# PERSONALIZED REAL ESTATE VALUATION SYSTEM WITH PRICE PREDICTION

## Team members

Jackson Jose &
Jijo John
Date: 01-Mar-2025 – 10-May-2025

## Summary

This project focuses on developing a robust machine learning system for accurately predicting real estate property values. Utilizing the Zillow Housing Dataset, the system employs an ensemble regression approach to capture complex relationships between property features, location, market trends, and seasonal effects. The project encompasses comprehensive data cleaning and preprocessing, exploratory data analysis (EDA) to uncover insights, training of multiple regression models, and the development of a Stacking Regressor for final predictions. Key deliverables include the ensemble model, feature importance analysis, and the groundwork for a future interactive Flask web application for property assessment. Evaluation metrics such as RMSE, MAE, and $R^2$ are used to demonstrate the model's predictive power and accuracy.

# 1 Introduction

Accurate real estate valuation is crucial for buyers, sellers, investors, and financial institutions, but traditional methods can be subjective and slow. Primary goal to develop a data-driven machine learning system for accurately predicting property sale prices.

Significance of the project is to demonstrates the efficacy of machine learning in providing more objective, efficient, and potentially more accurate property valuations compared to traditional approaches.

Core Methodology: Employs ensemble regression techniques, culminating in a Stacking Regressor, to capture complex market dynamics.

Key Project Phases: Includes meticulous data cleaning, insightful exploratory data analysis (EDA), training diverse models, and rigorous evaluation.

Target Outcome: Aims to deliver a robust predictive model, analyze key value drivers (feature importance), and lay groundwork for a practical valuation tool (e.g., Flask app).
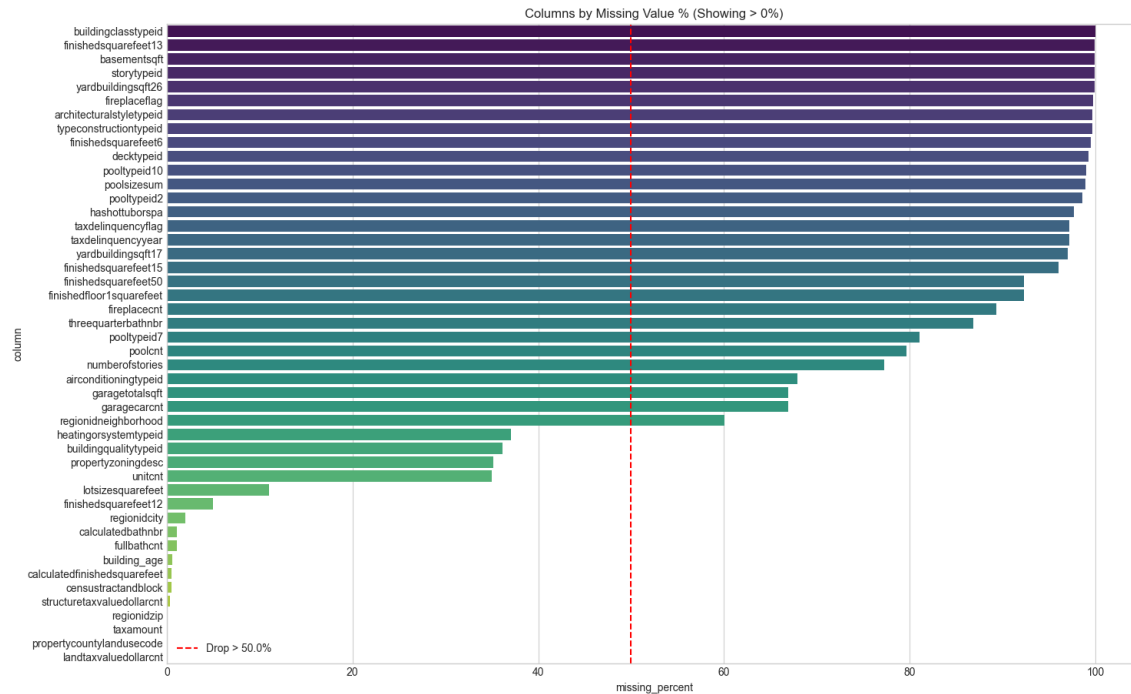
## 2. Dataset Overview

Dataset Used: Leverages the comprehensive Zillow Housing Dataset, rich in property features and transaction data.

