# Machine Learning 2 Project

Sergio Cárcamo Jara (466116)

# Overview

- Data Description
- Data Specificity
- Training/Test Data Division
- Comparison of Methods
- Summary and Conclusions
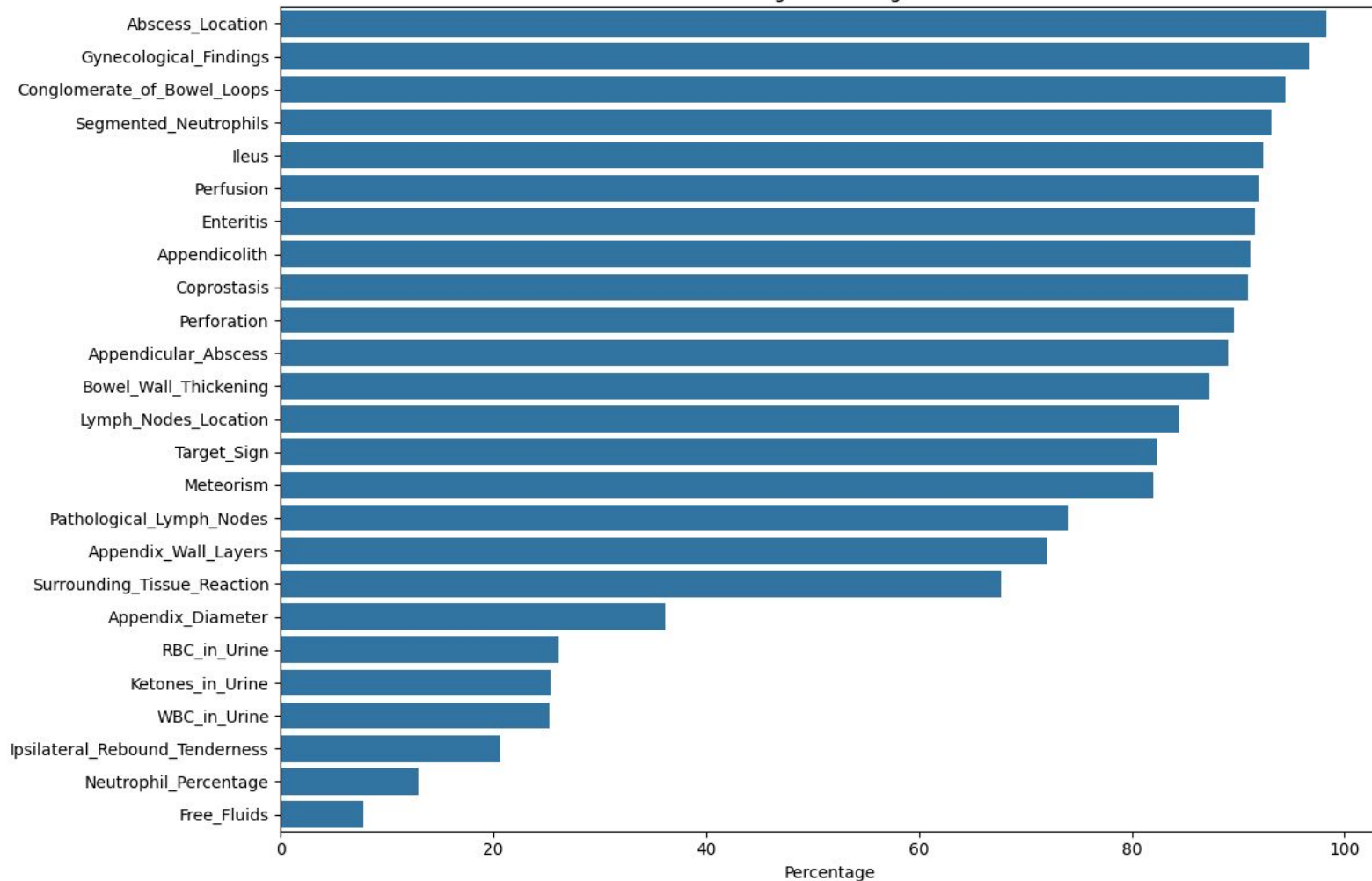- Questions

