# IMMO ELIZA
# DATA ANALYSIS

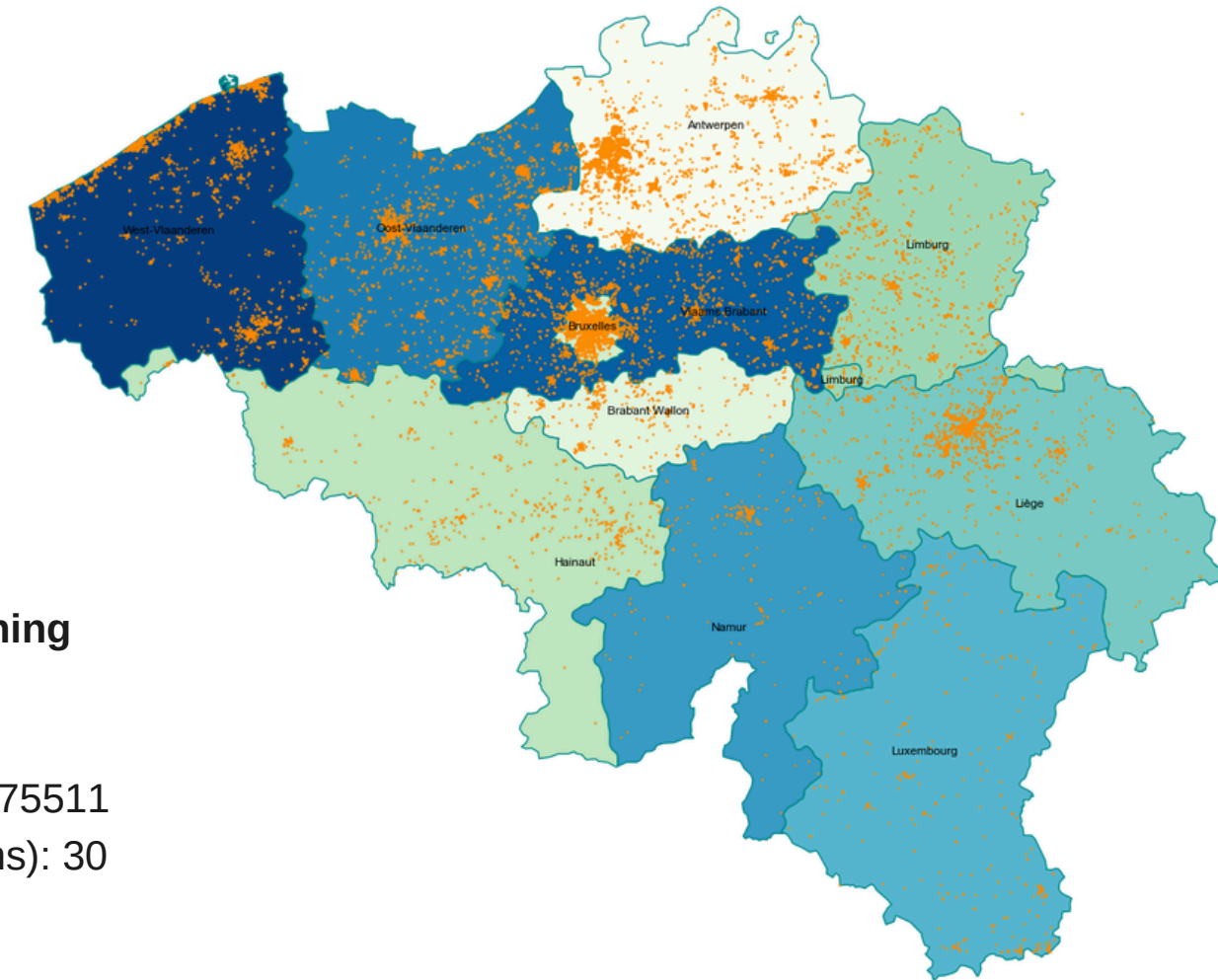TEAM 5

Yusra

Zelim

Rasmita

Muntadher

# OBJECTIVE

Perform an initial analysis of the scraped data to create visualizations and establish the foundation for the Machine Learning model.