## 3. Data Cleaning and Preprocessing

A rigorous data cleaning and preprocessing pipeline was implemented to prepare the data for modeling:
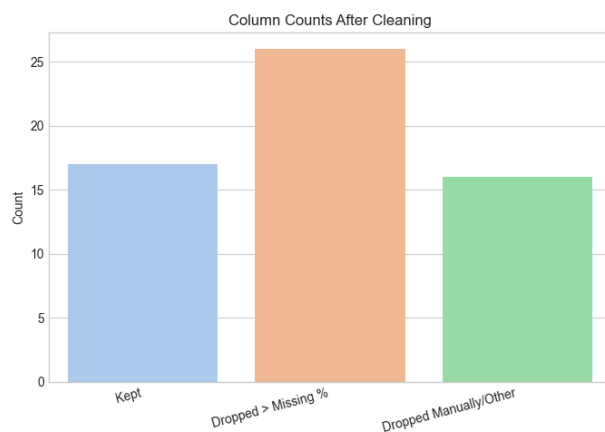
- Initial Data Loading & Merging:
- Dictionary Mapping & Target Variable Engineering:
- Handling Missing Values:
    1. An initial analysis of missing values across all features was conducted. The distribution of missing percentages is visualized below, highlighting features with substantial missing data.

Columns by Missing Value % (Showing > 0%)

- Columns with a missing value percentage greater than 50% were dropped to reduce noise and dimensionality.
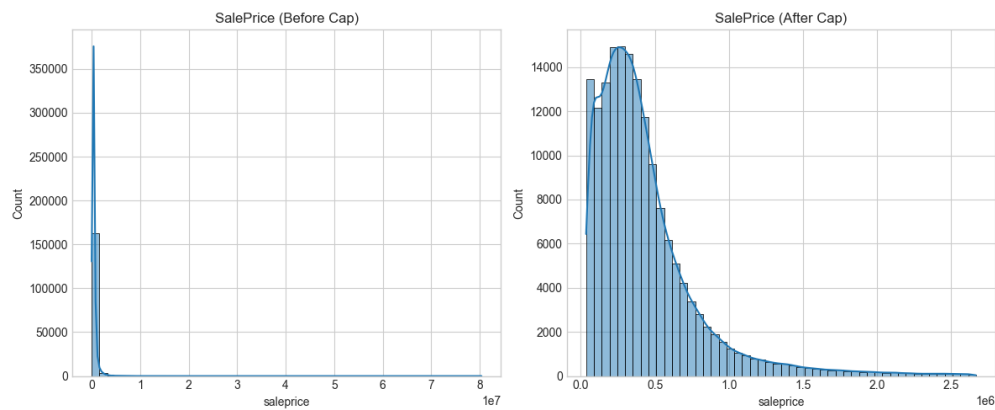
**Column Dropping (Manual/Redundancy/Leakage):**

- Identifier columns (`parcelid`), redundant location identifiers (while keeping `regionidzip` and `fips`), raw census data, and tax-related columns that could lead to data leakage were manually removed.
- The following chart summarizes the column reduction process:


Column Counts After Cleaning

(Caption: Figure 2: Summary of columns kept versus those dropped due to high missing percentages or manual selection.)

