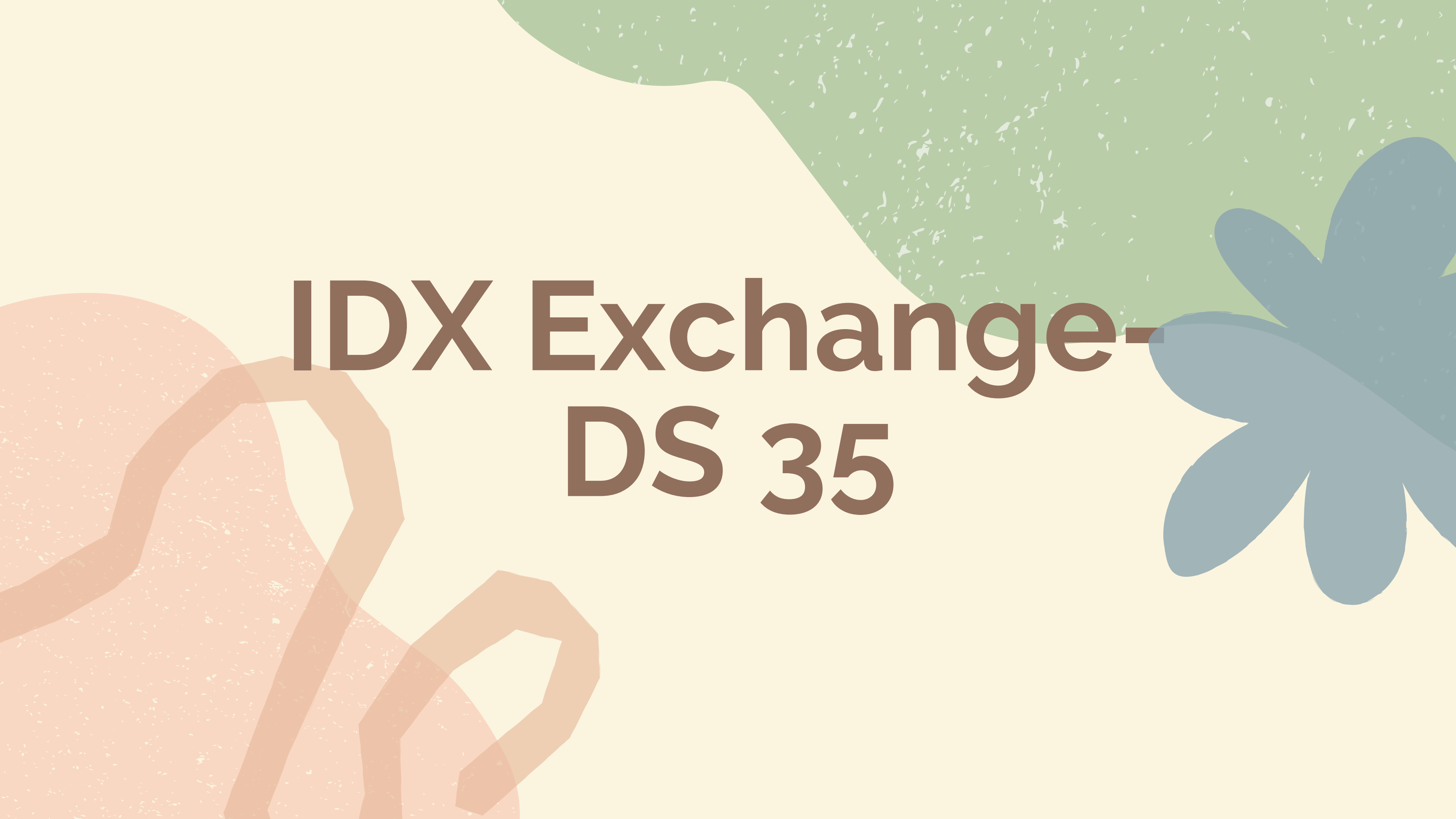


# IDX Exchange- DS 35



# Contents

1. Project Introduction
2. Data Overview
3. Our Models
4. Models Comparison
5. Future Steps / What We've Learned



# Introduction

**Goal:** Predicting California Property Close Prices

**Focus:** Single Family Residences & Residential

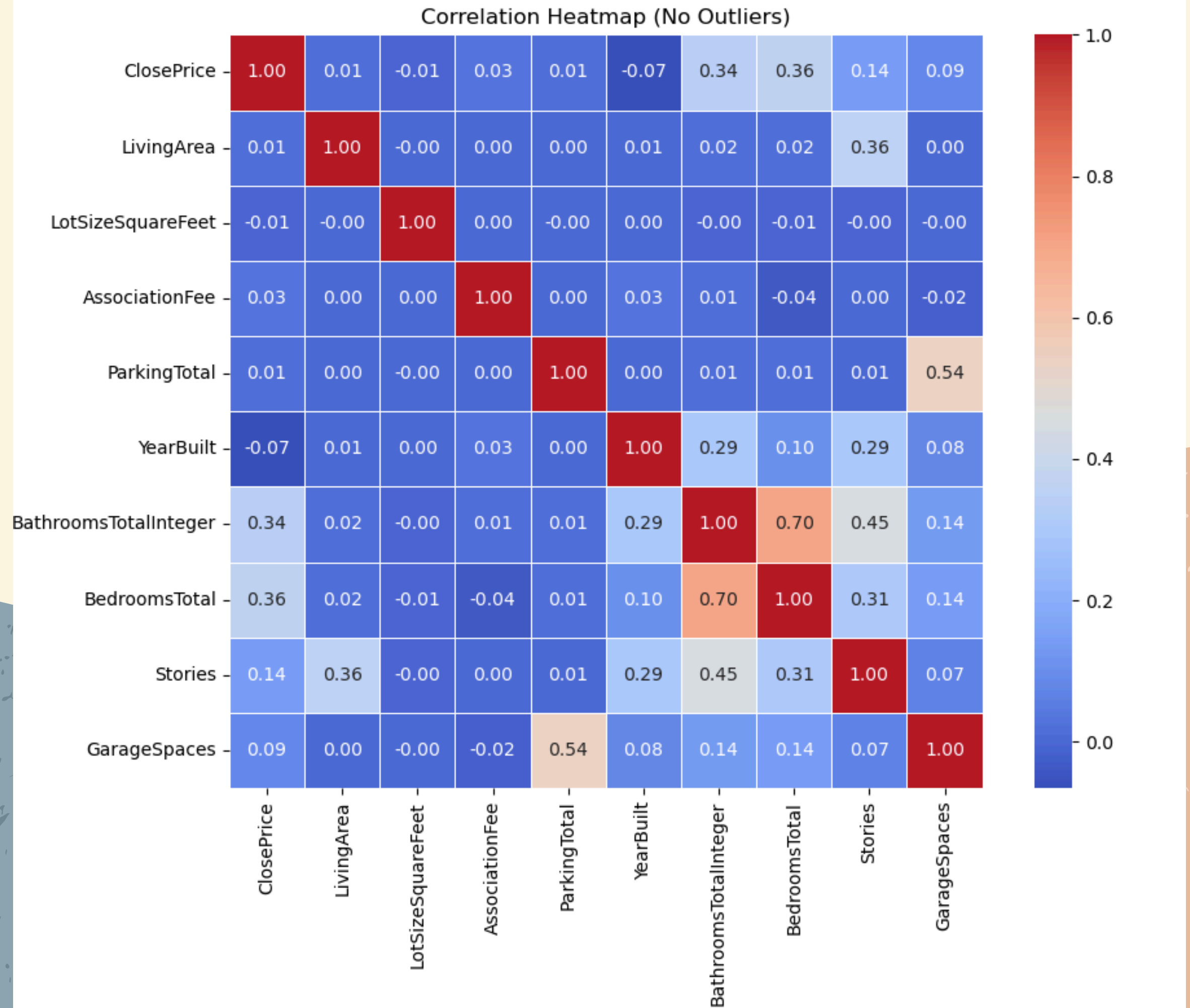
**Data:** Historical transaction data from CRMLs

