



HOUSE PRICE PREDICTION

Using Machine learning

PROJECT GUIDE: Prof. Suvarna Ranade

PROJECT BY–

- PRASANNA KAILAS SHINDE (ROLL NO 36)
- MRUNALI MANOHAR PETAKAR (ROLL NO 30)

Index :

03	GOAL OF PROJECT
04	DATA OVERVIEW
08	DATA CLEANING
10	EDA
13	PREPROCESSING & FEATURE SELECTION
14	MODEL BUILDING AND EVALUATION
18	PREDICTION AND FUTURE STEPS

goal / aim :

The goal of the project is to develop a robust and accurate machine learning model that can predict house prices based on various features related to real estate properties in Ames city. By leveraging advanced analytics and predictive modeling, our objective is to create a model that outperforms traditional methods, providing accurate and reliable predictions of house sale prices.



data overview :

- The dataset comprises real estate data from Ames city, including a comprehensive set of features related to residential properties.

- **Key Features:**

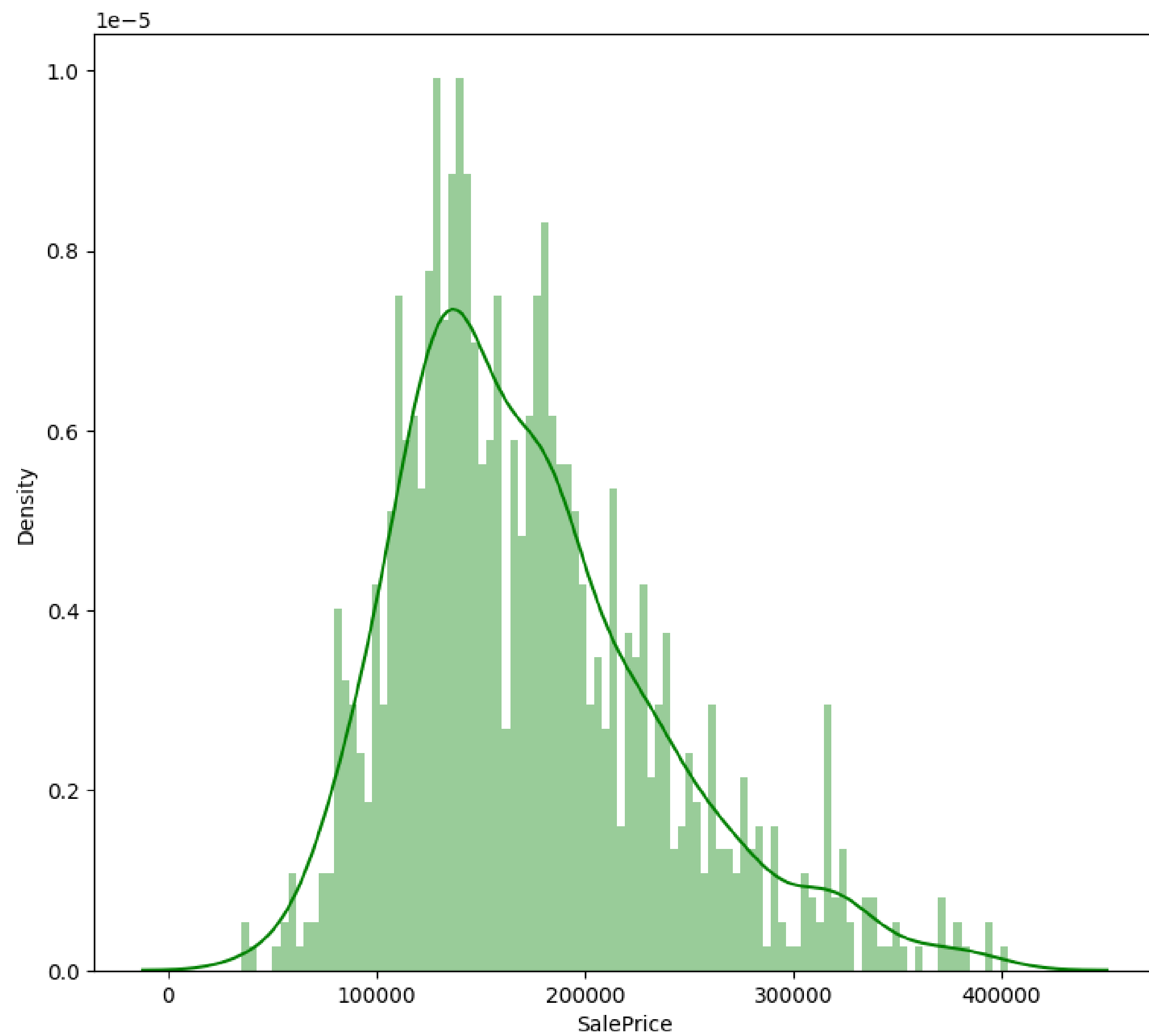
Property Details:

1. SalePrice: This is the target variable that the model aims to predict.
 2. LotArea: Indicates the size of the lot on which the property is situated.
 3. Utilities: Describes the types of utilities accessible at the property.
- and so on...

