

Preprocessing and Model Performance Report

1. Preprocessing Summary

The preprocessing pipeline included the following steps:

1. Handling Missing Values:

- Filled missing compound names with 'No compound'.

2. Encoding Categorical Variables:

- Binary encoding for Yes/No columns.
- Mapping finishing types from ordinal categories.
- Mapping view types into numerical categories.

3. Removing Unwanted Columns:

- Dropped listing_date and days_on_market.

4. Outlier Removal:

- Applied IQR filtering per compound to remove price, area, and distance outliers.

5. Removing Duplicate Rows:

- Ensured dataset contains unique apartment listings.

2. Model Performance Comparison

Linear Regression

Train R²: 0.7324

Train MAE: 395958.41

Test R²: 0.7448

Test MAE: 409566.5

LightGBM

Train R²: 0.7876

Train MAE: 351902

Test R²: 0.7436

Test MAE: 407062

Random Forest

Train R²: 0.8226

Train MAE: 320372

Test R²: 0.7371

Test MAE: 408382

XGBoost

Train R²: 0.8475

Train MAE: 295927

Test R²: 0.71

Test MAE: 431286

CatBoost

Train R²: 0.8751

Train MAE: 260120.62

Test R²: 0.8702

Test MAE: 284065.07