# OBSERVATIONS AND FEATURES

**Observations**: These are the rows in your dataset, representing each individual instance .
**Features**: These are the columns in your dataset, representing the attributes or characteristics of each observation. For example, features in a real estate dataset could include "location," "price," "number of bedrooms.
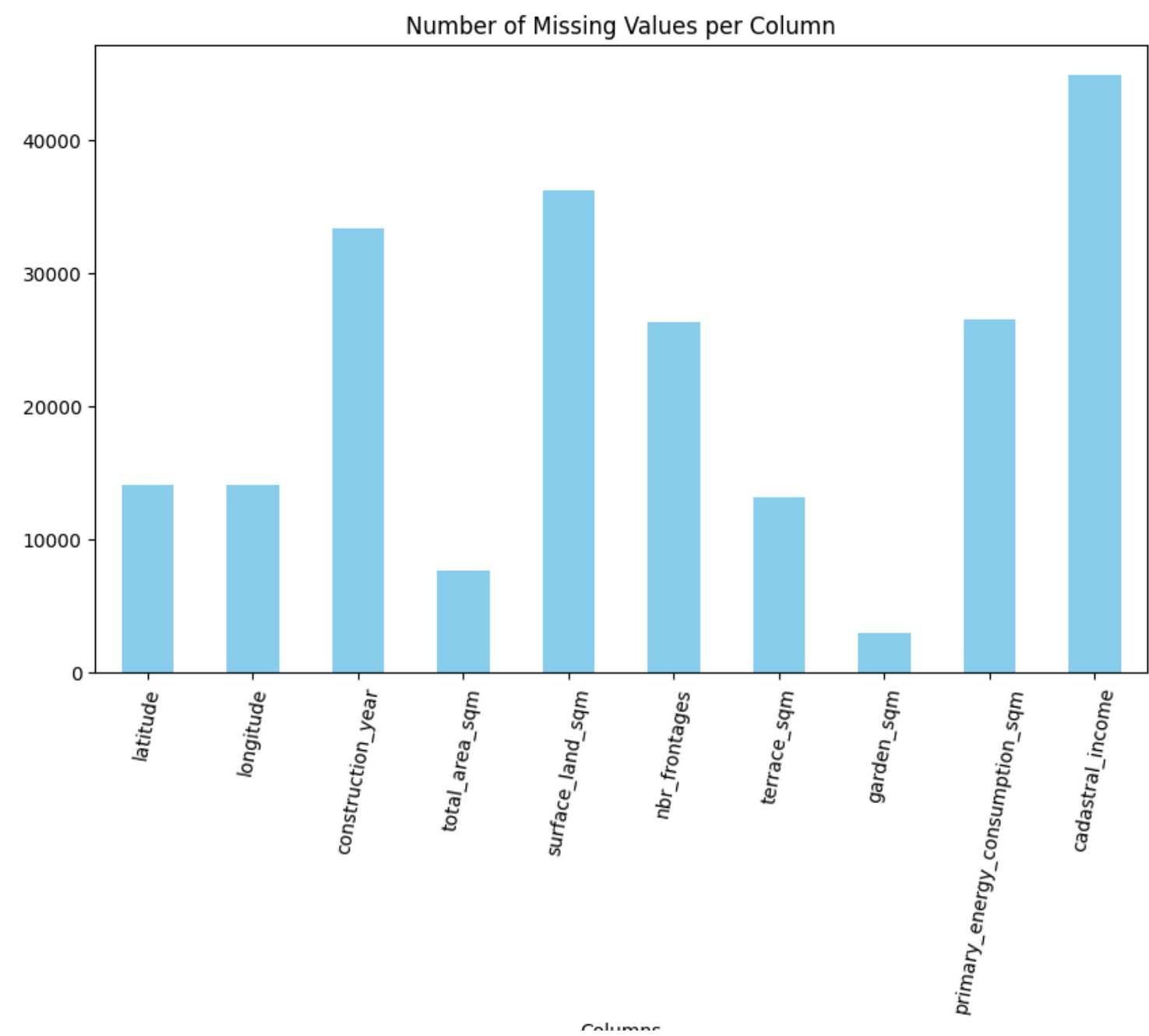


**Before Data Cleaning**

Number of (rows): 75511
Number of (columns): 30

**After Data Cleaning**

Number of (rows): 17258
Number of (columns): 23

## THE PROPORTION OF MISSING VALUES PER COLUMN.

- Construction year       44%
- Total area sqm          10.08%
- Surface land sqm        48.01%
- Frontages               34.89%
- Terrace sqm             17.40%
- Garden sqm              3.89%
- Latitude                18.67%
- Longitude               18.67%
- Primary Energy Consumption sqm 35.18%
- Cadastral Income 59.55%



Number of Missing Values per Column

# VARIABLES YOU WOULD  DELETE .

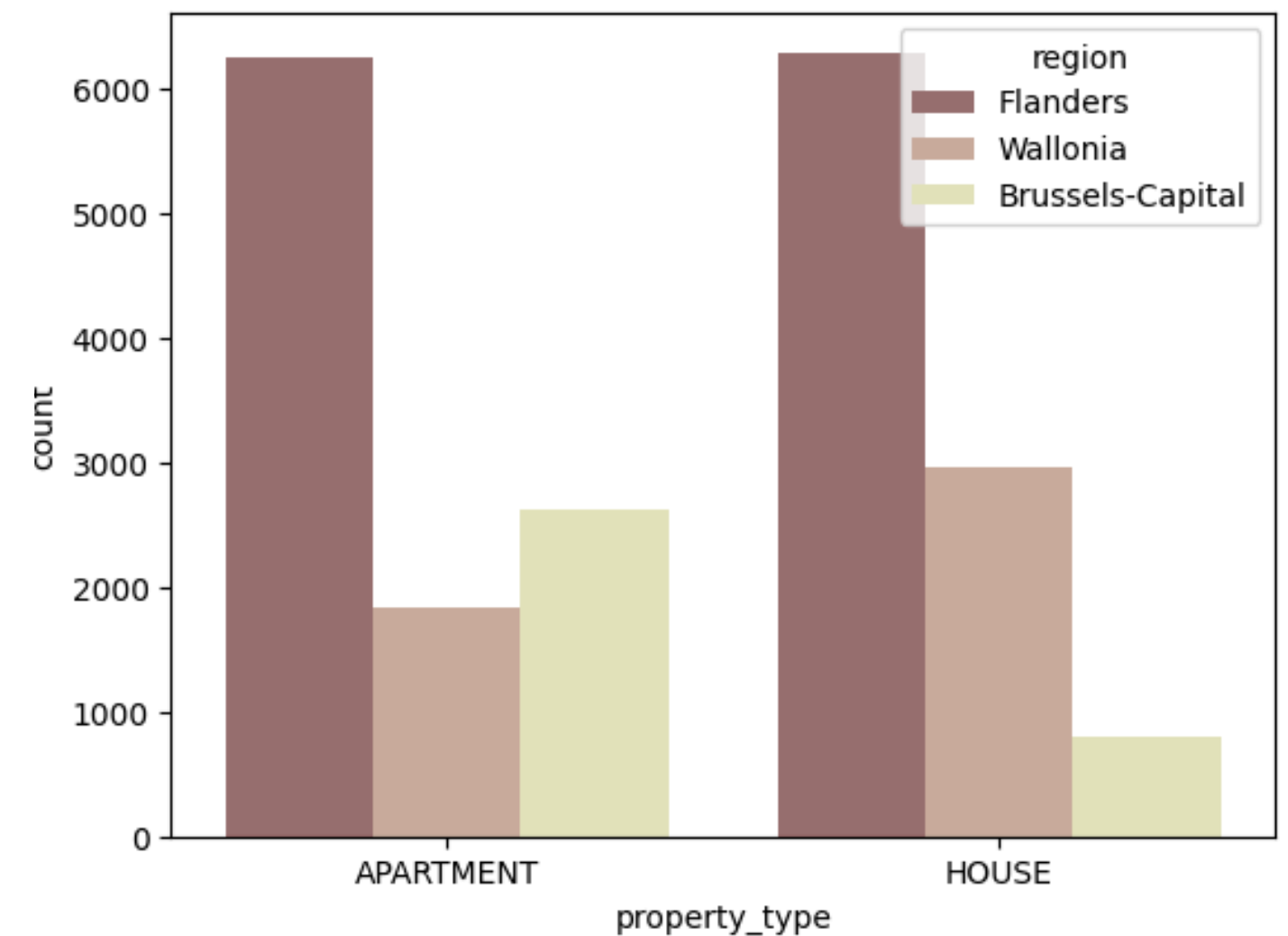**Variable Deletion in Real Estate Dataset for Immo Eliza**