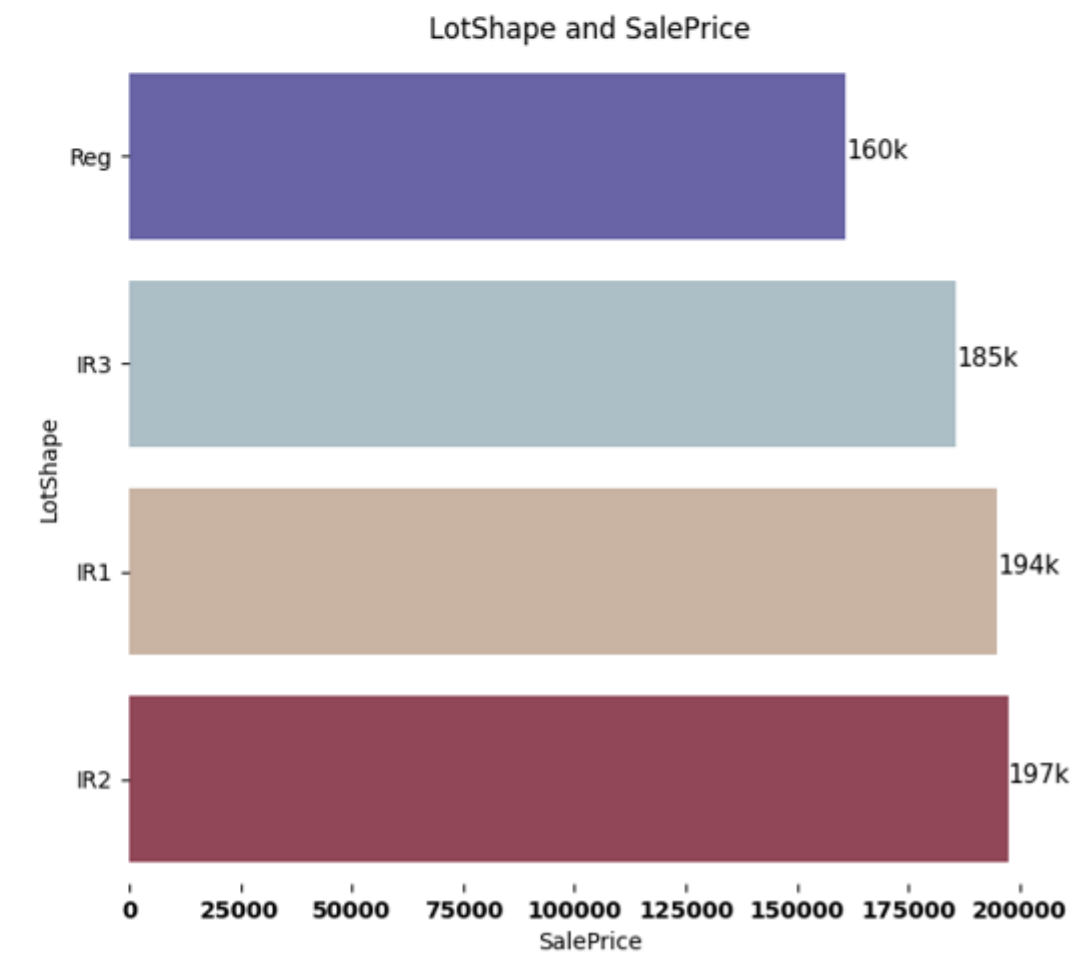
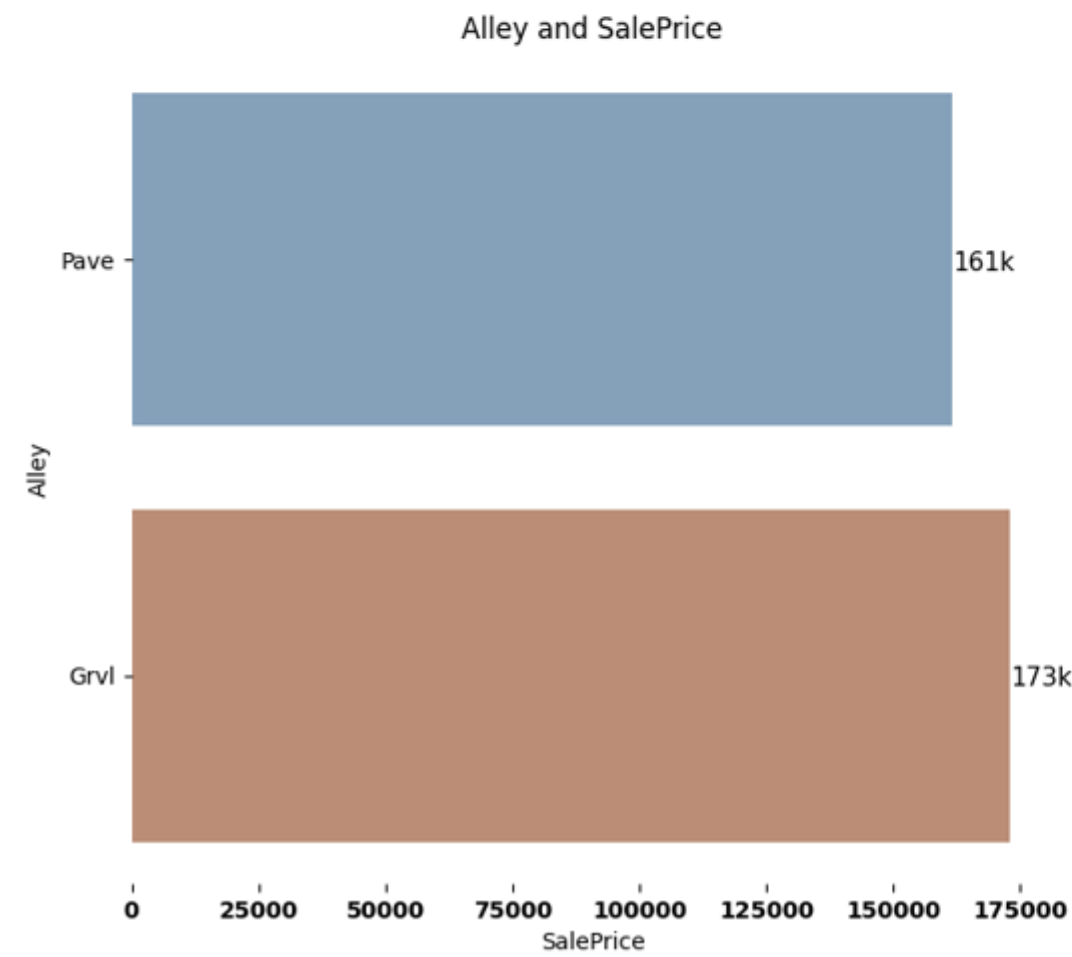
