

Satellite Imagery Based Property Valuation using Multimodal Regression

Name: Ashi Agrawal

Enrollment Number: 24118012

Date: 5 Jan 2026

1. Introduction

Property prices are influenced not only by intrinsic housing attributes but also by neighborhood characteristics such as infrastructure, accessibility, and environmental surroundings. Traditional tabular data often fails to capture this contextual information. This project incorporates satellite imagery alongside structured housing data to capture neighborhood-level visual context and improve property price prediction using a multimodal regression approach.

Modeling Overview

This project follows a multimodal regression approach for property valuation by combining structured housing attributes with satellite imagery. A tabular-only regression model is first trained to establish a baseline performance. In parallel, satellite images corresponding to property locations are processed using a pretrained convolutional neural network to extract neighborhood-level visual features. These image embeddings are concatenated with tabular features to form a multimodal feature representation, which is then used for price prediction. Model performance is evaluated by comparing the tabular-only and multimodal approaches.

2. Dataset Description

The dataset consists of residential property records containing structured attributes such as number of bedrooms, bathrooms, living area, lot size, construction grade, and geographical coordinates (latitude and longitude). The target variable for prediction is the market price of the property.

Satellite images corresponding to each property location were programmatically acquired using the latitude and longitude values. These images capture neighborhood-level characteristics that are not directly available in the tabular dataset.

3. Data Preprocessing and Feature Engineering

Data preprocessing involved removing duplicate entries and dropping non-essential columns such as identifiers and timestamps that do not contribute to prediction. The cleaned dataset was then used for feature engineering.

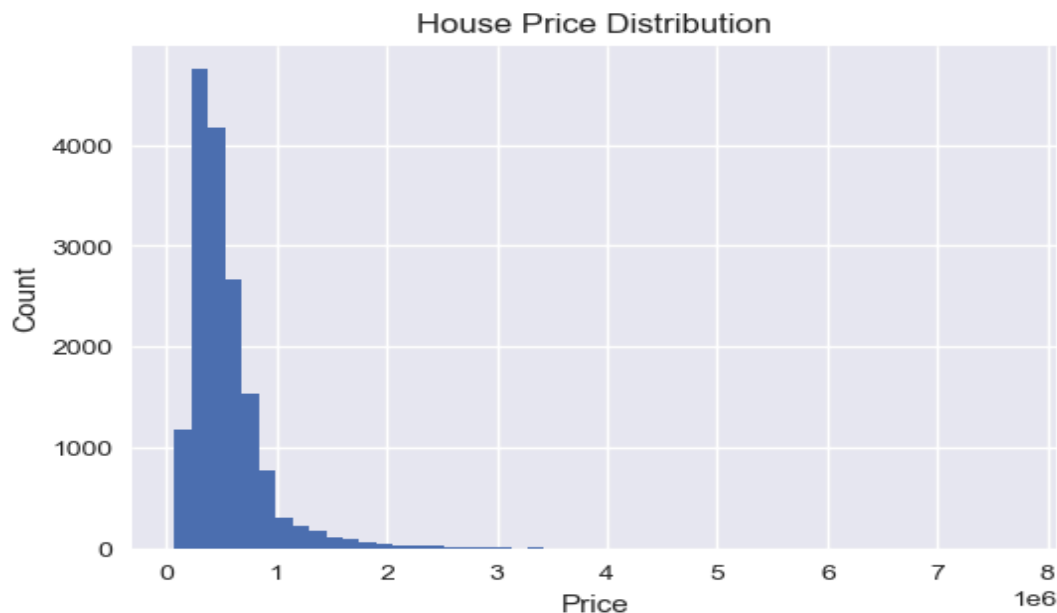
Several derived features were created to enhance the predictive power of the model, including:

- Property age calculated from the year of construction
- Basement area ratio relative to total living area
- Difference between property living area and neighborhood average
- A binary indicator for high-grade construction

These engineered features help incorporate domain knowledge and improve model interpretability.