**Outlier Handling:**

- The distribution of the target variable, sale price, was examined before and after outlier capping (e.g., at 1st and 99th percentiles) to ensure a more robust modeling process.
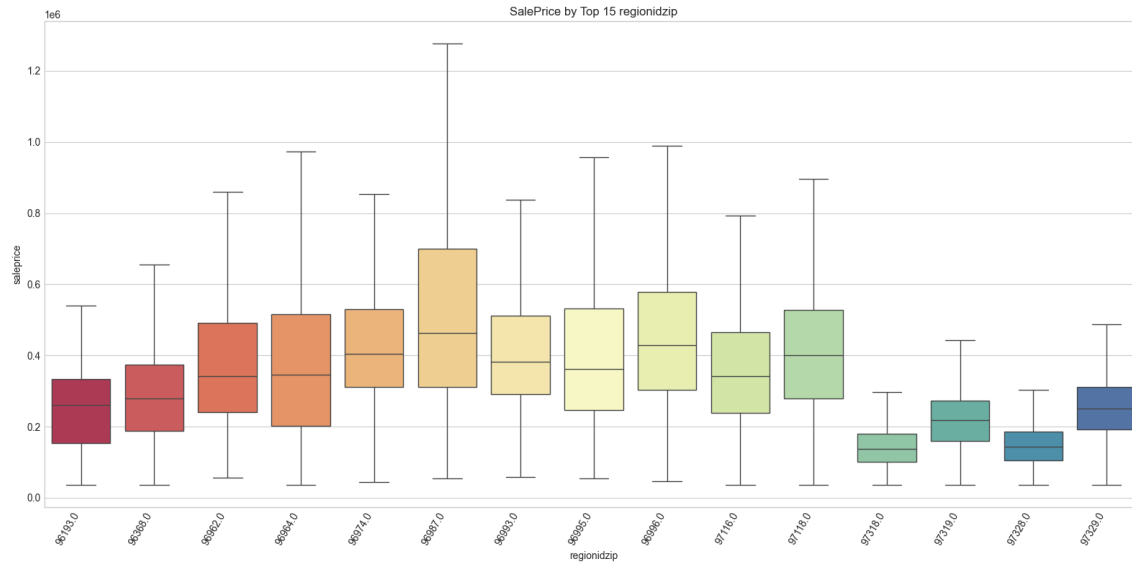


- Similar outlier treatment was applied to key numerical features like calculatedfinishedsquarefeet

- The final cleaned dataset was then used for EDA and model building.

## 4. Exploratory Data Analysis (EDA)

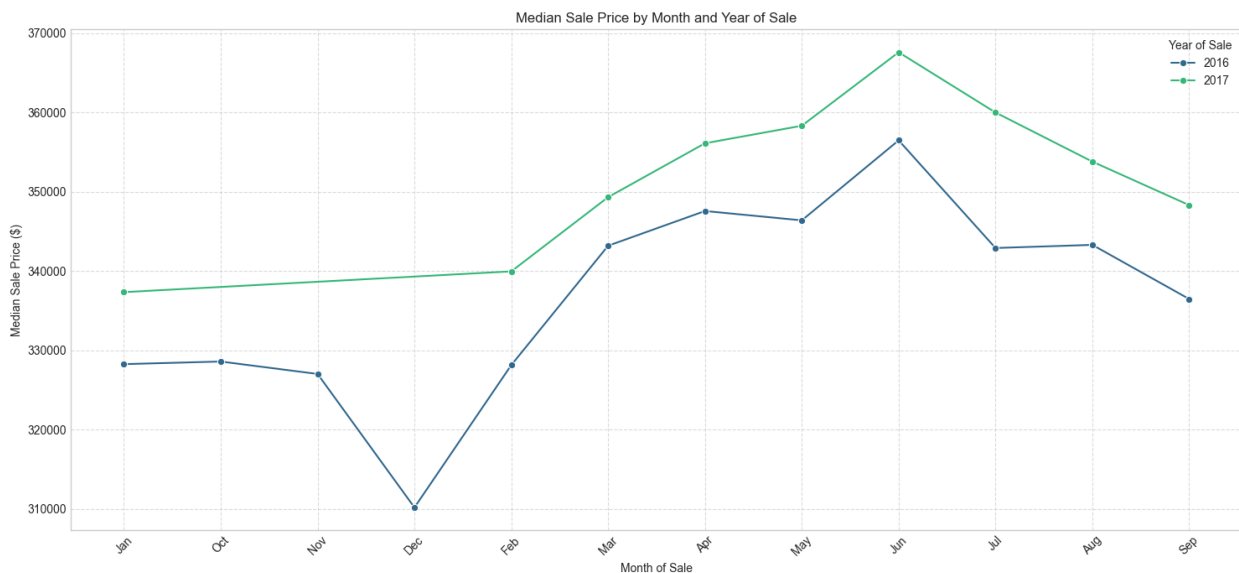EDA was performed on the cleaned dataset to uncover patterns...

