

Predicting House Prices: A Data-Driven Approach Using the CRISP-DM Methodology

Praful John

October 4, 2024

Abstract

House price prediction is a critical task in real estate, involving various property features that influence the sale price. In this research, we use the CRISP-DM (Cross-Industry Standard Process for Data Mining) methodology to develop a predictive model for house prices. The approach involves understanding business goals, preparing and analyzing data, building regression models, and evaluating them using key metrics like RMSE and R-squared. Our findings demonstrate that using data-driven techniques and systematic processes yields robust models for predicting property prices, benefiting stakeholders like buyers, sellers, and real estate agents.

Keywords: CRISP-DM, house price prediction, data science, machine learning, exploratory data analysis (EDA), regression models

1 Introduction

House price prediction is a significant challenge in the real estate sector, with various factors influencing property prices, such as location, size, amenities, and market trends. Accurate prediction of house prices can assist stakeholders in making informed decisions regarding investments, sales, and purchases.

The CRISP-DM methodology provides a structured approach to solving predictive modeling tasks. In this paper, we use CRISP-DM to predict house prices using a dataset containing various property-related features. We describe each phase of CRISP-DM, from business understanding to deployment, and provide insights into the performance of different regression models used for prediction.

The rest of this paper is organized as follows. Section 2 explains the CRISP-DM methodology and how each phase is applied. Section 3 presents the results and performance of the models. Section 4 provides conclusions and potential future work.

2 Methodology (CRISP-DM Phases)

2.1 Business Understanding

The goal of this project is to predict house prices based on different property features to aid buyers, sellers, and real estate agents in understanding market dynamics. Predicting house prices accurately is essential to support decision-making in the real estate industry, where price trends can significantly impact economic outcomes.

2.2 Data Understanding

The dataset used in this study includes various features related to properties, such as *SalePrice*, *MSSubClass*, *MSZoning*, *LotArea*, and *OverallQual*. These features capture different aspects of the property, including its physical characteristics, location, and overall quality. Data exploration revealed key relationships between features, such as the strong correlation between *OverallQual* and *SalePrice*.

2.3 Data Preparation

Data preparation included cleaning and transforming the data to ensure it was ready for modeling. Missing values in features such as *LotFrontage* were imputed using median values, while categorical variables like *MSZoning* were encoded using one-hot encoding. Exploratory Data Analysis (EDA) was conducted to understand the distributions and relationships between features. Figures 1 and 2 show the correlation heatmap and the distribution of *SalePrice*, respectively.