4. Exploratory Data Analysis (EDA)

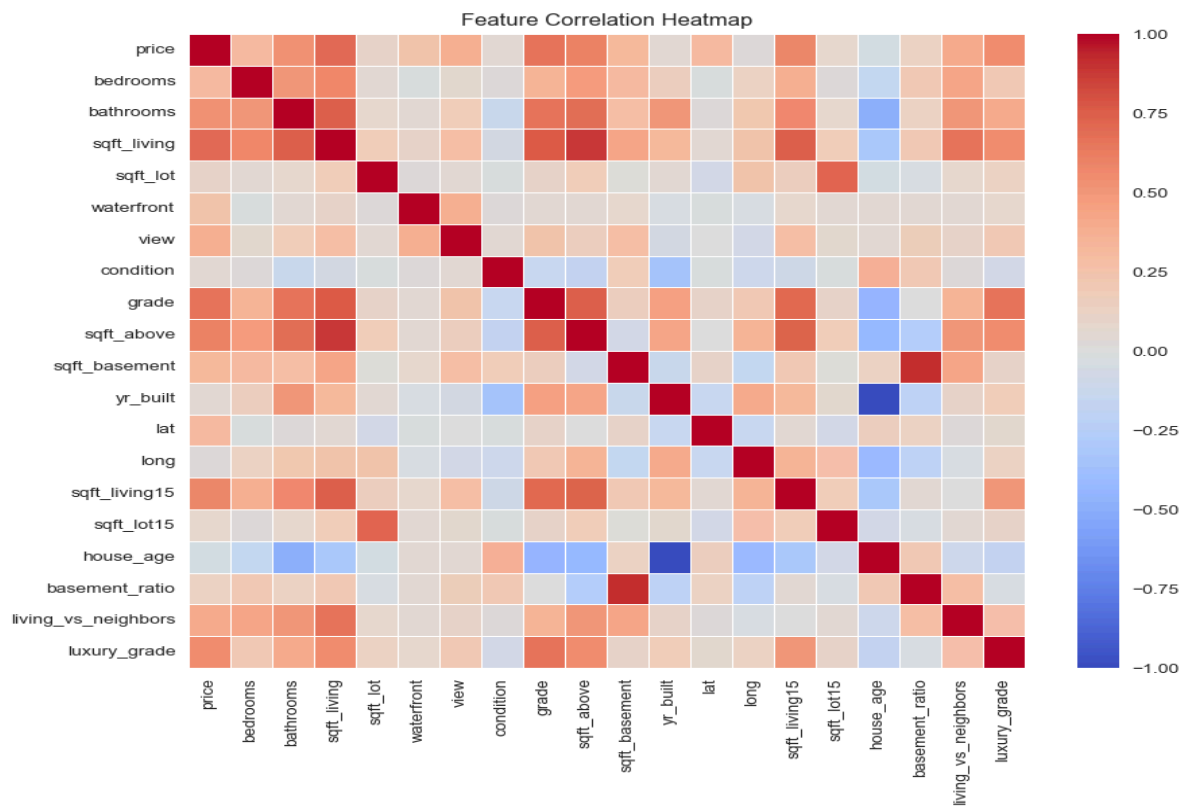
Exploratory data analysis was performed to understand relationships between features and the target variable. Distribution plots, scatter plots, and correlation heatmaps were used to analyze trends and dependencies in the data.



