

# Analiza zbioru danych mieszkania

## Wielowymiarowa analiza danych

Przemysław Peplinski, Wiktor Galewski, Mikołaj Zalewski

2026-01-18

## Spis treści

<b>1</b>	<b>Wstęp</b>	<b>2</b>
1.1	Autorzy . . . . .	2
1.2	Motywacja i cele . . . . .	2
1.3	Opis projektu . . . . .	2
<b>2</b>	<b>Opis danych</b>	<b>2</b>
<b>3</b>	<b>Czyszczenie danych</b>	<b>3</b>
3.1	Walidacja . . . . .	3
3.2	Obsługa braków wartości . . . . .	4
3.3	Obsługa obserwacji odstających . . . . .	4
<b>4</b>	<b>Analiza eksploracyjna</b>	<b>4</b>
4.1	Statystyki opisowe . . . . .	4
4.2	Analiza rozkładów . . . . .	5
4.3	Analiza korelacji . . . . .	8
4.4	Analiza braków danych . . . . .	9
4.5	Analiza obserwacji odstających . . . . .	10
<b>5</b>	<b>Obróbka zmiennych</b>	<b>11</b>
5.1	Kodowanie zmiennych kategorycznych . . . . .	11
5.2	Normalizacja zmiennych numerycznych . . . . .	11

<b>6</b>	<b>2 metody</b>	<b>11</b>
6.1	Analiza głównych składowych (PCA) . . . . .	11
6.2	Analiza korespondencji (CA) . . . . .	13
<b>7</b>	<b>Wizualizacja metod</b>	<b>13</b>
7.1	Wizualizacje PCA . . . . .	13
7.2	Wizualizacje CA . . . . .	13
<b>8</b>	<b>Wnioski</b>	<b>13</b>

# 1 Wstęp

## 1.1 Autorzy

Autorami są... Wkład w projekt prezentował się następująco...

## 1.2 Motywacja i cele

Coś tam...

## 1.3 Opis projektu

Coś tam...

# 2 Opis danych

Dane pochodzą.../Dane zostały zebrane...

```
## Rows: 106,169
## Columns: 19
## $ Link          <chr> "https://domy.pl/mieszkanie/gdansk-siedlce-malczews~
## $ `Data dodania` <dtm> 2021-02-02 10:24:00, 2021-02-02 11:08:14, 2021-02-~
## $ `Data modyfikacji` <dtm> NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, NA~
## $ Tytuł         <chr> "Trzypokojowe mieszkanie na sprzedaż:", "Gdańsk, Si~
## $ Opis          <chr> "Młyny Gdańskie Gdańsk, ul. Malczewskiego Młyny Gda~
## $ Cena          <dbl> NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, NA,~
## $ Powierzchnia  <dbl> 60.02, 89.63, 60.49, 64.89, 104.66, 72.04, 75.21, 6~
## $ `Cena za metr` <dbl> NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, NA,~
```

```
## $ `Liczba pokoi`      <dbl> 3, 4, 3, 3, 4, 4, 4, 3, 3, 3, 4, 3, 1, 2, 3, 3, 3, ~
## $ Piętro              <dbl> 4, 1, 0, 0, 1, 3, 1, 2, 3, 3, 2, 1, 3, 0, 1, 3, 4, ~
## $ `Liczba pięter`     <dbl> NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, ~
## $ `Rok budowy`        <dbl> 2023, 2023, 2023, 2023, 2023, 2023, 2023, 2023, 202~
## $ Ocena               <dbl> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, ~
## $ Dzielnica           <chr> "Siedlce", "Siedlce", "Siedlce", "Siedlce", "Siedlc~
## $ `Rodzaj budynku`    <chr> "Apartamentowiec", "Pozostałe", "Pozostałe", "Pozos~
## $ Rynek               <chr> "Pierwotny", "Pierwotny", "Pierwotny", "Pierwotny",~
## $ Ogłoszeniodawca     <chr> "Agencja", "Agencja", "Agencja", "Agencja", "Agencj~
## $ Portal              <chr> "Domy", "Morizon", "Morizon", "Morizon", "Morizon",~
## $ `Numer telefonu`    <chr> NA, "585055401", "585055401", "585055401", "5850554~
```

Opisać zmienne...

Wśród zmiennych niepotrzebnych do analizy można wymienić link, tytuł, opis oraz numer telefonu. Zmienne te zostaną zatem usunięte.