## Our Models:

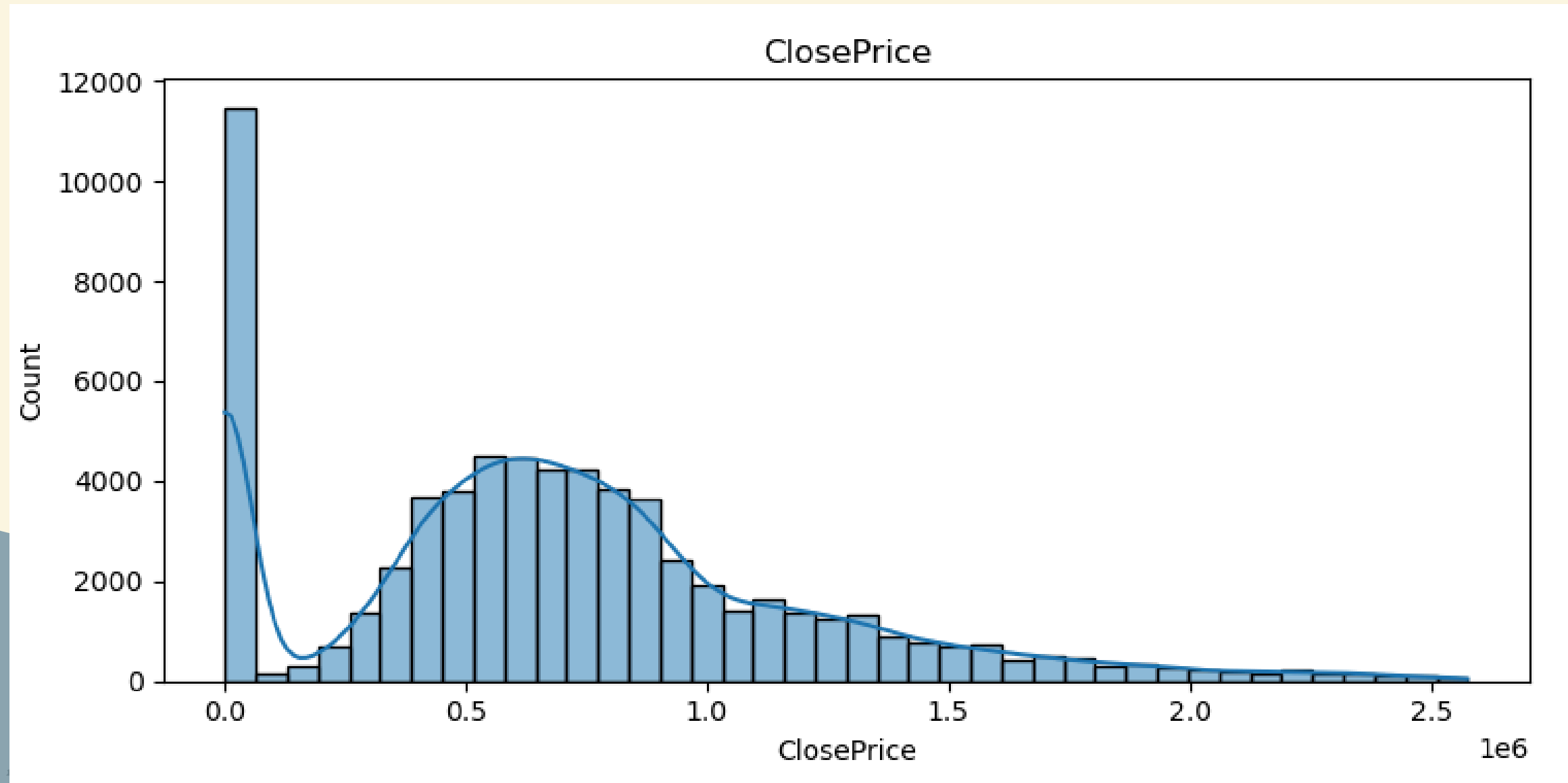
1. LightGBM
2. XGBoost
3. MLP Regression
4. Catboost