[FIGURE 1: Price distribution]



[FIGURE 2: Living area vs price]

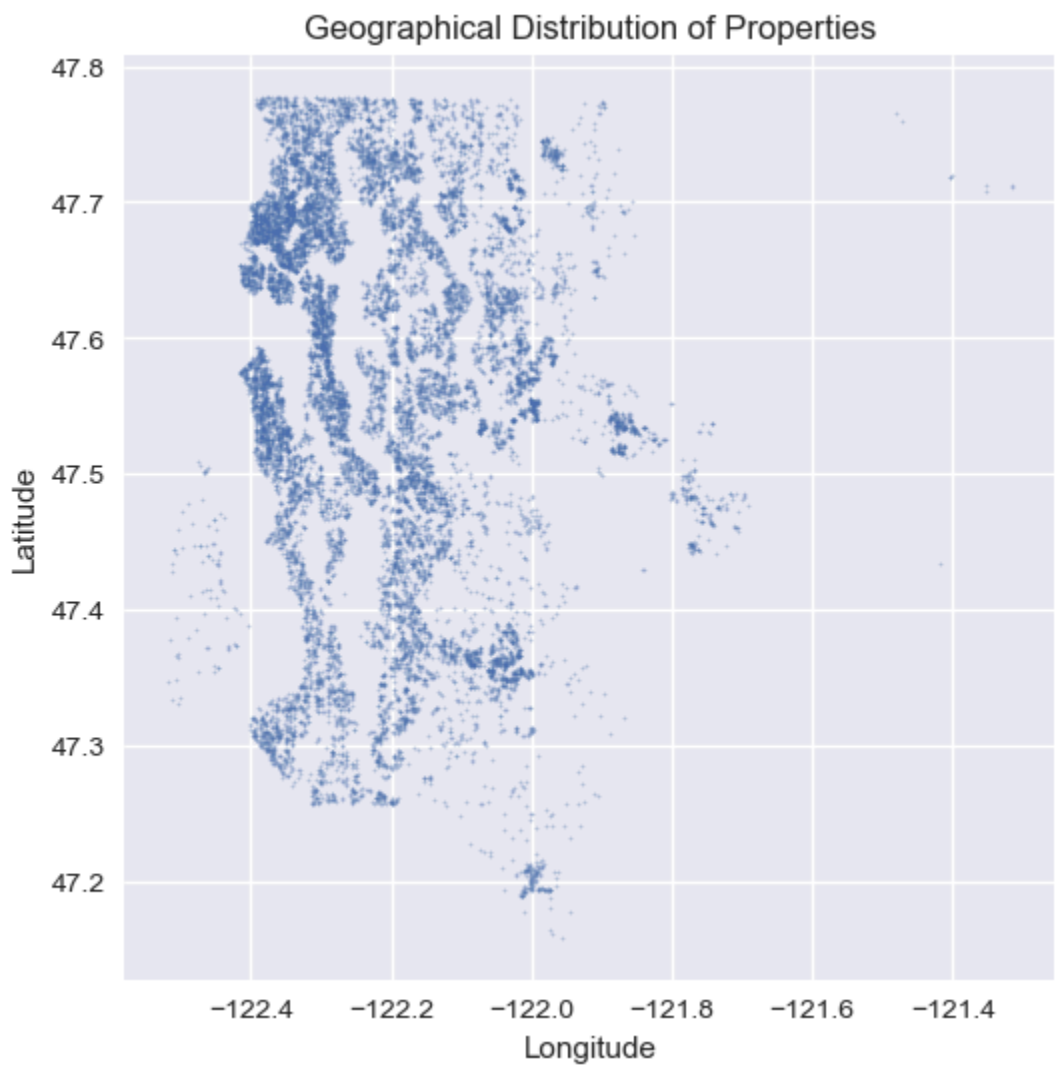


[FIGURE 3: Correlation heatmap]

The analysis showed strong relationships between property price and features such as living area, construction grade, and location-related attributes.

5. Geospatial Analysis

To study spatial patterns, a geospatial visualization was created using latitude and longitude values of the properties. This geospatial visualization of property locations highlights spatial clustering patterns corresponding to urban density. Properties in denser regions tend to exhibit higher prices, emphasizing the importance of spatial and neighborhood context in valuation tasks.



6. Satellite Image Acquisition

Satellite images were programmatically fetched using a static maps API based on the geographical coordinates of each property. A fixed zoom level was used to ensure consistent neighborhood-scale context across all images. Image downloading was fully automated through a Python script to ensure reproducibility.



