

# Statystyczne metody rozpoznawania obrazu

## Zadanie 4

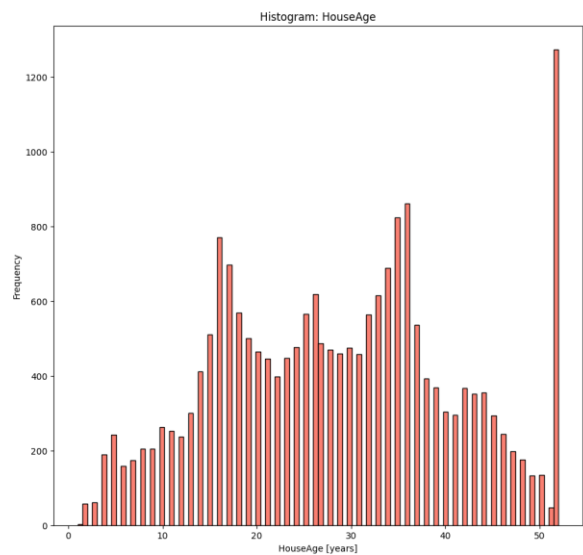
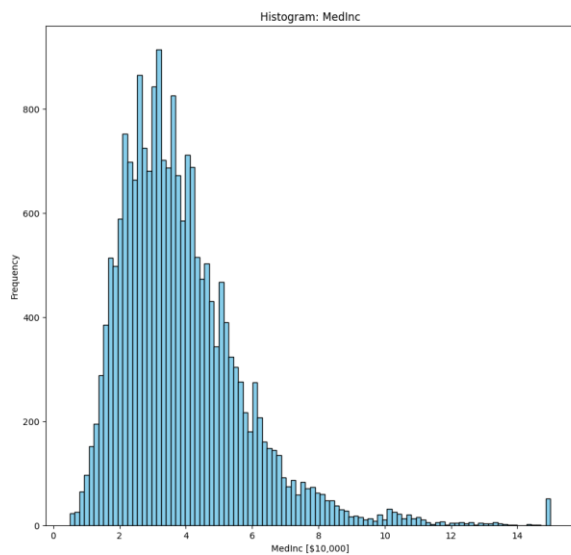
**Bartłomiej Gorzela**

### 1. Zidentyfikować wartości odstające w wybranym zbiorze danych.

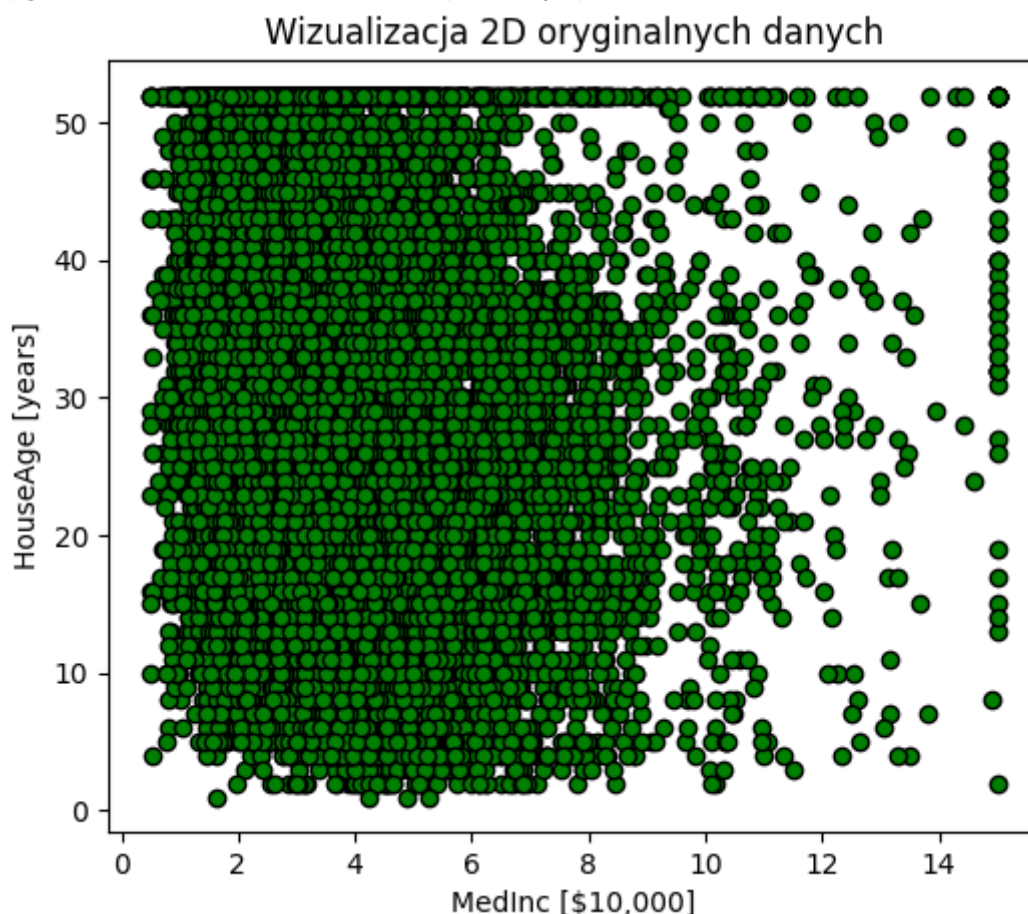
W celu wykonania zadania posłużyłem się przykładowym zbiorem danych 'California Housing Dataset' możliwym do zaimportowania bezpośrednio z biblioteki scikit-learn dostępnej w Python. Do identyfikacji wartości odstających zdecydowałem się na wybór dwóch atrybutów:

- *MedInc* (Median Income): przedstawia medianę dochodu gospodarstwa domowego.
- *HouseAge* (Median Age of Housing Units): średnia wieku wszystkich budynków mieszkalnych.

Wartości wybranych atrybutów zwizualizowałem za pomocą histogramów (argument 'bins' ustawiłem na wartość 100 w celu poprawy czytelności wizualizacji):



Oryginalne dane zwizualizowałem za pomocą wykresu 2D:

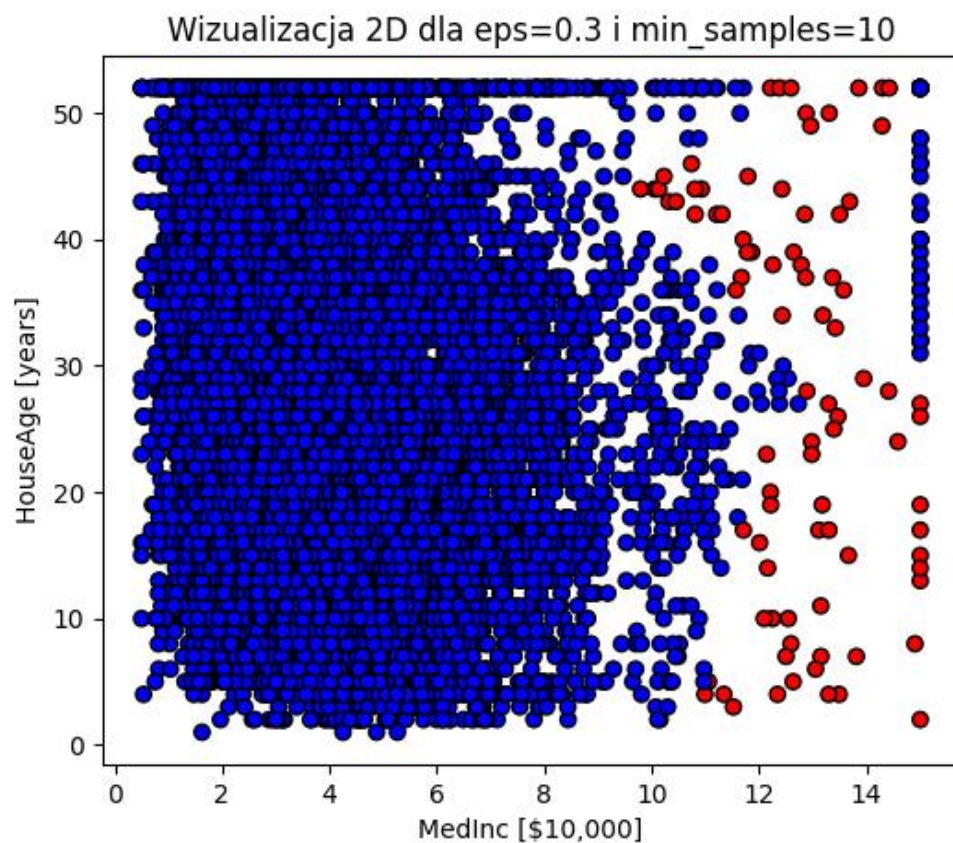


Kolejnym krokiem jaki wykonałem było zidentyfikowanie wartości odstających. Posłużyłem się w tym celu algorytmem DBSCAN (Density-Based Spatial Clustering of Applications with Noise), który umożliwia wykrywanie skupisk danych oraz jednocześnie oznaczanie punktów, które można uznać za wartości odstające. Algorytm jest dostępny w bibliotece scikit-learn.