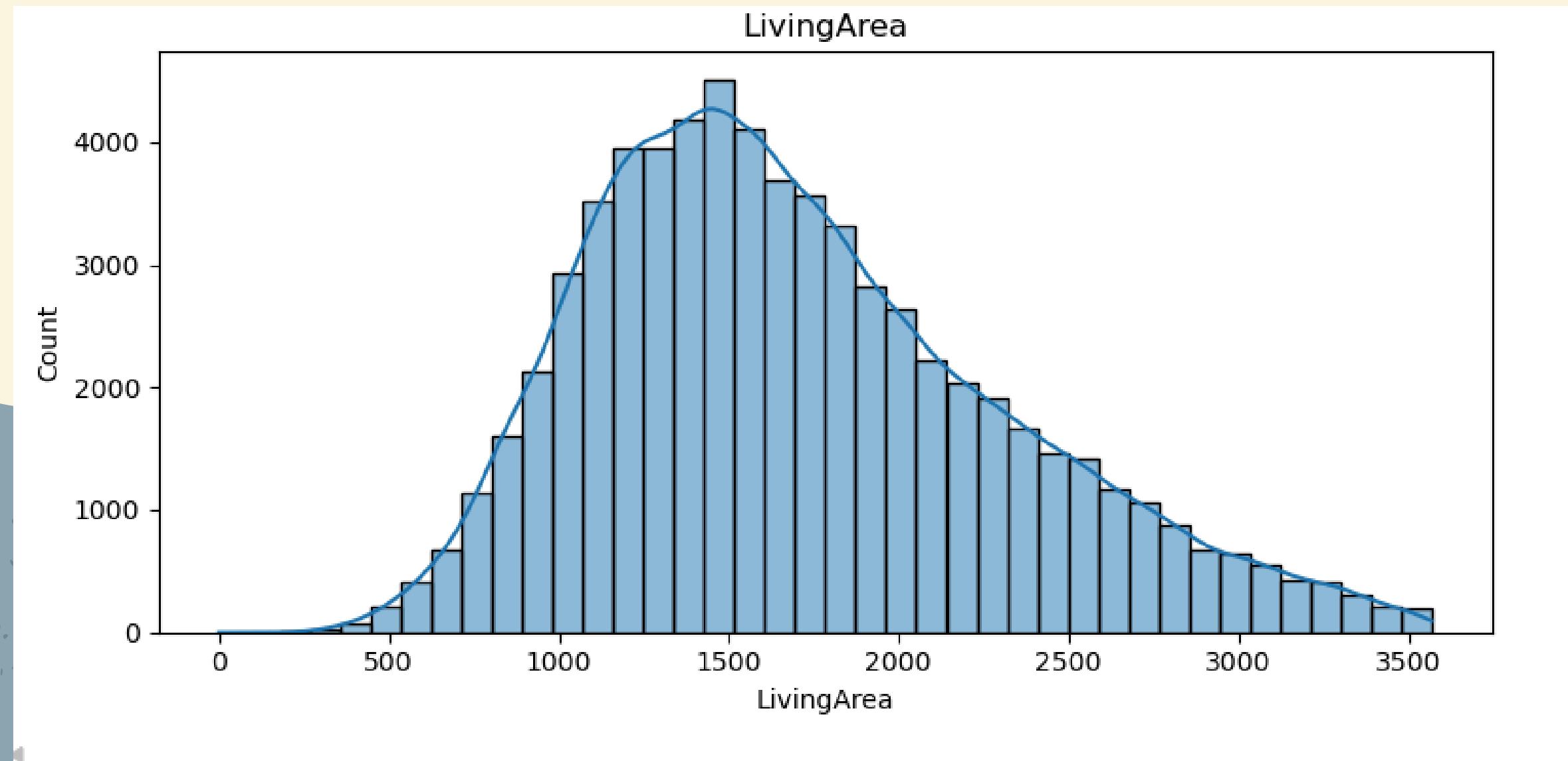
# Correlation table



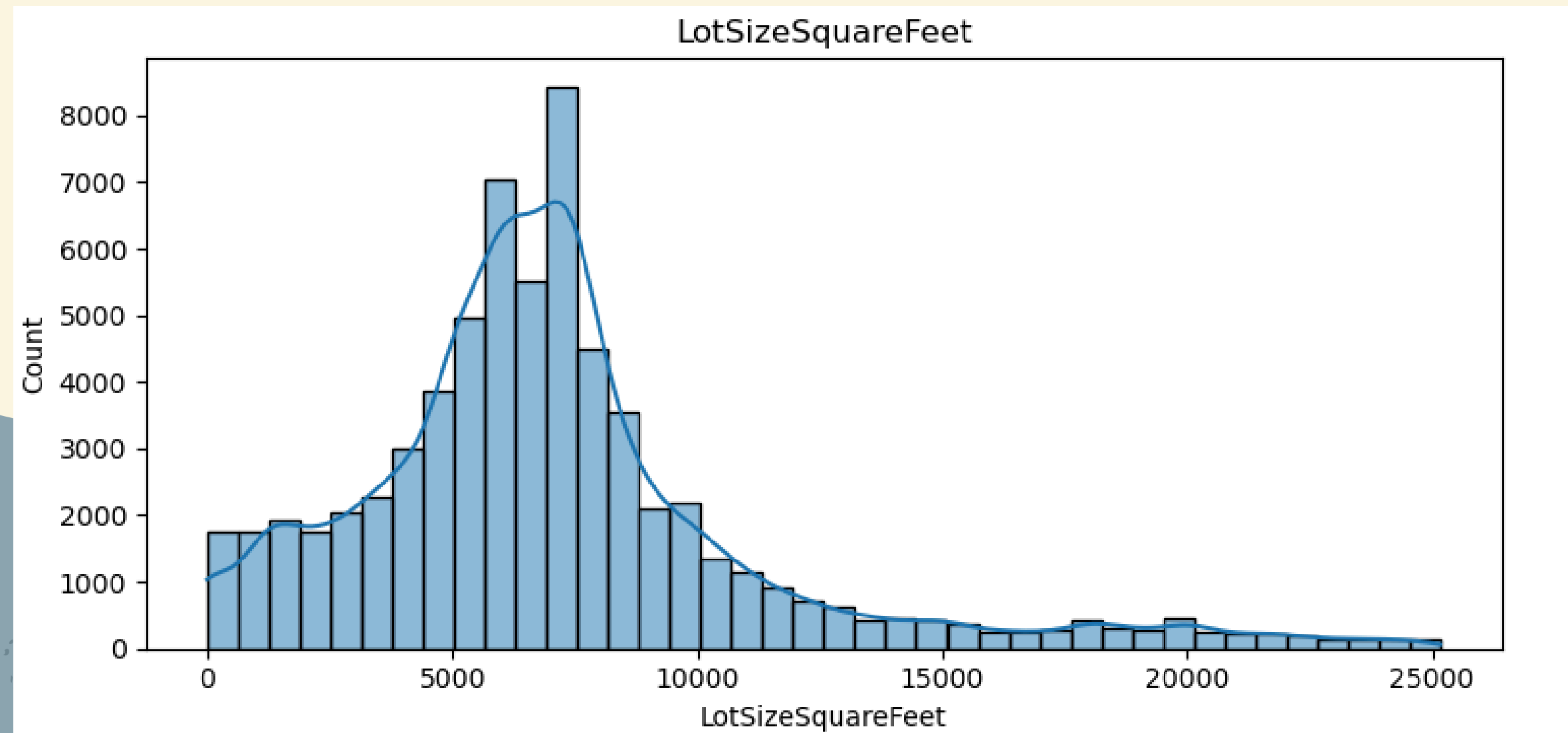
# Close price Distribution



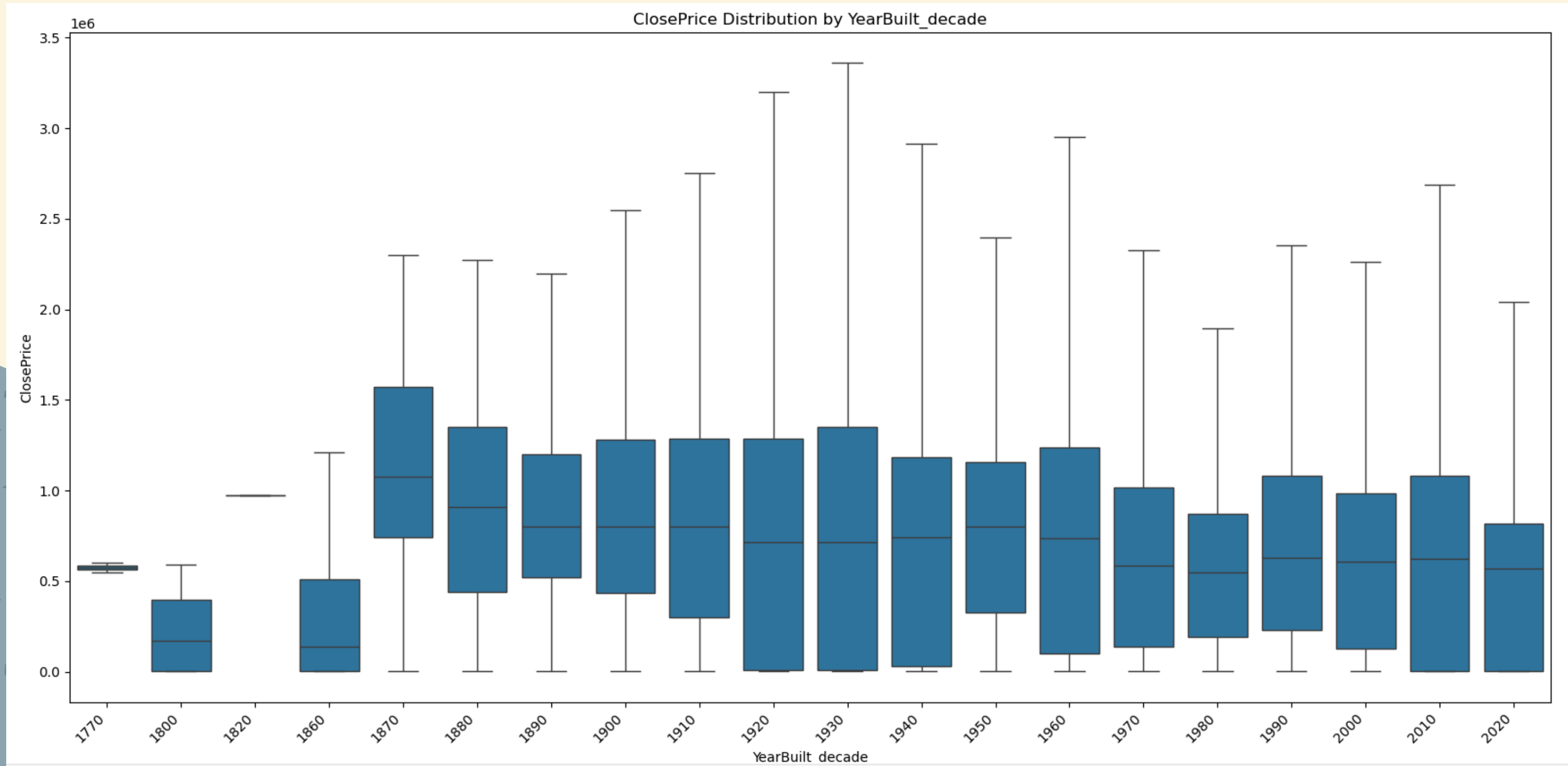
# Living Area Distribution



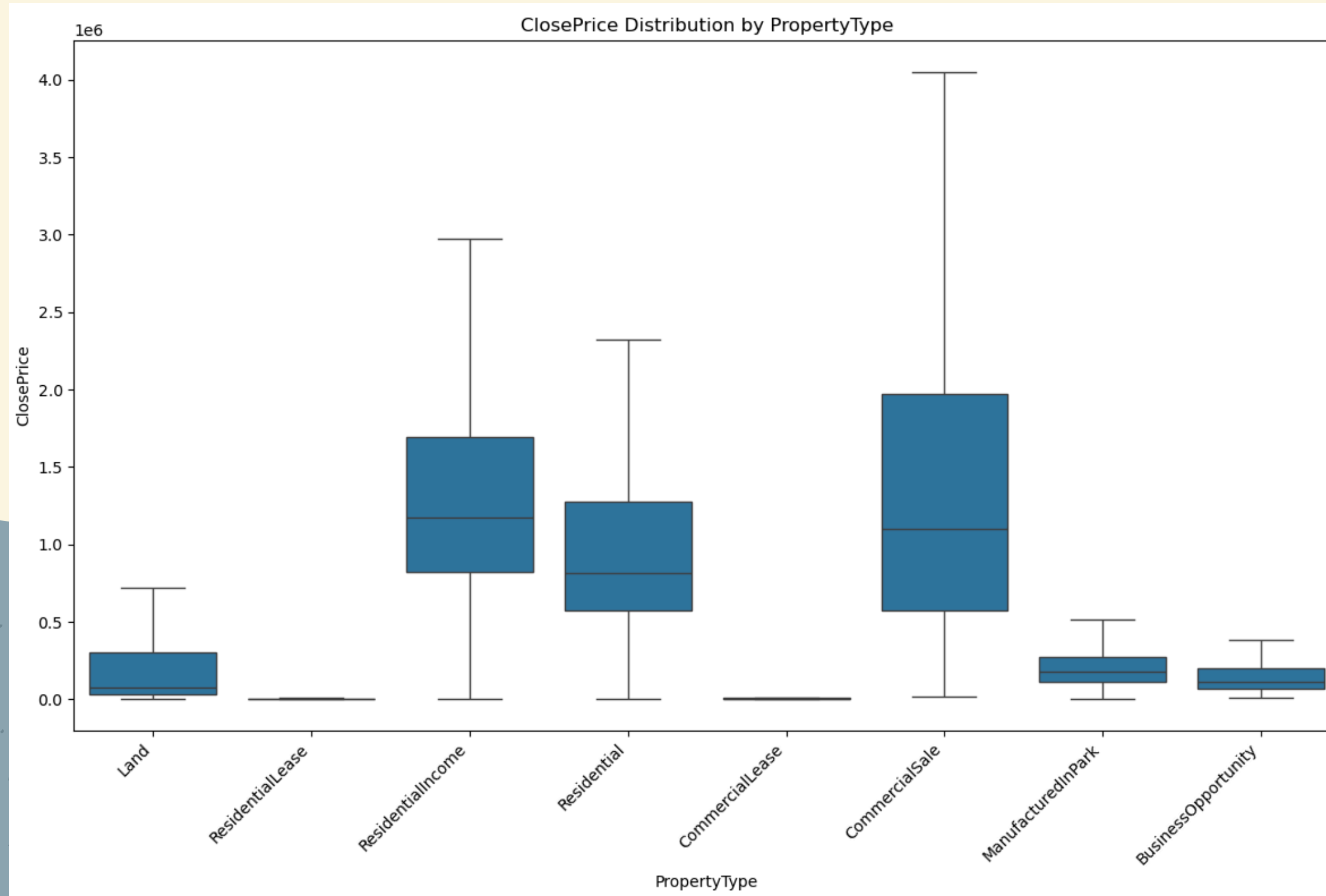
# Lot Size (ft<sup>2</sup>) Distribution



