## Introduction

The real estate market serves as a vital component of the economy, with housing prices influenced by a wide range of factors including property characteristics, location, and market dynamics. This project focuses on analyzing and predicting housing prices in Ames, Iowa, leveraging a robust dataset and machine learning techniques to explore the relationships between key property features and sale prices. By creating a detailed and accurate predictive model for Ames, this project also aims to build a generalized model for understanding housing markets in similar mid-sized cities across the United States. Such an approach ensures that insights derived from Ames can be reasonably applied to other markets with comparable economic and structural dynamics, providing a valuable framework for broader real estate analysis.

The dataset utilized for this project contains detailed information on 2,930 residential properties in Ames, encompassing 81 features. These features cover a broad spectrum, including physical attributes such as living area, basement size, and garage space; qualitative measures like overall construction quality; and locational characteristics such as neighborhood. The diversity and comprehensiveness of the dataset make it an excellent resource for building predictive models and uncovering critical insights into housing price determinants. The dataset's relatively controlled environment, with consistent pricing trends and a variety of property types, allows for a deep dive into localized price determinants while identifying patterns that can be generalized to similar mid-sized markets.

The primary goal of this project is to accurately predict housing prices using machine learning models while identifying the key factors that drive these prices. To achieve this, multiple machine learning techniques, including Multiple Linear Regression, K-Nearest Neighbors

(KNN), Decision Tree, and Random Forest, were applied. These models were chosen for their ability to capture a variety of relationships in the data, from linear trends to complex nonlinear interactions, and to rank the importance of features that influence housing prices.

To further enhance the predictive capabilities and deepen the analysis, several additional contextual and economic features were considered. These include neighborhood characteristics such as school quality ratings, crime rates, and proximity to parks, grocery stores, and hospitals, as these factors heavily influence buyer preferences and location desirability. Incorporating historical price trends allows the models to account for temporal fluctuations and broader economic factors, such as interest rates and inflation, providing a larger market context for current housing prices. Additionally, data on property renovations and maintenance history, such as kitchen upgrades, energy-efficient installations, or structural improvements, would help capture the added value these features bring to similarly sized properties.

Recognizing the importance of sustainability, the inclusion of energy-efficient and green features like solar panels, upgraded HVAC systems, and sustainable construction practices would provide insights into their influence on housing demand and pricing trends. Moreover, incorporating seasonal market dynamics and accessibility factors—like commuting times to employment hubs, public transportation availability, and major highway proximity—would allow the models to better predict location-based premiums and seasonal fluctuations in property values. Demographic and economic indicators, such as median household income, employment rates, and population growth, would further contextualize the housing market within broader affordability and demand trends.

This project highlights the value of data-driven analysis and predictive modeling in understanding and forecasting housing markets. By focusing on Ames, Iowa, it provides a detailed exploration of property valuation in a stable mid-sized market. At the same time, the project's framework can be generalized to predict housing prices in similar mid-sized cities across the United States, demonstrating its broader applicability. This analysis offers actionable insights for stakeholders, including homeowners, buyers, real estate professionals, and financial institutions, empowering them to make informed decisions and strategies. By leveraging machine learning models and incorporating diverse property, economic, and contextual features, this project underscores the power of data science in shaping real estate valuation and market understanding.

## Data Description

The Ames Housing dataset provides a comprehensive record of 2,930 residential property sales in Ames, Iowa, between 2006 and 2010, containing 81 features that cover structural, locational, and qualitative property attributes. The stable yet diverse housing market in Ames makes it an ideal case study for analyzing property valuation, particularly in mid-sized markets where trends are consistent and insights are generalizable. The dataset's features can be categorized into Property Characteristics, such as 1st Floor Square Footage (1st Flr SF) and Garage Area, which reflect living and storage spaces highly valued by buyers; Location Features, such as Neighborhood, which influence desirability, accessibility, and market value; Quality Indicators, including Overall Quality (Overall Qual), a 1-10 rating of material and construction standards; and Additional Amenities, such as Garage Cars and Total Basement Square Footage (Total Bsmt SF), which enhance functional utility and market appeal.