## 3 Czyszczenie danych

### 3.1 Walidacja

Tak

```
## Rows: 106,169
## Columns: 19
## $ Link                <chr> "https://domy.pl/mieszkanie/gdansk-siedlce-malczews~
## $ `Data dodania`      <dtm> 2021-02-02 10:24:00, 2021-02-02 11:08:14, 2021-02-~
## $ `Data modyfikacji`  <dtm> NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, ~
## $ Tytuł               <chr> "Trzypokojowe mieszkanie na sprzedaż:", "Gdańsk, Si~
## $ Opis                <chr> "Młyny Gdańskie Gdańsk, ul. Malczewskiego Młyny Gda~
## $ Cena                <dbl> NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, ~
## $ Powierzchnia        <dbl> 60.02, 89.63, 60.49, 64.89, 104.66, 72.04, 75.21, 6~
## $ `Cena za metr`      <dbl> NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, ~
## $ `Liczba pokoi`      <dbl> 3, 4, 3, 3, 4, 4, 4, 3, 3, 3, 4, 3, 1, 2, 3, 3, 3, ~
## $ Piętro              <dbl> 4, 1, 0, 0, 1, 3, 1, 2, 3, 3, 2, 1, 3, 0, 1, 3, 4, ~
## $ `Liczba pięter`     <dbl> NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, ~
## $ `Rok budowy`        <dbl> 2023, 2023, 2023, 2023, 2023, 2023, 2023, 2023, 202~
## $ Ocena               <dbl> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, ~
## $ Dzielnica           <chr> "Siedlce", "Siedlce", "Siedlce", "Siedlce", "Siedlc~
## $ `Rodzaj budynku`    <chr> "Apartamentowiec", "Pozostałe", "Pozostałe", "Pozos~
## $ Rynek               <chr> "Pierwotny", "Pierwotny", "Pierwotny", "Pierwotny",~
## $ Ogłoszeniodawca     <chr> "Agencja", "Agencja", "Agencja", "Agencja", "Agencj~
## $ Portal              <chr> "Domy", "Morizon", "Morizon", "Morizon", "Morizon",~
## $ `Numer telefonu`    <chr> NA, "585055401", "585055401", "585055401", "5850554~
```

## 3.2 Obsługa braków wartości

Kilka zmiennych do usunięcia, dla niektórych trzeba usunąć obserwacje, w innych imputować wg obliczeń.

```
## missForest iteration 1 in progress...done!
## estimated error(s): 1.379711e-05 0.03532471
## difference(s): 9.168541e-11 0.03261746
## time: 29.01 seconds
##
## missForest iteration 2 in progress...done!
## estimated error(s): 1.235244e-05 0.03504698
## difference(s): 5.99119e-13 0.004266807
## time: 28.19 seconds
##
## missForest iteration 3 in progress...done!
## estimated error(s): 1.242352e-05 0.03537198
## difference(s): 2.322088e-13 0.004550982
## time: 28.44 seconds
##
## missForest iteration 4 in progress...done!
## estimated error(s): 1.233419e-05 0.0350588
## difference(s): 2.511408e-13 0.004973066
## time: 29.28 seconds
```