**Sale Price Distribution Across Neighborhoods:**

- Boxplots show the Sale Price distribution across the top 15 (or 20, adjust based on your chart) most frequent ZIP codes. This reveals significant price variations by locality.
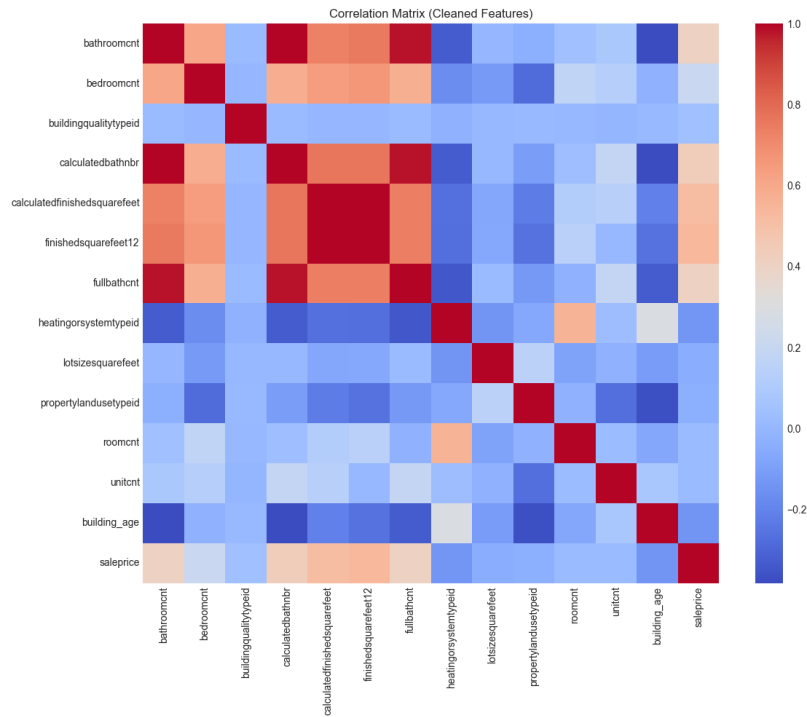
SalePrice by Top 15 regionidzip

- **Seasonal Trends in Property Prices:**

    o *As described previously, assuming seasonal_price_median_lineplot.png and seasonal_price_multiyear_lineplot.png were generated if data allowed. If not, this section might state that seasonal features were explored but showed minimal clear patterns or data was insufficient.*
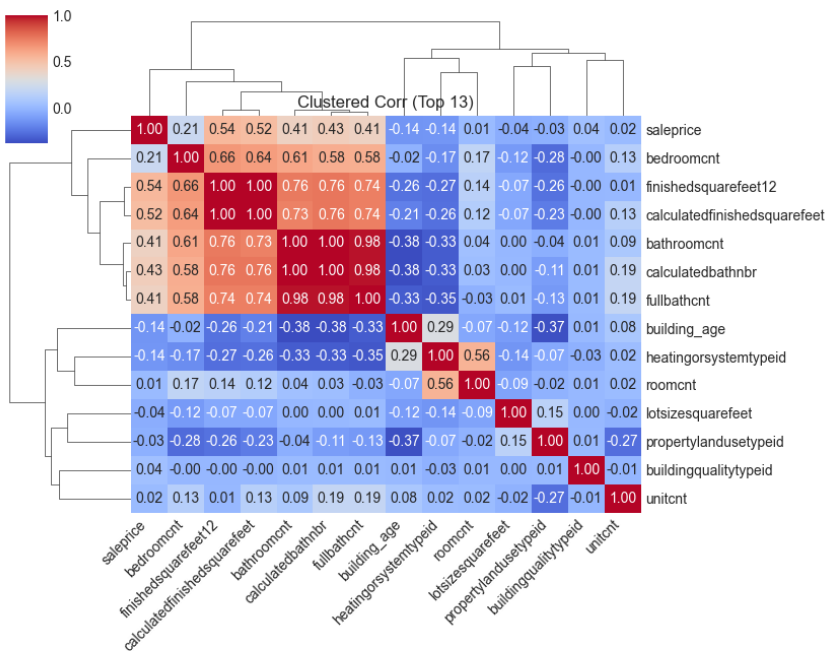

Median Sale Price by Month and Year of Sale

## Correlation Analysis:

- A correlation heatmap was generated for key numerical features and sale price from the cleaned dataset to identify linear relationships.