Six key predictors of SalePrice were identified based on high correlation values: Overall Qual, Garage Area, Garage Cars, 1st Flr SF, Total Bsmt SF, and Neighborhood. These features capture the structural and locational factors that most strongly influence property prices. For instance, Overall Qual emerged as the strongest predictor (correlation coefficient 0.80), underscoring the importance of material and construction quality, while Neighborhood highlighted locational disparities in pricing. Exploratory Data Analysis (EDA) provided further insights into these relationships. The Monthly Average Sale Prices graph showed significant price fluctuations between 2006 and 2010, peaking above $220,000 in early 2006 before dropping to approximately $140,000 by mid-2010, reflecting the prolonged impact of the 2008 financial crisis. Similarly, the Year-Over-Year Price Change graph revealed volatility, with sharp declines of -10% to -20% during 2008-2009, brief recoveries in 2009, and a negative trend persisting into 2010. Seasonal trends were evident in the Average Sale Price by Month chart, where prices peaked in January (~$200,000) and showed slight declines in spring, followed by recoveries during summer and a minor dip in November and December.

Scatter plots further illuminated the relationships between key features and SalePrice. The scatter plot of Garage Area (correlation 0.64) revealed that properties with garages between 400 and 800 sq. ft. sold for $100,000-$300,000, while those with larger garages exceeding 1,000 sq. ft. often surpassed $500,000. Total Bsmt SF (correlation 0.63) showed similar trends, with basements in the 500-2,000 sq. ft. range contributing to prices between $100,000 and $300,000, while larger basements were characteristic of luxury homes exceeding $600,000. Likewise, 1st Flr SF (correlation 0.62) highlighted that homes with floor spaces exceeding 2,500 sq. ft. consistently commanded prices above $400,000, particularly when paired with high Overall Qual ratings. Violin plots illustrated the concentration of lower-quality homes (rated 1-4) within the

$50,000-$150,000 range, while higher-quality homes (rated 8-10) predominantly exceeded $300,000. The analysis of SalePrice by Neighborhood further confirmed local disparities, with neighborhoods like North Ridge and Stone Brook averaging prices above $300,000, whereas others like Meadow Village and Briardale remained below $150,000.

To ensure robust analysis, missing data and outliers were carefully handled. Median imputation addressed missing values for key numerical variables like Garage Area, Total Bsmt SF, and 1st Flr SF, while frequent category imputation handled missing values in categorical variables like Neighborhood. Features with excessive missing data or limited relevance were excluded to streamline the analysis. These measures preserved data integrity and enhanced the predictive power of the models. The focus on six key predictors—Garage Cars, Garage Area, Total Bsmt SF, 1st Flr SF, Overall Qual, and Neighborhood—allowed the creation of accurate and interpretable models. These findings highlight the dominant role of quality, usable space, and location in property valuation, particularly in mid-sized, stable markets like Ames.

In conclusion, the Ames Housing dataset revealed valuable insights into the determinants of residential property values. Despite market challenges during 2006-2010, structural features such as 1st Flr SF and Garage Area, combined with location and quality indicators like Neighborhood and Overall Qual, consistently emerged as strong predictors of SalePrice. Monthly and year-over-year trends highlighted the impact of macroeconomic factors, while EDA confirmed robust relationships between key features and property values. These insights provide homeowners, buyers, and real estate professionals with actionable knowledge about the Ames housing market, offering broader implications for similar markets.

## Models and Methods

To predict housing sale prices in Ames, Iowa, I explored a variety of machine learning models, leveraging the comprehensive Ames Housing dataset. This dataset, with its 82 features capturing both numerical and categorical attributes, offered a rich foundation for building and testing predictive models. My goal was to compare different approaches to find the most effective one for this task while maintaining a focus on accuracy, interpretability, and reproducibility.