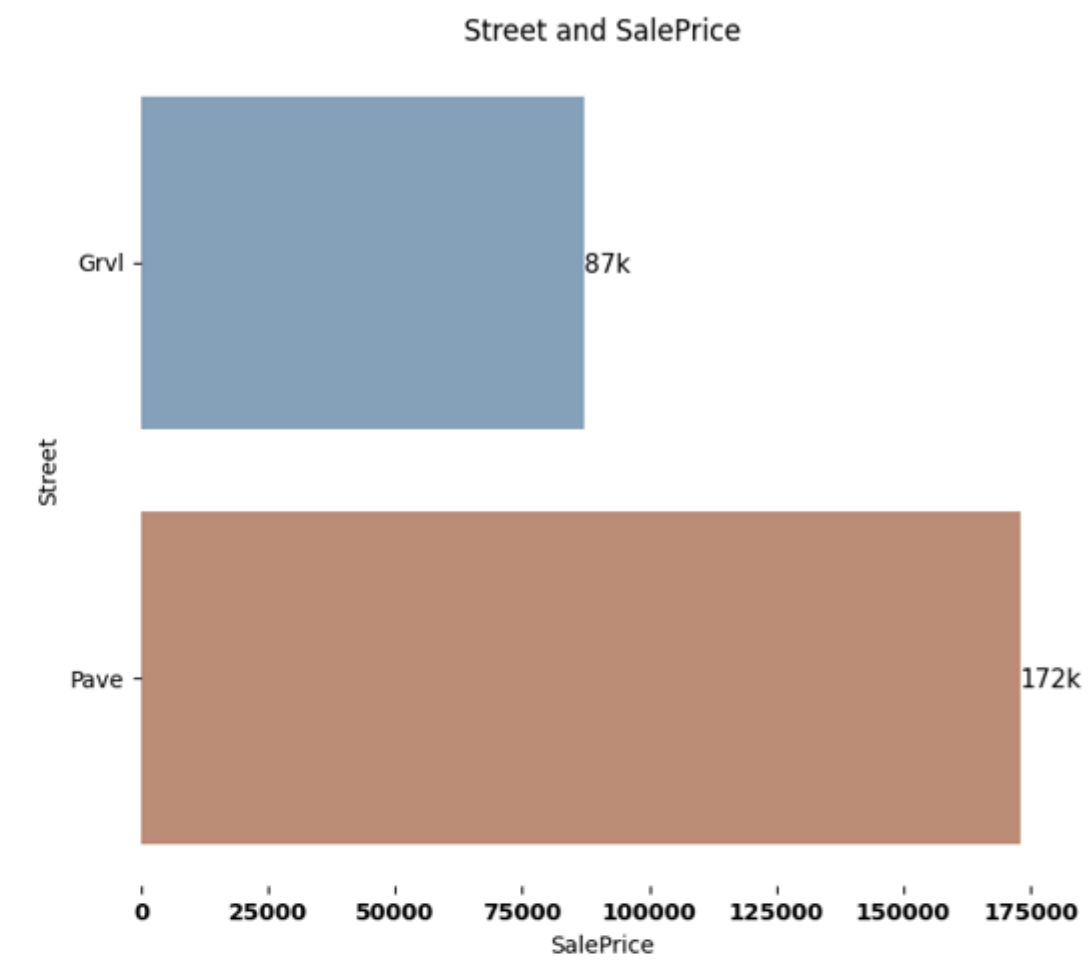
