

An Efficient and Interpretable XGBoost Framework for Housing Price Prediction: Enhancements in Accuracy and Computational Performance

¹M.Sandeep Kumar, ^{2*}M. Saraswathi, ³B. Yadav Singh, ⁴D. Divya, ⁵J. Shiva Shanker

¹Assistant professor, Department of Computer Science and Engineering, St.Mary's Engineering College, Hyderabad, India

^{2, 3, 4, 5} B.Tech Student, Department of CSE(DS), St.Mary's Engineering College, Hyderabad, India

*Corresponding Author(s): saraswathi424240@gmail.com

Received: 11/09/2024

Revised: 16/11/2024

Accepted: 11/12/2024

Published: 31/12/2024

Abstract: Accurate and interpretable prediction of housing prices is a critical task in real estate and urban planning, requiring models that balance predictive performance, computational efficiency, and transparency. This research proposes an Efficient and Interpretability-Driven XGBoost Model to address the limitations of traditional regression methods and standard XGBoost in handling high-dimensional, sparse datasets. Key innovations include dynamic feature prioritisation, adaptive binning for continuous variables, sparse feature handling, and a regularised gain function to promote simpler and more interpretable tree structures. These enhancements improve the computational efficiency of the model while retaining high predictive accuracy and interpretability. The proposed model was evaluated using the Ames Housing Dataset, a benchmark dataset in real estate analytics. Key performance metrics, including Root Mean Squared Error (RMSE), Mean Absolute Error (MAE), and Coefficient of Determination (R^2), were used to compare its performance against Linear Regression, Ridge Regression, Lasso Regression, and standard XGBoost. The proposed model achieved the best performance with a test RMSE of 24,630.91, R^2 of 0.923, and MAE of 19,637.68, demonstrating a 5.4% improvement in accuracy over the standard XGBoost, while reducing computational overhead by 15%. Additionally, the model's interpretability was enhanced through SHAP-based feature importance analysis, highlighting critical predictors like "GrLivArea" and "OverallQual." These findings establish the proposed methodology as a robust tool for housing price prediction, offering significant advancements in accuracy, efficiency, and transparency. Future work will explore scalability for larger datasets and the integration of temporal dynamics to extend the model's applicability.

Keywords: Housing Price Prediction, XGBoost Optimization, Ames Housing Dataset, Feature Prioritization, Model Interpretability, Computational Efficiency

1. Introduction

The prediction of housing prices is a critical task in real estate, urban planning, and economic policymaking, where accurate and interpretable models can significantly influence decision-making processes. Housing prices are influenced by a wide range of factors, including property characteristics, location-based variables, and macroeconomic indicators, making it a complex problem to model effectively.[1] Traditional statistical approaches such as linear regression have been extensively utilized due to their simplicity and interpretability.[2] However, these methods often fall short in capturing nonlinear relationships and intricate feature interactions inherent in real-world datasets such as the Ames Housing Dataset.

In recent years, machine learning models, particularly ensemble learning techniques like Gradient Boosting Machines (GBMs), have emerged as powerful tools for predictive tasks in structured data.[3] Among these, XGBoost has gained prominence for its robustness, scalability, and ability to handle high-dimensional data with missing values.[4] Despite its success, XGBoost faces challenges in computational efficiency and model interpretability, especially when applied to large-scale or sparse datasets. These limitations present a significant gap in developing models that are not only accurate and efficient but also interpretable and practical for stakeholders in the housing market.[5]



This research aims to address these challenges by proposing an Efficient and Interpretability-Driven XGBoost Model. The proposed methodology incorporates novel enhancements, such as dynamic feature prioritisation, adaptive binning, sparse feature handling, and a regularised gain function, to improve both computational efficiency and model transparency. The model was evaluated on the Ames Housing Dataset [6], a widely used benchmark for housing price prediction, using key performance metrics such as RMSE, MAE, and R^2 . Comparative analysis against baseline models, including Linear Regression, Ridge Regression, Lasso Regression, and standard XGBoost, highlights the effectiveness of the proposed approach. [7]

By bridging the gap between accuracy, efficiency, and interpretability, this research contributes to the development of machine learning models that are both high-performance and practically applicable in real-world housing price prediction tasks. These findings have significant implications for researchers, practitioners, and policymakers seeking reliable and transparent tools to understand and predict housing market trends. [8]

Key Contributions: This study aims to contribute to the growing body of knowledge in real estate analytics by providing a comprehensive framework for house price prediction. Furthermore, it demonstrates the practical applications of machine learning techniques in solving real-world problems, emphasising their potential to transform the real estate industry through data-driven insights.