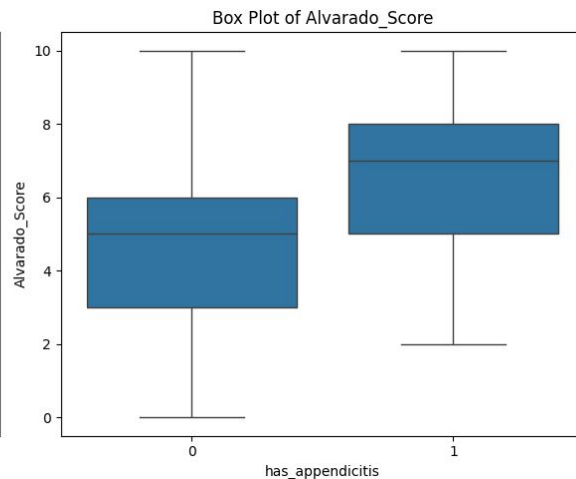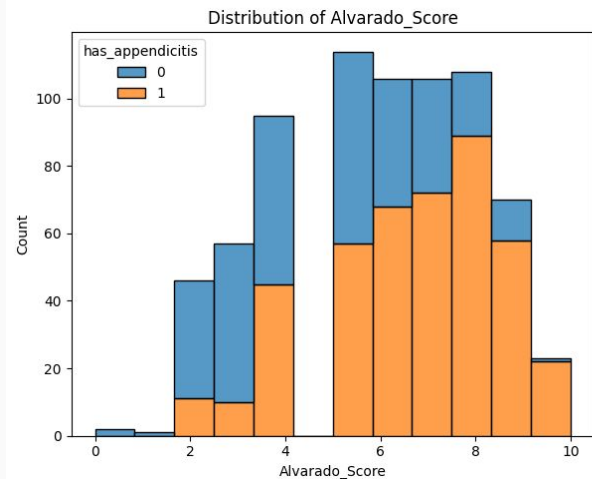# Pediatric Appendicitis

*Classification*

# Data description

Data corresponds to pediatric patients with suspected appendicitis admitted with abdominal pain to Children's Hospital St. Hedwig in Regensburg, Germany, between 2016 and 2021.
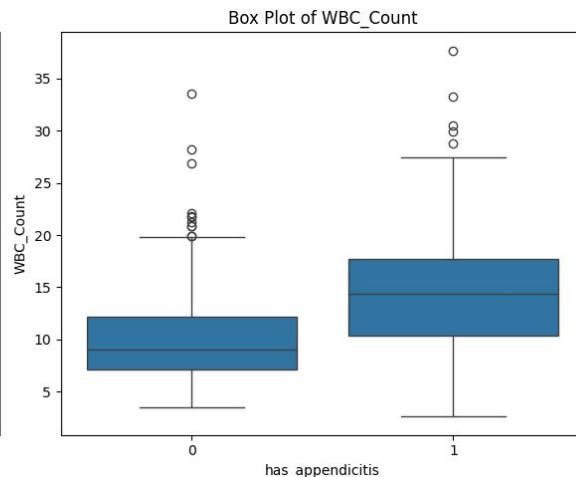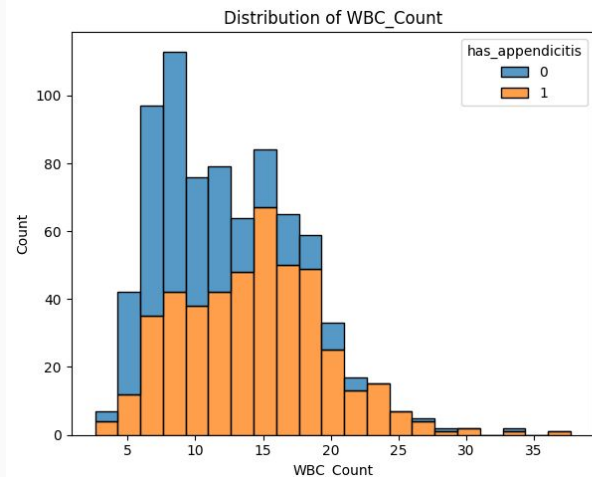
The dataset includes ultrasound images, laboratory, physical examination, scoring results and ultrasonographic findings extracted manually by the experts, and the target variable is diagnosis, to predict if the patient has appendicitis or not.

Percentage of Missing Values

The Alvarado score is a clinical scoring system used in the diagnosis of appendicitis. Based on symptoms and blood test results.

