# Forecasting Property Valuations in a Mid-Sized U.S. City: A SHAP-Gain Feature Selection and ElasticNet-Ensembled Approach with Optuna-Tuned XGBoost

Saisrijith Reddy Maramreddy

May 2025

## 1 Background on Property Valuation and Machine Learning

In the United States, property taxes are assessed based on the estimated fair market value of real estate, including both land and building components. This value is determined annually through a government-led assessment process, which may be a contested by property owners. Predicting these assessments accurately is a vital part of revenue forecasting for local governments and a meaningful application of machine learning.

### 1.1 Why Prediction Accuracy and Interpretability Matter

Property Assessment Values directly influence individual tax obligations, urban development decisions, and housing affordability analyses. However, public datasets are often noisy and incomplete-featuring missing renovation records, outdated area measurements, or abrupt shifts in land valuation. Therefore, a high-performing ML model must not only minimize prediction error but also provide interpretable insights that help stakeholders understand, trust and audit the predictions-especially when such predictions may inform public policy or fiscal planning.

## 2 Data Description and Objective?

The task is to predict 2019 assessed property values using historical records from 2015 to 2019. Each record includes includes detailed building features, land area and value, protested assessments, and neighborhood-level identifiers. The final 2019 target is defined as:

$$\text{TARGET} = \text{building\_value\_2019} + \text{land\_value\_2019}$$

This is a regression problem with approximately 600,000 training records and 400,000 test records. No external data was permitted.

## 3 Feature Engineering

Feature engineering plays a critical role in unlocking predictive signals from raw data, particularly in structured datasets involving temporal and categorical variables. Our goal was to construct meaningful, leakage-free features that improve model performance while maintaining generalization to unseen data. Before initiating any feature transformations, we performed a strict comparison of columns between the training and test datasets to ensure no information advantage was present. This safeguards against data leakage and ensures that engineered features reflect patterns learnable at inference time.

### 3.1 Logic-Driven Missing Value Handling and Imputation

Before applying general imputation strategies, we incorporatd feature-specific logic to better handle structural missingness and anomalies commonly found in public property records.

**Backfilling Year-Based Columns Across Feature Groups** To handle missing values in temporally structured features (e.g., `building_area`, `quality`, `full_bath`), we designed a consistent backfilling approach that uses older data to impute more recent years. Specifically, columns were ordered from newest to oldest (i.e., 2019, 2018, ..., 2015), and we applied `bfill(axis=1)` across these columns. This setup causes older values (e.g., from 2015 or 2016) to be used to fill in newer year columns (e.g., 2018 or 2019), effectively implementing a forward fill

in temporal logic. This approach assumes that older data reflects the property's original state more accurately, and helps prevent later-year anomalies or missing values from distorting long-term trends.

This logic was applied across multiple feature groups:

- **Residential count features:** `floors`, `full_bath`, `half_bath`, `bedrooms`, `total_rooms`

- **Area-based features:** `building_area`, `land_area`

- **Valuation features:** `building_value`, `land_value`, `assessed`

- **Categorical building attributes:** `foundation_type`, `grade`, `building_condition`, `quality`, `quality_description`, `physical_condition`, `exterior_walls`, `has_cooling`, `has_heat`



Figure 1: Backfiling of Categorical Features

**Zero-Aware Property Filtering**   In cases where `floor_area_total_2019 = 0`, we treated the property as non-residential or commercial and applied domain-specific logic to avoid inappropriate imputations or distortions in the modeling process.

- All related residential building features—such as `full_bath`, `total_rooms`, `garage_area`, and `porch_area`—were set to zero.

- We also zeroed out all `building_area_` variables and the corresponding `building_value_` variables to reflect the absence of a residential structure.

These records were retained in the dataset rather than dropped, as they likely represent a distinct class of properties where valuation is driven primarily by land characteristics. Explicitly identifying and treating these properties allowed the model to better separate residential and non-residential valuation patterns.



Figure 2: Zero-Aware Property Filtering

**Feature Drop Based on Sparse Signals**   For features that appeared largely irrelevant or unused across the dataset, we calculated the percentage of zero values in columns such as `mobile_home_area`, `deck_area`, and `porch_area`. If a column contained over 90% zeros, it was considered non-informative and dropped from the modeling pipeline to reduce dimensionality and noise.

