df_cleaned['Opłata miesięczna za pobyt'].mean())

print("3.4. Dane po uzupełnieniu braków w wybranych kolumnach (pierwszych 5):")
print(df_cleaned.head(5))
```

```
✓ 0/0
```

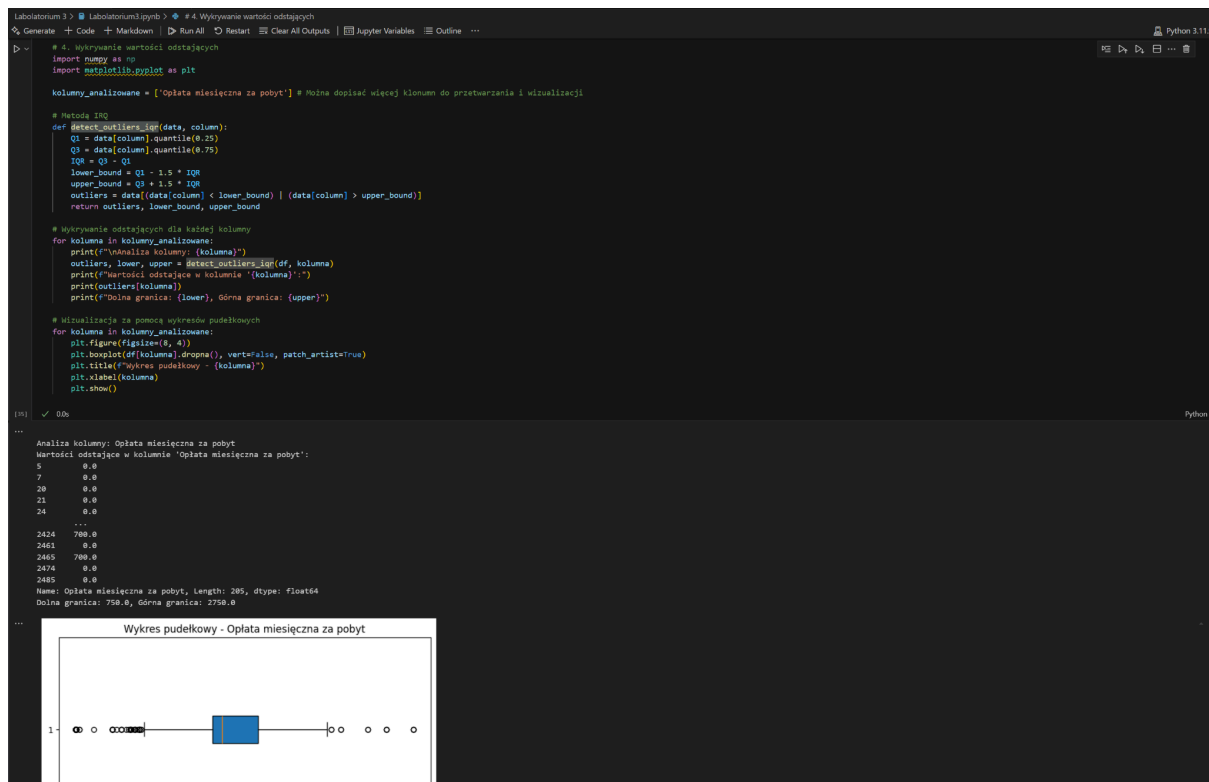
```
3.1. Liczba brakujących wartości w każdej kolumnie: Imię      0
Nazisko      0
Województwo   0
Powiat        0
Gmina         0
Kod pocztowy  0
Liczba miejsc 0
Liczba dzieci zapisanych      291
Liczba miejsc ze środków FERS 1401
Liczba dzieci zapisanych na miejsca z FERS 1401
Opłata miesięczna za pobyt      139
Opłata miesięczna na miejscach z KPO/FERS      0
Stawka godzinowa za każde godzinę powyżej 10 godzin 2152
Opłata godzinowa - podstawowa opłata ponoszona przez rodziców za pobyt dziecka (bez zniek i bez wyżywienia) 2398
Opłata za wyżywienie - miesięczna 2154
Opłata za wyżywienie - dzienna 485
dtype: int64