WBC: White blood cell count

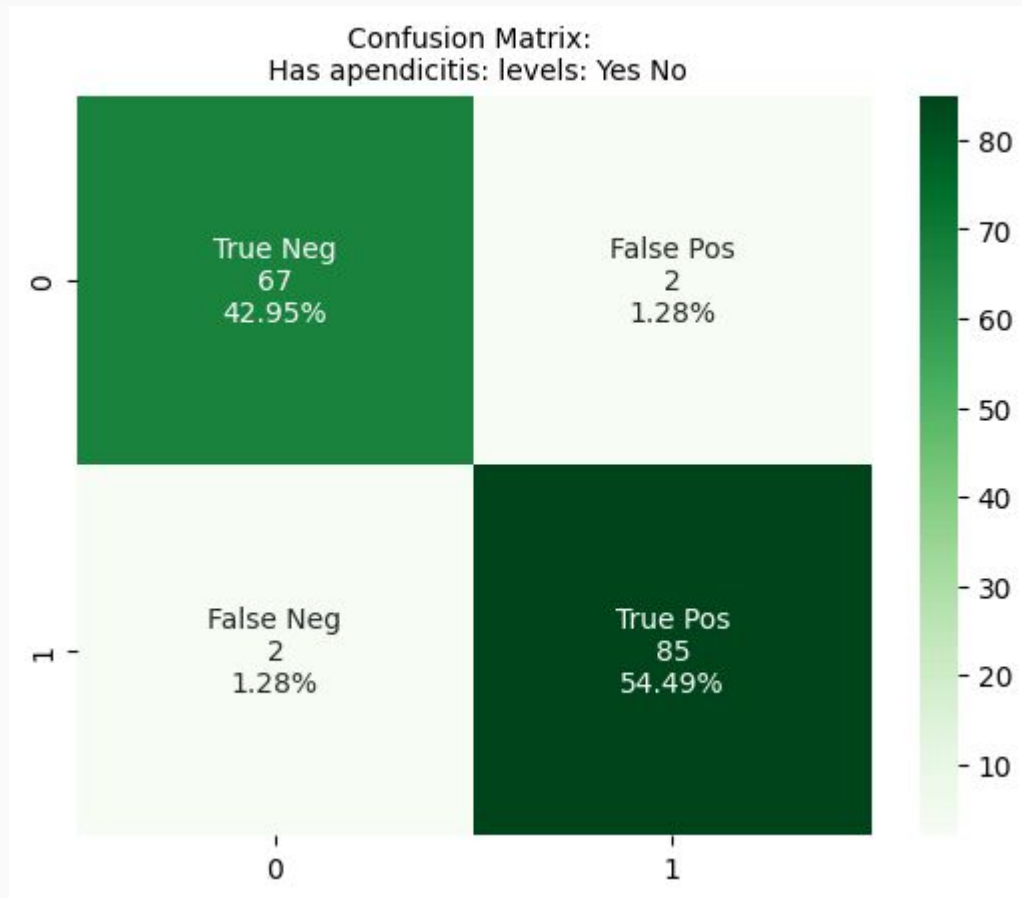| Alvarado score | |
|---|---|
| **Symptoms** | |
| Abdominal pain that migrates to the right iliac fossa | 1 |
| Anorexia (loss of appetite) or ketones in the urine | 1 |
| Nausea or vomiting | 1 |
| Tenderness in the right iliac fossa | 2 |
| **Signs** | |
| Rebound tenderness | 1 |
| Fever of 37.3 °C or more | 1 |
| **Laboratory** | |
| Leukocytosis > 10,000 | 2 |
| Neutrophilia > 70% | 1 |
| **TOTAL** | **10** |

# Training/Test Data Division

- Training/Test data split 80/20
- Randomized Search CV
    - parameters set for each model
    - 25 iterations
    - scoring F1
    - cv 5

# Comparison of Methods

| model | RandomForestClassifier | AdaBoostClassifier | XGBClassifier |
|---|---|---|---|
| accuracy | 0.929487 | 0.974359 | 0.948718 |
| precision | 0.94186 | 0.977011 | 0.964706 |
| recall | 0.931034 | 0.977011 | 0.942529 |
| F1 | 0.936416 | 0.977011 | 0.953488 |