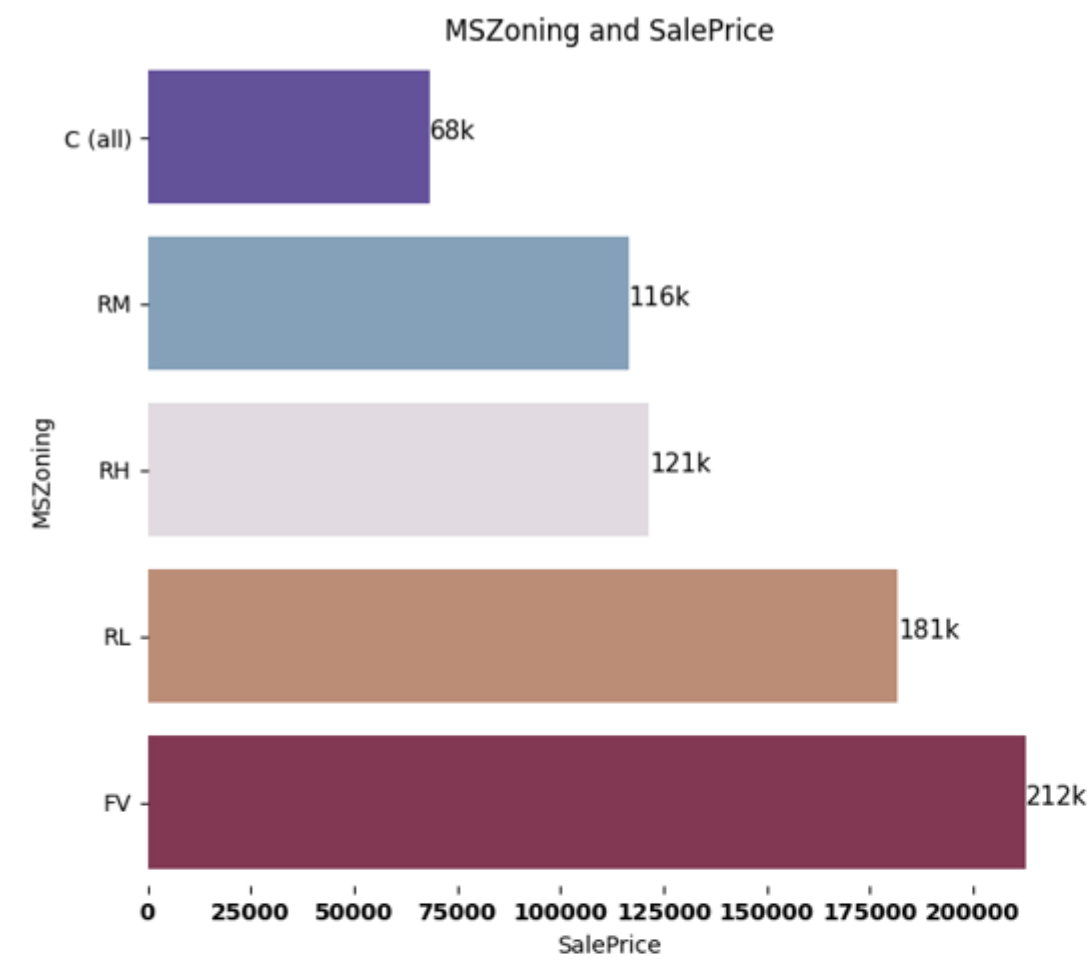
the link for dataset is given [here](#)

