# House Value Prediction for Strategic Real Estate Decision-Making

•••

**Team G1**
Anna Dominic
Karthikeyan Shanmugam
Manisha Goyal
Siri Desiraju

# Business Problem

**Objective:**

Develop a model to predict house prices in Ames, Iowa.

**Impact:**

- Enables data-driven real estate decisions in investments, sales, and development

- Enhances understanding of factors influencing property values

# Target Stakeholders

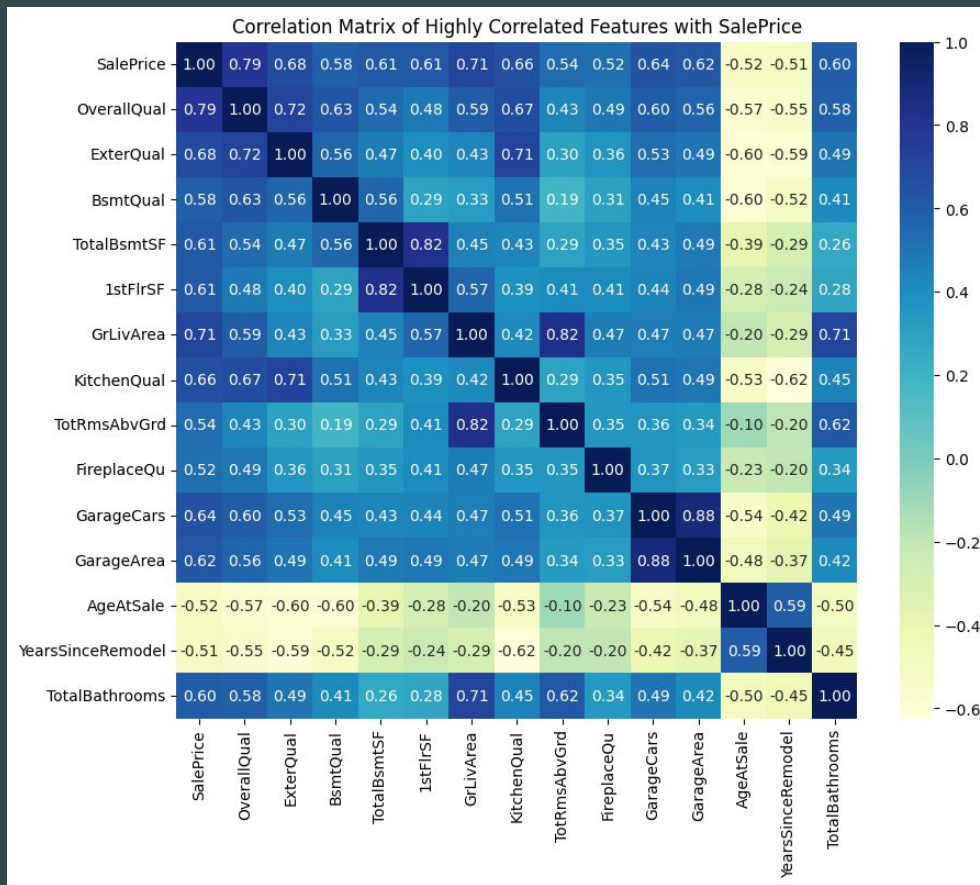| Who am I? | What is my goal? | How can the model help me? |
|---|---|---|
| Real Estate Agent | Set competitive prices to attract buyers & satisfy sellers | Use model predictions to enhance service quality and optimize listings |
| Real Estate Developer | Plan lucrative developments that captivate buyers | Use model insights to focus on profitable features and projects |
| Home Owner | Maximize home value with strategic remodeling | Use model for home value prediction and insights on remodel impact |
| Home Buyer | Find a home that meets personal needs at a fair price | Use model predictions to determine fair pricing and negotiate effectively |

# Understanding the Data

**Dataset:** House Prices of Ames, Iowa
(from Kaggle)

| >1400 | ~4 |
|---|---|
| properties | years |

| ~35 | ~46 |
|---|---|
| numeric | categorical |



Correlation Matrix of Highly Correlated Features with SalePrice

# Hypothesis Testing for Insights

Recently renovated properties sell at higher prices, controlling for the age of the property

OLS Regression Analysis

Market preferences reflect a higher valuation for certain dwelling types and styles that align with buyer demands