## 3.3 Obsługa obserwacji odstających

Jako tako zostawiamy

# 4 Analiza eksploracyjna

## 4.1 Statystyki opisowe

##	Data dodania	Data modyfikacji	Cena	Powierzchnia
##	Min. :18666	Min. :0.0000	Min. : 33000	Min. : 14.00
##	1st Qu.:19733	1st Qu.:0.0000	1st Qu.: 526925	1st Qu.: 40.00
##	Median :20058	Median :0.0000	Median : 669000	Median : 50.60
##	Mean :19980	Mean :0.3913	Mean : 794244	Mean : 54.38
##	3rd Qu.:20249	3rd Qu.:1.0000	3rd Qu.: 899000	3rd Qu.: 64.60
##	Max. :20437	Max. :1.0000	Max. :13000000	Max. :553.00
##				
##	Cena za metr	Liczba pokoi	Piętro	Rok budowy

```

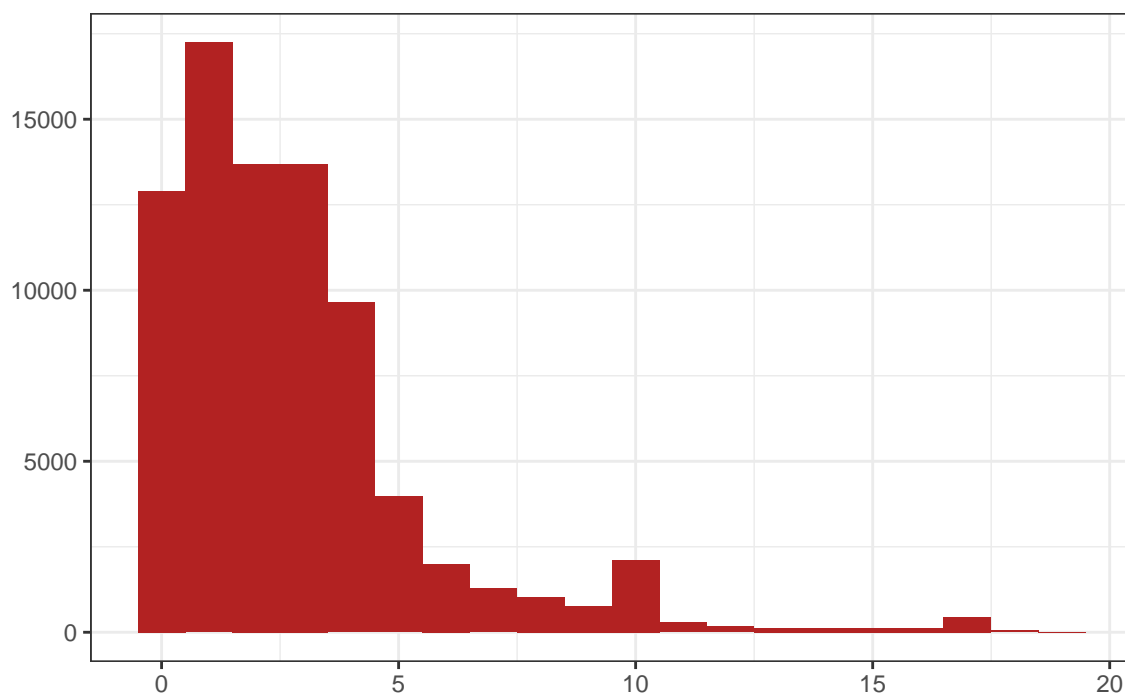
## Min.      : 1121    Min.      : 1.000    Min.      : 0.000    Min.      :1780
## 1st Qu.:11200    1st Qu.: 2.000    1st Qu.: 1.000    1st Qu.:1986
## Median :13414    Median : 2.000    Median : 2.000    Median :2020
## Mean    :14664    Mean    : 2.535    Mean    : 2.786    Mean     :2004
## 3rd Qu.:16667    3rd Qu.: 3.000    3rd Qu.: 4.000    3rd Qu.:2024
## Max.     :62298    Max.     :24.000    Max.     :19.000    Max.     :2028
##
##                               Dzielnica                Rodzaj budynku        Rynek
## Śródmieście           :11064    Apartamentowiec:24873    Min.      :0.0000
## Ujeścisko-Łostowice:10920    Blok                :47654    1st Qu.:0.0000
## Jasień                : 9296    Kamienica           : 7236    Median :0.0000
## Wrzeszcz              : 5576                                Mean    :0.3485
## Przymorze             : 5052                                3rd Qu.:1.0000
## Siedlce               : 4915                                Max.     :1.0000
## (Other)               :32940
## Ogłoszeniodawca                Portal
## Min.      :0.0000    Otodom                :25920
## 1st Qu.:1.0000    Trojmiasto            :16197
## Median :1.0000    Okolica               : 7993
## Mean    :0.8734    Nieruchomosci-online: 7119
## 3rd Qu.:1.0000    Allegro               : 6505
## Max.     :1.0000    Adresowo              : 3257
##                               (Other)                :12772