**Multi-Level Median and Mode Imputation**  After applying logic-based pruning, we used a three-level median imputation strategy for continuous features (e.g., `assessed_value_2017`, `building_area_2017`) based on the following hierarchy:

- **Level 1:** neighborhood-level median

- **Level 2:** region-level median

- **Level 3:** global median (fallback)

For categorical variables such as `foundation_type` or `building_condition`, we applied single-level mode imputation using the most frequent category within the training data. While this approach is less localized, it provided a simple and stable method for handling missing values in features with low cardinality.



Figure 3: Multi-Level Median Imputation

## 3.2   Neighborhood and Region-Level Statistical Features

To capture localized pricing dynamics and identify anomalies in property assessments, we engineered a suite of statistical features using the 2018 assessed values as a proxy for prior valuation context. We first computed neighborhood-level metrics including the mean, median, standard deviation, and interquartile range (IQR) of `assessed_2018`, grouped by `neighborhood`. Similarly, region-level statistics were computed using the `region` variable.

Using the merged and imputed stats, we computed derived features such as:

- `assess_minus_neigh_mean`: the raw deviation of a property's 2018 assessed value from its neighborhood mean

- `assess_ratio_neigh_mean`: a normalized ratio of a property's value to its local average

- `z_score_assess_neigh`: a z-score based on neighborhood-level variation

- Corresponding region-level counterparts: `assess_minus_region_mean`, `assess_ratio_region_mean`, and `z_score_assess_region`

These features helped contextualize each property's assessed value relative to other properties within the same neighborhood or region. They proved useful in capturing outliers and potentially undervalued homes that deviated from local valuation patterns.

## 3.3   Frequency Encoding of High-Cardinality Geographic Variables

To convert high-cardinality categorical variables into numerical features while preserving signal strength, we applied frequency encoding to four key geographic identifiers: `neighborhood`, `region`, `zone`, and `subneighborhood`. This encoding strategy served two purposes. First, it allowed the model to retain information about how common or rare a spatial unit was. A frequently occurring neighborhood (i.e., one with high frequency) likely has more properties, which implies greater residential or commercial development in that area. Second, areas with more properties are also likely to have more consistent and well-understood assessment patterns—the government has "seen" more properties there, which may reduce valuation volatility. These areas are more visible or prioritized in municipal processes. Finally, frequency encoding avoids the dimensional explosion caused by one-hot encoding, which is especially problematic for variables with high cardinality like `neighborhood` or `subneighborhood`.

Figure 4: Frequency Encoding

## 3.4 Boolean and Ordinal Encoding

**Boolean Encoding:** We focused on three boolean variables: `has_cooling`, `has_heat`, and `protested`. These features were encoded across all five years using binary values (0/1).

**Ordinal Encoding:** For ordinal features such as `quality`, `quality_description`, `grade`, `building_condition`, and `physical_condition`, we performed domain-informed cleaning and then applied ordinal encoding based on defined category hierarchies. Prior to encoding, raw values were standardized through column-specific replacements. For instance, extreme or ambiguous values like `X`, `None`, or overly granular subgrades (e.g., `X-`, `E+`) were either mapped to more interpretable categories or treated as missing. Some detailed conditions like `Unsound` and `Very Poor` were collapsed into broader categories such as `Poor`. Unknown values were handled gracefully by assigning an encoded fallback of `-1`. This process ensured that ordinal information was preserved in a numerically meaningful way, allowing models to leverage the ordered nature of these features without exploding dimensionality as one-hot encoding would.



Figure 5: Ordinal Encoding

## 3.5 Target Encoding of Nominal Categorical Variables (2015–2019)

For certain nominal features that lacked ordinal structure but exhibited high cardinality, such as `foundation_type` and `exterior_walls`, we applied target encoding across all years from 2015 to 2019. This encoding replaces each category with a smoothed version of the mean target value (`assessed_2018`) observed for that category.

To avoid overfitting and data leakage, we implemented a 5-fold cross-validated target encoding procedure. For each fold, the mean target value was computed from the training portion and mapped to the validation fold. We used a smoothing parameter of 10 to balance the influence of the global mean versus the category-specific mean, especially for infrequent categories.