During this review, certain variables were identified as <u>redundant</u> or <u>irrelevan</u>t for our specific objectives and were therefore removed. Also variables with a high percentage of missing values (e.g., more than 40%).
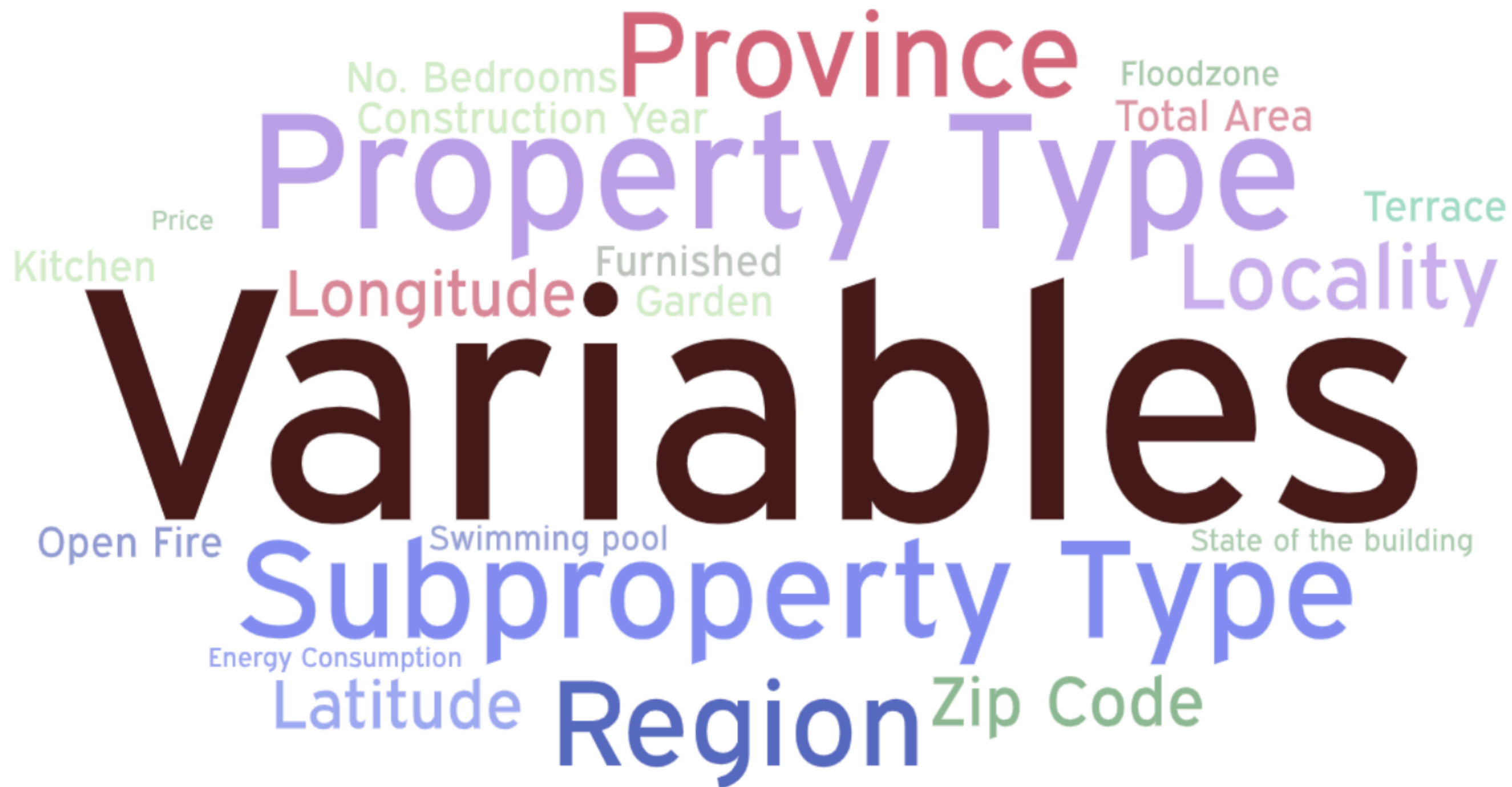
1. Surface land sqm
2. Frontages
3. Cadastral Home
4. Terrace FL
5. Garden FL
6. Primary Energy Consumption
7. Heating type

## REMOVING THESE VARIABLES IS EXPECTED TO

- **Improve Model Simplicity:** Fewer variables reduce model complexity, which can lead to better interpretability and faster computation.