1460
rows

81
columns

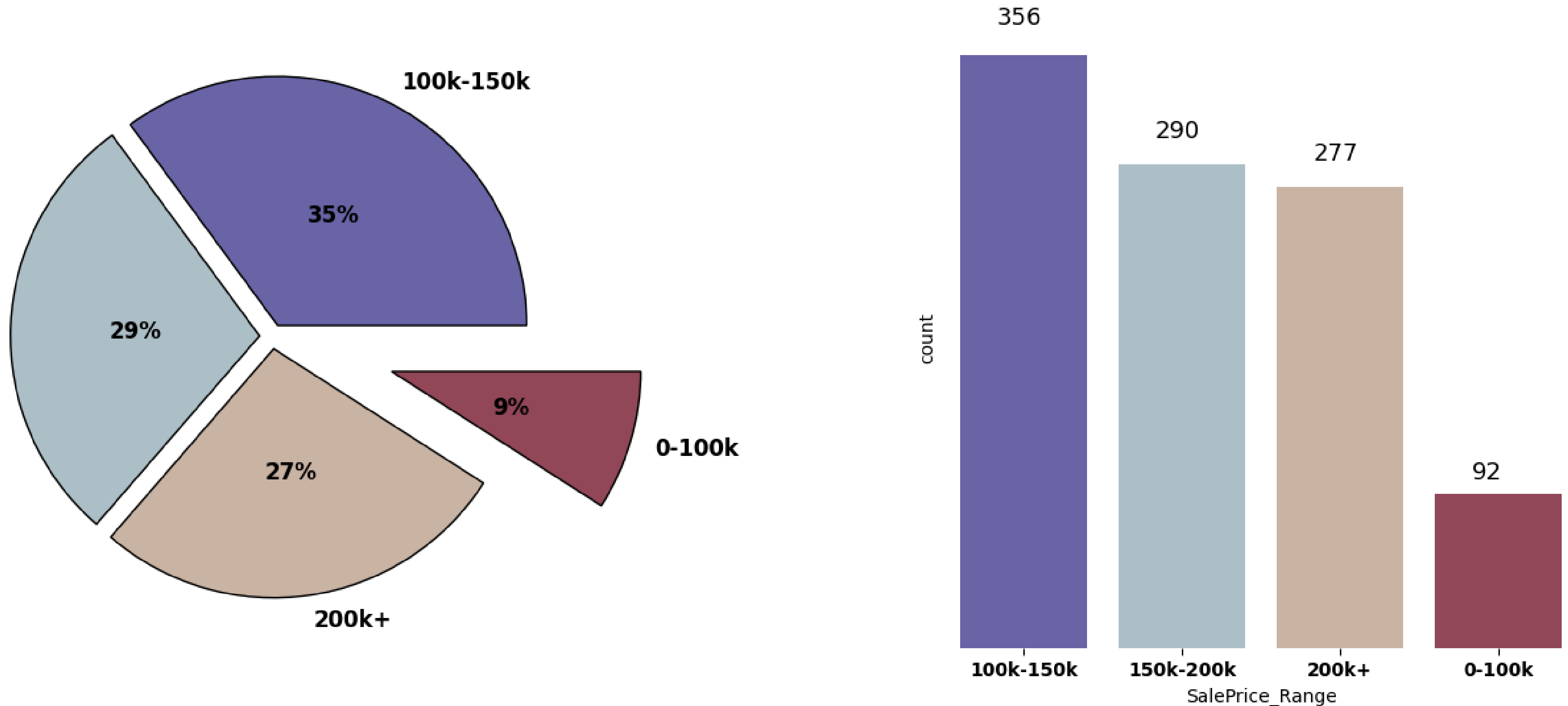


Saleprice is positively distributed



We can observe the price distributions across categorical variables through this visualization. It's evident how the prices are influenced by certain unique features of houses.

► SalePrice Distribution ◀



We created a temporary categorical variable to check the distribution of the target variable. As seen, our dataset predominantly contains houses from the middle and upper segments.

data exploration & cleaning:

first of all we imported the important libraries like pandas, numpy, scikit-learn, matplotlib, seaborn etc.

1. Handling Missing Values:

For numerical features: Imputation with mean.