This method enabled us to capture predictive signal from nominal features without creating high-dimensional one-hot encodings or imposing artificial ordinal structure.



Figure 6: Target Encoding

## 3.6 Quantile Binning of Features

To enhance robustness and reduce sensitivity to outliers, we converted few continuous features into categorical bins using quantile-based binning. Growth metrics such as `land_value_growth`, `building_value_growth`, and `assessed_growth` were binned into four quantiles, with thresholds computed only on the training data to prevent information leakage. If quantile binning failed due to low cardinality (e.g., repeated values), we defaulted to equal-width binning. All binned variables were explicitly cast as categorical to ensure compatibility with tree-based models.

Additionally, we binned `year_built_final` into five quantiles to capture generational differences in construction periods. This replaced raw year values with interpretable ordinal categories. Original continuous features were removed after binning to avoid redundancy and reduce multicollinearity.



Figure 7: Quantile Binning of Growth Features

## 3.7 Rare Frequency Suppression in Spatial Encodings

Following frequency encoding of high-cardinality spatial variables (`region`, `neighborhood`, `zone`, `subneighborhood`), we applied a rare-value suppression step to mitigate the noise introduced by sparsely represented categories. For each frequency-encoded column, we identified values that occurred in less than 0.1% of the training data and replaced them with a neutral value of zero in both the training and test sets. This was done using thresholds derived solely from the training distribution to prevent data leakage.

The intuition behind this strategy is that extremely rare spatial groupings may not provide reliable or generalizable signals to the model. Treating them as a common fallback class (i.e., assigning them a frequency of zero) improves model stability and reduces overfitting to idiosyncratic, low-support locations. This transformation preserves the informativeness of frequent categories while smoothing out sparse tail behavior in the feature space.

## 3.8 Log Transformation and Distribution Smoothing for Ridge Regression

To satisfy linear model assumptions and reduce skew-related distortion in Ridge regression, we applied a targeted log transformation to select continuous features. Specifically, we identified variables related to building size, land area, and valuation (e.g., `building_value_2019`, `land_area_2018`, `neigh_assess_mean`) whose skewness exceeded a threshold of 2.0 in the training set. For these features, we applied a `log1p` transformation, which effectively stabilized variance, compressed long-tailed distributions, and improved linear fit potential.

This transformation was particularly useful for the Ridge regression model, which benefits from normally distributed inputs and is sensitive to extreme values. By selectively applying `log1p` only to features with high skew, we preserved model interpretability while enhancing numerical stability and predictive performance.



Figure 8: Log Transformation for Ridge Regression

## 3.9 Adaptive Quantile Clipping for Tree-Based Models

To further control the influence of extreme values in tree-based models, we implemented an adaptive quantile clipping strategy informed by skewness severity. Using a precomputed skewness report, we categorized numeric features (excluding binary and target-encoded variables) into two groups: *ultra-skewed* (skewness > 100) and *moderately-skewed* ($2 <$ skewness $\leq 100$). Features were considered only if they had more than ten unique values and were not binary.

For ultra-skewed features, we applied clipping at the 0.5th and 99.5th percentiles. For moderately skewed features, we clipped at the 0.1st and 99.9th percentiles. All thresholds were derived solely from the training data and applied to both training and test sets to ensure leakage-free transformations. This clipping procedure helped suppress extreme values that might otherwise dominate decision paths or split criteria in tree-based learners.

These transformations were specifically designed for use with XGBoost and LightGBM, where reducing the influence of outliers improves model generalization and enhances interpretability in leaf-based decision structures.



Figure 9: Adaptive Quantile Clipping for Tree-Based Models

## 3.10 Interaction Features for Linear and Nonlinear Models

To enrich model expressiveness, we engineered a comprehensive set of interaction features used across both Ridge regression and tree-based models (XGBoost and LightGBM). These included multiplicative and ratio-based terms such as `grade_quality_index`, `value_per_age`, `area_x_quality`, and `assess_to_neigh_mean`, capturing relationships between physical dimensions, valuation, quality, and neighborhood context. For Ridge regression, these features acted as implicit basis expansions—effectively enabling the linear model to capture non-additive effects by introducing new combinations of input features.