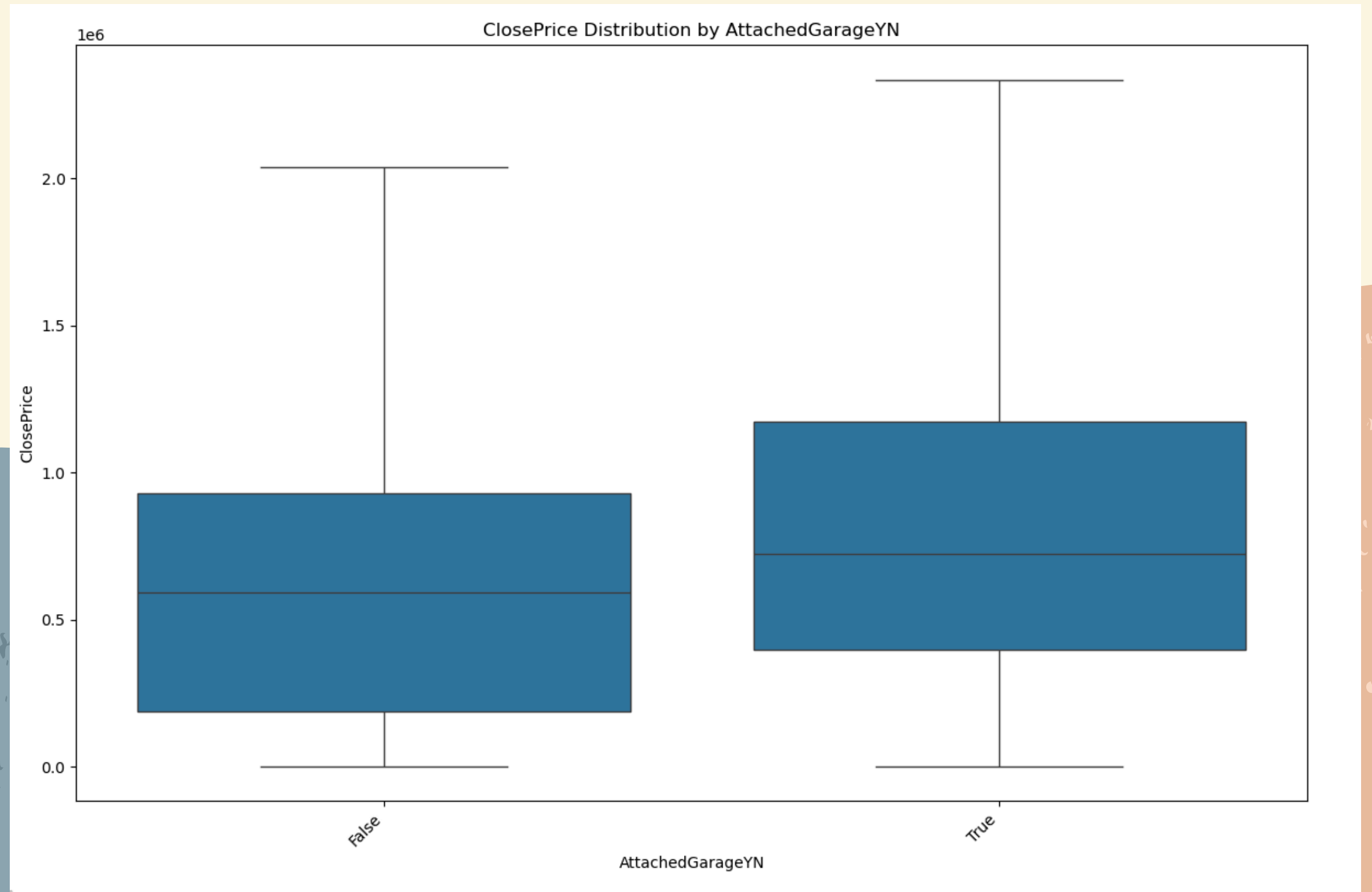
# Close Price by YearBuild



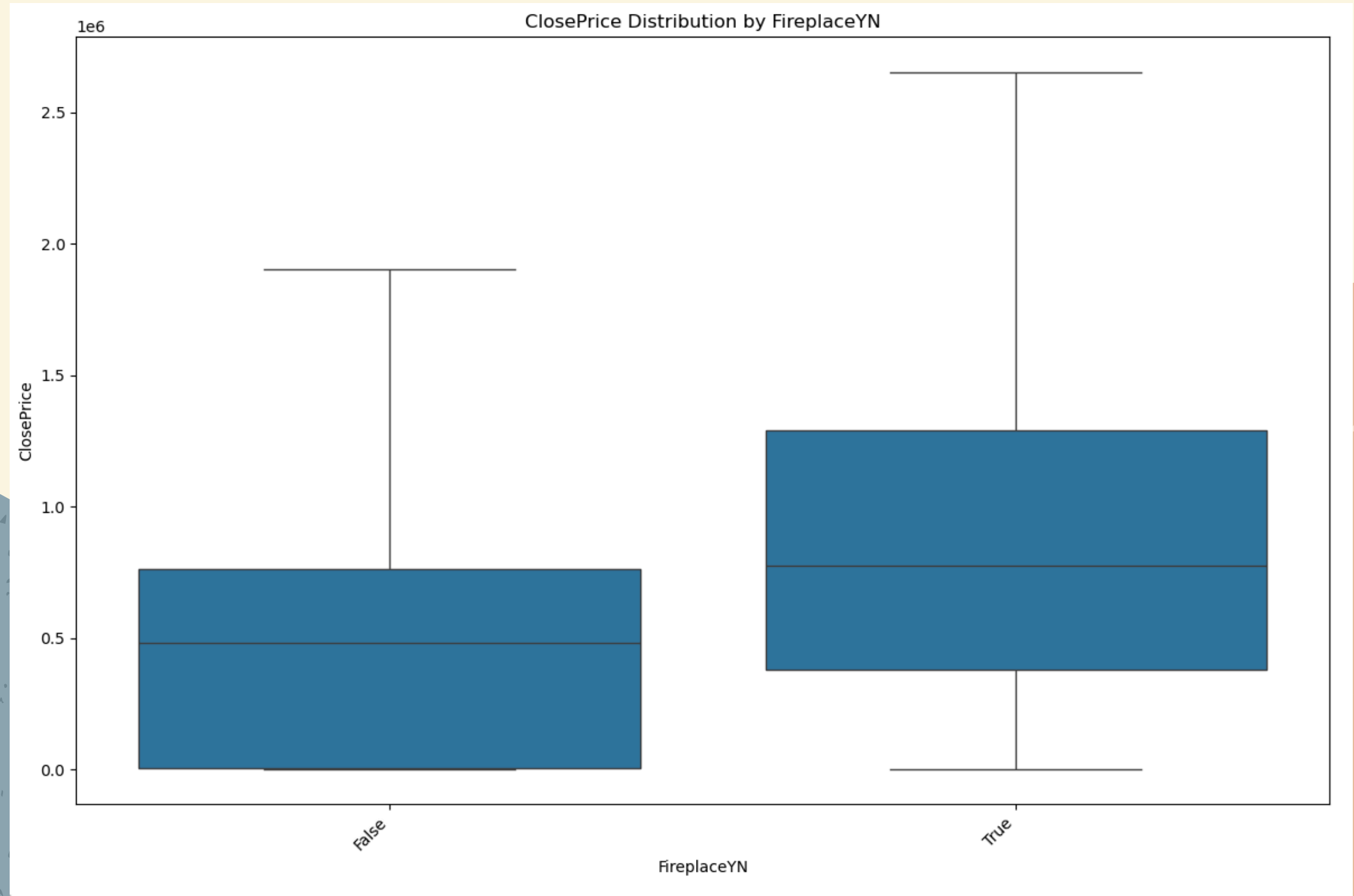
# Close Price by Property type



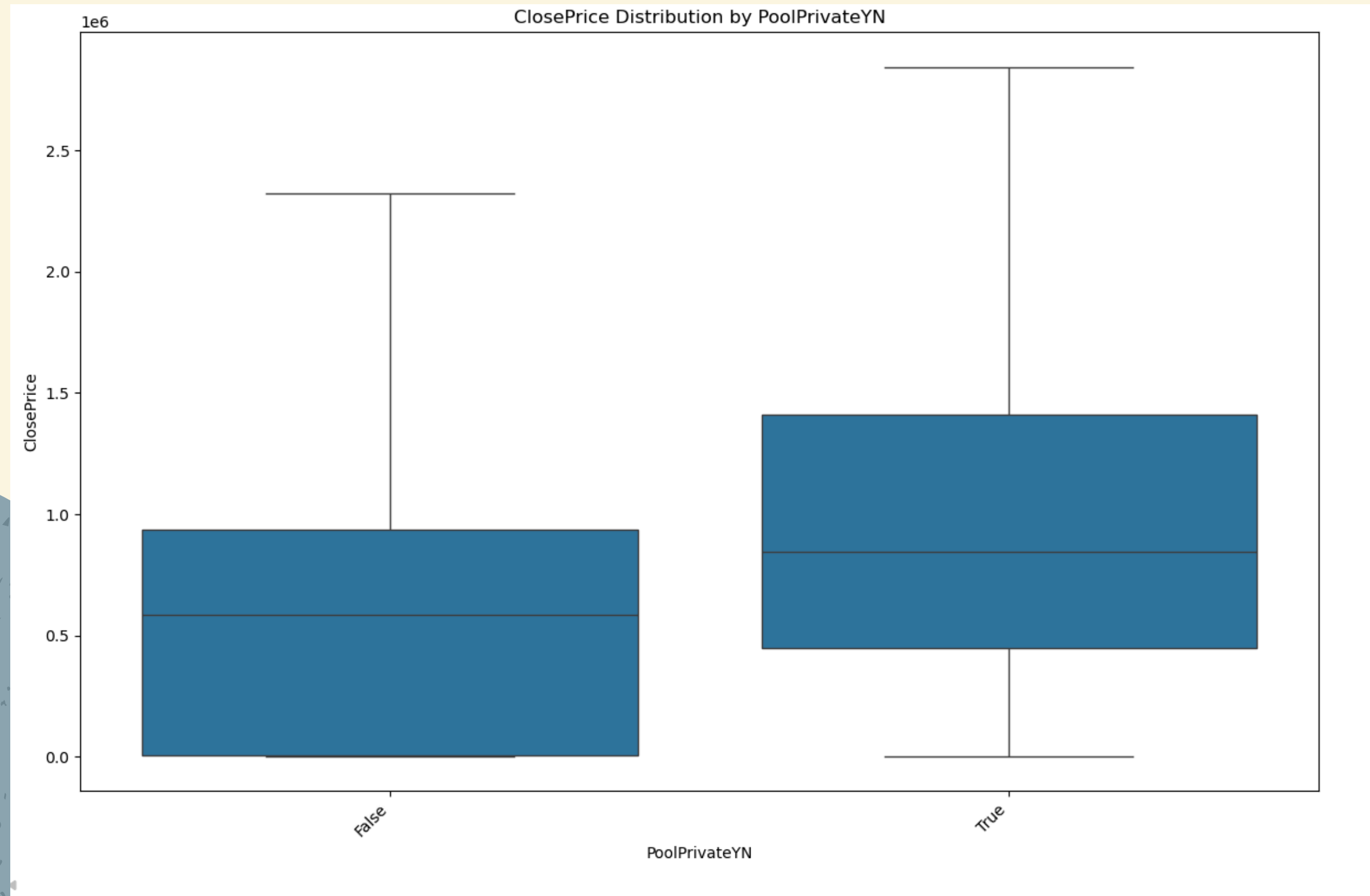
# Close Price by Attached Garage



# Close Price by Fireplace



# Close Price by Private Pool



# What is LightGBM?

Overview	How It Works	Advantages
<ul style="list-style-type: none"><li>• Light Gradient Boosting Machine</li><li>• Decision Tree-based Algorithms</li><li>• High-performance Gradient Boosting Framework</li><li>• Developed by Microsoft</li></ul>	<ul style="list-style-type: none"><li>• Boosting rounds</li><li>• Small decision trees</li><li>• Correcting errors</li><li>• Learning patterns</li></ul>	<ul style="list-style-type: none"><li>• Speed</li><li>• Efficiency</li><li>• Scalability</li><li>• Large datasets</li></ul>