- **Focus on Key Features:** Concentrating on variables that have a more direct impact on property prices will help refine the accuracy of predictions.

# QUANTITATIVE

- Price
- Total Area
- No. of bedrooms
- Terrace(sqm)
- Garden(sqm)
- Primary Energy Consumption per sqm

# QUALITATIVE

## Nominal

- ID
- Property Type
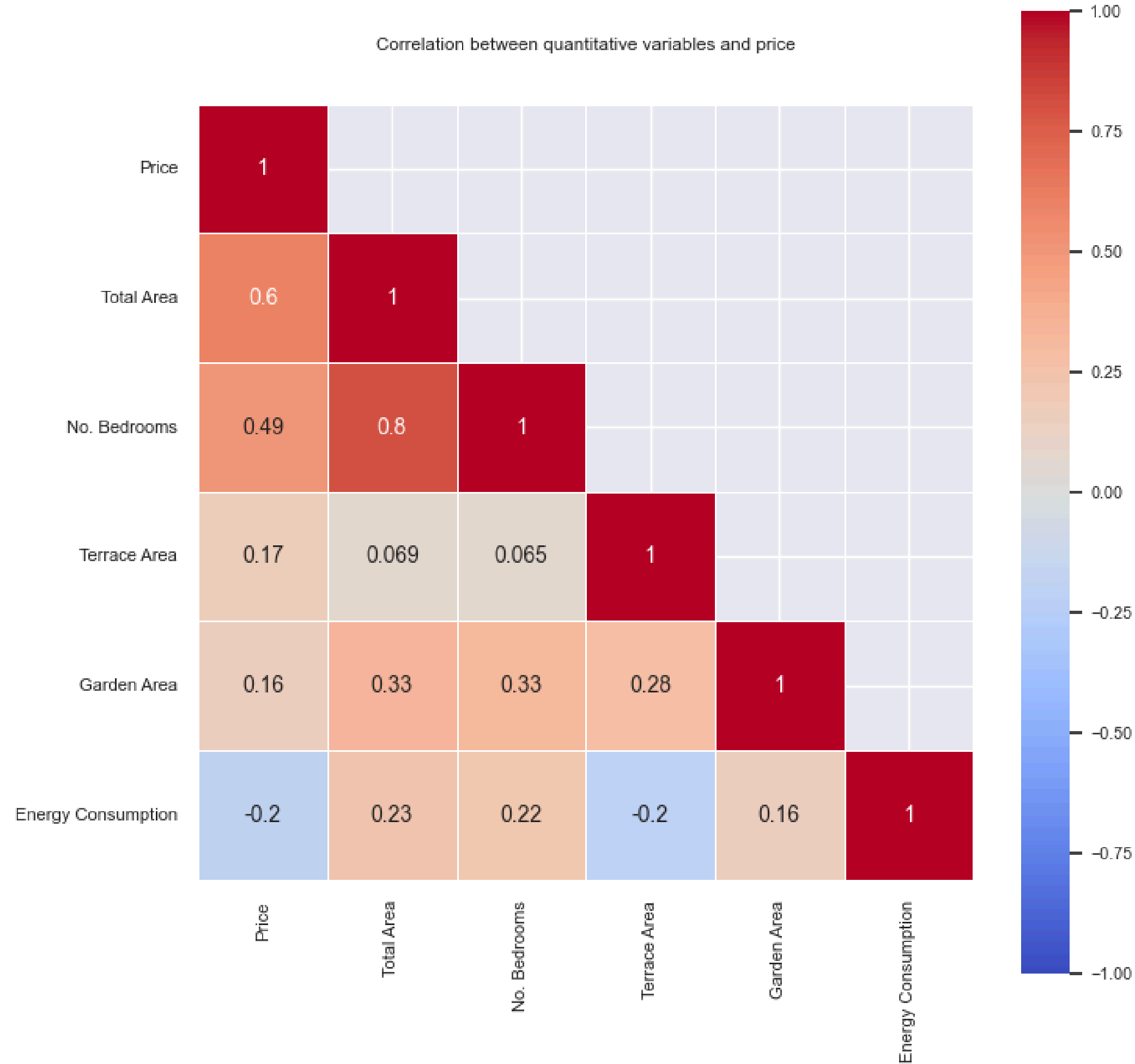- Sub property type
- Location data