3.2. Procent brakujących danych w każdej kolumnie:
Imię      0.0
Nazisko    0.0
Województwo  0.0
Powiat      0.0
Gmina       0.0
dtype: float64

3.3. Kolumny usunięte z powodu dużej liczby braków (> 50%):
...
1      NaN
2      20.0
3      20.0
4      NaN
Output is truncated. View as a scrollable element or open in a text editor. Adjust cell output settings...
```

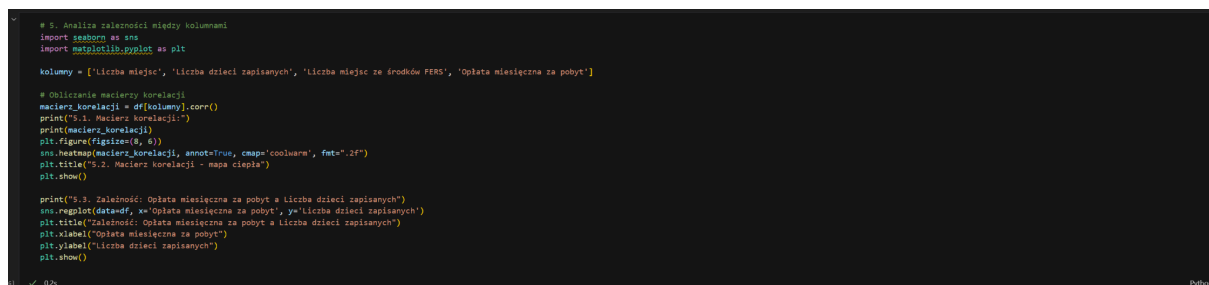
Rys. 3. Fragment kodu i wynik jego wykonania

Zidentyfikowano kolumny zawierające brakujące wartości oraz obliczono ich procentowy udział. Użyto różnych metod uzupełniania braków w danych, usunięcia kolumn z dużą liczbą braków, co pozwoliło przygotować dane do dalszej analizy.