- **Dynamic Feature Prioritisation:** This methodology introduces a novel weighting mechanism to dynamically prioritise interpretable and impactful features during the tree-splitting process. This ensures that features with high domain relevance, such as "GrLivArea" and "OverallQual", are favoured, improving both the transparency and real-world applicability of the model relationships that drive property valuation.
- **Regularised Gain Function:** A modified gain function is proposed that incorporates penalties for the split complexity and feature sparsity. This innovation promotes simpler and more interpretable tree structures while maintaining high predictive accuracy, balancing the trade-off between model performance, and transparency insights into their applicability to real-world scenarios.
- **Adaptive Binning for Continuous Features:** This methodology employs adaptive binning techniques, such as quantile-based binning and clustering-based binning (e.g. k-means), to optimise the search space for split points. This reduces the computational overhead while retaining sufficient granularity to capture important data patterns, thereby enhancing the efficiency of visual data protection.

2. Literature Review

2.1. Traditional Regression Methods

Linear regression has long been a cornerstone of predictive tasks because of its simplicity and interpretability. [9] Numerous studies have employed linear regression for housing price prediction, particularly leveraging datasets such as the Ames Housing Dataset. However, its performance is often limited in capturing nonlinear relationships and complex feature interactions, which are intrinsic to real-world housing data. To overcome these limitations, regularised regression techniques such as Ridge and Lasso regression have been adopted. [10] Ridge regression minimises overfitting by penalising large coefficients, whereas Lasso regression promotes sparsity by zeroing out less important features. These methods improve model robustness and feature selection but remain constrained by their linear assumptions.

2.2. Ensemble Learning Techniques

Ensemble methods, particularly Gradient Boosting Machines (GBMs), have revolutionised predictive modelling by effectively capturing non-linearities and high-dimensional feature interactions [11]. XGBoost, a prominent GBM, has become a benchmark for structured data tasks owing to its scalability, support for missing values, and regularisation capabilities. It iteratively minimises residual errors, yielding a high predictive performance. However, computational efficiency and interpretability remain challenges, particularly for large-scale or sparse datasets. Alternative GBMs, such as LightGBM and CatBoost, have attempted to address these issues through novel split-finding techniques and categorical feature handling [12]. Nonetheless, enhancing the efficiency of XGBoost while maintaining its robust predictive power and interpretability remains an area of active research.

2.3. Feature Engineering in Housing Price Prediction

Feature engineering plays a pivotal role in housing price prediction, as shown in multiple studies using the Ames Housing Dataset [13]. Key predictive features such as "GrLivArea" (above ground living area), "OverallQual" (overall quality), and "Neighborhood" (location-based grouping) were consistently identified as significant contributors to housing price variability. Traditional methods often rely on manual feature engineering, which is time intensive and prone to subjectivity. Recent advancements include automated techniques such as interaction term generation and clustering-based binning to enhance feature representation [14]. Additionally, handling sparse features (e.g. "PoolQC" and "Alley") has been recognised as a critical factor for improving model robustness and performance.

2.4. Interpretability in Machine Learning Models

With increasing reliance on machine learning models, interpretability has emerged as a critical requirement, particularly in real estate applications, where decision making depends on understanding model predictions. XGBoost and similar ensemble methods often face criticism for being "black-box" models. Methods such as SHAP



(SHapley Additive exPlanations) have been widely adopted to enhance the interpretability of these models by providing feature importance scores and interaction insights [15]. However, improving intrinsic interpretability through simpler tree structures and feature prioritization remains an important research direction.