Additionally, we created a specialized set of log-transformed interaction terms—such as `log_area_x_grade`, `log_assess_x_age`, and `log_value_diff`—used exclusively in the Ridge pipeline. These features helped linearize multiplicative relationships and reduce skew, improving fit under Ridge's sensitivity to input distribution. Log-based interactions were excluded from tree models, which are inherently robust to skew and insensitive to monotonic transformations like log, as they rely only on the relative ordering of feature values when making splits.

# 4 Model Development and Tuning

## 4.1 Ridge Regression with Cross-Validation

We implemented a Ridge regression model using `RidgeCV` to automatically select the regularization strength $\alpha$ through nested cross-validation. A 3-fold outer loop was used for estimating out-of-fold (OOF) performance, while each inner fold evaluated a grid of $\alpha$ values ranging from $10^{-3}$ to $10^{2}$ on a logarithmic scale. Input features were standardized within a `Pipeline` using `StandardScaler` to ensure scale-invariant regression coefficients.

The model selected a different optimal $\alpha$ for each fold, reflecting local variance in validation behavior:

- Fold 1 RMSE: 42,050.33—Best $\alpha$: 2.1544

- Fold 2 RMSE: 41,036.52—Best $\alpha$: 27.8256

- Fold 3 RMSE: 40,619.40—Best $\alpha$: 0.5995

The final out-of-fold RMSE across all folds was **41,239.79**, with an average best $\alpha$ of approximately **10.1932**. This indicates that moderate regularization consistently improved generalization across different training splits. We

saved both OOF predictions and test forecasts as NumPy arrays—`ridgecv_oof_preds.npy` and `ridgecv_test_preds.npy`—for later use in model ensembling. These stored outputs served as reliable building blocks for downstream blending and stacking strategies.

## 4.2 Tree-Based Models with Optuna and SHAP-Gain Feature Selection

To capture nonlinear interactions and leverage automatic handling of missing values and categorical splits, we trained two gradient boosting models: LightGBM and XGBoost. Both models followed a structured pipeline consisting of hyperparameter optimization using Optuna, followed by SHAP- and gain-based feature selection, and a final retraining on the selected features.

**Step 1: Hyperparameter Tuning with Optuna.** For each model, we defined an Optuna objective that trained 3-fold cross-validated models using early stopping. We explored hyperparameter ranges tailored to each algorithm, with LightGBM using a native pruning callback and XGBoost leveraging `XGBoostPruningCallback`. During tuning, we stored the best out-of-fold (OOF) predictions across all trials to later use in ensembling.



Figure 10: XGBOOST Hyperparameter Tuning with Optuna

**Step 2: SHAP and Gain-Based Feature Selection.** After tuning, we trained new LightGBM and XGBoost models using the best parameters on each fold of the training data. For each fold, we computed SHAP importance values and LightGBM/XGBoost gain importances. We retained features that collectively accounted for 95% of total importance in either SHAP or gain, and constructed a union of these high-signal features across all folds. This union was used to define the final reduced feature space for retraining.



Figure 11: SHAP and Gain-Based Feature Selection

**Step 3: Final Model Training and Inference.** Each final model was retrained on the full training set using only the selected features, with early stopping enabled to prevent overfitting. Predictions on the test set were

generated using the best iteration count. All model outputs—including OOF predictions and test forecasts—were saved for submission and later use in ensemble blending.

This hybrid approach—combining Optuna-based tuning with SHAP-driven interpretability—allowed us to retain only high-impact features, thereby improving generalization and reducing overfitting without sacrificing performance. The best out-of-fold RMSE achieved was **40,925.29** with XGBoost and **41,641.42** with LightGBM, confirming the robustness of both pipelines.

# 5 Ensembling Strategy

## 5.1 Weighted Model Blending with Optuna

To consolidate the strengths of our top-performing base models—XGBoost, RidgeCV, and LightGBM—we employed a weighted blending strategy optimized using `Optuna`. This approach directly searched for the optimal linear combination of model predictions that minimized RMSE on a holdout set.