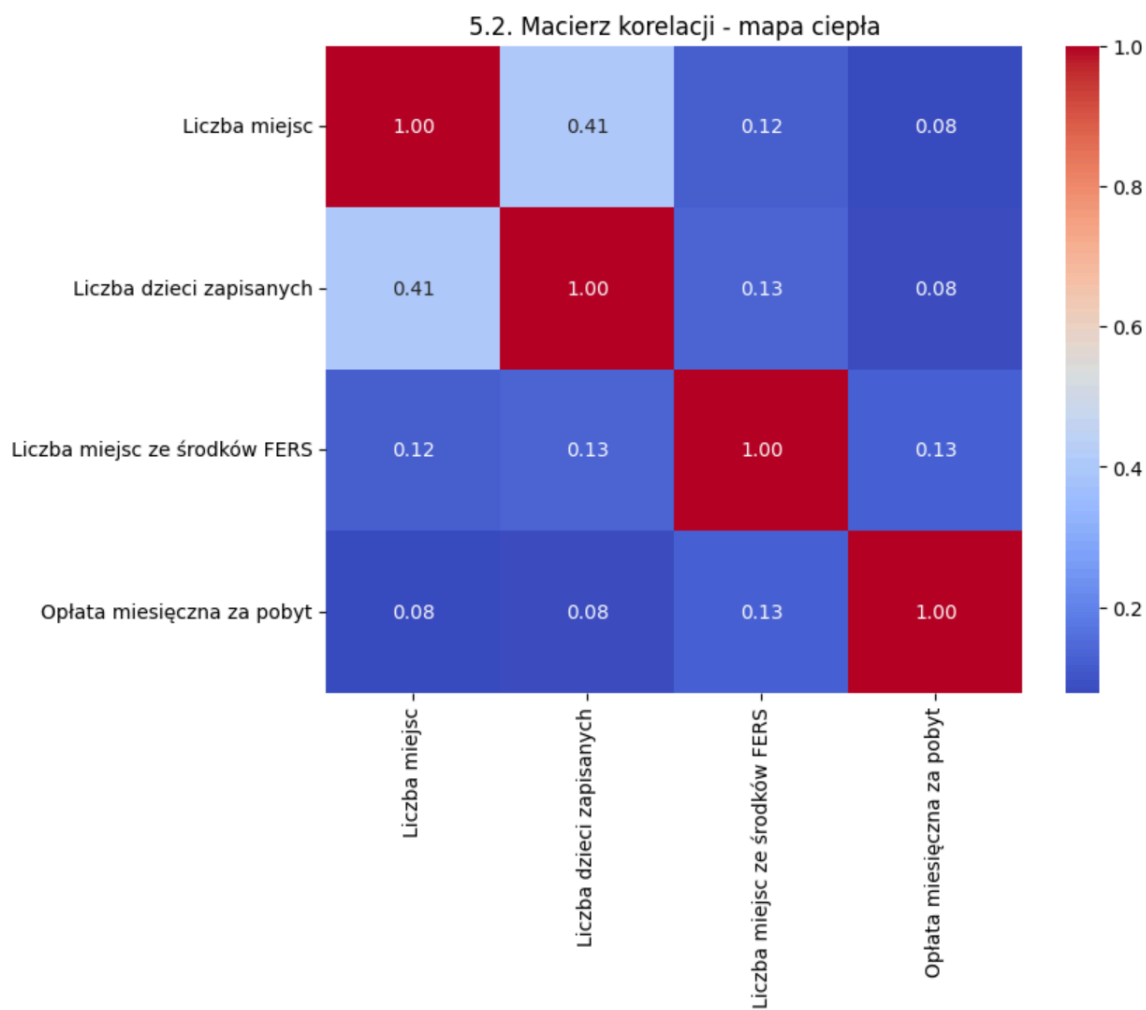


Rys. 4. Fragment kodu i wynik jego wykonania

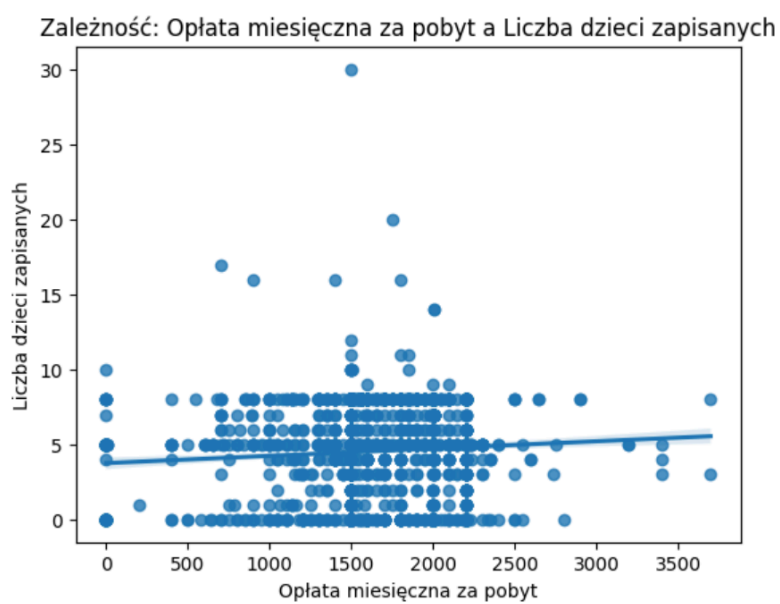
Za pomocą metody IQR wykryto wartości odstające w danych. Wizualizacja tych wartości za pomocą wykresów pudełkowych ułatwiła zrozumienie ich rozmieszczenia.



Rys. 5. Fragment kodu



Rys. 6. Macierz korelacji - mapa ciepła



Rys. 7. Macierz korelacji - mapa ciepła

Obliczono korelacje między wybranymi kolumnami liczbowymi oraz przeprowadzono analizę wizualną tych zależności. Użyto wykresów punktowych i heatmap, które pomogły w zrozumieniu powiązań pomiędzy danymi.