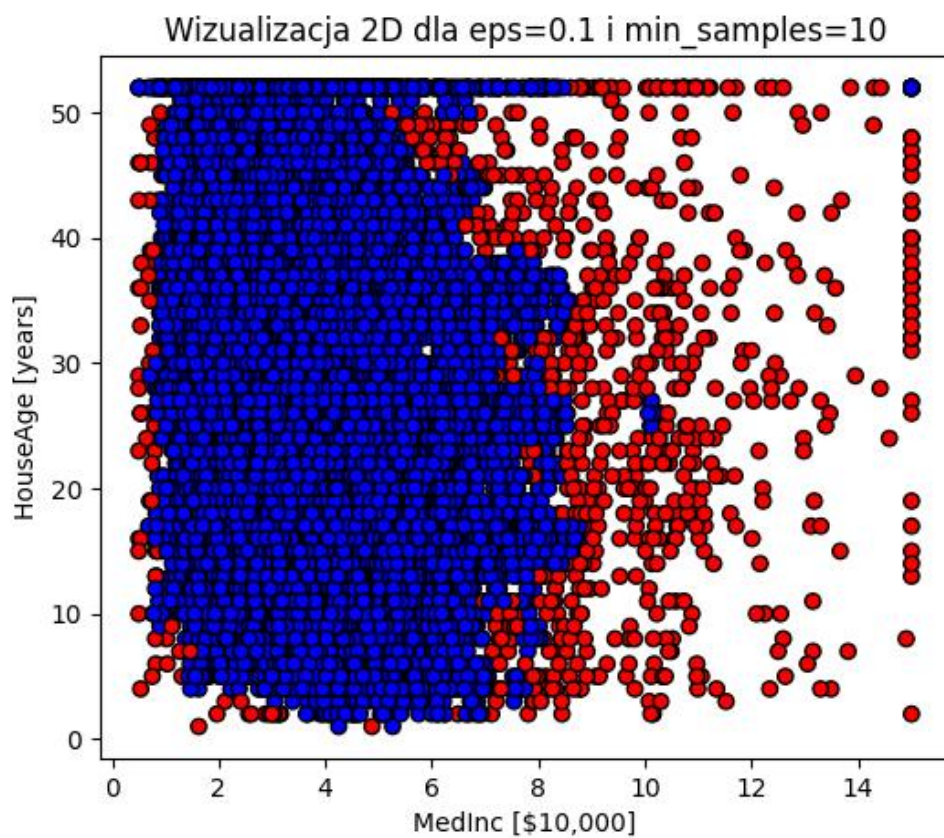
Algorytm wymagał przeprowadzenia standaryzacji danych, czego dokonałem przy użyciu obiektu `StandardScaler()`. Po przeskalowaniu danych, zastosowałem algorytm DBSCAN z parametrami: **eps=0.3** oraz **min\_samples=10**. Oznacza to, że punkt jest uznany za rdzeniowy jeśli ma co najmniej 10 sąsiadów w promieniu 0.3 jednostki (w przestrzeni standaryzowanej).

W wyniku działania DBSCAN każdemu punktowi została przypisana etykieta klastra (`dbscan_label`), przy czym wartość -1 oznacza punkt uznany za odstający.

Poddane działaniu algorytmu atrybuty ponownie zwizualizowałem za pomocą wykresu 2D, zaznaczając kolorem czerwonym punkty uznane za wartości odstające.