# AdaBoost Classifier



Confusion Matrix:
Has apendicitis: levels: Yes No

|  | 0 | 1 |
|---|---|---|
| 0 | True Neg 67 42.95% | False Pos 2 1.28% |
| 1 | False Neg 2 1.28% | True Pos 85 54.49% |

# Summary and Conclusions

- Chosen model is AdaBoost Classifier (with Decision Tree Classifier as base estimator)
- Find variables that can be removed without sacrificing performance (such as ultrasound images)
- Improve imputation methods

- Best hyper parameters
  - algorithm: SAMME
  - learning_rate: 0.63435404
  - n_estimators: 111
  - imputer: mean
  - estimator criterion: entropy
  - estimator max_depth: 7
  - estimator min samples leaf: 3
  - estimator min samples split: 3

# Real estate prices

*Regression*

# Data description

For the regression problem, the data was scraped from otodom from september to december 2024. The scope is properties for sale in the Warsaw area, including flats, excludes houses.

The dataset includes data on the property itself, and the building if the property is a flat. The target variable is the price.

# Data description

Property characteristics:
- estate: FLAT, HOUSE
- **area_m2**
- rooms_number
- floor_number
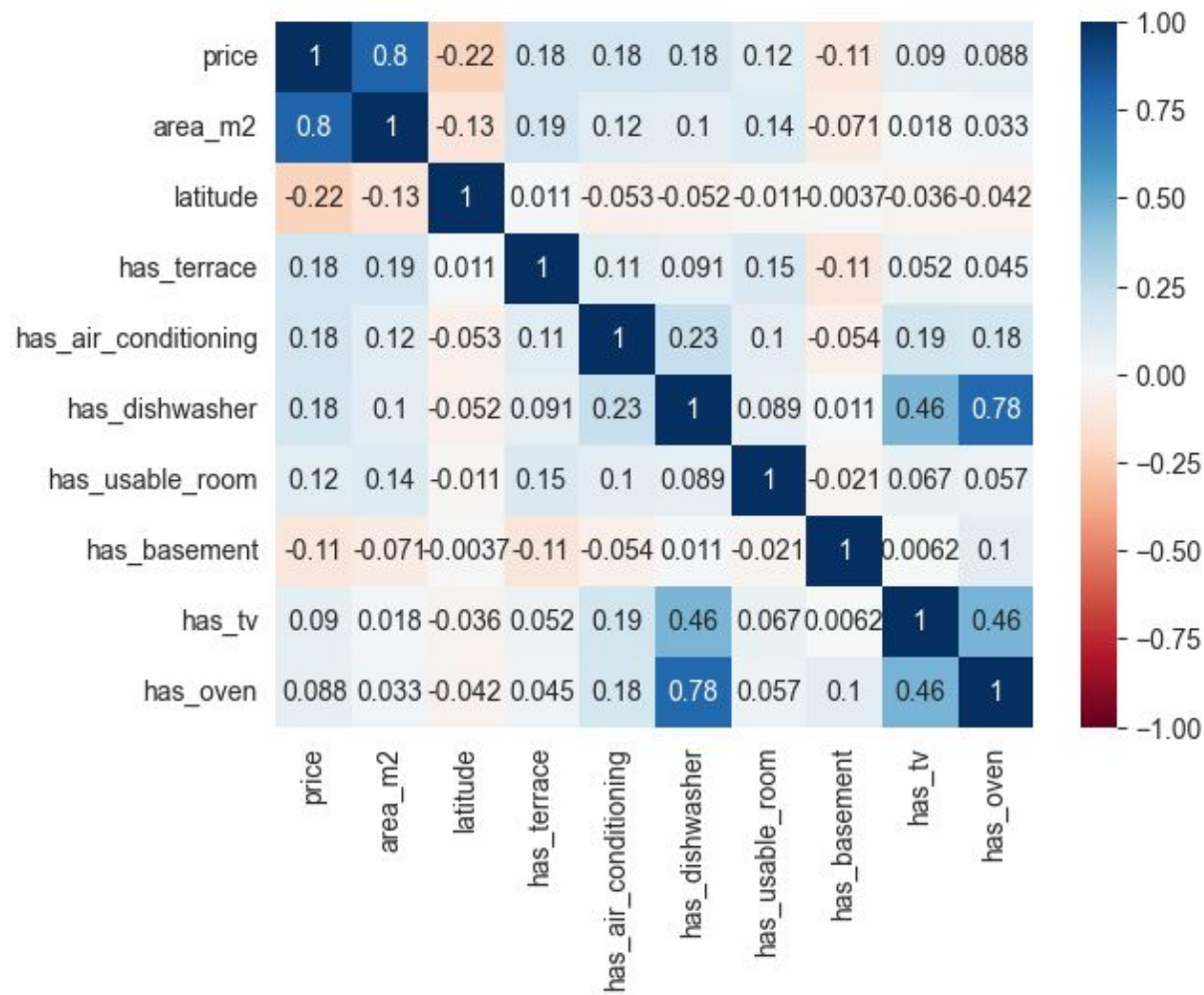- windows_type
- heating
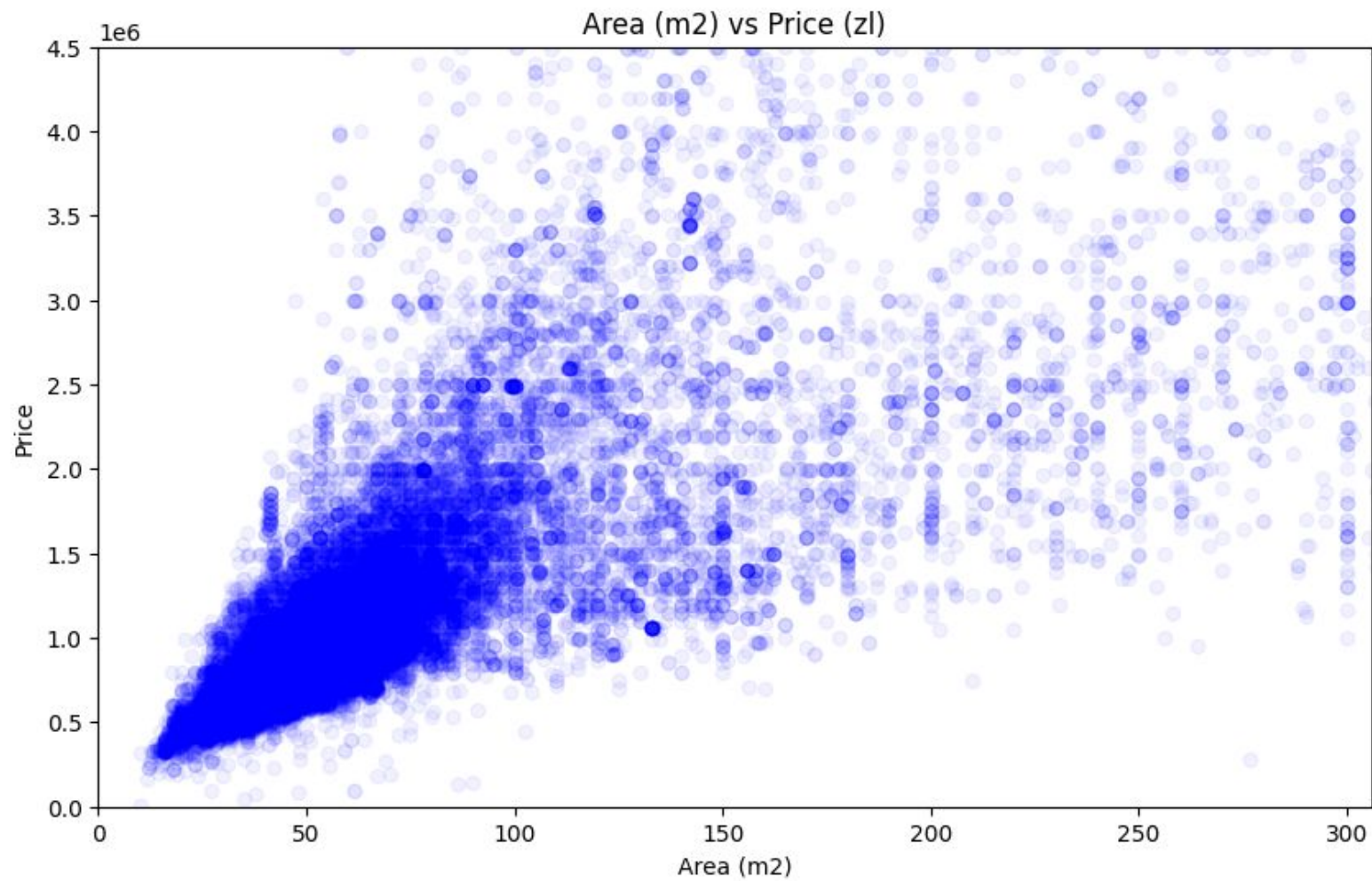- **price**