Correlation Matrix (Cleaned Features)

- A clustered heatmap was also created for the top 13 (or as per your chart) most correlated features with saleprice (or a relevant subset), reordering features based on similarity to highlight groups of correlated variables.


Clustered Corr (Top 13)

## 5. Modeling Approach

- Split data into train/test sets (80/20 split)
- Trained multiple models: Ridge, Lasso, Random Forest, Gradient Boosting, LightGBM
- Combined models using Stacking Regressor for improved performance
- Evaluated using RMSE, MAE, and R² Score

## 6. Model Evaluation Summary

- **Performance Metrics:** Models were evaluated using RMSE, MAE, and $R^2$ Score. The table below summarizes the performance of the key models tested on the held-out test set.
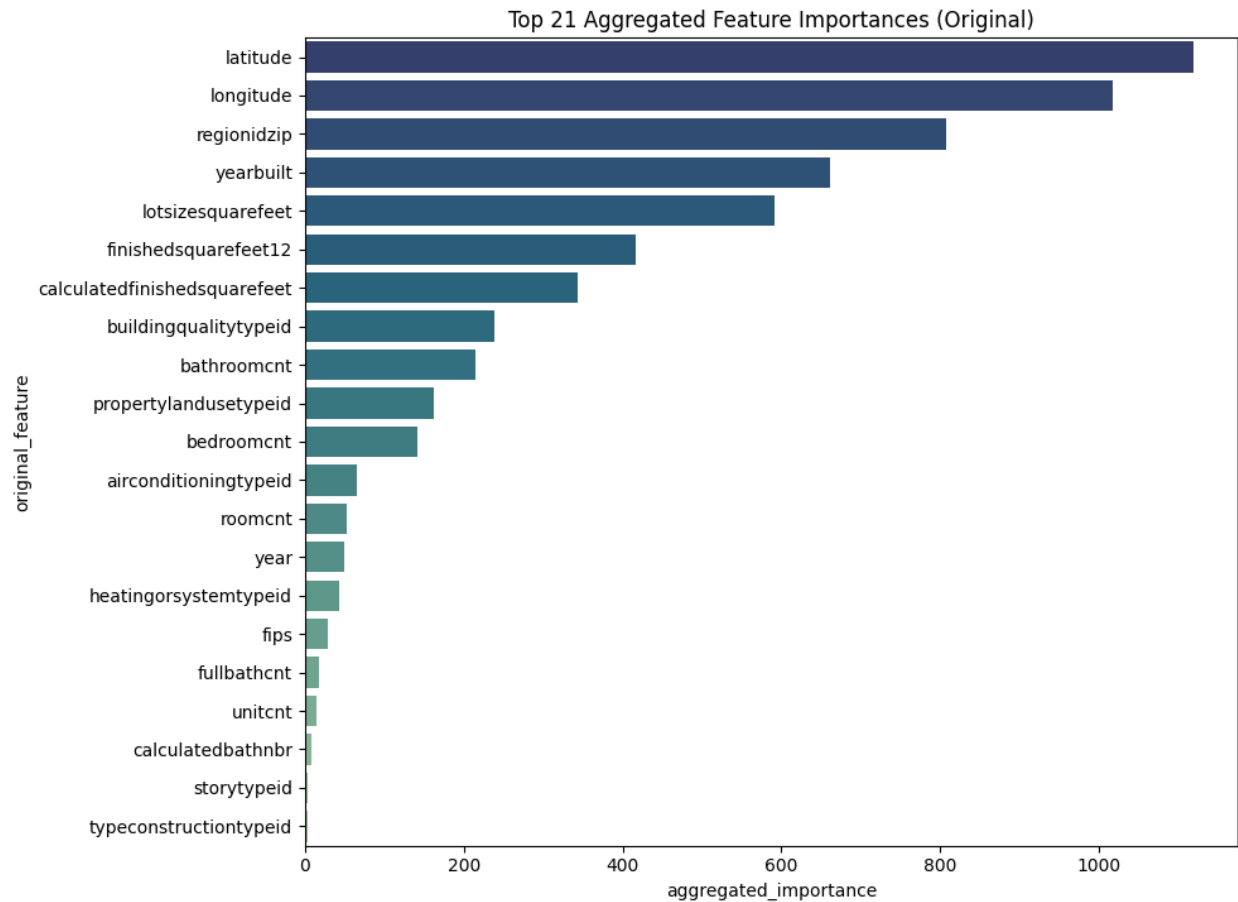
- **Table 1: Model Performance Comparison**