I began by setting up my development environment, using essential Python libraries like NumPy and Pandas for data manipulation, Matplotlib and Seaborn for visualizations, and Scikit-learn for model building and evaluation. This foundational setup allowed for a streamlined workflow from preprocessing to model evaluation.

For modeling, I selected several regression techniques that represent a range of strengths and methodologies. I started with Linear Regression for its simplicity and ease of interpretation. Next, I added the K-Nearest Neighbors (KNN) Regressor, a non-parametric model that relies on feature similarity to make predictions. I also implemented the Decision Tree Regressor, known for its clear decision rules, and the Random Forest Regressor, a robust ensemble method that combines multiple decision trees to improve accuracy and reduce overfitting. This diverse selection allowed me to compare models with different complexities and capabilities.

Preprocessing was a critical step due to the dataset's complexity and diversity. Using Scikit-learn's tools, I built a preprocessing pipeline that standardized numerical features with StandardScaler, encoded categorical features like neighborhood using OneHotEncoder, and handled missing values with SimpleImputer. To manage these transformations efficiently, I used Pipeline and ColumnTransformer, which ensured a clean and reproducible workflow.

To evaluate model performance, I adopted a rigorous approach using cross-validation. This technique allowed me to assess how each model performed on multiple subsets of the data, providing a more reliable measure than a single train-test split. For optimization, I applied GridSearchCV, which systematically tested various hyperparameter combinations to find the best configuration for each model. Key evaluation metrics included mean squared error (MSE) to quantify prediction error and R-squared ($R^2$) to measure the proportion of variance explained by the models.

Understanding model behavior was another priority. I used permutation importance to identify the features that had the most significant impact on predictions, enhancing interpretability. For decision tree models, I visualized their structure using plot_tree, which provided a clear representation of how predictions were made. These tools made the models not just accurate but also understandable.

By combining simple models like Linear Regression with more advanced approaches like Random Forests, I could balance interpretability and predictive power. Random Forests, in particular, stood out for their ability to handle complex feature interactions while maintaining robustness. Overall, this systematic approach—from preprocessing and feature engineering to model selection and interpretation—enabled me to develop a reliable framework for predicting housing prices, offering valuable insights into the factors driving property values in Ames.

## Results and Interpretation

This analysis of the Ames housing dataset implemented and compared multiple predictive models to determine their ability to predict housing prices. The models included a baseline

predictor, multiple regression, k-nearest neighbors (KNN), decision tree, and random forest. Performance evaluation was conducted using metrics such as **Mean Squared Error (MSE)**, **Root Mean Squared Error (RMSE)**, **R² Score**, and **cross-validation scores**. These metrics provided insights into model accuracy, generalizability, and their ability to handle complex relationships in the data.

The **baseline model**, which predicts the mean sale price for all observations, resulted in a **high MSE of approximately 637,970,549.8**. This simple model set a benchmark for evaluating the performance of more advanced models. Its substantial error highlighted the limitations of ignoring predictor variables and demonstrated the importance of using sophisticated techniques to capture the underlying patterns in housing prices.

The **multiple regression model** incorporated both categorical and numerical features, such as Garage Area, Total Basement SF, and Neighborhood. Through preprocessing steps like **one-hot encoding** and **scaling**, the model effectively captured linear relationships within the data. It achieved a **training MSE of 1,297,575,175.16** and a **testing MSE of 1,527,669,934.49**, with an **R² score of 0.857** for training and **0.8095** for testing. Cross-validation produced an average R² score of **0.7726**, indicating strong performance and generalizability. However, its reliance on linear assumptions limited its ability to capture more complex, nonlinear patterns, particularly for high-value properties where residual errors were higher.