## Ordinal

- Construction year
- State of building
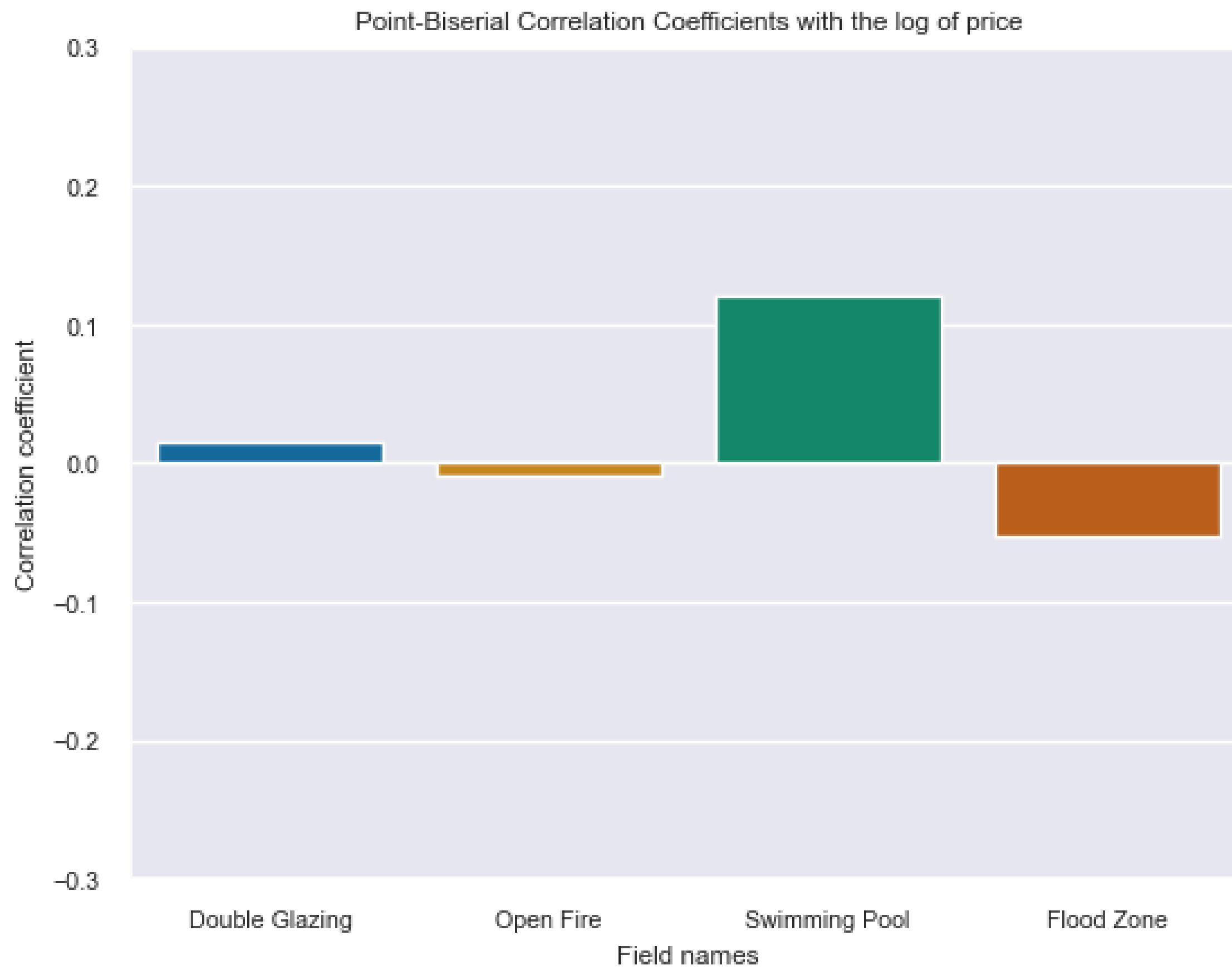- Kitchen level

## Binary

- Swimming pool
- Double glazing
- Furnished
- Open Fire
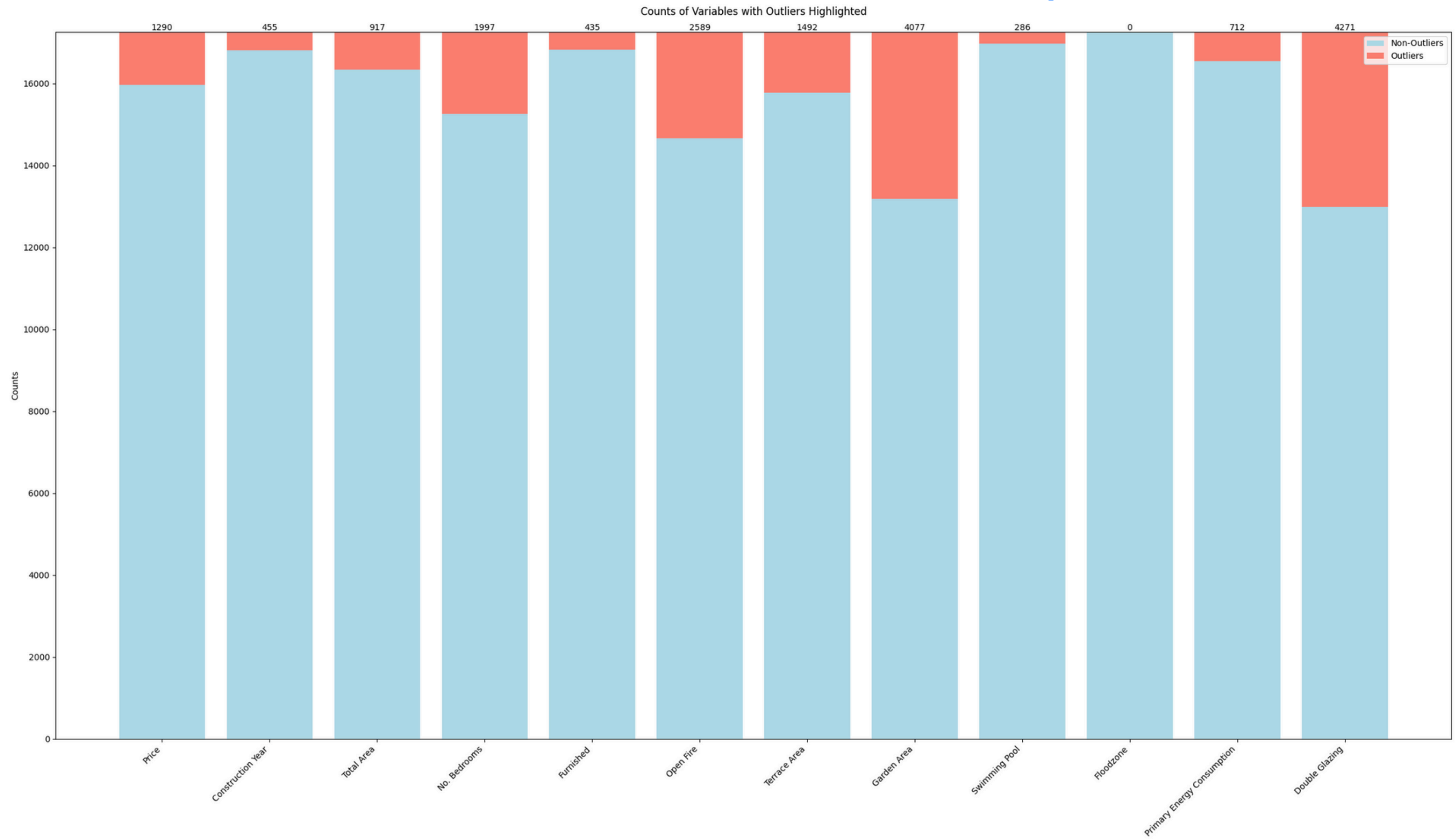- Flood zone