# Feature Engineering & Spatial Intelligence

KNN Neighborhood Features	Domain-Specific Feature Engineering
Added location intelligence beyond city/ZIP → Used latitude/longitude	Created real-estate ratios
Generated neighborhood price statistics (mean, median, quartiles)	Captured real estate behavior (size, age, discounting, tax burden)
<b>IMPACT:</b>	Boosted $R^2$ and reduced MAPE

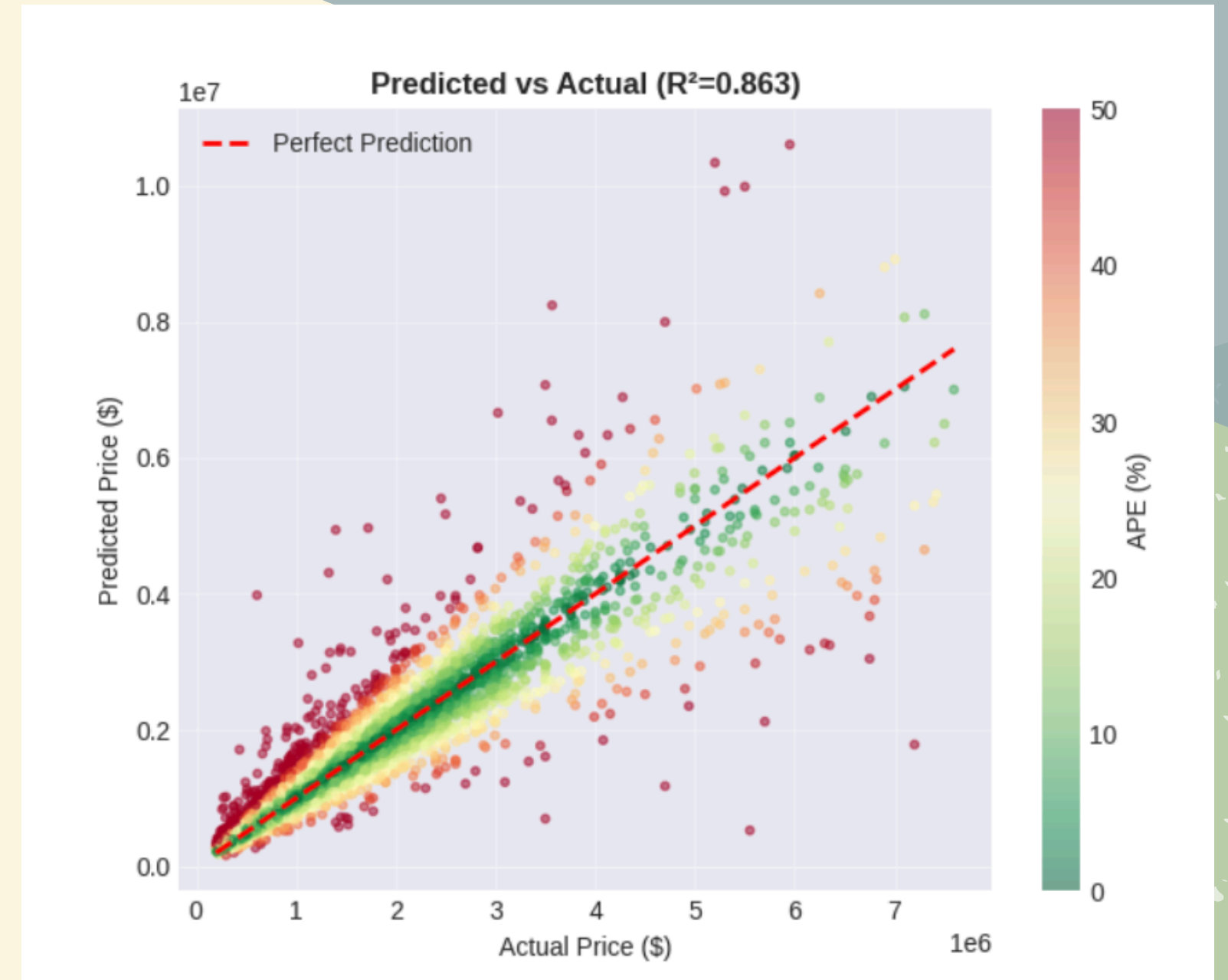
# Hyperparameter Tuning with Optuna

- Performed an automated search to minimize validation MAPE
- Used early stopping to prevent overfitting
- Fined-tune key parameters:
  - num\_leaves, max\_depth, min\_child\_samples, min\_gain\_to\_split
  - learning\_rate, feature\_fraction, bagging\_fraction
  - reg\_alpha, reg\_lambda

→ Achieved the lowest MAPE.

# Results : Key Metrics

- $R^2 = 0.86$  → strong signal capture in noisy real-estate data
- MAPE  $\approx 12\%$  → reliable valuation accuracy.
- Half of all homes within  $\sim 8\%$  error (MdAPE) → model is consistently accurate across most properties, not just on average.



# Why chose XGBoost?

Overview	How It Works	Advantages
<ul style="list-style-type: none"><li>• Tree-based Gradient Boosting Model</li><li>• High-performance XGBoost Regressor</li><li>• Designed for structured/tabular prediction tasks</li><li>• Widely used in industry and Kaggle competitions</li></ul>	<ul style="list-style-type: none"><li>• Ensemble of decision trees built sequentially</li><li>• Each tree corrects residual errors from previous trees</li><li>• Uses gradient boosting to optimize loss</li><li>• L1/L2 regularization to reduce overfitting</li><li>• Handles mixed features and missing values efficiently</li></ul>	<ul style="list-style-type: none"><li>• Strong performance on non-linear, high-dimensional tabular data</li><li>• Captures complex feature interactions naturally</li><li>• Robust with engineered features (geo, text, categorical, etc.)</li><li>• Built-in regularization improves generalization</li><li>• Provides feature importance for interpretability</li></ul>

# Data Wrangling

## Removed Extreme Outliers

- Identified unrealistic prices (multi-million anomalies).
- Such extremes distort model learning and increase variance.
- Applied percentile rules to trim top & bottom extreme values.
- Result: more stable, more generalizable XGBoost performance.

## KNNImputer for Numerical Features

- Real estate data exhibits local similarity (area, bedrooms, location).
- KNNImputer uses nearest neighbors to fill missing numeric values.
- Preserves natural variability vs. flattening via mean/median imputation.
- More realistic feature distribution → better model quality.

## “Missing” Category for Categorical Features

- Categorical missingness can be informative on its own.
- Added explicit “Missing” category instead of mode imputation or dropping rows.
- Enables model to learn patterns associated with missing values.
- Well-aligned with XGBoost’s strong categorical handling

# Feature Engineering

## Numeric Scaling

- Numeric features were processed using median imputation to handle missing values.
- After imputation, features were standardized using StandardScaler to normalize scales and reduce variance differences.
- This prepares numeric inputs for consistent handling by the XGBoost model and prevents features with larger ranges from dominating splits.

## Text Vectorization

- Textual features (Flooring and Levels) were transformed using TF-IDF vectorization (max 100 features each).
- Converts raw text into numeric vectors representing token importance, allowing the model to leverage descriptive property information effectively.
- Ensures that repeated important terms contribute more to feature representation, enhancing predictive power.

## Geospatial Encoding

- Categorical location features (State, County, City) were encoded using hierarchical target encoding, replacing each category with the mean ClosePrice for that group.
- Within each City, properties were grouped using KMeans clustering on (Latitude, Longitude) to generate City\_Cluster labels.
- This captures both local pricing patterns and neighborhood-level spatial clusters, helping the model distinguish areas with systematically different prices.

# Model Tuning

- Performed automated hyperparameter optimization with Optuna to minimize MdAPE (Median Absolute Percentage Error)
- Integrated early stopping in every trial to prevent overfitting
- Optimized key XGBoost parameters impacting tree structure, sampling, regularization, and learning rate

Achieved the lowest MdAPE.

## Basis Results:

- $R^2$ : 0.8038
- MAPE: 22.77%
- MdAPE: 14.97%

## Best Results:

- $R^2$ : 0.8111
- MAPE: 22.21%
- MdAPE: 14.36%

## Conclusion:

There is no significant improvement. It could be the limitation of the model and the included features.

Next Step:

## Exploring Features

Ratios:

- Price per sqft =  $\text{ClosePrice} / \text{LivingArea}$
- Lot-to-living ratio =  $\text{LotSizeArea} / \text{LivingArea}$

Further analysis on price by region:

- Include median and variance for each district.

Also, it could be the limitation of explainability on XGBoost.

We can explore on other advanced model as well such as CatBoost and MLP. Or mix models by ensemble.

# MLP Introduction

What is MLP Regression?

Why choose MLP Regression?