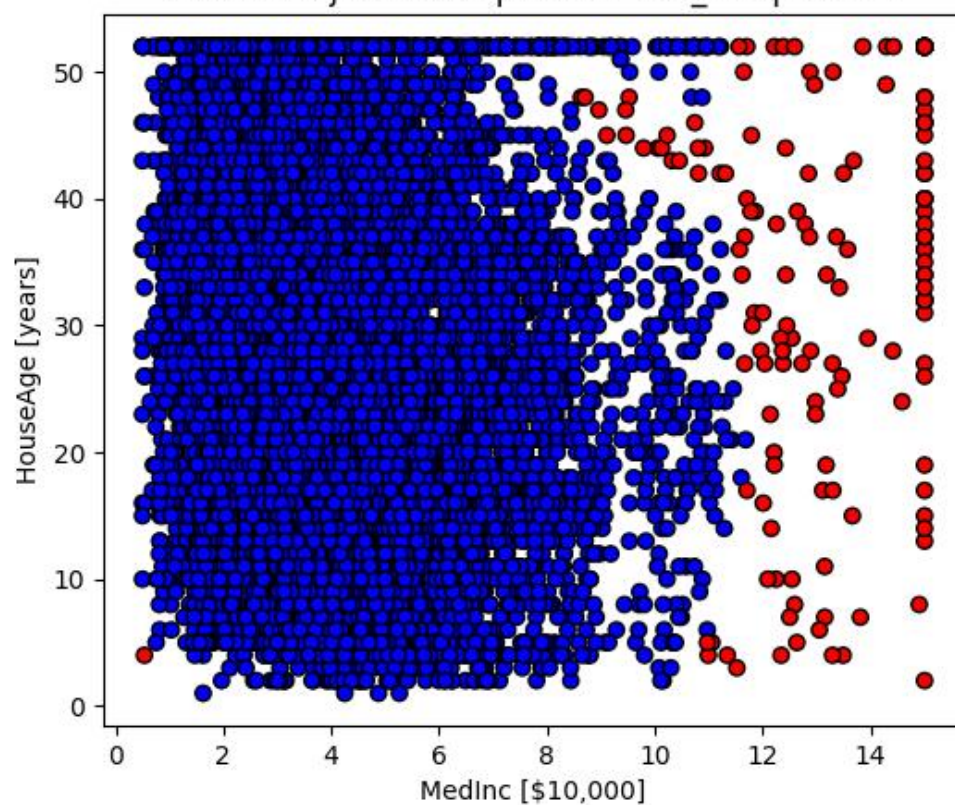


Ten sam krok powtórzyłem dla modyfikując wartości 'eps' oraz 'min\_samples'. Wyniki znajdują się poniżej:

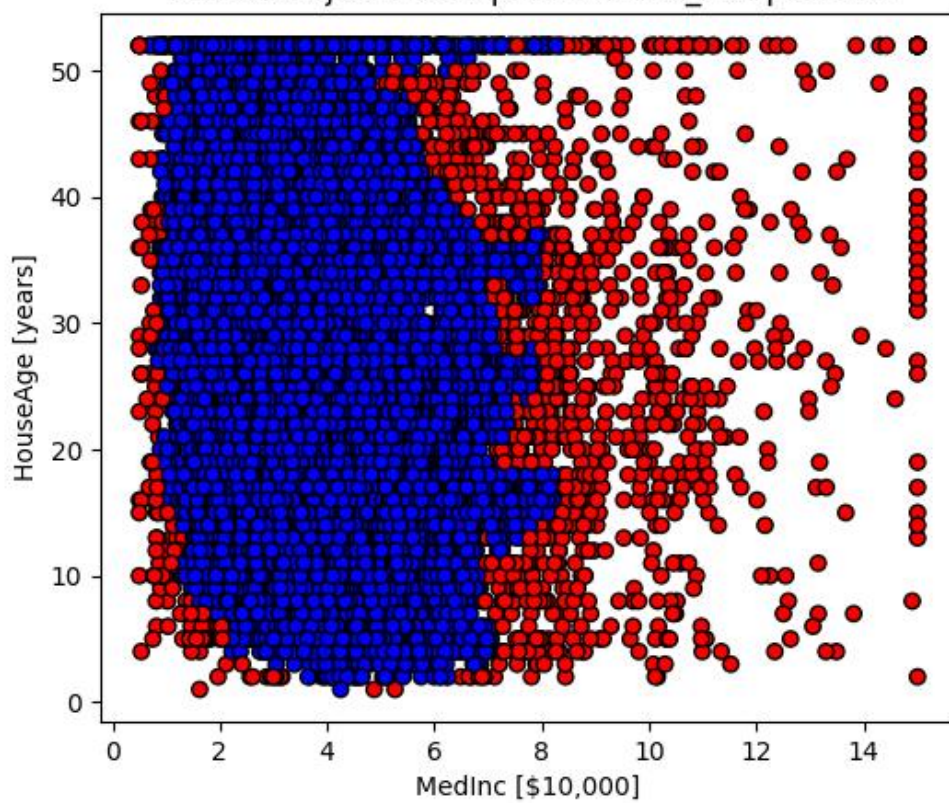




Wizualizacja 2D dla  $\text{eps}=0.3$  i  $\text{min\_samples}=15$



Wizualizacja 2D dla  $\text{eps}=0.1$  i  $\text{min\_samples}=15$



Kod źródłowy:

Bartłomiej\_Gorzela\_Zadanie\_4.py U

StatisticalPatternRecognitionMethods > Task4 > Bartłomiej\_Gorzela\_Zadanie\_4.py > ...

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 from sklearn import datasets
4 from sklearn.preprocessing import StandardScaler
5 from sklearn.cluster import DBSCAN
6
7 california = datasets.fetch_california_housing()
8 california_data = california.data
9 california_df = pd.DataFrame(california_data, columns=california.feature_names)
10
11 # Wizualizacja danych za pomocą histogramów
12 plt.figure(figsize=(12, 5))
13 plt.subplot(1, 2, 1)
14 plt.hist(california_df['MedInc'], bins=100, color='skyblue', edgecolor='black')
15 plt.title(f'Histogram: MedInc')
16 plt.xlabel('MedInc [$10,000]')
17 plt.ylabel('Frequency')
18 plt.subplot(1, 2, 2)
19 plt.hist(california_df['HouseAge'], bins=100, color='salmon', edgecolor='black')
20 plt.title(f'Histogram: HouseAge')
21 plt.xlabel('HouseAge [years]')
22 plt.ylabel('Frequency')
23 plt.tight_layout()
24 plt.show()
25
26 # Wizualizacja 2D oryginalnych danych
27 plt.figure(figsize=(6, 5))
28 plt.scatter(california_df['MedInc'], california_df['HouseAge'], color='green', edgecolor='k')
29 plt.xlabel('MedInc [$10,000]')
30 plt.ylabel('HouseAge [years]')
31 plt.title('Wizualizacja 2D oryginalnych danych')
32 plt.show()
33
34 # DBSCAN
35 attributes_for_scaling = california_df[['MedInc', 'HouseAge']]
36 scaler = StandardScaler()
37 attributes_scaled = scaler.fit_transform(attributes_for_scaling)
38 dbscan = DBSCAN(eps=0.3, min_samples=10)
39 # dbscan = DBSCAN(eps=0.1, min_samples=10)
40 # dbscan = DBSCAN(eps=0.3, min_samples=15)
41 # dbscan = DBSCAN(eps=0.1, min_samples=15)
42
43 labels = dbscan.fit_predict(attributes_scaled)
44 california_df['dbscan_label'] = labels
45
46 # Wizualizacja 2D z etykietami DBSCAN
47 plt.figure(figsize=(6, 5))
48 colors = ['red' if label == -1 else 'blue' for label in labels]
49 plt.scatter(california_df['MedInc'], california_df['HouseAge'], c=colors, edgecolor='k')
50 plt.xlabel('MedInc [$10,000]')
51 plt.ylabel('HouseAge [years]')
52 plt.title('Wizualizacja 2D dla eps=0.3 i min_samples=10')
53 # plt.title('Wizualizacja 2D dla eps=0.1 i min_samples=10')
54 # plt.title('Wizualizacja 2D dla eps=0.3 i min_samples=15')
55 # plt.title('Wizualizacja 2D dla eps=0.1 i min_samples=15')
56 plt.show()
```