```

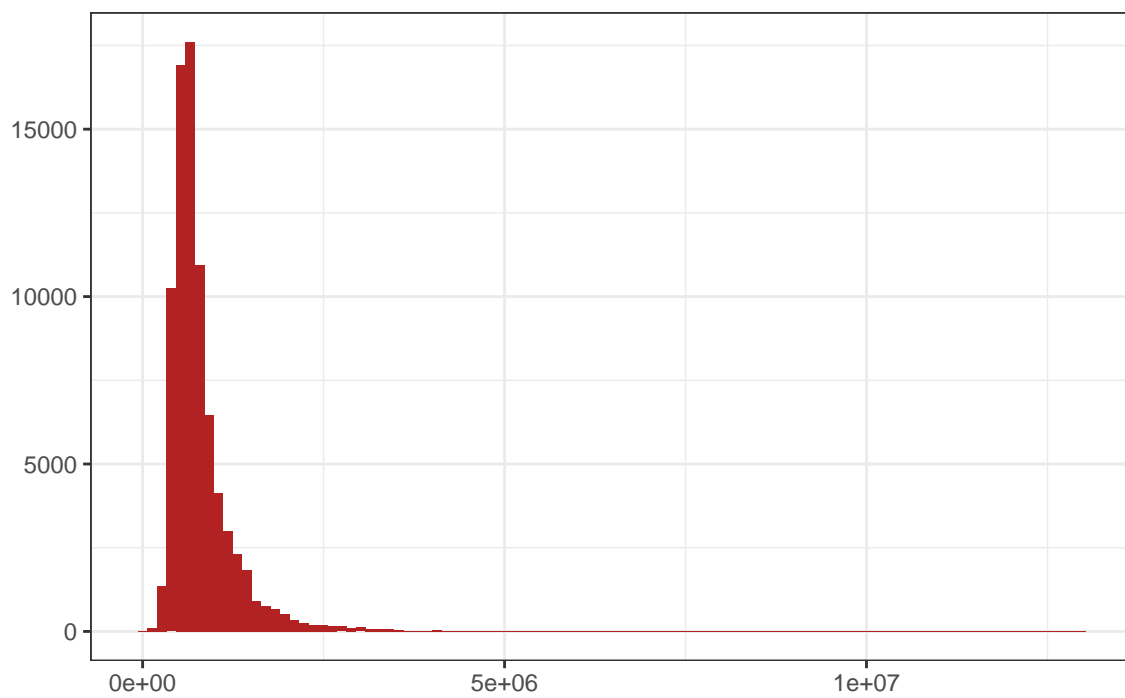
## 4.2 Analiza rozkładów

Histogramy zmiennych numerycznych, wykresy kolumnowe zmiennych kategorycznych

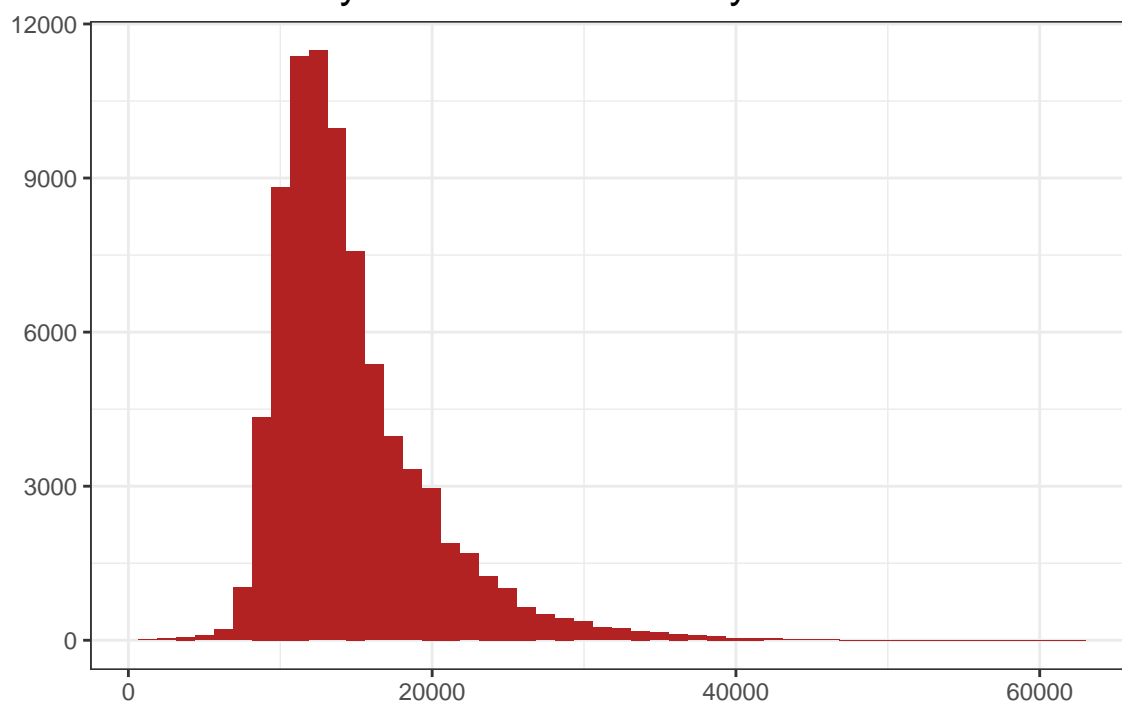