Definition	Advantage	Disadvantage
<ul style="list-style-type: none"><li>• Multilayer Perceptron, a kind of <b>neuron network</b> with multiple layers</li><li>• Input → layer 1 → layer 2 → ... → output</li><li>• Nonlinear activation function</li><li>• Can be used for classification and regression tasks</li></ul>	<ul style="list-style-type: none"><li>• Flexible design, suitable for ANY function</li><li>• Strong accuracy with large datasets</li><li>• Regression, Classification, Pattern recognition, ...</li><li>• Housing price are often non-linear</li></ul>	<ul style="list-style-type: none"><li>• Large data &amp; demanding computation resources</li><li>• Highly hyper-parameter choices sensitive</li><li>• “Black-box” model</li><li>• Local minimum risk</li></ul>

# Feature Engineering

## Outliers

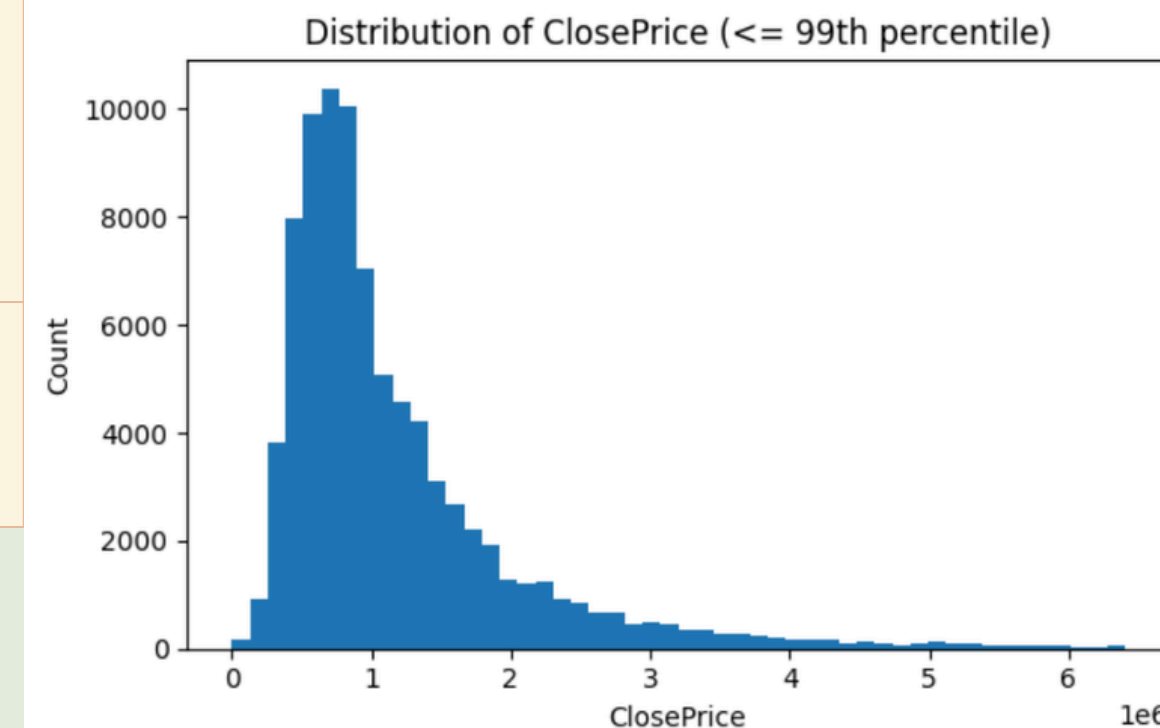
- Filter both **Price** and **Price Per Square Feet**
- **Upper bound:** min(Tukey bound, quantile bound) - stable results
- Price < \$10000

## Feature Building

- **Additional variables:** Property Age, Baths per Bed, Lot per Living,  $\log(\text{Living Area})^2$ ,  $\log(\text{Living Area}) \times \text{Total Bathrooms}$
- **Out-of-fold target encoding with smoothing:** deal with geographic features (City, Zip, ...)
- **Multi-scale RBF** (radial basis function): transfer longitude and latitude variables to distances to related centers
- **Out-of-fold KNN:** the impact of prices in neighborhood

## Final Input

- Total number of features: **188**



# MLP Model

**Model structure:** Input → Dense(512) → Dense(256) → Dense(128) → Output(1)

Output(1) is **log(Price)**

Model Overview	Training	Post Training
<ul style="list-style-type: none"><li>• Dense -&gt; ReLU</li><li>• BatchNorm -&gt; Dropout -&gt; L2</li></ul>	<ul style="list-style-type: none"><li>• Early stopping + LR scheduler</li></ul>	<ul style="list-style-type: none"><li>• Smearing Correction</li></ul>
<p><b>Dense + ReLU</b> brings:</p> <ol style="list-style-type: none"><li>1. Non-linearity;</li><li>2. Model stability;</li><li>3. Computation efficiency;</li><li>4. Approximation to any function.</li></ol> <p><b>BatchNorm</b> - layer normalization, brings:</p> <ol style="list-style-type: none"><li>1. Faster convergence;</li><li>2. Overfitting prevention - <b>layer level</b>.</li></ol> <p><b>Dropout</b> brings:</p> <ol style="list-style-type: none"><li>1. Overfitting prevention - <b>neuron</b> level.</li></ol> <p><b>L2</b> brings:</p> <ol style="list-style-type: none"><li>1. Overfitting prevention - <b>weight</b> level.</li></ol>	<p><b>Early stopping</b> brings:</p> <ol style="list-style-type: none"><li>1. Best performance checkpoint;</li><li>2. Overfitting prevention.</li></ol> <p><b>Learning rate</b> scheduler brings:</p> <ol style="list-style-type: none"><li>1. Fine-tune around the minimum.</li></ol>	<p><b>Smearing correction</b> brings:</p> <ol style="list-style-type: none"><li>1. Downward bias caused by exponentiation log-predictions</li></ol>

# Results

- The model performs well on the cleaned dataset
- BUT ...

<b>R<sup>2</sup></b>	87.15%
<b>MAPE</b>	12.79%
<b>MdAPE</b>	8.72%

Model performance on Full dataset:

<b>R<sup>2</sup></b>	0.32%
<b>MAPE</b>	13.2%
<b>MdAPE</b>	8.83%

Price Segment:

<b>0k - 500k</b>	<b>R<sup>2</sup></b>	-0.58
<b>500k-1000k</b>	<b>R<sup>2</sup></b>	0.33
<b>1000k-infk</b>	<b>R<sup>2</sup></b>	0.75

- **Model Interpretation**

Can we interpret machine learning models?

**SHAP Method** – Complex machine learning models interpretation

SHapley Additive exPlanations, a game-theoretic approach to explain the output of **any machine learning model** by calculating the **marginal contribution** of a feature to the prediction.

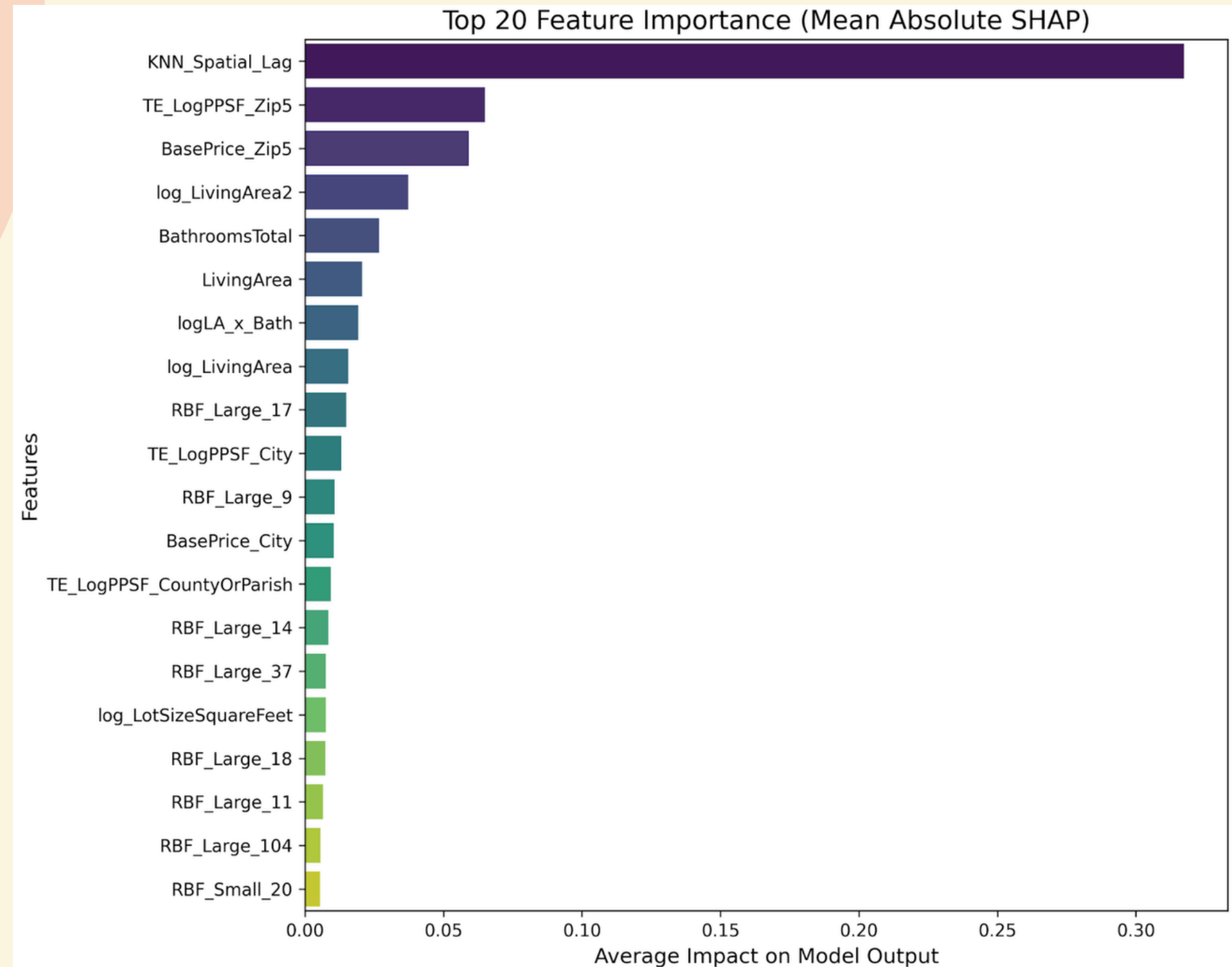
**Advantage:**

1. Consistency
2. Directionality

# Results

This model prioritizes location over physical attributes.