```
# 6. Przetwarzanie danych
df['Suma miejsc i dzieci'] = df['liczba miejsc'] + df['liczba dzieci zapisanych']

print("6.1. Nowa kolumna 'Suma miejsc i dzieci':")
print(df.head())

print("6.2. Grupowanie po województwie i obliczanie średnich:")
grouped_by_województwo = df.groupby('województwo')
print(grouped_by_województwo.head())

print("6.3. Sortowanie po kolumnie 'Opłata miesięczna za pobyt':")
sorted_df = df.sort_values(by='Opłata miesięczna za pobyt', ascending=True)
print(sorted_df.head())

print("6.4. Sortowanie po 'województwo' i 'liczba dzieci zapisanych':")
sorted_multi = df.sort_values(by=['województwo', 'liczba dzieci zapisanych'], ascending=[True, False])
print(sorted_multi.head())
```

6.1. Nowa kolumna 'Suma miejsc i dzieci':

	Imię	Nazwisko	województwo	Powiat	Gmina
0	Irabella	Wojcik	MAZOWIECKIE	Warszawa	Warszawa - Praga-Północ
1	Paulina	Pomorska	MAZOWIECKIE	Warszawa	Warszawa - Mokotów
2	Anna	Spiechow	MAZOWIECKIE	otwocki	Otwock
3	Angelika	Rogaczewska	PODKARPACKIE	Gdańsk	Gdańsk
4	Paulina	Adamus	ŁÓDZKIE	Łódź	Łódź-Śródmieście

Kod pocztowy Liczba miejsc Liczba dzieci zapisanych \

	0	1	2	3	4
0	01-484	8	5.0		
1	02-621	8	8.0		
2	05-480	8	8.0		
3	80-044	5	5.0		
4	90-266	5	8.0		

Liczba miejsc ze środków FERS Liczba dzieci zapisanych na miejsca z FERS \

	0	1	2	3	4
0	0.0	0.0	0.0		
1	0.0	0.0	0.0		
2	NaN	NaN	NaN		
3	NaN	NaN	NaN		
4	0.0	0.0	0.0		

Opłata miesięczna za pobyt Opłata miesięczna na miejscach z KPO/FERS \

	0	1
0	NaN	0.0 z1
1	1950.0	0.0 z1
...		
3387	20.0	
621	15.0	
1310	18.0	
52	16.0	

Output is truncated. View as a [scrollable element](#) or open in a [text editor](#). Adjust cell output [settings](#).

Rys. 8. Fragment kodu i wynik jego wykonania

Stworzono nową kolumnę na podstawie dwóch istniejących. Następnie dane zostały posortowane według wybranych kolumn, a także grupowane po regionie (województwo), co pozwoliło na dalsze analizy.

3. Wnioski

W trakcie przetwarzania danych zauważono, że wiele kolumn zawiera brakujące wartości. Obsługa tych braków poprzez uzupełnianie, usuwanie kolumn czy konwersję pozwala przygotować dane do dalszej analizy.

Wartości odstające często wpływają na wyniki analiz, dlatego ich identyfikacja i wizualizacja za pomocą wykresów pudełkowych jest niezbędna do poprawnej interpretacji danych.

Przeprowadzenie analizy korelacji oraz wizualizacja tych zależności za pomocą mapy ciepła pozwala na zrozumienie, jak poszczególne cechy danych są ze sobą powiązane, co może prowadzić do bardziej trafnych decyzji analitycznych.