The **K-nearest neighbors (KNN) model** provided a proximity-based approach to price prediction by tuning the number of neighbors (k = 3) and weights. This model achieved a **testing MSE of 1,321,604,517.26**, an RMSE of **36,491.17**, and an **R² score of 0.8321**. Cross-validation yielded an average R² score of **0.8031**. The KNN model excelled in identifying localized trends,

such as neighborhood-specific variations in housing prices. However, it faced challenges when dealing with high-dimensional data and rare feature combinations, resulting in slight performance degradation compared to the ensemble methods. Despite this, KNN demonstrated its strength in modeling regional price dynamics.

The **decision tree model** introduced hierarchical splits that made the model highly interpretable. It achieved a **training MSE of 1,471,466,553.52** and a **testing MSE of 1,889,473,553.00**, with an **R² score of 0.7643**. Cross-validation yielded an average R² score of **0.7529**, slightly lower than KNN and regression. While decision trees provided clarity in understanding feature importance, they struggled to generalize due to their tendency to overfit, particularly on training data. The residual analysis showed errors concentrated for high-value properties, reflecting the limitations of binary splits in capturing complex interactions.

The **random forest model**, an ensemble of decision trees, delivered the best overall performance. By combining multiple trees and reducing overfitting, the model achieved a **training MSE of 382,989,351.36** and a **testing MSE of 1,210,127,955.51**, with an **R² score of 0.8491**. Cross-validation further validated its robustness, producing an average R² score of **0.8175**, the highest among all models. Feature importance analysis revealed that **Overall Quality**, **1st Floor SF**, and **Garage Area** were key predictors of housing prices, aligning with domain knowledge. The random forest model effectively captured nonlinear relationships and generalized well across the test set. However, residual analysis indicated slightly higher errors for extreme property values, suggesting the potential for further refinement through advanced techniques like gradient boosting or XGBoost.

The evaluation of the best-performing model was based on a combination of testing MSE, $R^2$ scores, and cross-validation results. The random forest model outperformed all others, achieving the lowest test MSE and the highest $R^2$ score, indicating its superior accuracy and generalizability. While multiple regression performed well, its linear assumptions limited its ability to model nonlinear patterns. KNN provided localized insights but struggled with high-dimensional data, and the decision tree model, though interpretable, fell short in predictive accuracy compared to the ensemble method.

## Conclusion

The Random Forest model clearly emerged as the strongest performer among all the predictive approaches applied to the Ames housing dataset. It offered the best balance between predictive accuracy, feature importance insights, and generalization. By aggregating results from multiple decision trees, the Random Forest effectively minimized both training and testing errors, achieving a high $R^2$ score of 0.85 on the test set and the lowest Mean Squared Error (MSE) among all the models. This performance highlights its capability to capture both linear and nonlinear relationships in the data while maintaining resilience to overfitting, which was evident in the consistency between training and testing results.

In contrast, the **Multiple Regression model**, while effective in capturing linear relationships and achieving strong results with an $R^2$ score of 0.81, showed limitations in handling the nonlinear dependencies and outliers within the dataset. Despite its solid performance, particularly with scaled numerical features and one-hot encoded categorical variables, the residual patterns indicated that the model struggled to predict high-value properties accurately.

The **K-Nearest Neighbors (KNN) model** also provided valuable insights by capturing localized patterns and neighborhood-specific pricing trends. With a tuned parameter $k=3k = 3k=3$, the model achieved a competitive $R^2$ score of 0.83. However, its reliance on proximity-based methods limited its ability to generalize across high-dimensional data, and computational inefficiencies posed challenges for larger datasets. While KNN performed well in specific regions of the data, it fell short compared to ensemble methods like Random Forest.

The **Decision Tree model** delivered excellent interpretability, offering clear visualizations of feature importance and decision rules. Features such as **Overall Quality** and **Garage Area** emerged as significant predictors, reinforcing their importance in determining housing prices. However, the model's simplicity came at the cost of accuracy, as its $R^2$ score of 0.76 indicated a weaker ability to generalize compared to the ensemble methods. The binary splitting nature of decision trees limited their capacity to capture complex interactions, leading to higher errors for extreme property values.