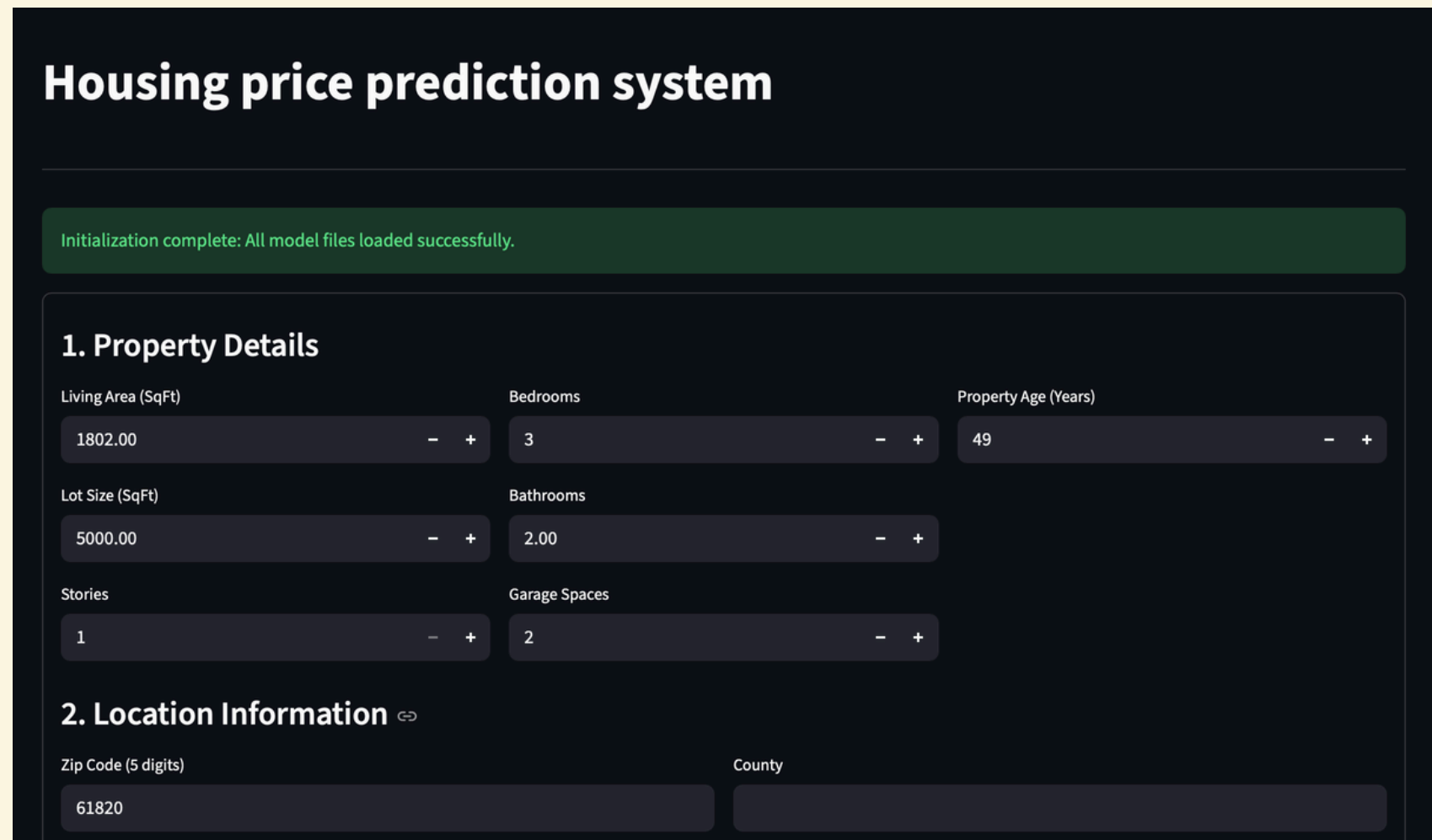
- **Dominant Feature:** Spatial Lag
- **Macro-Location:** Zip Code
- **Physical Attributes:** Living area, Bedrooms, ...



- This is a high performance model on cleaned data and high-end properties but with limitations on generalization.
- Lower-priced homes have higher variance in property condition that our current data doesn't capture.

# Future steps

- **Segmented modeling:** Develop separate models for the Luxury Segment (>1M) and the Economy Segment (<500k).
- **Unstructured data integration:** Incorporate NLP on listing descriptions or Computer Vision on property photos.
- **Robust loss functions:** Experiment with Huber Loss or Quantile Regression instead of standard MSE.
- **Streamlit App:** A real useable model.



The screenshot shows a web application titled "Housing price prediction system". At the top, a green notification bar states "Initialization complete: All model files loaded successfully." Below this, the interface is divided into two main sections: "1. Property Details" and "2. Location Information".

**1. Property Details**

Living Area (SqFt)	Bedrooms	Property Age (Years)
1802.00	3	49

Lot Size (SqFt)	Bathrooms
5000.00	2.00

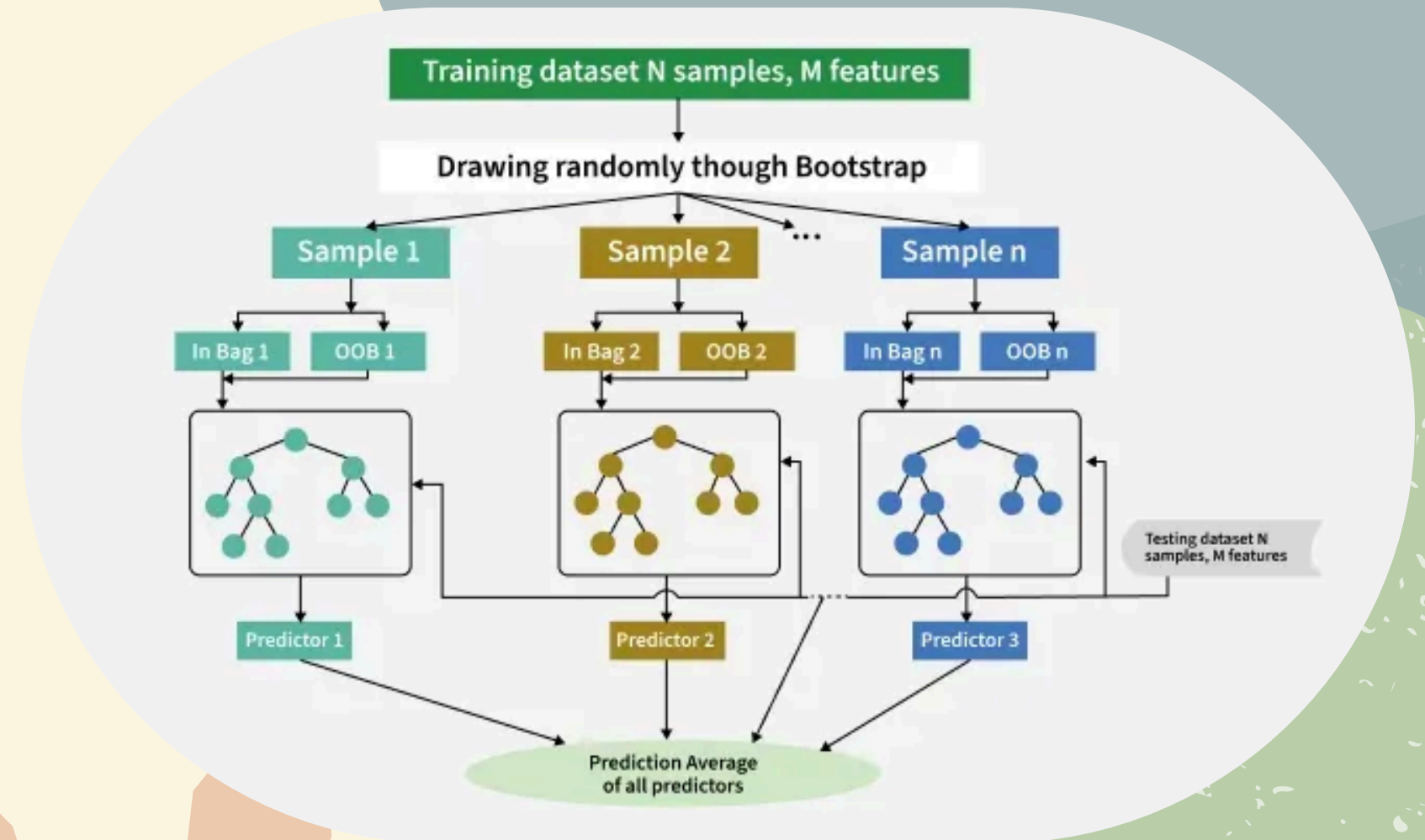
Stories	Garage Spaces
1	2

**2. Location Information**

Zip Code (5 digits)	County
61820	

# What is CatBoost?

- Gradient Boosting framework by Yandex
- Optimized for tabular data with categorical features
- Handles categorical variables automatically
- Strong performance with minimal preprocessing