[SAMPLE PHOTOS]

7. CNN-Based Image Feature Extraction

A pretrained ResNet-18 convolutional neural network was used as a feature extractor. The final classification layer was removed, and each satellite image was converted into a 512-dimensional embedding vector. These embeddings provide a compact numerical representation of neighborhood-level visual context, enabling seamless integration with tabular data.

8. Multimodal Regression Model

8.1 Tabular-Only Baseline

As a baseline, a Random Forest Regressor was trained using only the structured tabular features. This model establishes a reference performance level without any visual information.

8.2 Multimodal Model

To build the multimodal model, image embeddings were concatenated with tabular features. The combined feature vector was then used to train the same regression model, allowing a fair comparison between tabular-only and multimodal approaches.

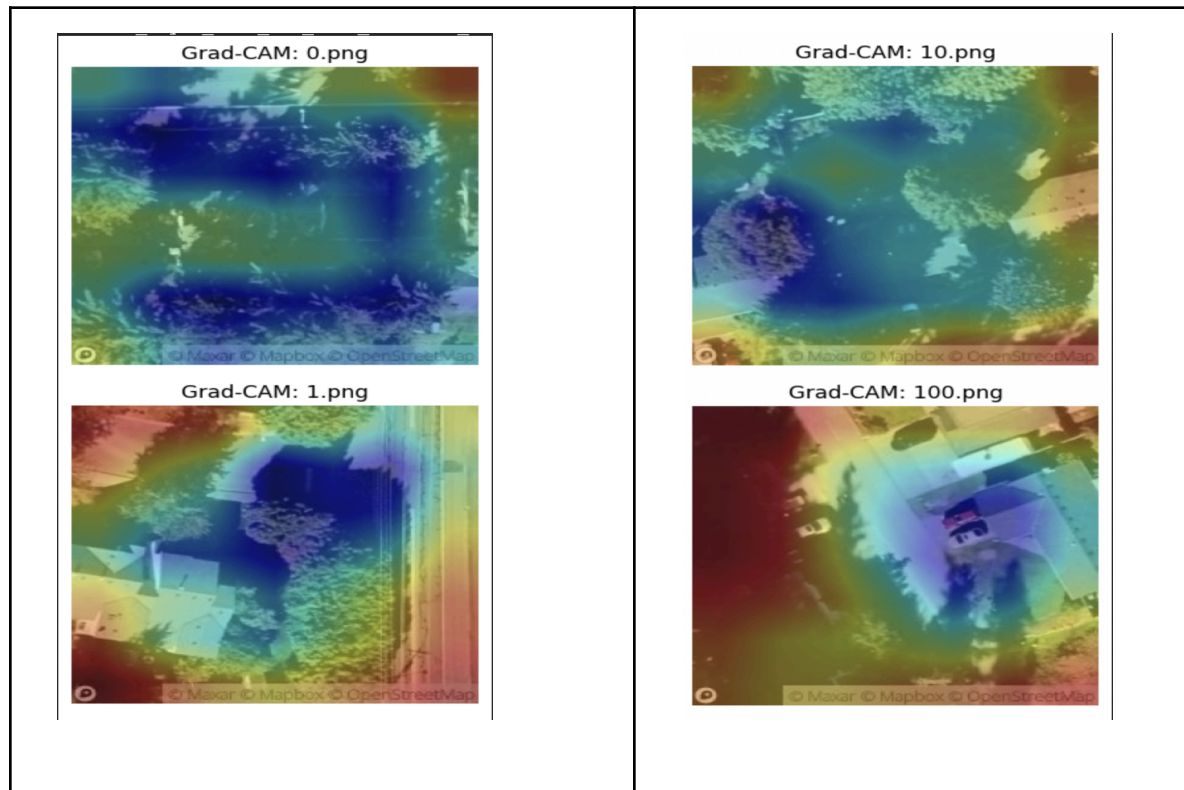
9. Model Performance Comparison

Model performance was evaluated using Root Mean Squared Error (RMSE) and R^2 score on a validation set.