# CONVERSION METHODS

- Numeric Encoding
- Label Encoding
- One hot encoding
- Frequency encoding

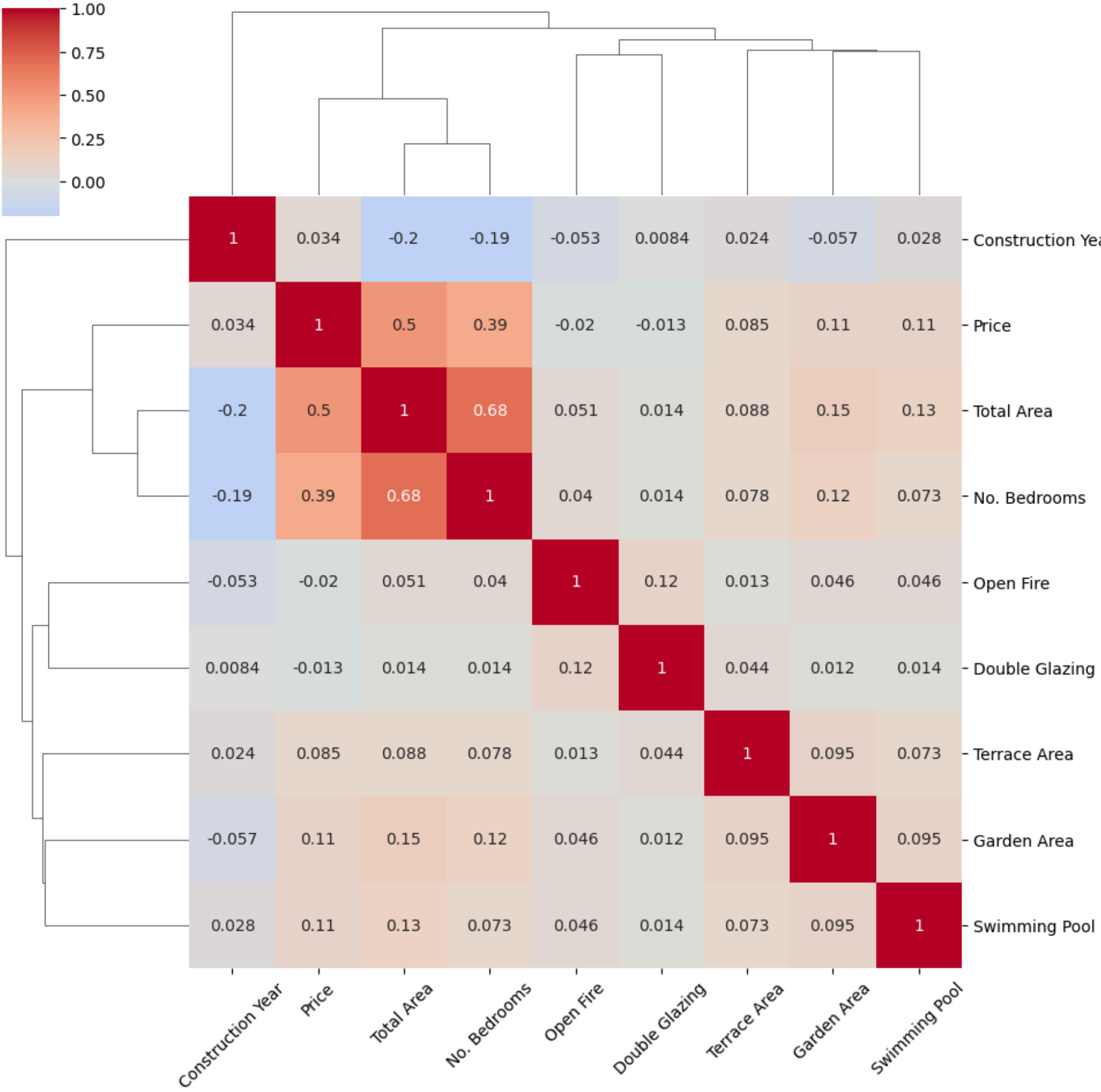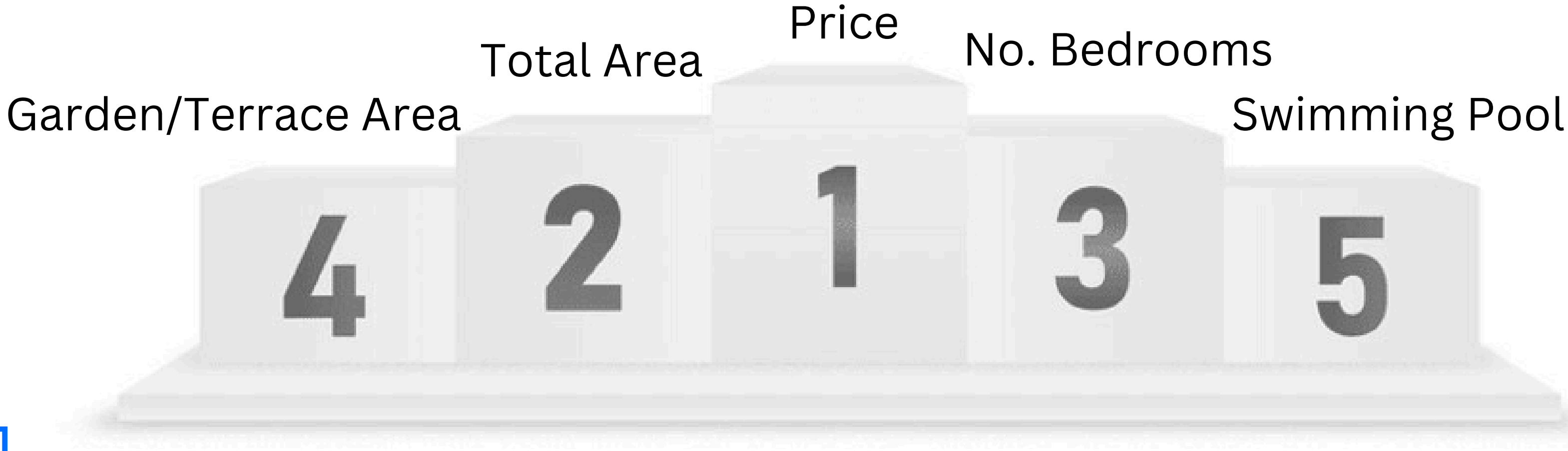Correlation between quantitative variables and price