Location
- **district**
- latitude
- longitude

Building characteristics:
- building_year
- **building_age**
- building_type
- building_floors_num
- construction_status
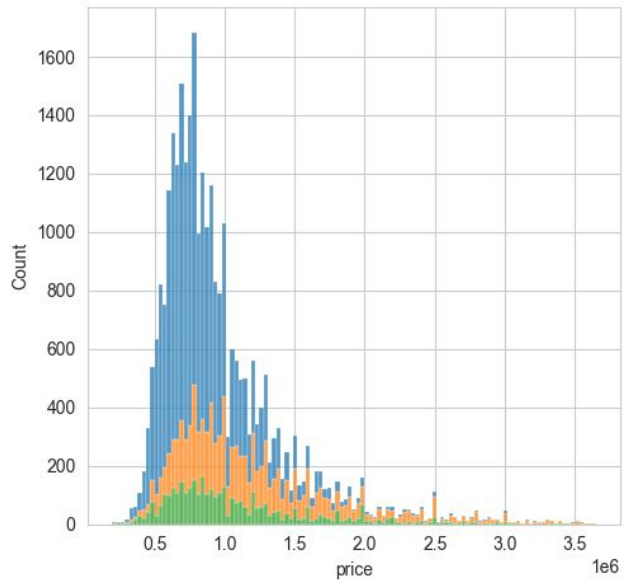- building_material
- building_ownership

# Data description: amenities

- has_lift
- has_internet
- has_furniture
- has_air_conditioning
- has_tv
- has_oven
- has_stove
- has_dishwasher
- has_fridge
- has_washing_machine
- has_separate_kitchen