We first constructed a meta-training set consisting of out-of-fold (OOF) predictions from each base model. A corresponding test matrix was constructed from each model's final test predictions. The blending weights were constrained to be non-negative and normalized to sum to one.

An `Optuna` study was run for 100 trials, where each trial proposed a new set of blending weights and evaluated their performance via RMSE on the holdout split. The final optimized weights were:

- XGBoost: $w_0 = 25.98\%$

- RidgeCV: $w_1 = 33.53\%$

- LightGBM: $w_2 = 40.49\%$

These weights were then used to produce a final blended prediction for the test set. The resulting predictions achieved an RMSE of **36,239.91** on the holdout set—outperforming all individual base models and demonstrating the value of combining linear and tree-based perspectives. This ensembling strategy offered a robust and interpretable improvement by leveraging complementary model strengths, particularly balancing Ridge's stability with the non-linear flexibility of the tree-based models.

## 5.2 Stacked Ensembling with ElasticNetCV

To complement our Optuna-based weighted average ensemble, we implemented a stacked generalization approach using `ElasticNetCV` as a meta-learner. This method treats out-of-fold (OOF) predictions from the base models—XGBoost, RidgeCV, and LightGBM—as features in a second-level regression model. By learning how to optimally combine base predictions, the meta-model can capture nonlinear inter-model relationships while applying regularization to prevent overfitting.

**Meta-Model Training.** We concatenated the OOF predictions into a 3-column meta-feature matrix and used it to fit an `ElasticNetCV` model wrapped in a `StandardScaler` pipeline. The meta-model searched over a grid of $\ell_1$ ratios and $\alpha$ values, using 3-fold cross-validation to identify the optimal regularization configuration.

**Holdout Evaluation.** For evaluation, we trained the meta-learner on an 80% split and evaluated on a 20% holdout. The resulting RMSE was **36,344.64**, closely aligned with the Optuna-weighted blend. The selected $\alpha$ was **0.01**, indicating strong regularization and robust coefficient shrinkage.

**Final Test Predictions.** The final model was retrained on the full meta-feature set and used to predict the test set.This stacked approach provided a robust and regularized alternative to linear averaging, automatically downweighting weaker models while maintaining interpretability and reproducibility. On the competition leaderboard, the `ElasticNetCV` ensemble—combining Ridge, XGBoost, and LightGBM predictions—secured the **top rank** with a **private leaderboard RMSE of 36,021**, just 2 points above the public leaderboard score of 36,019. Interestingly, while an Optuna-weighted linear blend achieved a lower public RMSE, it ranked below the ElasticNet ensemble on the final evaluation, underscoring the latter's generalization strength.

Figure 12: ElasticNetCV ensemble

# 6 Residual Error Analysis

To evaluate model robustness, we conducted residual analysis across key dimensions. Three consistent patterns emerged:

**1. Underprediction for High-Value Properties** Residuals increased with actual property value, particularly above \$5M. This indicates systematic underestimation of high-end properties, likely due to their rarity and unique characteristics.

**2. Volatility in Protested Properties** Properties with multiple protests between 2015–2018 exhibited larger residual variance, despite median residuals near zero. This suggests that frequently disputed properties are harder to predict, potentially due to unrecorded structural changes or historical valuation disagreements not captured in the dataset.

**3. Poor Generalization in Sparse Neighborhoods** Residuals were more variable and extreme in neighborhoods with low frequency in the training data. This suggests weaker model generalization in underrepresented regions. Frequency-aware encodings or collapsing sparse neighborhoods into broader groups may enhance stability and reduce prediction noise in these areas.



(a) Residuals vs Actual Value     (b) Residuals vs Protest Count     (c) Residuals vs Neighborhood Frequency

Figure 13: Residual Analysis Across Value, Protest Count, and Neighborhood Frequency

# 7 SHAP-Based Interpretability and Insights

**Top Predictive Features** SHAP analysis revealed that features like `assessed_2018`, `building_value_2018`, and `land_value_2018` were primary drivers of the model's predictions. Structural attributes such as `building_area_2019`, `floor_area_x_grade`, and `grade_2019` also carried strong explanatory power. These results confirmed the value of engineering ratio and interaction terms that encode economic density, build quality, and age-adjusted valuation. Neighborhood-level variables, especially `neigh_assess_std` and frequency encodings, further demonstrated the importance of local context in real estate assessment.