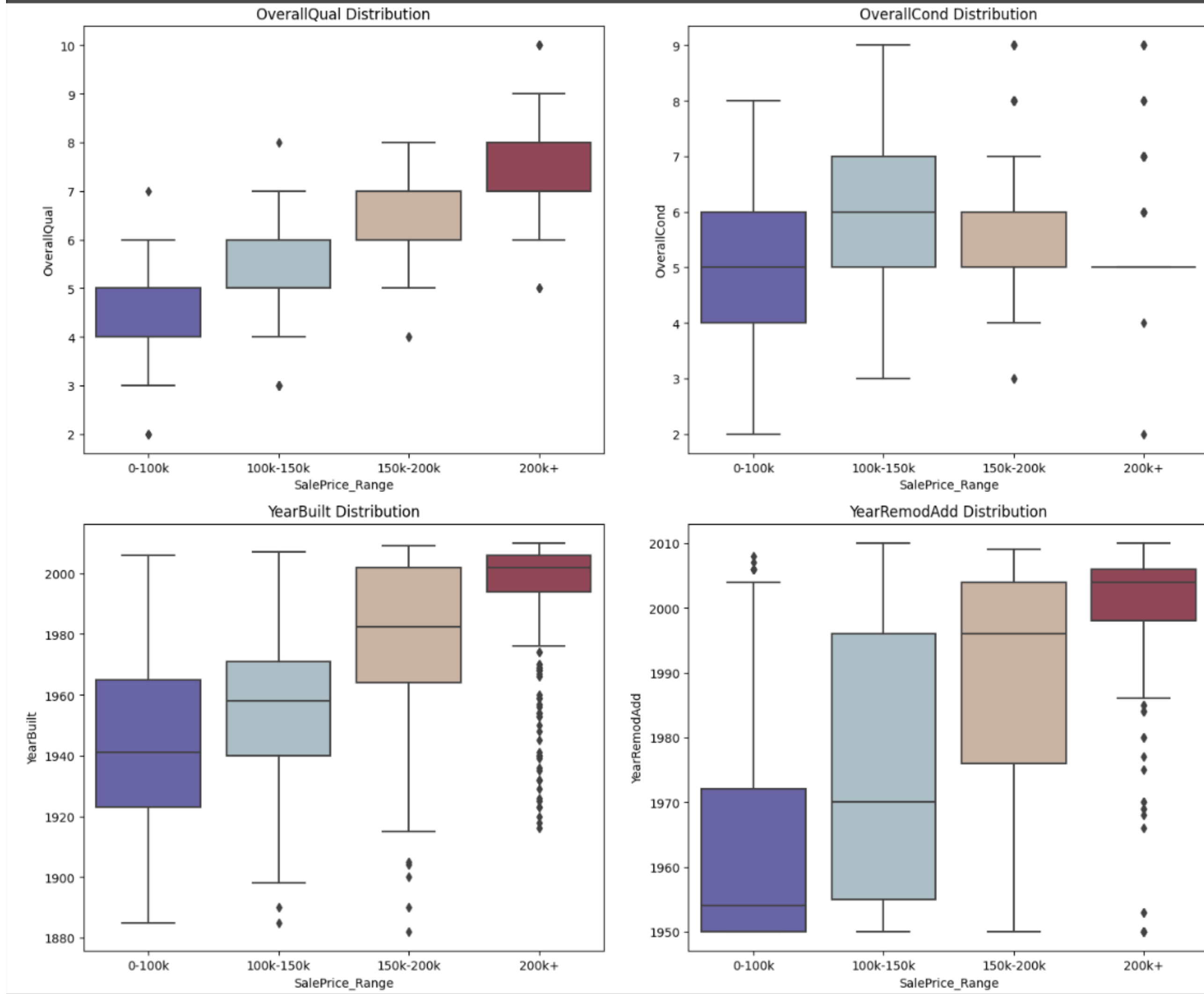
For categorical features: Imputation with the mode.

2. Outlier Detection and Removal:

- Divided columns into categorical (cat)(43) and continuous (con)(38) variables.
- Used Z-score (scipy.stats.zscore) to identify outliers in continuous variables.
- Detected outliers using a threshold (Z-score > 3 or Z-score < -3).

Handling Strategy:

- Removed rows containing outliers.
- Reindexed the dataset for consistency.
- 1015 rows × 24 columns