- has_garage
- has_usable_room
- has_terrace
- has_balcony
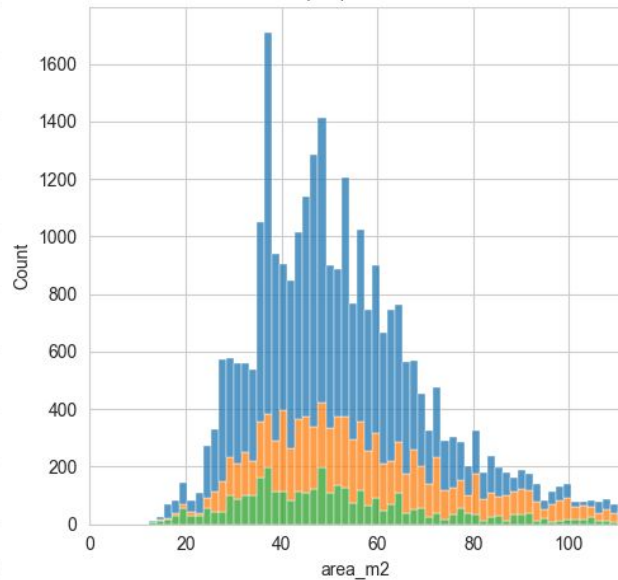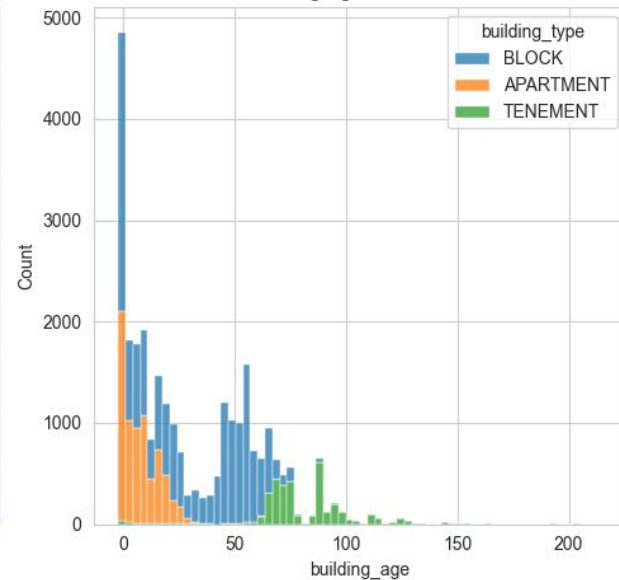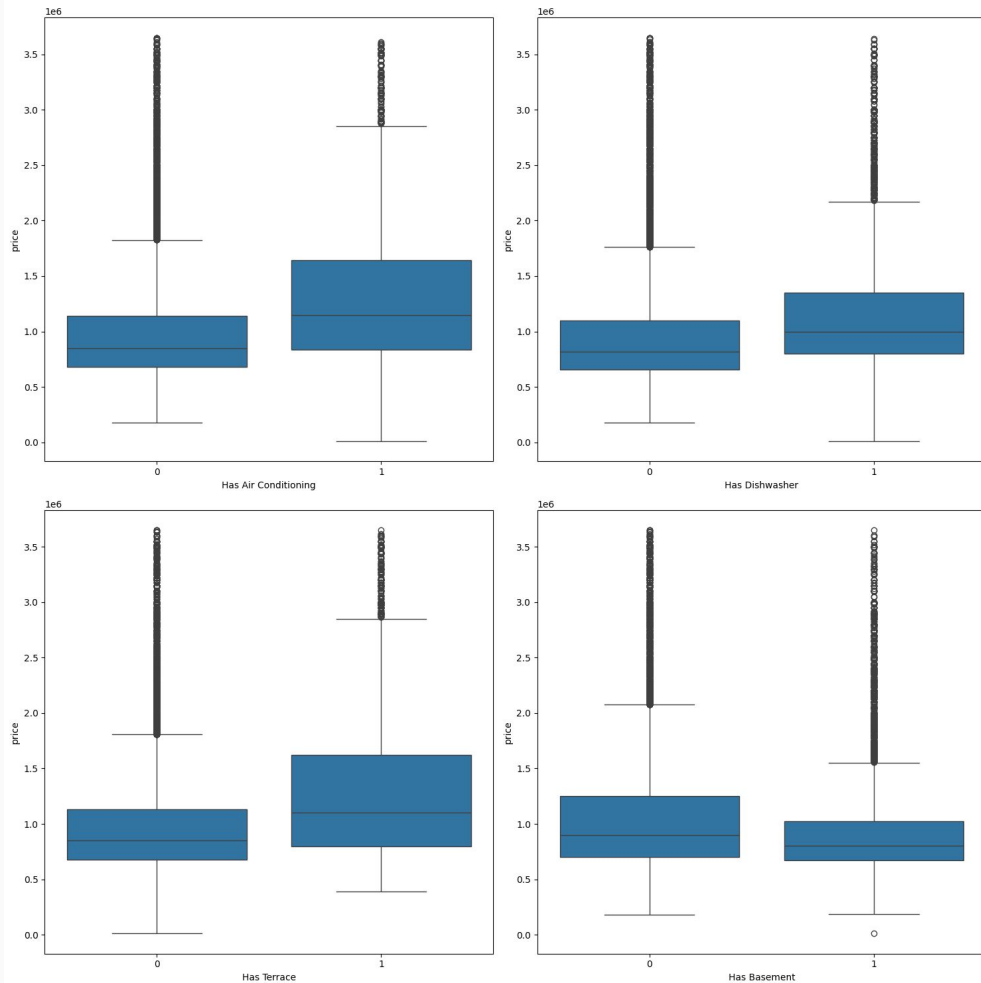- has_attic
- has_basement
- has_garden
- has_pool