# Advantages

Automatic Handling of Categorical Features

Symmetric Tree Structure

Supports GPU Training and Missing Values

# Disadvantages

Slower on Very Large Datasets

Fewer Hyperparameter Tuning Options

Limited Interpretability

# Data Preprocessing

Filtering	Cleaning	Pruning	Target Transformation
<ul style="list-style-type: none"><li>Retained only Residential and SingleFamilyResidence properties</li></ul>	<ul style="list-style-type: none"><li>Converted string-based numeric fields to numeric types</li></ul>	<ul style="list-style-type: none"><li>Removed listing-related fields</li><li>Kept only samples with sale prices within the 0.5%–99.5% quantile range</li></ul>	<ul style="list-style-type: none"><li>Applied <math>\log_{1p}(\text{ClosePrice}) \rightarrow y_{\log}</math> to reduce skewness in price distribution</li></ul>

# Model Iteration

## Difference:

- $R^2 \uparrow 0.62$
- MAPE  $\downarrow 6\%$
- MdAPE  $\downarrow 5\%$

## v1: Baseline

- CatBoost with external data (Income, Schools).
- Result: Poor performance due to missing external data.
  - $R^2 \approx 0.25$
  - MAPE  $\approx 17\%$
  - MdAPE  $\approx 12\%$

## v2: Spatial

- Engineered "Mean Neighbor Price (20 nearest)".
- Result: Massive accuracy boost.
  - $R^2 \approx 0.86$
  - MAPE  $\approx 11\%$
  - MdAPE  $\approx 8\%$

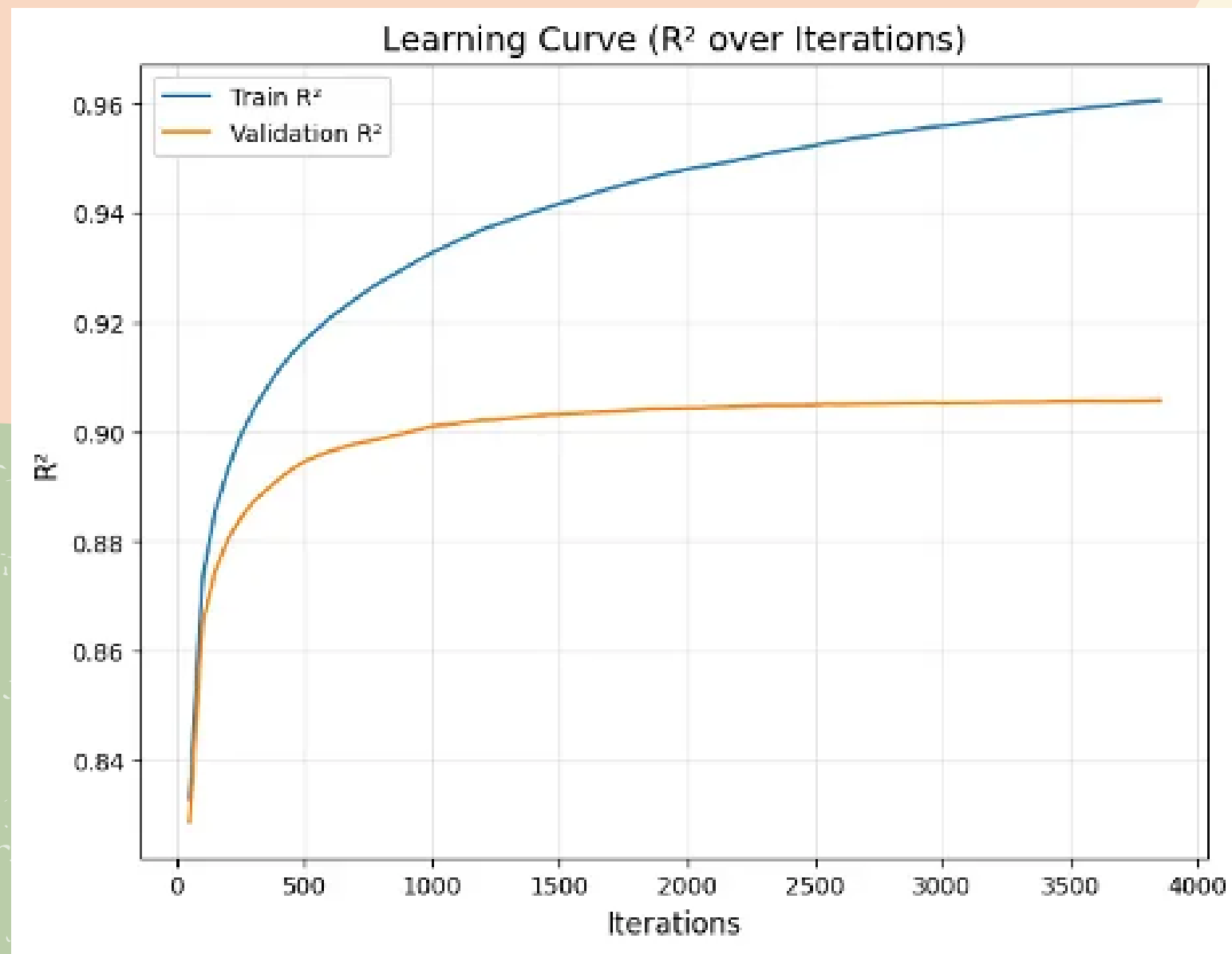
## v3: Remote-Aware

- Hybrid approach separating remote vs. urban areas.
- Result: Limited gain due to sparse remote data.

## v2.1: Final

- Refined v2 with stricter California-only filtering.
- Result: Most robust performance.
  - $R^2 \approx 0.87$
  - MAPE  $\approx 11\%$
  - MdAPE  $\approx 7\%$

# Model Diagnostics



## Residual Patterns

- Residuals are mostly centered around zero, indicating no strong systematic bias.
- High-priced properties exhibit a funnel-shaped pattern, where prediction variance increases with price.
- A small number of extreme outliers are observed (approximately  $\pm 2M$  to  $\pm 6M$ ), suggesting difficulty capturing rare or highly unique homes.

## Learning Curves (RMSE)

Training and validation RMSE decrease together, showing no signs of overfitting.

## Learning Curves ( $R^2$ )

Validation  $R^2$  increases steadily and remains close to the training curve, indicating good generalization and stable learning behavior.

# Next Steps / Future Work

- **Increase historical data for remote areas**
  - Helps reduce errors caused by small sample size and distinct distribution patterns in remote regions.
- **Build a dedicated model for high-priced properties**
  - Since luxury homes have unique characteristics, a specialized model may capture their pricing patterns more accurately.

# Models Performance Comparison

(Training: 2025.01–2025.08; Testing: 2025.09)

Model	$R^2$	MAPE	MdAPE
<b>MLP</b>	87.15%	12.79%	8.72%
<b>LightGBM</b>	85.74%	12.14%	7.75%
<b>XGBoost</b>	81.11%	22.21%	14.36%
<b>CatBoost</b>	86.82%	11.29%	7.71%

# Advantages & Disadvantages

Model	Advantages	Disadvantages
MLP	<ul style="list-style-type: none"><li>• Powerful nonlinear learning</li><li>• highly flexible</li></ul>	<ul style="list-style-type: none"><li>• Needs large data- easily overfits</li><li>• low interpretability</li></ul>
LightGBM	<ul style="list-style-type: none"><li>• Very fast</li><li>• scalable</li><li>• handles missing values</li></ul>	<ul style="list-style-type: none"><li>• Sensitive to tuning</li><li>• weak categorical support</li><li>• higher overfitting risk</li></ul>
XGBoost	<ul style="list-style-type: none"><li>• Strong tabular performance</li><li>• good regularization</li><li>• handles missing values</li></ul>	<ul style="list-style-type: none"><li>• Heavy training &amp; tuning cost</li><li>• risk of overfitting with poor features</li></ul>
CatBoost	<ul style="list-style-type: none"><li>• Best for categorical data</li><li>• strong defaults</li><li>• stable</li></ul>	<ul style="list-style-type: none"><li>• Slower on very large data</li><li>• fewer tuning options</li></ul>

# Common Feature Engineering

- **Geospatial / Neighborhood Features**

KNN + coordinates to create neighborhood price/distance features  
→ boosts  $R^2$  and lowers MAPE.

- **Target Log Transformation**

Apply  $\log_{1p}$  to ClosePrice → reduces skew and stabilizes training.

- **Outlier Filtering**

Use price bounds or percentile/IQR rules  
→ remove extremes and improve model stability.

# Future Steps

Data Enhancements	Modeling Enhancements	Validation Improvements
<ul style="list-style-type: none"><li>• Expand dataset (Sep–Dec 2025) for seasonality</li><li>• Add macroeconomic signals (rates, CPI, unemployment)</li><li>• Enrich location intelligence (schools, commute distance)</li></ul>	<ul style="list-style-type: none"><li>• Experiment with ensemble stacking (LightGBM + XGBoost + LSTM).</li><li>• Apply deeper explainability methods (SHAP, PDPs).</li><li>• Deploy the model as an interactive Streamlit/dashboard app.</li></ul>	<ul style="list-style-type: none"><li>• Implement time-based cross-validation for realistic forecasting.</li><li>• Generate regional error heatmaps.</li><li>• Identify geographic overpricing/underpricing clusters.</li></ul>

# What We've Learned



- Recognized that data cleaning drives model quality.
- Observed that feature engineering often outperformed algorithm changes.
- Improved model accuracy using neighborhood + ratio features.
- Strengthened teamwork through structured collaboration.

Thank you!