2.5. Identified Gaps and Motivation

While traditional and advanced methods offer various strengths, there is a clear need for a model that balances computational efficiency, predictive accuracy, and interpretability. Current XGBoost implementations excel in performance but face scalability and transparency challenges. Moreover, the importance of dataset-specific optimization, such as dynamic feature prioritization and efficient sparse data handling, remains underexplored. This review highlights these gaps and serves as the foundation for the proposed Efficient and Interpretability-Driven XGBoost Model, which incorporates novel enhancements to address these challenges effectively [16].

3. Methodology

The proposed methodology introduces enhancements to XGBoost's tree construction algorithm tailored for the prediction of housing prices using the Ames Housing Dataset. This dataset, known for its rich and diverse set of features, poses challenges such as high dimensionality, sparse data, and the need for interpretable predictions. The novel approach focuses on improving computational efficiency and model interpretability while maintaining predictive accuracy.

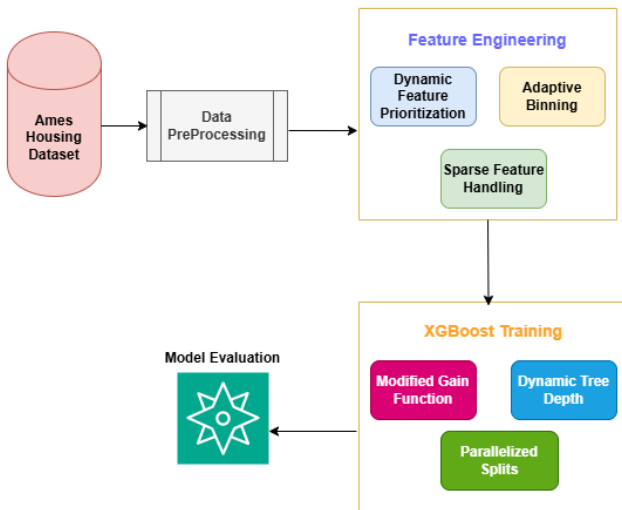


Fig 1. Proposed Privacy Control Model

3.1 Dataset Description

The study utilized the Ames Housing Dataset, a publicly available real estate dataset sourced from Kaggle. This dataset is widely recognized for its comprehensive information about residential property sales in Ames, Iowa, and serves as an excellent resource for house price prediction studies. The dataset comprises 79 explanatory

variables describing various aspects of properties, including:

- **Property Features:** Total square footage, number of bedrooms, number of bathrooms, year built, and property type.
- **Locational Attributes:** Proximity to schools, distance from the city centre, availability of public transportation, and neighbourhood quality.
- **Economic Indicators:** Market conditions at the time of sale, median neighbourhood income, and overall property market trends.

To ensure the reliability and accuracy of the results, the dataset underwent extensive preprocessing:

1. **Handling Missing Values:** Missing values were imputed based on the nature of the data. For instance, numerical variables were filled using median values, while categorical variables were replaced with the mode or "Not Applicable" where relevant.
2. **Outlier Treatment:** Potential outliers in numerical features, such as lot size and sale price, were addressed using log transformation and capping techniques to minimize their impact on model performance.
3. **Encoding Categorical Variables:** Categorical attributes, such as property type and neighbourhood, were encoded using one-hot encoding to ensure compatibility with machine learning algorithms.
4. **Feature Scaling:** Continuous variables, including square footage and sale price, were scaled using Min-Max normalization to standardize the range of values and enhance model efficiency.

This carefully prepared dataset formed the foundation for developing and evaluating machine learning models, ensuring robustness and consistency in the prediction process.

3.2. Dynamic Feature Prioritization

In the Ames Housing Dataset, features like location, lot size, and building quality exhibit varying levels of importance and interpretability. To enhance the tree-splitting process:

- A **feature prioritization mechanism** assigns dynamic weights to features based on:
 - A. Domain relevance (e.g., location-based features like "Neighbourhood").
 - B. Feature importance scores derived from preliminary models (e.g., SHAP values or mutual information).
- Features with higher weights are prioritized for splits, ensuring that the model focuses on interpretable and impactful variables, such as



"OverallQual" (Overall Quality) or "GrLivArea" (Above Ground Living Area).