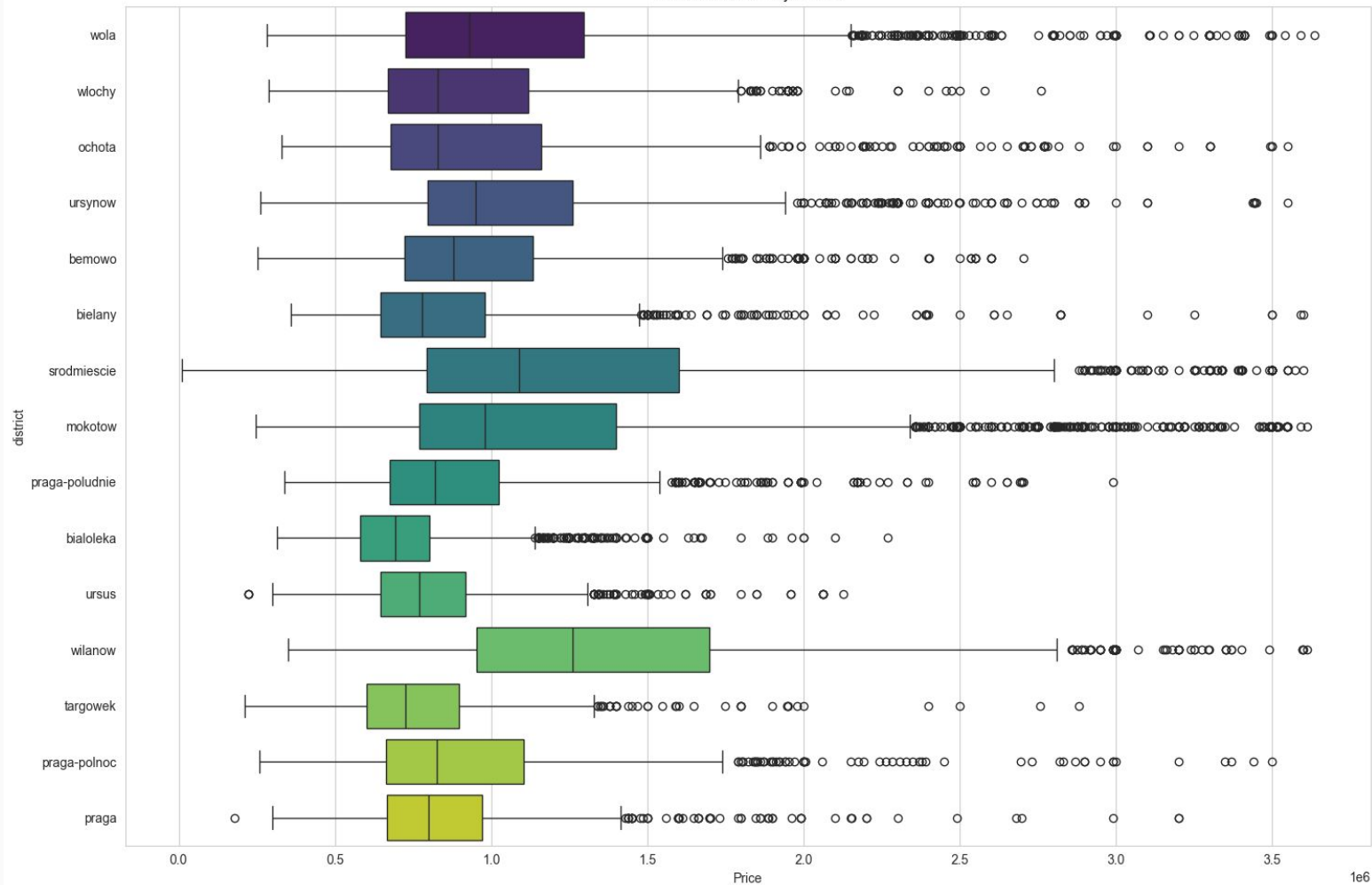Area (m2) vs Price (zl)

* tenement: kamienica

Price Distribution by Property Features

## Highlights

- Air conditioning: correlated with newer and furnished flats
- Dishwasher: correlated with furnished properties
- Terrace: correlated with larger and newer properties
- Basement: correlated with older properties

Price Distribution by District

# Training/Test Data Division

For splitting the data, I used a 80/20 split.

Top 30 features selected with highest correlation.

Randomized Search CV used for hyperparameter tuning.

- 50 iterations
- metric neg_root_mean_squared_error
- cv 5

# Comparison of Methods

| model | Decision Tree Regressor | XGBoost Regressor | Random Forest Regressor | Ensemble Learning* |
|---|---|---|---|---|
| RMSE | 209411 | 120715 | 147654 | 131948 |
| R2 | 0.80 | 0.932 | 0.90 | 0.92 |
| MAPE | 0.14 | 0.068 | 0.09 | 0.079 |

*Ensemble Learning: XGBoost, Random Forest Regressor, GradientBoostingRegressor

# Summary and Conclusions

- Selected model is XGBoost Regressor due to lower RMSE and MAPE
- Improve feature selection methods
- Improve imputing methods
- Adjust ensemble pipeline
- Comparison of performance between training and test seems to indicate overfitting

# Thank you!