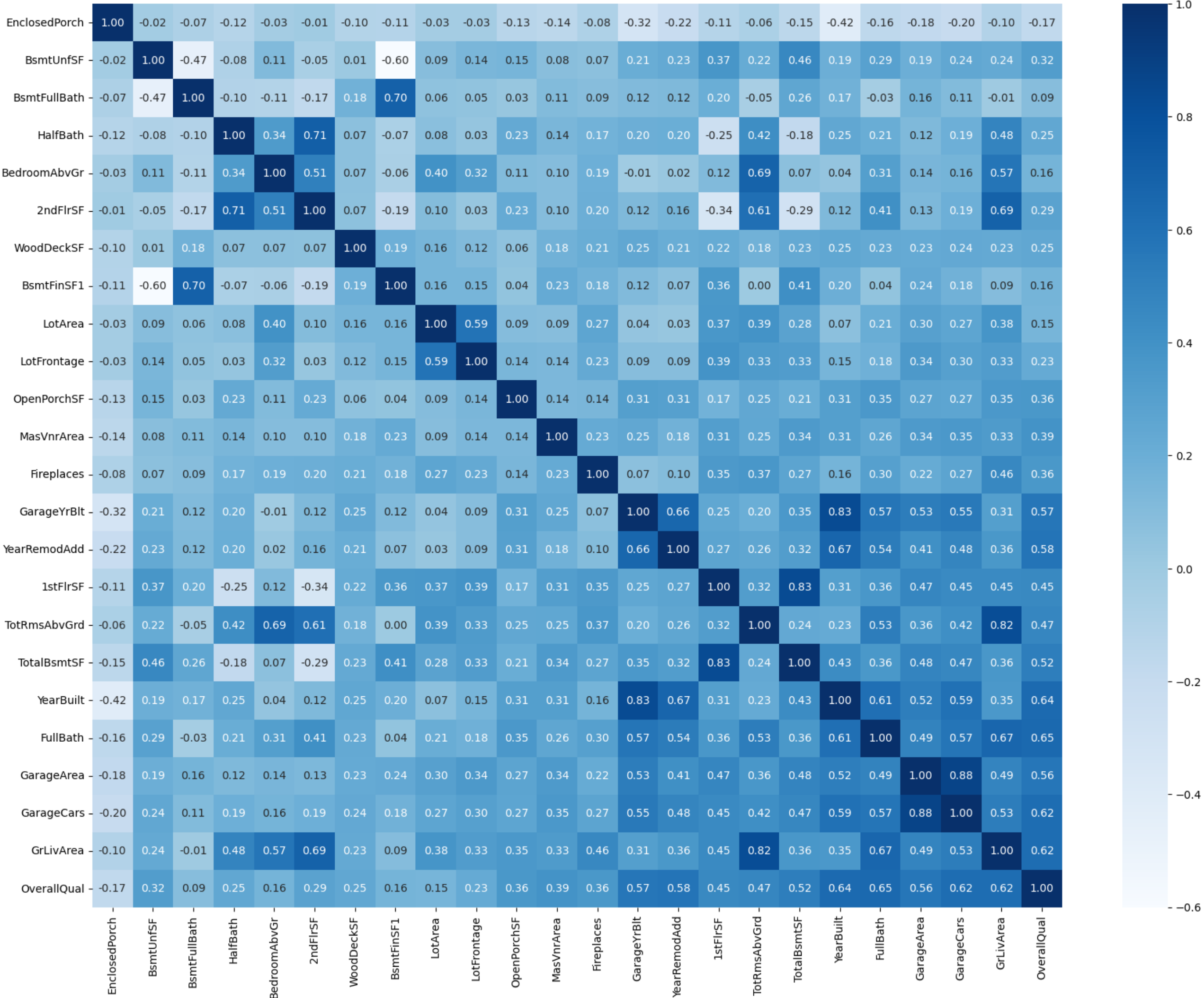


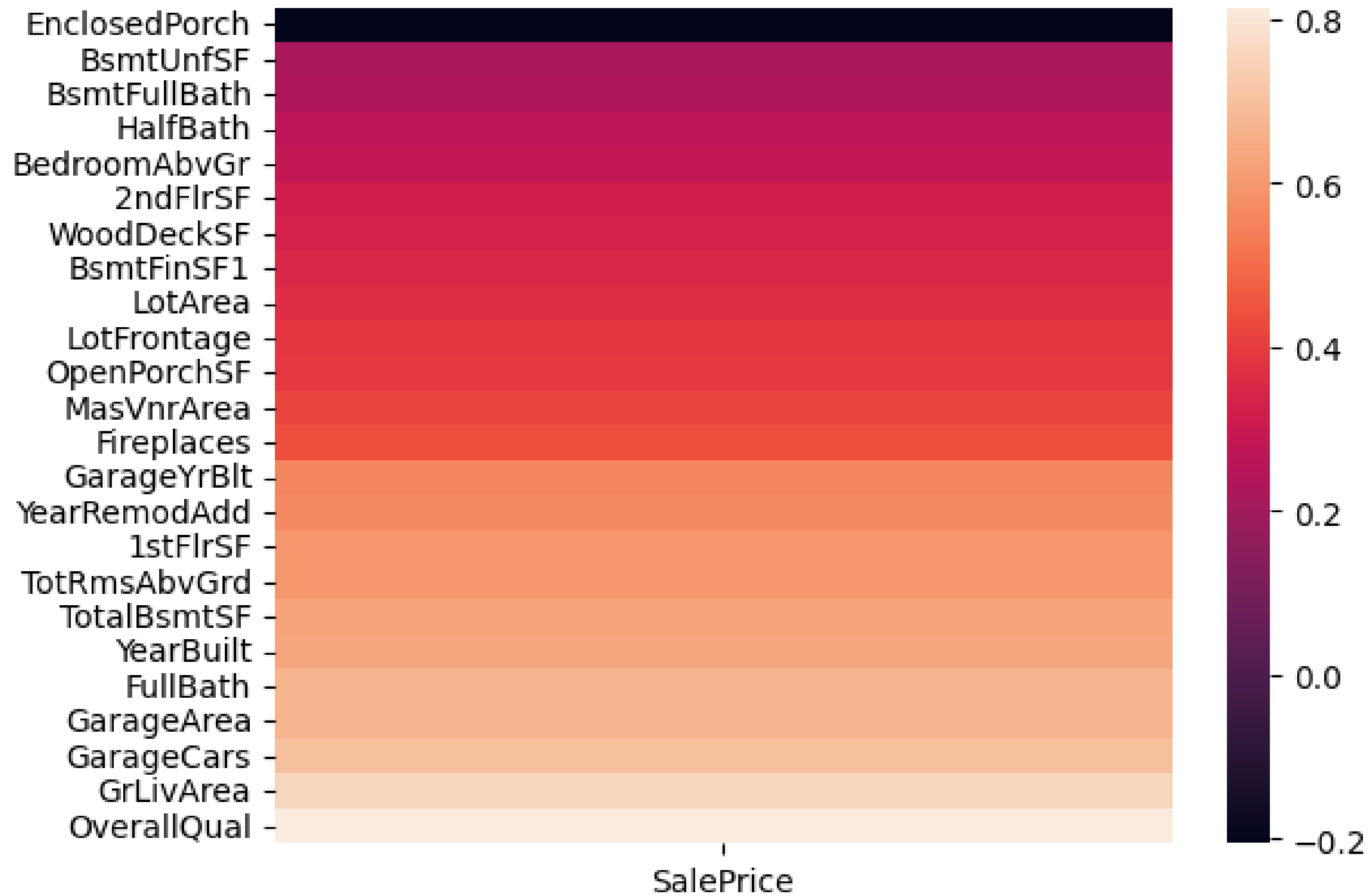
EDA (exploratory data analysis)

Exploratory Data Analysis (EDA) plays a crucial role in understanding the relationships within the dataset and identifying patterns that can guide the feature selection process.