Ultimately, the Random Forest model stood out as the best predictor due to its ability to balance accuracy, generalizability, and robustness. Its feature importance analysis highlighted critical predictors such as **Overall Quality**, **1st Floor SF**, and **Garage Area**, aligning with domain knowledge and reinforcing confidence in its results. While all models demonstrated improvements over the baseline predictor, Random Forest's superior performance across evaluation metrics makes it the most reliable model for understanding and predicting housing prices in the Ames dataset.

## Next Steps

To enhance the predictive capabilities of my models and gain deeper insights into the Ames housing market, incorporating additional contextual, economic, and structural features would be invaluable. These enhancements would provide a more comprehensive understanding of the diverse factors influencing housing prices, enabling the models to better capture market intricacies and improve their predictive accuracy. While this model is based on Ames, its methodologies and findings could be reasonably extended to predict and model housing trends in most midsize cities across the USA, as they share similar market dynamics, housing diversity, and economic conditions.

A critical area of improvement would be the inclusion of detailed neighborhood characteristics. Data on school quality ratings, crime rates, and proximity to public amenities like parks, grocery stores, and hospitals would provide a more nuanced understanding of how location impacts property values. These factors often drive buyer preferences, and their integration into the analysis would allow the models to better differentiate between similar properties located in different neighborhoods, accounting for locational desirability and quality of life.

Another key addition would be historical price trends. By incorporating data on how property values have evolved over time across neighborhoods, the models could account for temporal fluctuations and broader economic conditions such as interest rate changes and inflation. Situating current housing prices within this larger market context would enable the models to predict long-term value trajectories more effectively, providing a forward-looking perspective to buyers and real estate professionals.

Details on property renovations and maintenance history could also significantly enhance the models. Information about recent upgrades to critical areas, such as kitchens and bathrooms, or the addition of energy-efficient appliances, would help quantify the value added by these improvements. This data would allow the models to capture differences in buyer interest and pricing for similarly sized homes that vary in condition and quality.

The growing emphasis on sustainability highlights the importance of incorporating features such as solar panels, energy-efficient HVAC systems, and other green upgrades. These features not only reflect current buyer preferences but also influence property pricing due to long-term cost savings and environmental benefits. Including data on these sustainable features would enable the models to evaluate their impact on property values, adding a modern and environmentally conscious dimension to the analysis.

Seasonal market dynamics represent another critical area worth exploring. Housing prices often fluctuate based on the time of year, with spring typically witnessing increased buyer activity and competition compared to winter months. Incorporating this information would allow the models to better predict seasonal impacts on property valuations and buyer behavior, further improving accuracy.

Accessibility factors, such as commuting times to employment hubs, proximity to major highways, and the availability of public transportation, would also enrich the analysis. Buyers often prioritize convenience and reduced travel times, particularly in midsize cities, where commuting dynamics play a significant role. Including these location-based elements would enable the models to assess premiums associated with transportation accessibility.

Finally, incorporating demographic and economic indicators like median household income, population growth, and employment rates would provide valuable market context. These indicators would help the models understand broader economic conditions that influence property demand and affordability, particularly in specific neighborhoods. By analyzing these factors alongside structural and locational features, the models could offer a more holistic view of market pressures and their effects on housing prices.

While this enhanced model is rooted in Ames, its robust methodology and focus on critical predictors make it highly applicable to most midsize cities in the USA. These markets share a combination of stable growth, diverse housing inventory, and consistent buyer preferences, enabling the model to generalize effectively. By accounting for key variables like property characteristics, location dynamics, and economic conditions, this approach provides actionable insights into housing market trends, empowering homeowners, buyers, and real estate professionals to make informed decisions across similar markets.