Model	RMSE	R^2
Tabular Only	130984.8027	0.863
Tabular + Satellite Images	148720.2242	0.824

The tabular-only model achieved better quantitative performance than the multimodal model on the validation set. This is likely because structured features such as living area, construction grade, and location already capture strong predictive signals. In contrast, the satellite image embeddings introduce high-dimensional visual features that may contain noise or redundant information not directly aligned with the price prediction task. Additionally, the CNN feature extractor was pretrained on a generic image dataset and not fine-tuned for real estate valuation, which may limit the usefulness of the extracted visual features. These factors together explain the reduced R^2 observed in the multimodal setting.

10. Explainability Using Grad-CAM

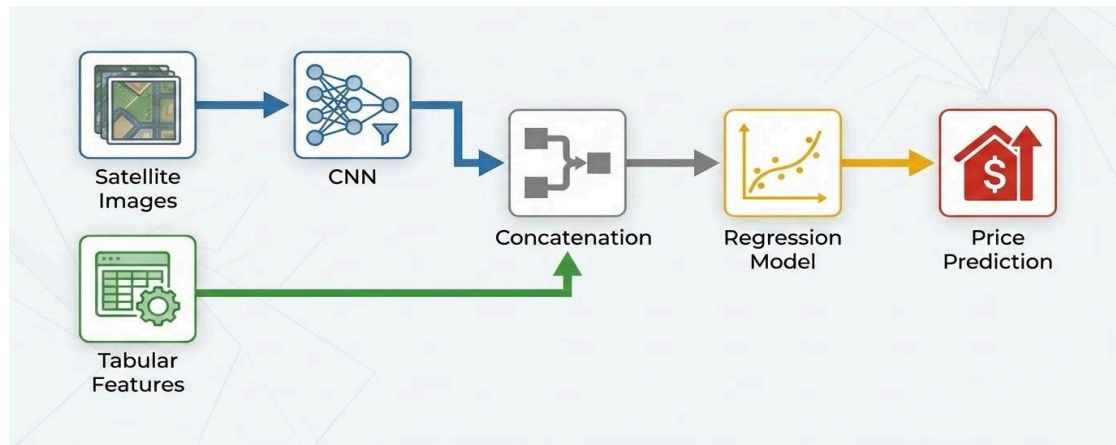


Grad-CAM was applied to the CNN feature extractor to visualize salient regions in satellite imagery. The heatmaps show attention to road networks, built-up areas, and green spaces, indicating that neighborhood context is captured by the visual features.

11. Financial and Visual Insights from Satellite Imagery

Analysis of satellite imagery reveals several visually derived factors with direct financial implications for property valuation. Grad-CAM visualizations indicate that the CNN focuses on road connectivity, urban density, and green cover. Properties near well-connected road networks and dense infrastructure typically command higher prices due to accessibility advantages. Similarly, green spaces and open areas contribute positively to property value by improving environmental quality and livability. In some cases, proximity to water bodies or premium urban clusters is also highlighted, reflecting established real estate pricing trends. These insights demonstrate that satellite imagery captures economically meaningful neighborhood characteristics not fully represented by tabular attributes alone.

12. Architecture Diagram



13. Conclusion

This project investigated a multimodal approach to property valuation by combining structured housing data with satellite imagery. While the tabular-only model achieved stronger validation performance, the multimodal framework demonstrates how neighborhood-level visual context can be integrated into traditional valuation pipelines. The results highlight that incorporating image-based features requires careful feature alignment and task-specific tuning to be effective. Overall, the project provides a practical and balanced view of the potential and challenges of multimodal learning for real estate analytics.

14. Future Work

Future improvements could include the use of higher-resolution satellite imagery, temporal satellite data, or end-to-end deep learning models that jointly learn visual and tabular representations.