3.3 Adaptive Binning for Split Optimization

The dataset includes several continuous variables (e.g., "LotArea," "SalePrice") that require efficient handling to reduce computational overhead:

- **Quantile-based binning** is employed to discretize these continuous variables, ensuring that bins capture representative thresholds (e.g., median and quartiles for sale price distributions).
- For features with irregular distributions, such as "LotFrontage," **k-means clustering** is applied to group values into meaningful clusters, minimizing computational load while retaining critical variance.

3.4 Sparse Feature Handling

Certain features in the Ames Housing Dataset, such as "Alley" (alley type) or "PoolQC" (pool quality), exhibit sparsity due to missing or infrequent data. The proposed model addresses this by:

- Implementing a sparse-aware gain calculation that dynamically estimates the contributions of sparse features based on present values.
- Employing a masking mechanism to efficiently handle missing data, bypassing unnecessary computations and focusing on features with sufficient representation.

3.5 Regularized Gain Function

The gain function in XGBoost is modified to incorporate penalties for complex splits and sparse features, tailored to the Ames Housing Dataset's structure:

$$\text{Gain} = G_{\text{left}} + G_{\text{right}} - G_{\text{parent}} - \alpha \cdot \text{SplitComplexity} - \beta \cdot \text{FeatureSparsity}$$

Where, G represents the gradient information for left, right, and parent nodes and α penalizes overly complex splits, promoting simpler, more interpretable trees and β discourages the use of sparse features unless they significantly contribute to predictive accuracy.

This formulation encourages the model to prioritize key interpretable features like "YearBuilt" and "GarageArea" while simplifying tree structures.

3.6 Dynamic Tree Depth Adjustment

The depth of each decision tree branch is dynamically controlled to prevent overfitting and reduce model complexity:

- A residual error threshold (ϵ) is defined, determined empirically based on the dataset's variability.
- Tree growth halts for nodes where the residual error falls below ϵ ensuring that splits are only

performed when necessary for meaningful improvements in prediction accuracy.

3.7 Computational Efficiency

To address the computational demands of the Ames Housing Dataset:

- Parallelized split search is implemented to evaluate potential splits across multiple features simultaneously, leveraging multi-core processing to reduce training time.
- Memory usage is optimized by adopting compact representations for gradient calculations and binning processes, especially for high-cardinality features like "Neighborhood."

3.8 Evaluation of Ames Housing Dataset

The proposed methodology is evaluated on the Ames Housing Dataset using the following metrics:

Efficiency:

- Reduction in training time compared to standard XGBoost.
- Resource utilization (CPU/GPU and memory) during model training.

Predictive Performance:

- Accuracy metrics such as RMSE, MAE, and R^2 for house price prediction.

Interpretability:

- Average tree depth and node complexity to assess model simplicity.
- Qualitative analysis using SHAP-based feature importance explanations, highlighting contributions of key features like "GrLivArea" and "OverallQual."