Rozkład pieter mieszkán



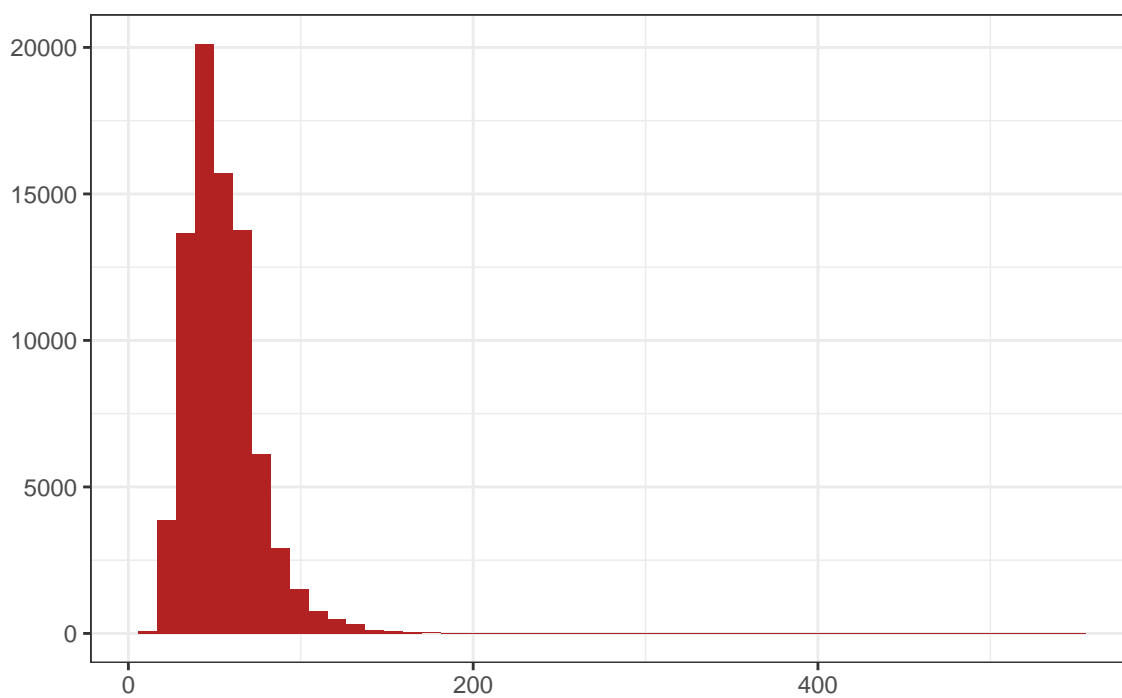
Rozklad cen mieszkán w zł

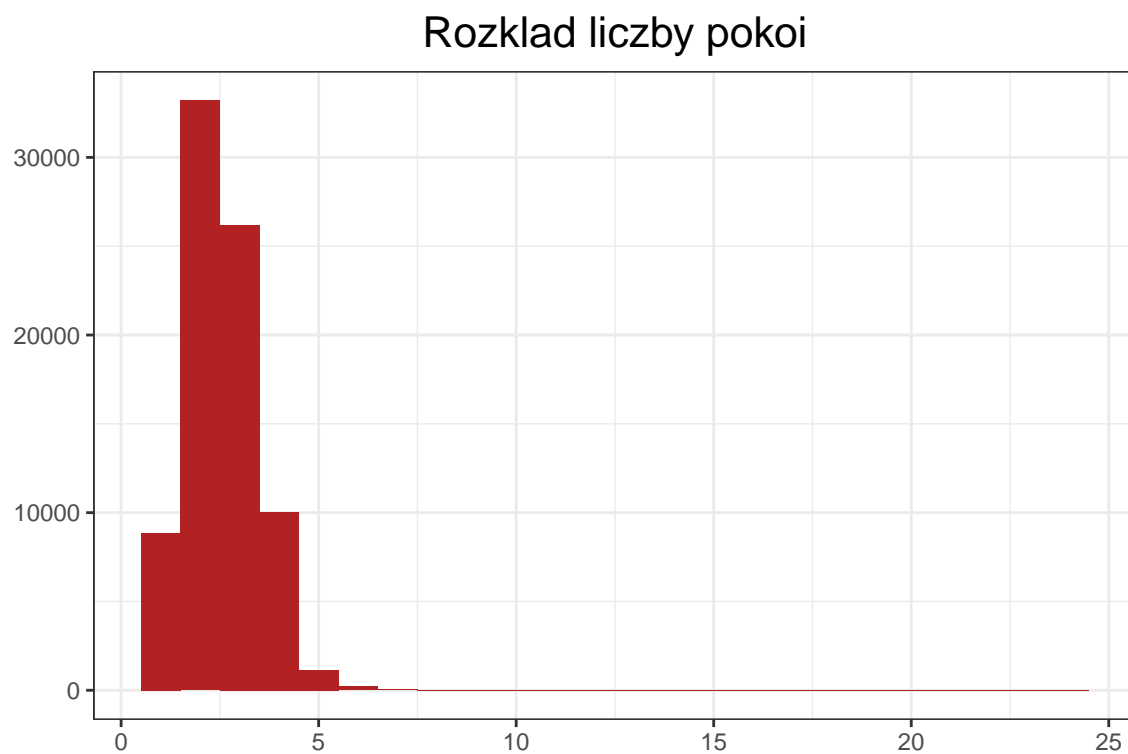