Method:Spearmans

Point-Biserial Correlation Coefficients with the log of price

Method: Point biserial

Counts of Variables with Outliers Highlighted

Clustermap of correlated variables

|  | Construction Year | Price | Total Area | No. Bedrooms | Open Fire | Double Glazing | Terrace Area | Garden Area | Swimming Pool |
|---|---|---|---|---|---|---|---|---|---|
| Construction Year | 1 | 0.034 | -0.2 | -0.19 | -0.053 | 0.0084 | 0.024 | -0.057 | 0.028 |
| Price | 0.034 | 1 | 0.5 | 0.39 | -0.02 | -0.013 | 0.085 | 0.11 | 0.11 |
| Total Area | -0.2 | 0.5 | 1 | 0.68 | 0.051 | 0.014 | 0.088 | 0.15 | 0.13 |
| No. Bedrooms | -0.19 | 0.39 | 0.68 | 1 | 0.04 | 0.014 | 0.078 | 0.12 | 0.073 |
| Open Fire | -0.053 | -0.02 | 0.051 | 0.04 | 1 | 0.12 | 0.013 | 0.046 | 0.046 |
| Double Glazing | 0.0084 | -0.013 | 0.014 | 0.014 | 0.12 | 1 | 0.044 | 0.012 | 0.014 |
| Terrace Area | 0.024 | 0.085 | 0.088 | 0.078 | 0.013 | 0.044 | 1 | 0.095 | 0.073 |
| Garden Area | -0.057 | 0.11 | 0.15 | 0.12 | 0.046 | 0.012 | 0.095 | 1 | 0.095 |
| Swimming Pool | 0.028 | 0.11 | 0.13 | 0.073 | 0.046 | 0.014 | 0.073 | 0.095 | 1 |