**Low-Impact Features** Although SHAP-ranked bottom 30 features showed limited average contribution, they were retained via the 95% SHAP + Gain union due to their potential complementary value. Examples include early-year indicators like `quality_description_2015`, `fireplaces_2016`, and spatial ratios like `porch_ratio`.

Their inclusion likely enhanced generalization by supporting edge cases, and their low impact helped confirm that more aggressive feature pruning would have offered little gain.
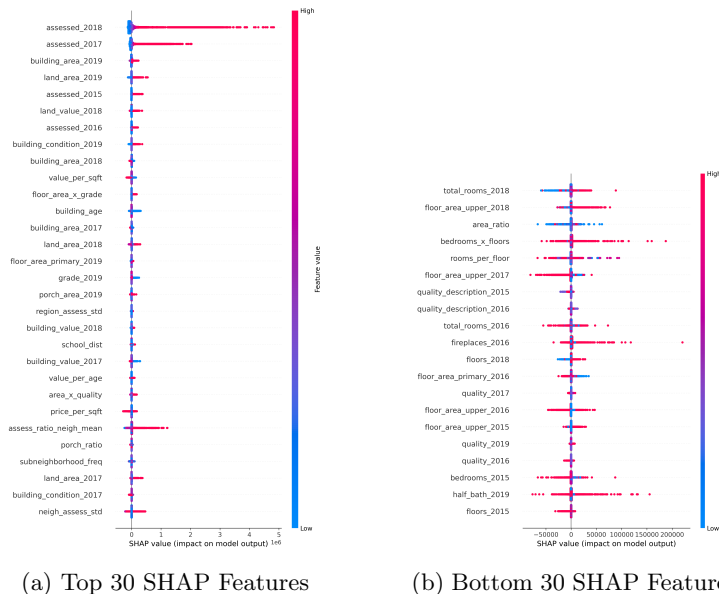


(a) Top 30 SHAP Features　　　　(b) Bottom 30 SHAP Features

Figure 14: SHAP Feature Importance: Top and Bottom Ranked Variables

# 8　Reflections and Lessons Learned

Our modeling pipeline blended SHAP-driven interpretability with ensemble-based prediction to achieve both transparency and predictive strength. Key lessons emerged across three dimensions:

- **What Worked Well:**
  - The SHAP + Gain union was highly effective in denoising the feature space and avoiding overfitting, retaining only high-impact predictors across folds.
  - The ElasticNet ensemble combined linear and non-linear model strengths to capture nuanced patterns in both well-represented and sparse regions.
  - Optuna tuning reduced manual trial-and-error and consistently improved generalization across Ridge, LGBM, and XGBoost pipelines.

- **What Could Improve:**
  - We did not model temporal dependencies across yearly features—treating them as time-series sequences or using transformer architectures may better account for evolving valuation dynamics over time.
  - Geographical variation was only partially captured using frequency encoding. Future work could explore residual stacking or blended stacking techniques with subgroup-aware features (e.g., neighborhood frequency, protest history) to correct localized error patterns more explicitly.

Overall, the final ElasticNet ensemble—integrating Ridge, SHAP-informed LGBM, and Optuna-tuned XG-Boost—delivered a strong performance on the private leaderboard. Future extensions may benefit from incorporating time-aware modeling, subgroup-specific residual correction, or spatially informed representations to further improve accuracy and fairness in municipal assessment systems.

# References

[1] Weylandt, M. (2025). *STA 9890: Housing Price Prediction Competition Guidelines*. Retrieved from `https://michael-weylandt.com/STA9890/competition.html#final-report-50-points`

# Appendix: Jupyter Notebooks

A ZIP archive containing all code notebooks is included with this submission:

- `Ridge(log transformed features).ipynb`

- `XGBOOST & LGBM(with clipping).ipynb`

- `RMSE OPTIMIZATION.ipynb`

To reproduce results or explore the full training pipelines, extract and run these notebooks in a Python 3.9+ environment with the listed dependencies.