Rozkład ceny za metr kwadratowy mieszkań w zł/m<sup>2</sup>



Rozkład powierzchni mieszkań w m<sup>2</sup>

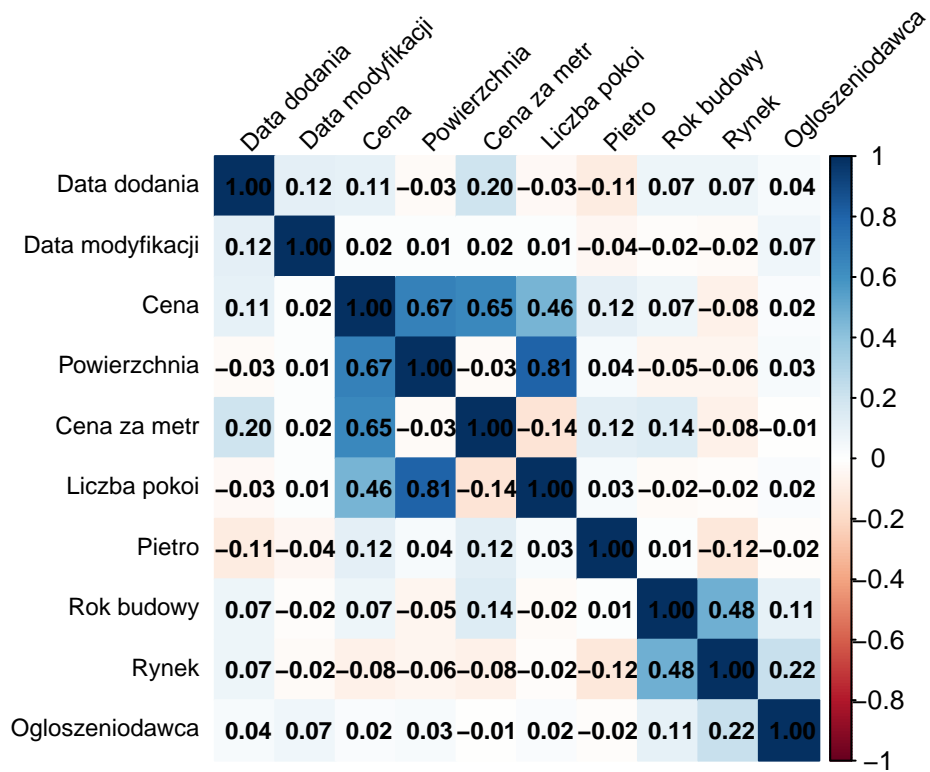




## 4.3 Analiza korelacji

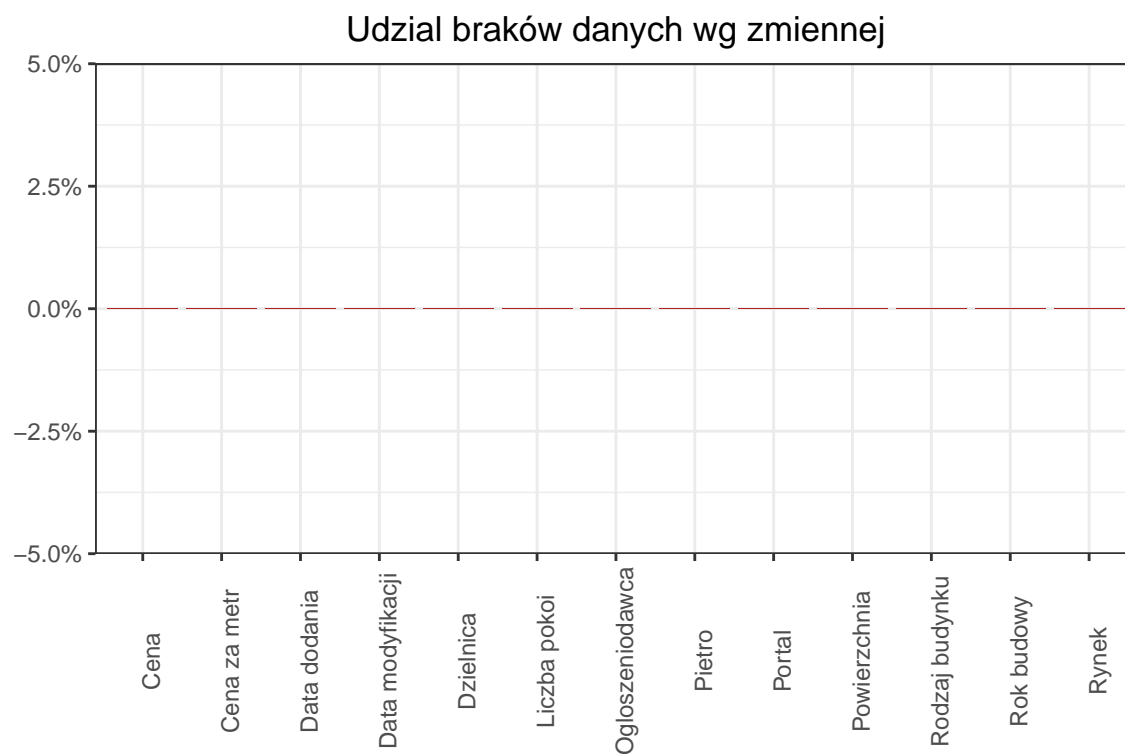