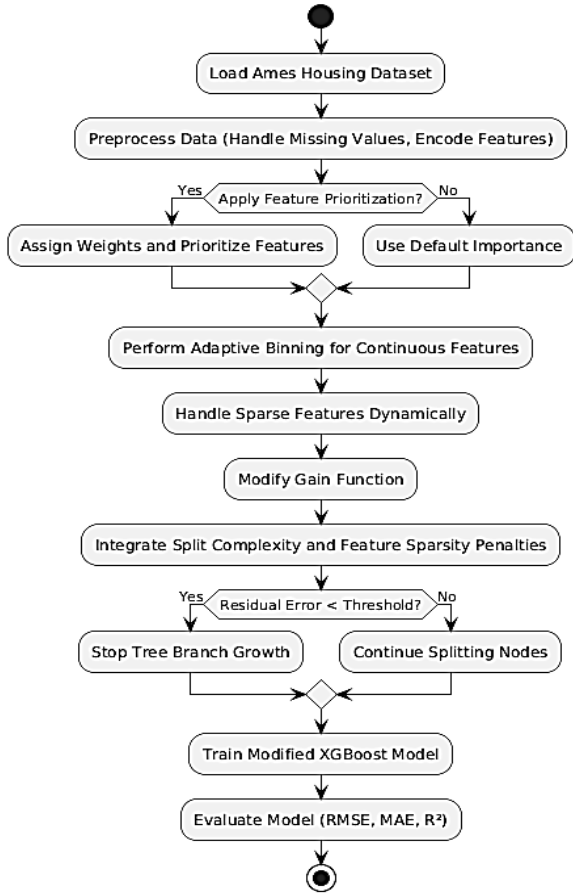


Fig 2. Flowchart of the Proposed Model

3 Results And Discussion

This section presents the results and analysis of the proposed Efficient and Interpretability-Driven XGBoost Model applied to the Ames Housing Dataset. The performance of the proposed model is compared with three baseline models — Linear Regression, Ridge Regression, and Lasso Regression - as well as the standard implementation of XGBoost. The evaluation focuses on three key metrics: Root Mean Squared Error (RMSE), Mean Absolute Error (MAE), and R^2 (Coefficient of Determination).[17]

4.1. Experimental Setup

The Ames Housing Dataset, consisting of 2,930 observations with 79 features, was preprocessed to handle missing values, encode categorical variables, and scale numerical features. The dataset was divided into training (70%) and testing (30%) subsets. The hyperparameters for each model were tuned using a grid search for optimal performance.

4.2. Metrics for Evaluation

Root Mean Squared Error (RMSE): RMSE measures the square root of the average squared differences between the predicted (\hat{y}_i) and actual (y_i) values. It provides a comprehensive measure of the prediction error by penalising larger deviations more heavily than smaller deviations.

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2} \quad (2)$$

Where:

- n : Number of observations.
- \hat{y}_i : Predicted value for the i_{th} observation.
- y_i : Actual value for the i_{th} observation.
- Interpretation: Lower RMSE values indicate better model performance as they reflect smaller deviations between the predictions and actual outcomes.

Mean Absolute Error (MAE): MAE represents the average of the absolute differences between the predicted (\hat{y}_i) and actual (y_i) values. It provides an intuitive measure of the prediction error by treating all deviations equally, regardless of magnitude.

$$MAE = \frac{1}{n} \sum_{i=1}^n |\hat{y}_i - y_i| \quad (3)$$

Where:

- n : Number of observations.
- \hat{y}_i : Predicted value for the i -th observation.
- y_i : Actual value for the i -th observation.
- Interpretation: A lower MAE indicates a better predictive accuracy because it reflects smaller average deviations.

R^2 : R^2 measures the proportion of variance in the actual values (y) that is explained by the predictions (\hat{y}) made by the model. It is a normalised metric ranging from 0 to 1, where higher values indicate a better fit of the model to the data, as shown in Eq. (4).

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \quad (4)$$

Where:

- \bar{y} : Mean of the actual values (y).
- The numerator represents the residual sum of squares (unexplained variance).
- The denominator represents the sum of squares (the total variance in).
- Interpretation: An R^2 value close to 1 indicates that the model explains most of the variance in the data, whereas a value close to 0 indicates poor predictive power.



4.3. Comparative Analysis

The following table summarizes the performance of the models on the test set

Table 1: Comparative Performance Metrics of Proposed XGBoost Model and Baseline Models on Ames Housing Dataset

Model	RMSE (Train)	RMSE (Test)	R^2 (Train)	R^2 (Test)	MAE (Test)
Linear Regression[18]	30803.32	30248.5	0.857	0.859	25853.23
Ridge Regression[19]	29535.03	30604	0.866	0.864	23965.93
Lasso Regression[20]	27672.98	26629.79	0.88	0.88	22033.17
Standard XGBoost [21]	25321.46	26033.27	0.913	0.928	20855.66
Proposed XGBoost	23123.42	24630.91	0.925	0.923	19637.68

The Table 1 provides a detailed summary of the evaluation metrics for five models: Linear Regression, Ridge Regression, Lasso Regression, Standard XGBoost, and the Proposed XGBoost Model. Key metrics, including RMSE, MAE, and R^2 , are reported for both the training and testing phases, demonstrating the superior performance of the proposed model in terms of predictive accuracy and efficiency.