# Top 5 variables

Garden/Terrace Area

Total Area

Price

No. Bedrooms

Swimming Pool
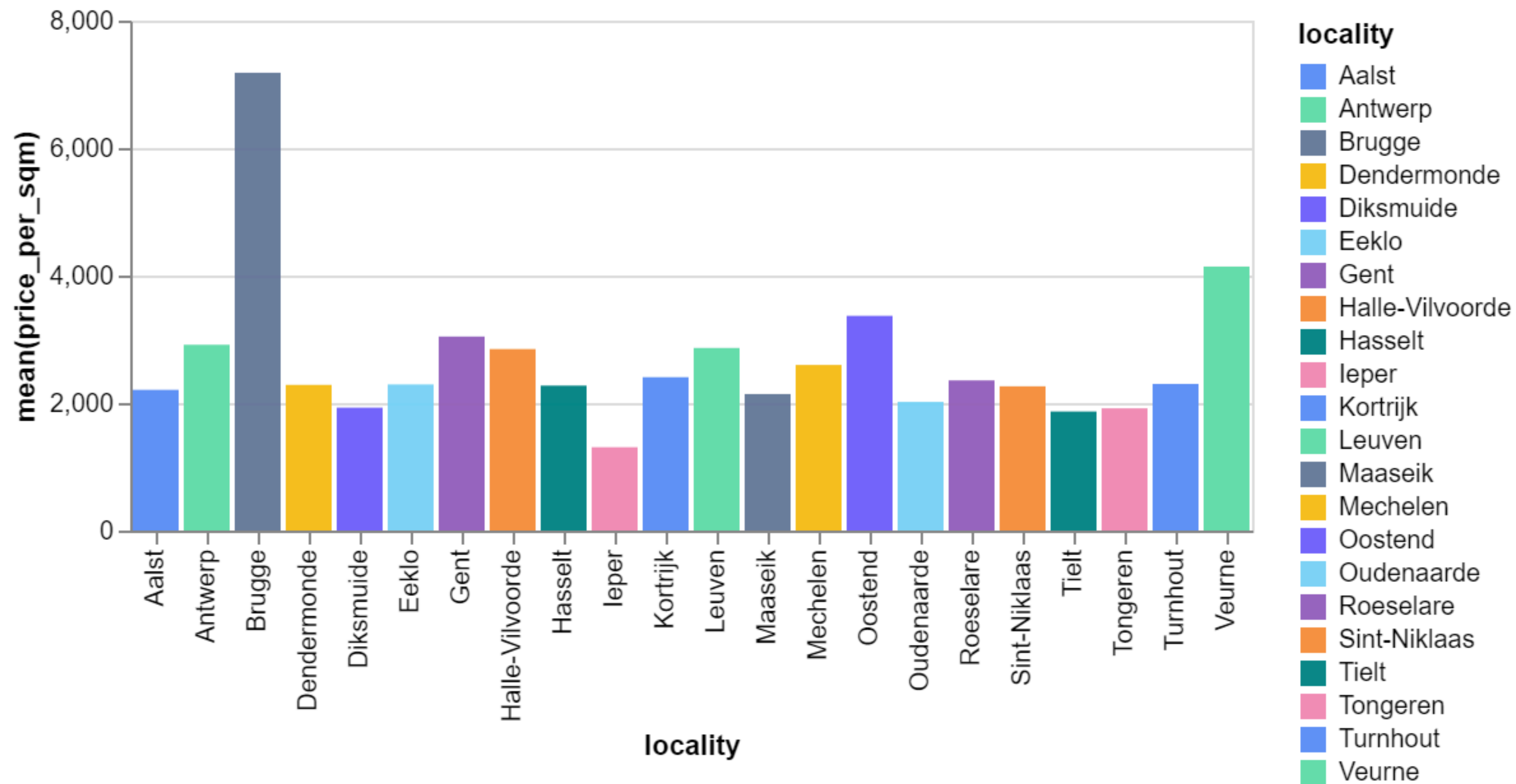
4

2

1

3

5

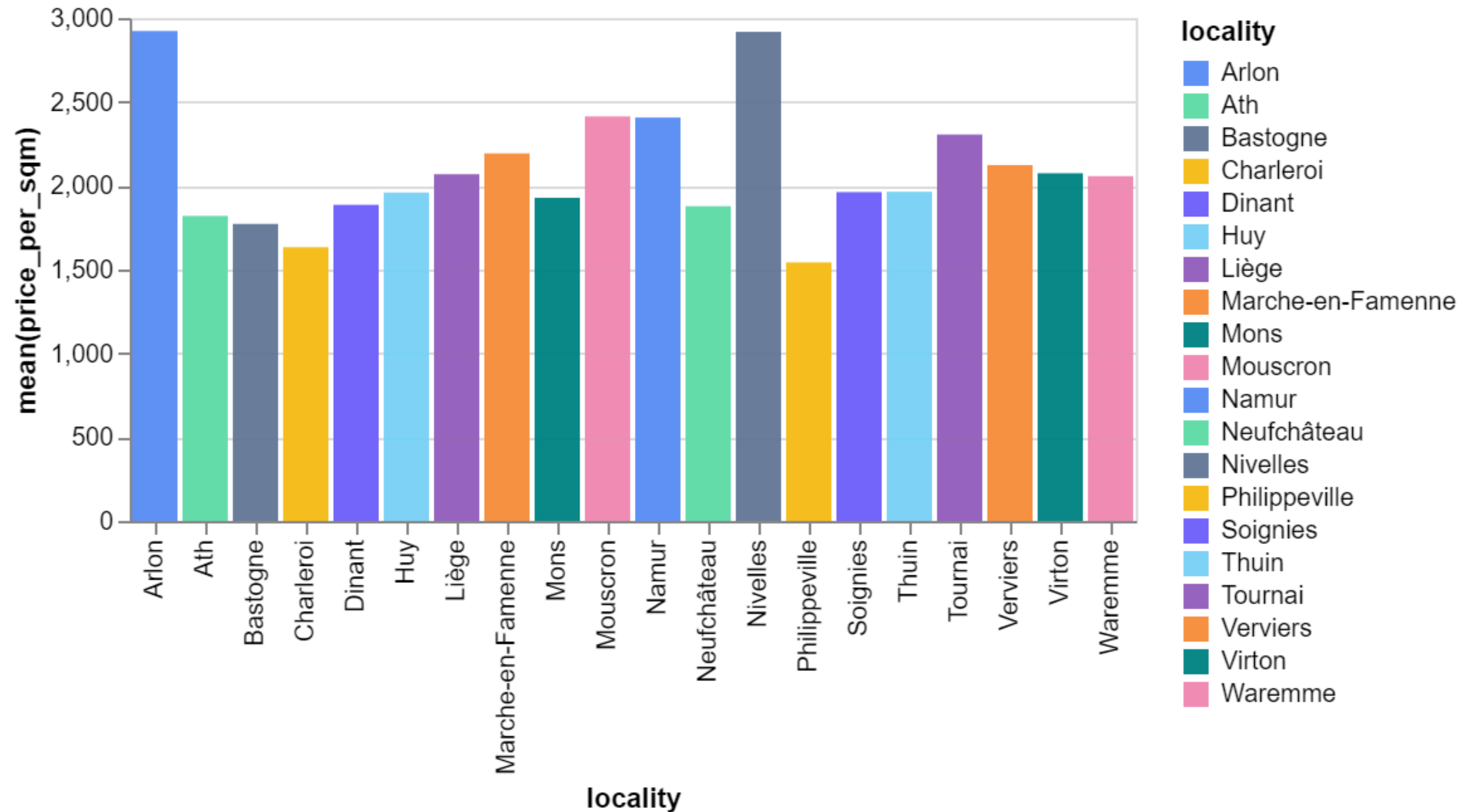Distribution of Properties by Total Area (sqm)

# Average Prices per sqm in Flanders

# Average Prices per sqm in Wallonia

# Conclusion

- **Price** is most correlated with the values of **total area** and **number of bedrooms**.
- Ordinal variables in our dataset have low correlation with price.
- **Garden area** has the most outliers.
- **Brugge** is the most expensive and **Ieper** is the least expensive locality in Flanders.
- **Nivelles** is the most expensive and **Philippeville** is the least expensive locality in Wallonia.
- **West Flanders** is the most expensive and **Hainaut** is the least expensive province in all of Belgium.

# THANK YOU