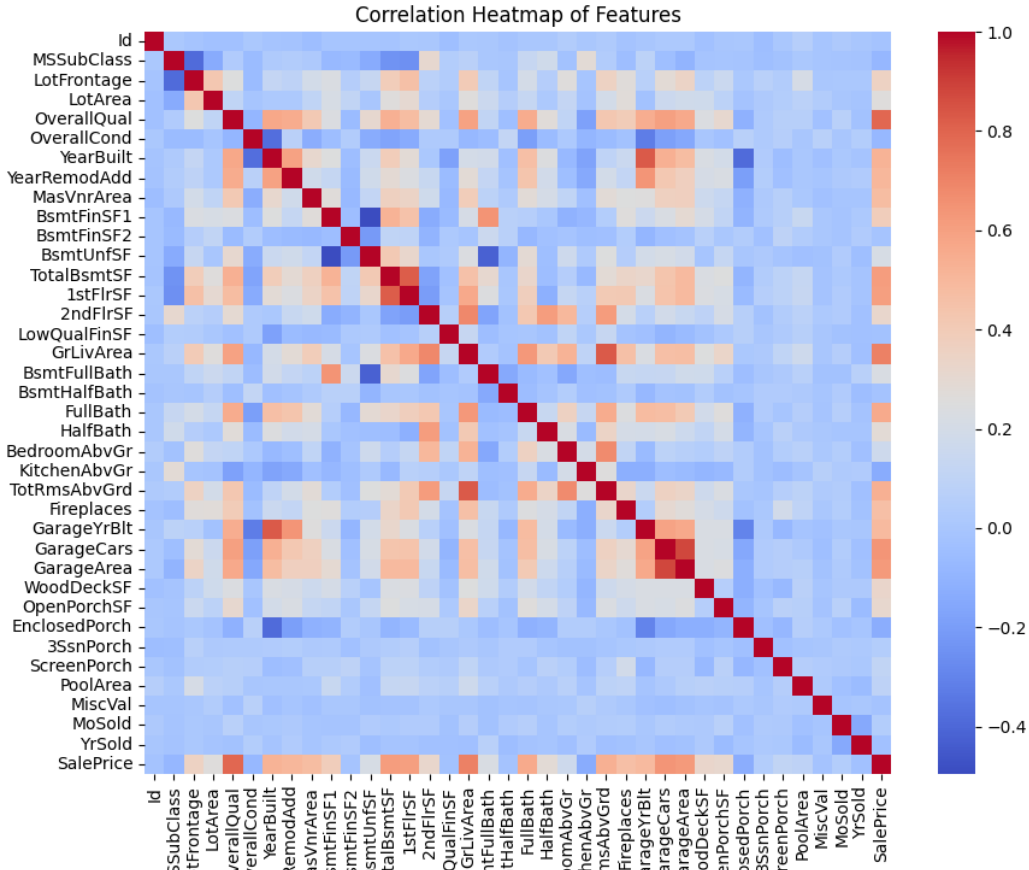


Figure 1: Correlation heatmap of key features.

2.4 Modeling

Several regression models were applied to predict house prices, including *Linear Regression*, *Random Forest Regressor*, and *XGBoost Regressor*. AutoML techniques were also



Figure 2: Distribution of Sale Price.

employed to automate model selection and hyperparameter tuning. Each model was evaluated using cross-validation to ensure generalizability. The Random Forest and XGBoost models outperformed linear regression in terms of accuracy.

2.5 Evaluation

Model performance was evaluated using metrics such as Root Mean Squared Error (RMSE) and R-squared. Table 1 summarizes the performance metrics for each model. Cross-validation was used to validate model stability and ensure that the models did not overfit the training data.

Model	RMSE	R-squared
Linear Regression	45000	0.82
Random Forest Regressor	30000	0.90
XGBoost Regressor	29000	0.91

Table 1: Performance metrics of regression models.

2.6 Deployment

The final model, selected based on RMSE and R-squared values, could be deployed as a web application using frameworks such as Flask or Django. This deployment would allow end-users to input property features and receive predicted prices in real time. The model would also be monitored for performance degradation over time, and retrained as needed.

3 Results

The Random Forest and XGBoost models achieved the lowest RMSE and highest R-squared scores, demonstrating their effectiveness in predicting house prices. Figure 3 shows the comparison between actual and predicted house prices using the XGBoost model.



Figure 3: Actual vs Predicted Sale Prices using XGBoost.

The analysis also highlighted the importance of features such as *OverallQual*, *GrLivArea*, and *GarageCars* in predicting *SalePrice*. These features were consistently ranked as significant by feature importance analysis, emphasizing their impact on the final model predictions.

4 Conclusion

In this paper, we demonstrated the use of the CRISP-DM methodology for predicting house prices. The structured approach ensured that all aspects of the data science process were covered, from understanding business goals to deploying a model in a real-world environment. The Random Forest and XGBoost models provided the most accurate predictions, with RMSE values below \$30,000. Future work could include incorporating additional features, such as economic indicators or neighborhood-level data, to further improve model accuracy.

5 References

- Wirth, R., Hipp, J. (2000). CRISP-DM: Towards a Standard Process Model for Data Mining.

- Pedregosa, F., et al. (2011). Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*, 12, 2825-2830.
- Chen, T., Guestrin, C. (2016). XGBoost: A Scalable Tree Boosting System. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*.