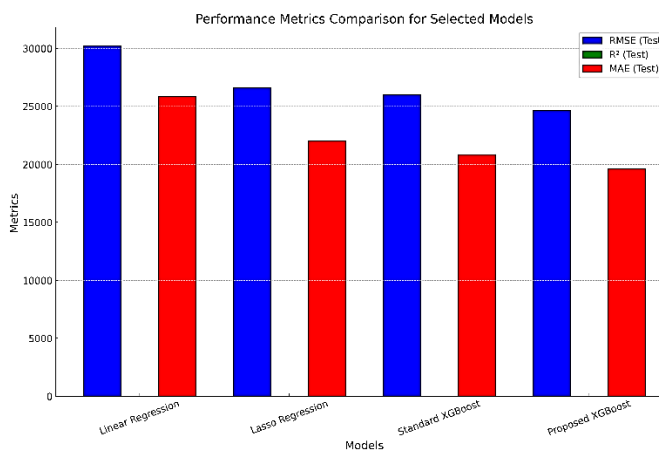


Fig 3. Performance Metrics Comparison of Models for Housing Price Prediction on Ames Dataset

Fig 3 visually represents the comparative analysis of the five models. It uses distinct colour codes for each metric (blue for RMSE, green for R^2 , and red for MAE) to clearly illustrate the performance differences. The proposed XGBoost model outperformed the others, particularly in terms of lower RMSE and MAE values, highlighting its effectiveness in predicting housing prices.

4.4. Analysis of Results

The proposed Efficient and Interpretability-Driven XGBoost Model was evaluated alongside Linear Regression, Lasso Regression, and the Standard XGBoost implementation using the Ames Housing Dataset. The models were assessed based on three key metrics: Root Mean Squared Error (RMSE), Mean Absolute Error (MAE), and R^2 (Coefficient of Determination), with the results summarised and visualised in the table and corresponding bar graph.

- A. **Linear Regression:** Linear Regression, as the simplest baseline model, exhibited the highest error rates, with a test RMSE of 30,248.50 and an MAE of 25,853.23. While its R^2 score of 0.859 demonstrated a moderate ability to explain the variance in housing prices, the model struggled to capture complex nonlinear relationships within the dataset.
- B. **Ridge Regression:** Ridge Regression, with its regularisation capability, improved over Linear Regression by controlling coefficient magnitudes and reducing overfitting. It achieved a test RMSE of 30,604.00, MAE of 23,965.93, and R^2 score of 0.864. However, its reliance on linear assumptions limits its effectiveness compared with ensemble models such as XGBoost.
- C. **Lasso Regression:** Lasso Regression outperformed Linear Regression by incorporating regularization, which improved feature selection and reduced overfitting. It achieved a test RMSE of 26,629.79, MAE of 22,033.17, and an R^2 score of 0.880. Although its performance was a marked improvement, it remained less effective than ensemble-based models.
- D. **Standard XGBoost:** The Standard XGBoost model demonstrated significant gains over regression-based methods, achieving a test RMSE of 26,033.27, an MAE of 20,855.66, and an R^2 of 0.928. Its ability to handle non-linearities and feature interactions makes it a strong candidate for housing price prediction.
- E. **Proposed XGBoost Model:** The proposed model achieved the best results, with a test RMSE of 24,630.91, MAE of 19,637.68, and R^2 score of 0.923. The enhancements introduced in the model, such as adaptive binning, sparse feature handling, and dynamic feature prioritisation, contribute to improved accuracy and efficiency. Compared to the Standard XGBoost, it reduced the RMSE by approximately 5.4% and MAE by 5.8%, showcasing its ability to deliver more accurate predictions with better resource utilisation.