Macierz korelacji





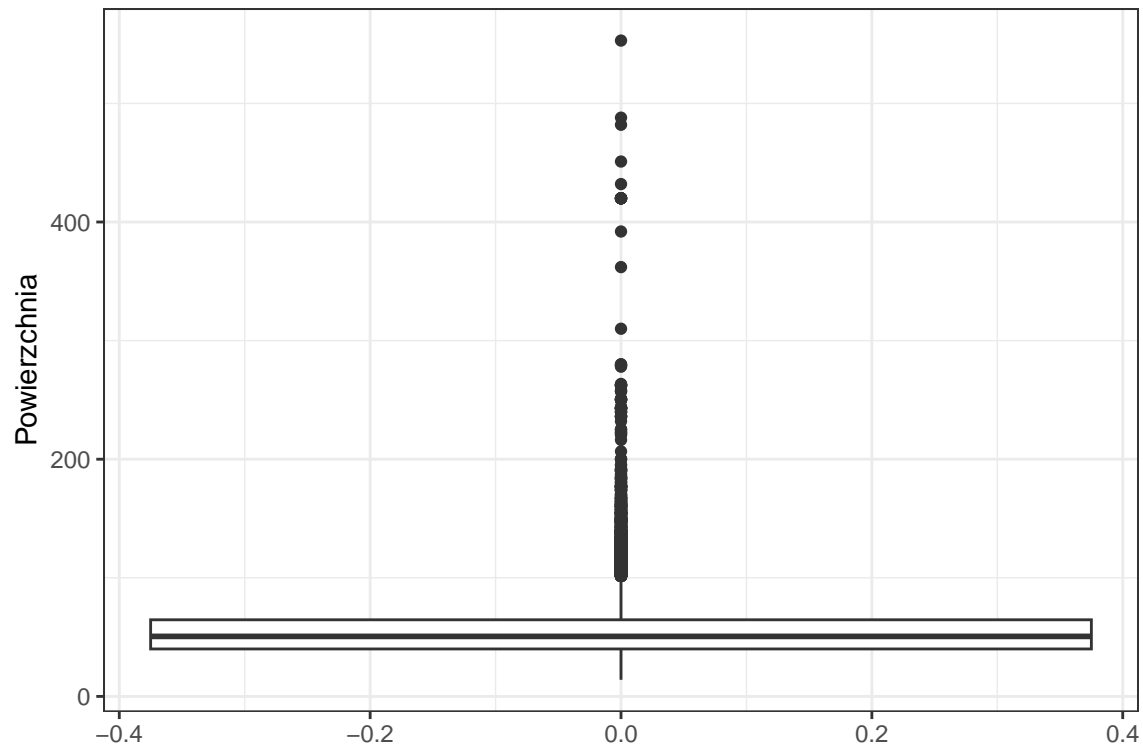
## 4.4 Analiza braków danych

Krótkie podsumowanie braków danych (wykres słupkowy?)