ANOVA & Tukey's HSD Post-hoc Analysis

# Data Processing

| Missing Values | Encoding Categorical | Feature Engineering |
|---|---|---|
| Imputed moderate to low missing values and preserved meaningful 'NA's | Ordinal features mapped to numeric scales | Developed new features like 'AgeAtSale' |
| High missing values evaluated for retention or exclusion | Nominal variables transformed via one-hot encoding | Normalized feature values using Min-Max scaling |

After data preparation, there were **213 features**

# Selecting Features for the Model

| Feature Selection | Number of Features | R² |
|---|---|---|
| Lasso Regression | 88 | 0.873 |
| Recursive Feature Elimination | 5 | 0.664 |
| Principal Component Analysis | 107 | 0.872 |

Feature selection methods evaluated on baseline linear regression model

Lasso Regression has the highest R² value and is more interpretable than PCA

# Model Development Strategy

Employed both interpretable (e.g., Random Forest) and complex (e.g., XGBoost) models to balance understandability with predictive power

Conducted comprehensive feature analysis to determine model efficacy using both full and reduced feature sets
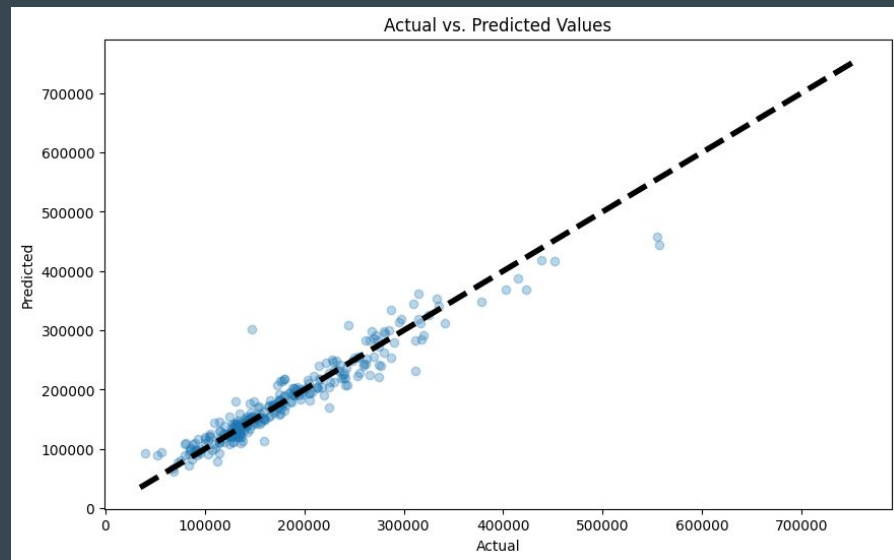
Utilized Grid Search with Cross-Validation to fine-tune model parameters, ensuring optimal performance
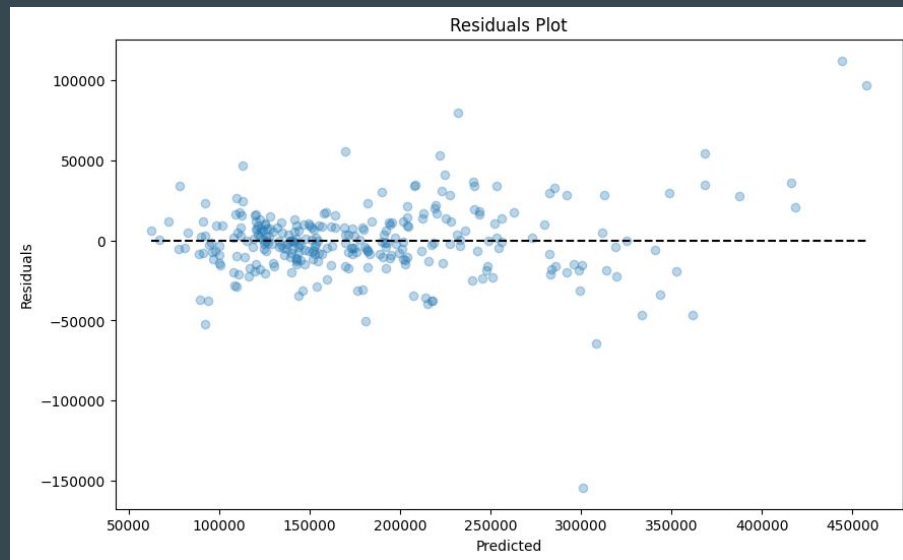
# Model Results

| Model | Feature Selection | R² | RMSE |
|---|---|---|---|
| Lasso Regression | Lasso | 0.872 | 27854.35 |
| Random Forest | - | 0.891 | 25769.59 |
| XGBoost | - | 0.918 | 22326.72 |
| SVM | - | 0.550 | 52399.20 |
| Neural Network | Lasso | 0.871 | 28105.04 |
| CatBoost | - | 0.826 | 25937.00 |
| Ensemble Models (Stacking Method) | - | 0.917 | 22446.62 |