- Determined which features are strongly correlated with the target variable (**SalePrice**).
- Selected features(24 features) with correlations above and below threshold.
- A correlation heatmap visually represents the correlation matrix.
- Features with stronger correlations are highlighted using a color gradient.

Heatmap:





Preprocessing and Feature Selection

- Preprocessing and feature selection are critical steps to ensure that the data is in a suitable format for machine learning models and to identify the most relevant features for accurate predictions.

1. Scaling of Continuous Features:

Ensured that continuous features are on a similar scale to prevent certain features from dominating others.

Utilized MinMaxScaler to scale continuous features between 0 and 1.

2. Dealing with categorical and continuous data:

- We created new dataframe (X1) containing all important continuous features(total 24).
- Then we converted categorical columns into a format suitable for machine learning algorithms.
- Employed `pd.get_dummies()` to create binary columns for each category within a variable.

- Created a new DataFrame (X2) containing these one-hot encoded categorical features(total 224).
- Then joined Categorical and Continuous Data.
- This consolidated dataset (Xnew) is used for subsequent steps in the machine learning process.

Benefits:

The effective handling of categorical data ensures that all features are in a format suitable for machine learning algorithms.

One-hot encoding enables the incorporation of categorical information without introducing ordinal relationships that could impact model performance.

3. Divide the dataset into training and testing sets(75%-25%) for model training and evaluation.

4. Feature Selection using OLS Regression:

Utilized Ordinary Least Squares (OLS) regression and adjusted R-squared values to iteratively drop less significant features.

Model building & evaluation

Model building and evaluation represent the core of the predictive analysis, involving the construction of machine learning models and assessing their performance on both training and testing datasets.

1. **Linear Regression Model:**

Constructed a baseline linear regression model.

Training:

Used the LinearRegression model from scikit-learn to train on the training set (xtrain, ytrain).

Testing Evaluation:

Evaluated the model on the testing set (xtest, ytest).

used Mean Absolute Error (MAE), Mean Squared Error (MSE), R-squared, and Adjusted R-squared value to check accuracy.

Accuracy - R_squared: 0.89

Adj_R2: 0.83

2. Ridge and Lasso Regression Models

- Ridge and Lasso introduce a penalty term to the linear regression objective function, discouraging the model from assigning excessively large weights to any particular feature. This helps prevent overfitting by promoting simpler models.
 - There are many features, and some are believed to be less relevant, Lasso can be used to automatically select a subset of the most important features.
- Used the Ridge and Lasso models from scikit-learn to train on the training set.
- Ridge and Lasso can improve a model's generalization performance by preventing it from fitting the noise in the training data too closely.

Accuracy - Ridge(R_squared: 0.92 , Adj_R2: 0.87)

Lasso(R_squared: 0.90 , Adj_R2: 0.84)

3. Comparison and Model Selection

Compared performance metrics (MAE, MSE, R-squared, Adjusted R-squared) across linear, Ridge, and Lasso models.

- After thorough model training and evaluation, the Ridge regression model has been selected as the final model for predicting house prices based on a combination of performance and interpretability factors.

- Performance:
 - Evaluation Metrics: Comparative evaluation of performance metrics, including Mean Absolute Error (MAE), Mean Squared Error (MSE), R-squared, and Adjusted R-squared, across linear, Ridge, and Lasso models.
 - Balanced Performance: Ridge regression demonstrated a balance between bias and variance, minimizing overfitting while maintaining a good fit to the data.
- Interpretability:
 - Feature Coefficients: Ridge regression, being a regularized linear model, assigns coefficients to features. The interpretability of these coefficients allows for a better understanding of the impact each feature has on the predicted house prices.
 - Model Simplicity: Ridge regression strikes a balance between simplicity and complexity, making it interpretable for stakeholders without sacrificing predictive performance.
- Conclusion: The Ridge regression model is selected for its overall balanced performance and the interpretability it provides, aligning well with the goals of accurate prediction and model transparency.

Predictions and Future Steps:

- The final Ridge regression model was applied to the test set for predicting house prices, and results were prepared for submission.
- The model's predictions are ready for real-world application.
- By exploring these avenues for improvement, our team is committed to staying at the forefront of data science, ensuring our models are not only accurate but also adaptive to emerging challenges and opportunities in the real estate domain. We look forward to the continuous evolution of our predictive capabilities.
- The Github link.

thank you!