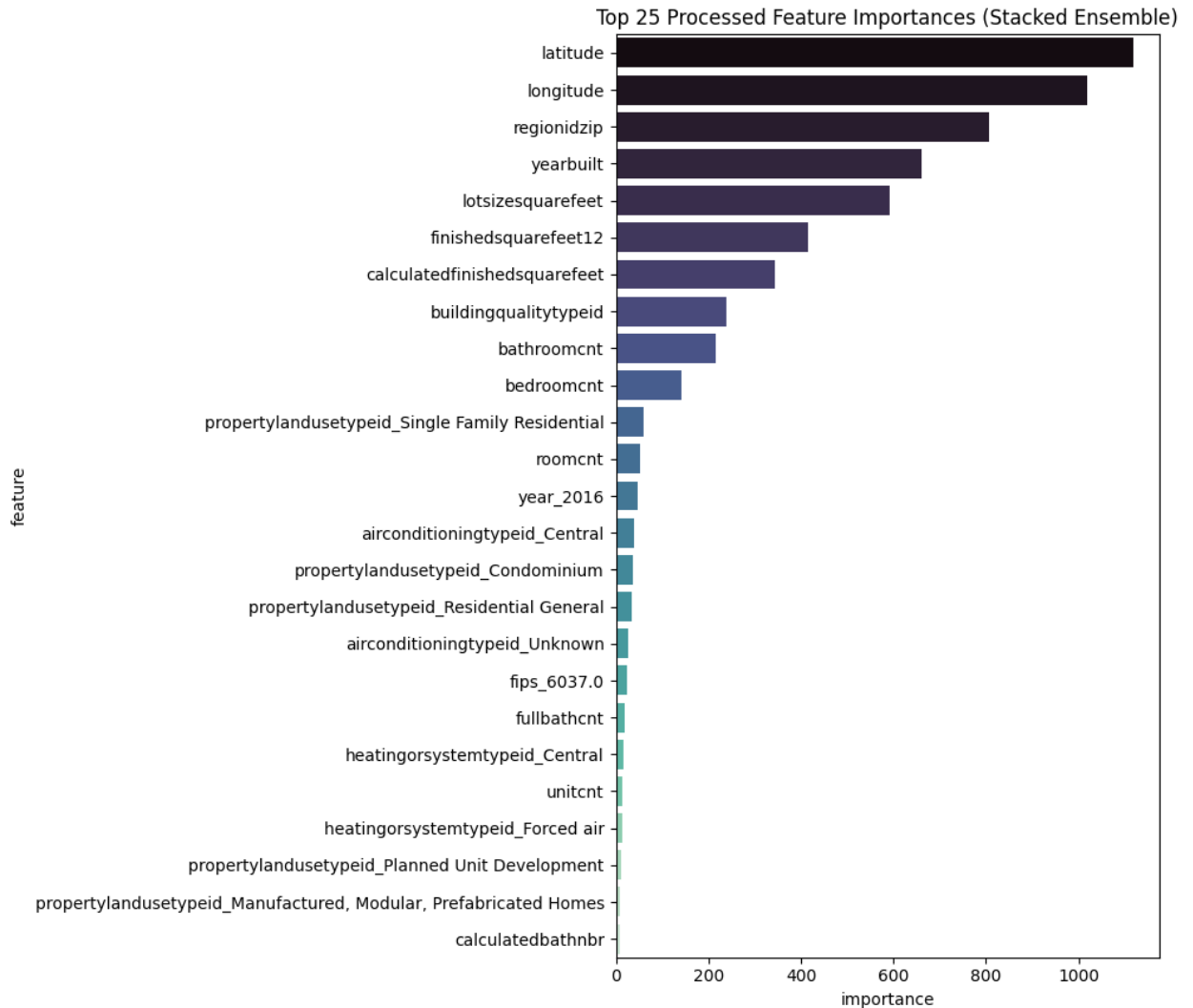| Model | RMSE | MAE | R2 | Training Time (s) |
|---|---|---|---|---|
| Ridge | 276846.6443 | 188488.9682 | 0.358018 | 0.10938 |
| Lasso | 276842.0568 | 188486.7022 | 0.358039 | 71.763522 |
| Linear Regression | 276839.628 | 188484.7048 | 0.358051 | 0.652342 |
| Random Forest | 238581.7946 | 157308.2173 | 0.523219 | 38.579632 |
| LightGBM | 238404.455 | 160702.7364 | 0.523928 | 1.514389 |
| Stacked Ensemble | 236202.6629 | 157769.7973 | 0.532681 | 180.655839 |