# Evaluating the Model



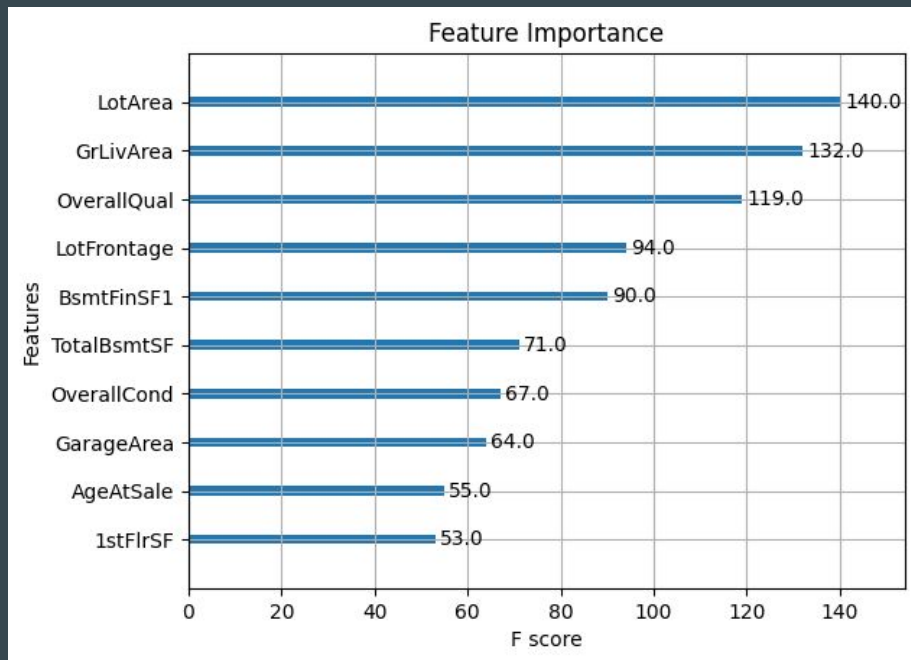Actual vs. Predicted Values

Residuals Plot
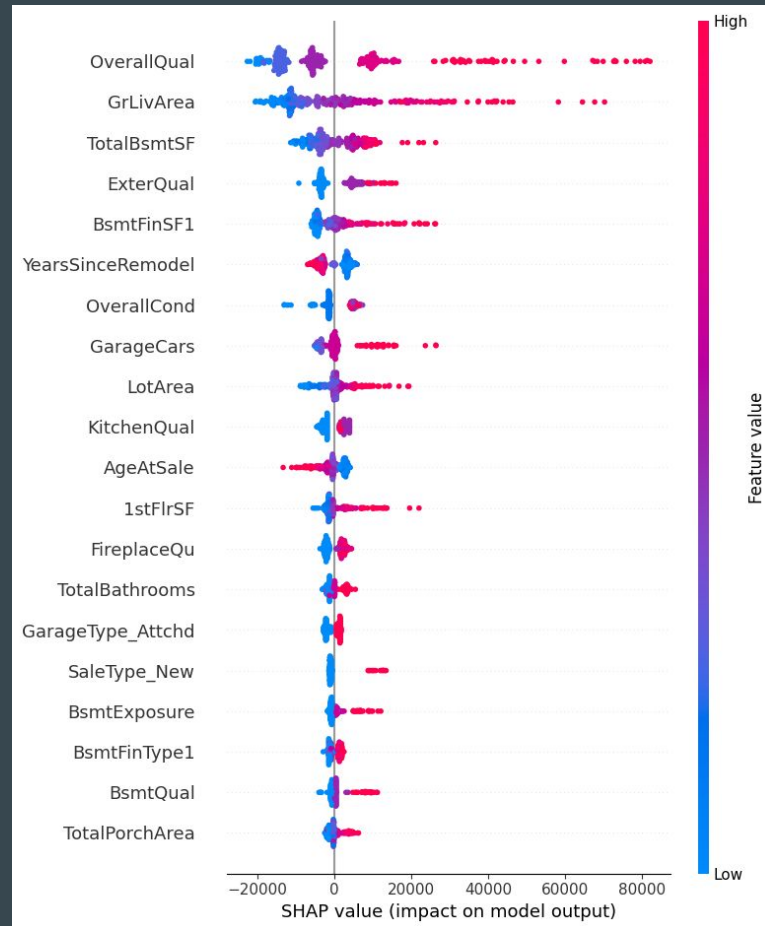
**Actual vs Predicted Values**