## 4.5 Analiza obserwacji odstających

Zrobić boxploty



## 5 Obróbka zmiennych

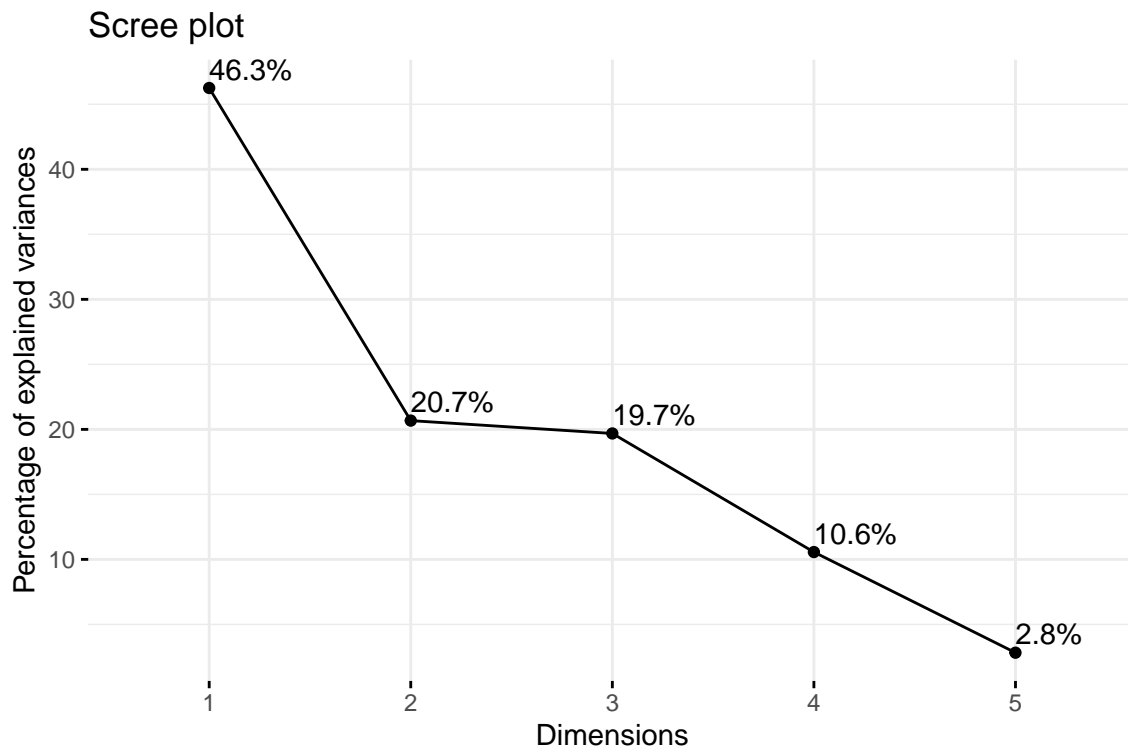
### 5.1 Kodowanie zmiennych kategoriycznych

### 5.2 Normalizacja zmiennych numerycznych

## 6 2 metody

### 6.1 Analiza głównych składowych (PCA)

```
## $chisq
## [1] 167.9269
##
## $p.value
## [1] 7.451025e-31
##
## $df
## [1] 10
```



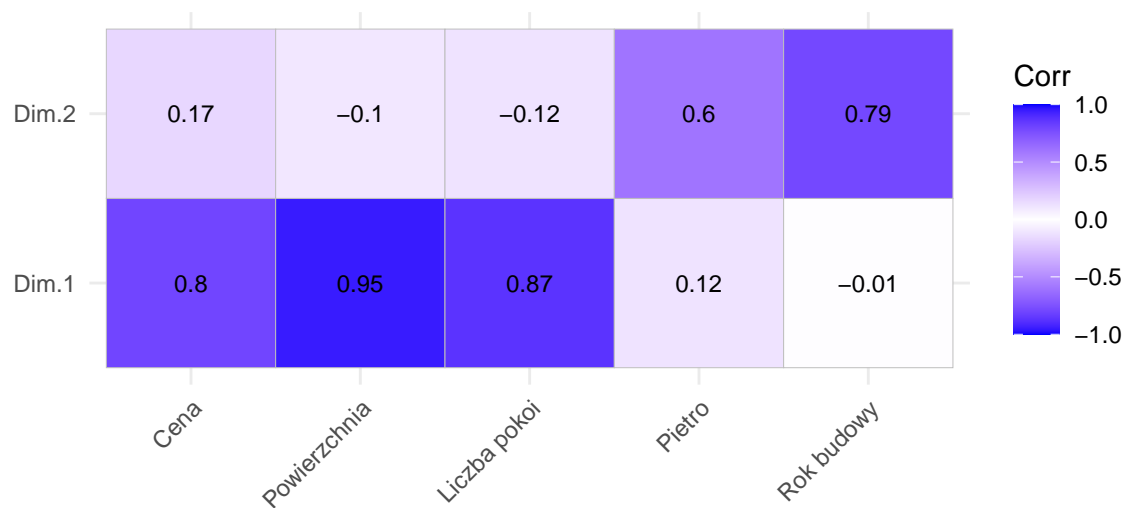
##	eigenvalue	variance.percent	cumulative.variance.percent
## Dim.1	2.3127535	46.255070	46.25507
## Dim.2	1.0336234	20.672469	66.92754
## Dim.3	0.9842327	19.684654	86.61219
## Dim.4	0.5281741	10.563482	97.17567
## Dim.5	0.1412163	2.824326	100.00000

Korzystając z kryterium Keisera można określić, że mieszkania można określić za pomocą dwóch zmiennych, ponieważ dla nich wartości własne są większe od 1. Taki sam wniosek można wyciągnąć na podstawie oceny wykresu osypiska, gdyż po drugim wymiarze zmiana wariancji jest znacznie mniejsza.

## 6.2 Analiza korespondencji (CA)

# 7 Wizualizacja metod

## 7.1 Wizualizacje PCA



## 7.2 Wizualizacje CA

# 8 Wnioski