**Best Performing Model:** Stacked Ensemble

**Feature Importance:**

- Feature importances were extracted to understand the key drivers of property value. The "Top 16 Aggregated Importances" chart shows the influence of the original input features (after aggregating importances from OHE components where applicable). Location identifiers (regionidzip, fips) and property size (finishedsquarefeet12, calculatedfinishedsquarefeet) often rank highly.

Top 21 Aggregated Feature Importances (Original)

- The "Top 25 Processed Importances" chart details the importances of features as seen by the model after preprocessing (including individual One-Hot Encoded categories). This often highlights specific ZIP codes or FIPS areas if the model was, for example, Ridge, where OHE'd features get individual coefficients.

Top 25 Processed Feature Importances (Stacked Ensemble)

## 7. Artifact Saving

To ensure reproducibility, enable deployment, and facilitate future analysis, key components of the developed system were serialized and saved. This included:

- **Final Model Pipeline (final_model_pipeline.joblib):** The complete, fitted pipeline containing the optimal preprocessor and the best-performing trained ensemble model (e.g., Stacked Regressor). This is the primary artifact for making new predictions.

- **Fitted Preprocessor (preprocessor.joblib):** The preprocessor (ColumnTransformer) instance, fitted on the training data, saved separately for potential independent use or inspection.
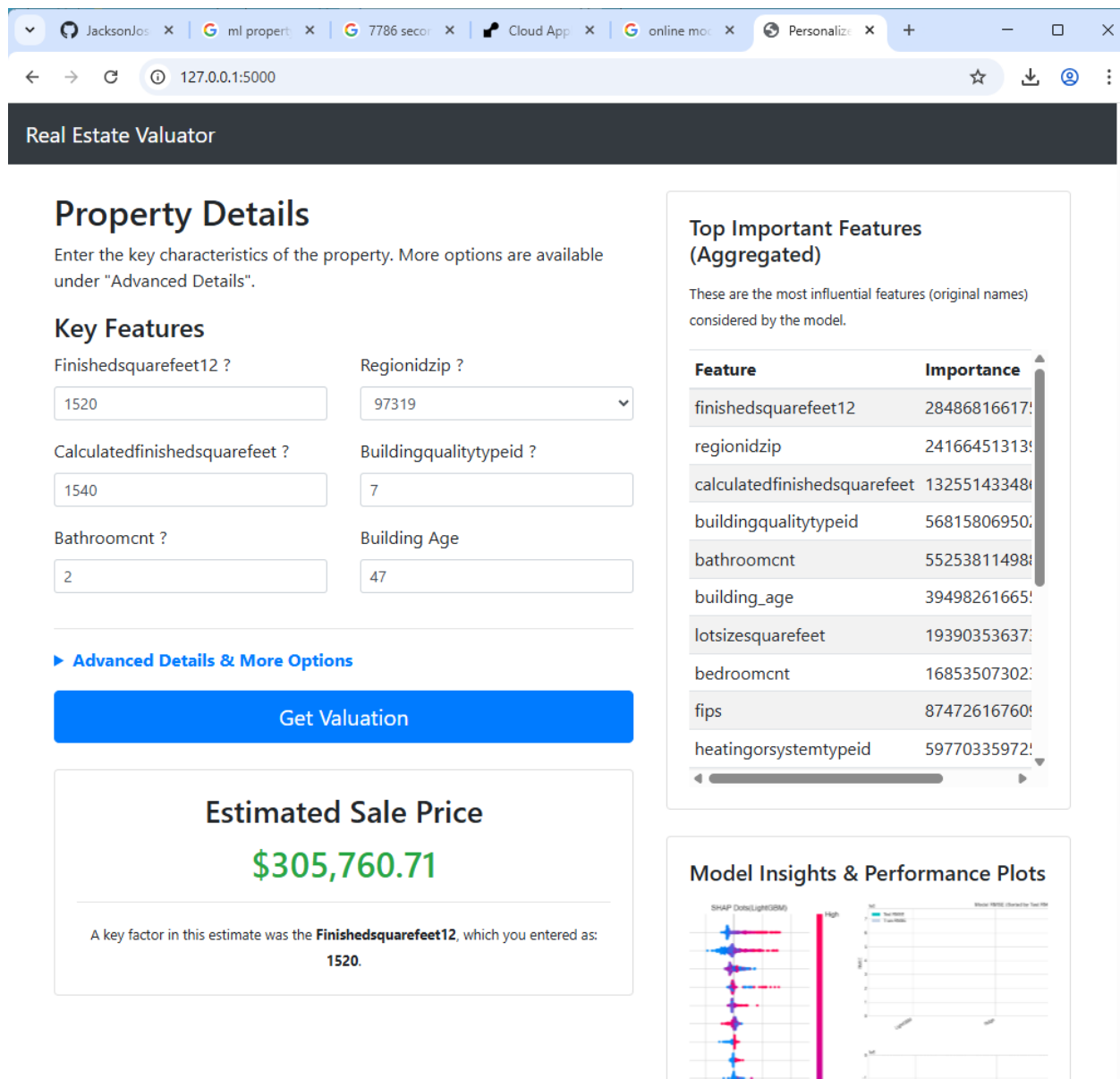
- **Feature Definitions & UI Data:**

o   expected_columns_before_preprocessing.joblib: List of raw input feature names the model expects.