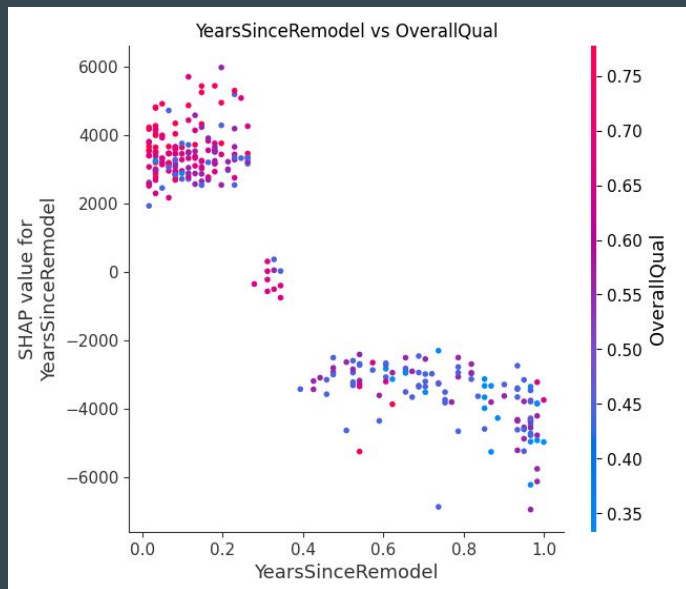**Residuals Plot**

# Identifying Key Features
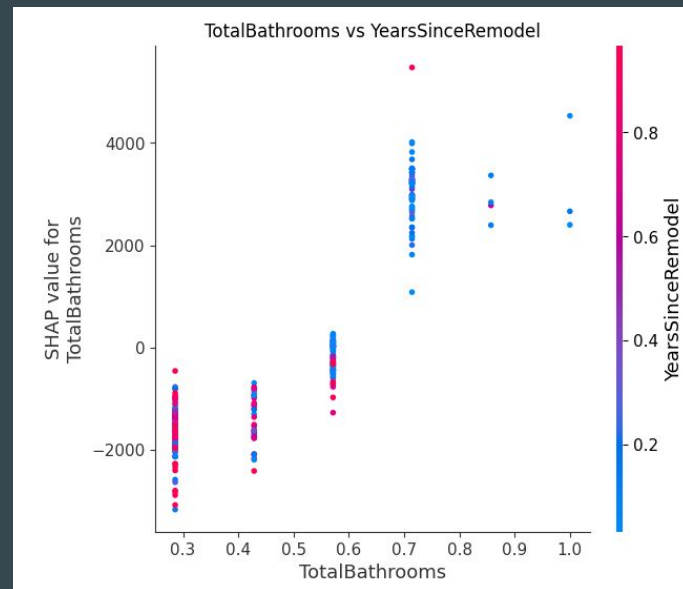


XGBoost Feature Importance

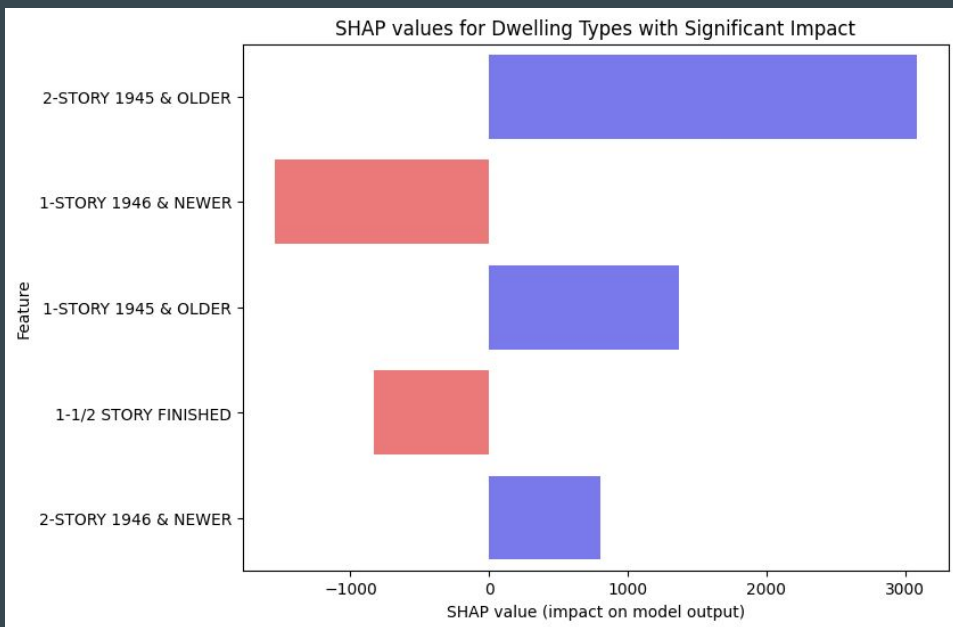SHAP Explainer

# Identifying Key Features



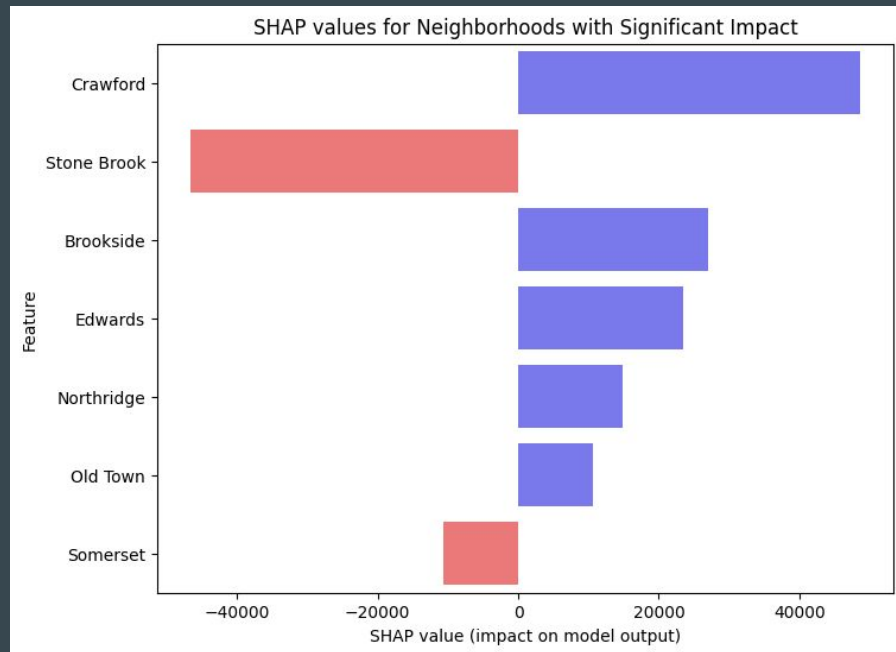Factors like recent renovations drive up 'Overall Quality' of the house

Recent renovations have significant impact by improving 'Number of Bathrooms'

# Identifying Key Features



SHAP Explainer - Dwelling Types

SHAP Explainer - Neighborhoods

# Recommendations to Stakeholders

| Real Estate Agent |
|---|
| - Evaluate house prices with given features<br>- Guide premium buyers to Crawford neighborhood and budget buyers to Stone Brook neighborhood |

| Real Estate Developer |
|---|
| - Design dwelling types like '2-Story'<br>- Prioritize living space quality and garage size |

# Recommendations to Stakeholders

| Home Owner |
| --- |
| - Assess house value<br>- Evaluate targeted renovations to increase value:<br>    - Repaint and remodel interior and exterior finish<br>    - Ensure basement is finished<br>    - Remodel kitchen |

| Home Buyer |
| --- |
| - Assess house value based on features and location<br>- Assess possible features and locations given a budget |

# Extending Scope to other Cities

**Motivation**: Increase business value and broaden applicability

**Methodology**: Leverage the best-performing model (XGBoost) to make predictions across diverse markets

| City/State | R² | RMSE | Important Features |
|---|---|---|---|
| New York | 0.589 | 3218761.26 | Property area, Zip code, Number of bathrooms |
| California | 0.707 | 61945.62 | Median household income, Distance to the ocean, Number of bedrooms |
| Seattle | 0.577 | 344297.69 | Apartment size, Zip code, Number of bathrooms |

# Future Work

**Enhance Model Relevance and Accuracy**

- **Cross-City Validation:** Explore the potential for better results with more comprehensive datasets

- **Local Economic Indicators:** Integrate metrics like employment rates and income levels to predict market shifts

- **Consumer Behavior:** Use surveys to capture buyer preferences and predict desirability trends

- **Real-Time Market Data:** Enhance predictions with live data on listings and sales