o   final_feature_names_post_ohe.joblib: List of feature names after One-Hot Encoding, used for interpreting processed feature importances.

o   feature_importance_table_data.joblib: Data (top N features and scores) for the UI's detailed feature importance table.

o   field_order_preference.joblib: List of original feature names, sorted by aggregated importance, to guide form field order in the UI.

o   Supporting dictionaries (field_descriptions.joblib, id_value_mappings.joblib) and original feature type lists from data processing were also saved for the Flask application.

- **Visualizations:** All relevant plots generated during EDA and feature importance analysis were saved as image files (e.g., in static/plots/) for inclusion in reports and the web application.

This systematic saving of artifacts ensures that the model can be easily loaded and utilized in other environments, such as the planned Flask web application.

## 8. Conclusion & Future Work

**Conclusion:** This project successfully developed an ensemble regression model for predicting real estate prices. The model captures influences from location, size, and other property attributes. The EDA provided valuable insights into data characteristics and price drivers.

**Flask Web Application:** The generated artifacts pave the way for deploying an interactive Flask application.

**Docker:** Deploy Docker Image and Configure Service

**Git :** https://github.com/JacksonJosekwt/SUBMISSION2

**Future Work:**

Include more granular location data, advanced feature engineering, hyperparameter tuning, alternative encoding, time series analysis, SHAP, regular retraining, confidence intervals.