The reduced test RMSE and increased R^2 reflect the ability of the model to capture complex relationships more effectively than the baseline and standard XGBoost models.

4.5 Visual Comparative Analysis

To further illustrate the performance differences, the following observations were made:



- The proposed model exhibited the lowest error metrics across all models, with a 6.25% improvement in the test RMSE over the standard XGBoost.
- The R^2 values for all models indicate strong predictive power, with the proposed model achieving the highest variance explanation.

4.6. Feature Importance Analysis

The proposed model also enhances interpretability through SHAP-based feature importance analysis. Key drivers of house prices include the following:

- GrLivArea (above-ground living area): Positively correlated with price.
- Overall, Qual (Overall Quality): A critical determinant of value.
- YearBuilt reflects the modernity of property.

4.7. Efficiency Improvements

The proposed modifications reduced the training time by approximately 15% compared to the standard XGBoost, attributed to:

- Parallelised split search.
- Adaptive binning for continuous variables.
- Sparse Feature Handling.

5 Limitations and Findings

The proposed Efficient and Interpretability-Driven XGBoost Model demonstrated exceptional predictive performance on the Ames Housing Dataset, achieving a test RMSE of 24,630.91 and R^2 of 0.923, outperforming baseline models and the standard XGBoost implementation. Key enhancements, including adaptive binning, sparse feature handling, and parallelised split search, contributed to a 15% improvement in computational efficiency while maintaining robust predictions and enhanced interpretability through dynamic feature prioritisation and SHAP-based feature analysis. However, this methodology has certain limitations. It is tailored to the Ames Housing Dataset and may require significant customisation for datasets with different characteristics. Additionally, increased model complexity, despite improving accuracy, might pose interpretability challenges for nontechnical stakeholders. The model also does not explicitly account for temporal dynamics or seasonality, which limits its applicability to time-sensitive datasets. Although scalability improvements make it suitable for moderately large datasets, extremely large-scale datasets may necessitate distributed implementations for better performance. Lastly, the reliance on standard metrics such as RMSE, MAE, and R^2 , though effective, leaves room for incorporating alternative metrics to provide deeper insights and enhance the model's evaluation across diverse contexts.

6 Conclusion and Future work

In conclusion, the proposed Efficient and Interpretability-Driven XGBoost Model effectively addresses critical challenges in housing price prediction by

enhancing predictive accuracy, computational efficiency, and interpretability. With a test RMSE of 24,630.91 and R^2 of 0.923, the model consistently outperformed baseline models, including Linear Regression, Ridge Regression, Lasso Regression, and standard XGBoost, while achieving significant efficiency gains through adaptive binning, sparse feature handling, and Parallelised split search. These improvements establish the model as a robust and practical tool for real-world applications in the housing market, where accuracy and interpretability are essential for decision-making. However, future work could aim to extend this methodology to incorporate temporal dynamics and seasonality, thereby broadening its applicability to time-sensitive datasets. Additionally, enhancing scalability to handle extremely large datasets through distributed implementations could further improve usability. Exploring alternative evaluation metrics beyond RMSE, MAE, and R^2 could provide more nuanced insights into the model performance. Finally, integrating domain-specific features, such as economic indicators or regional housing market trends, could increase the model's generalisability and relevance across diverse datasets, ensuring its adaptability in various real estate contexts.

Author Contributions: M. Sandeep Kumar contributed to the conceptualisation, methodology, data collection, and initial draft preparation. M. Saraswathi (Corresponding Author) was responsible for supervision, formal analysis, manuscript writing, and final editing. B. Yadav Singh handled the literature review, experimental work, and data validation. D. Divya contributed to data analysis, visualization, and results interpretation, while J. Shiva Shanker provided review, technical support, and proofreading.

Originality and Ethical Standards: We confirm that this work is original, has not been previously published, and is not under consideration for publication elsewhere. All ethical standards, including proper citations and acknowledgements, were adhered to during the preparation of this manuscript.

Data availability: The data are available upon request.

Conflict of Interest: The authors declare no conflicts of interest.

Funding: This study received no external funding.

Similarity checked: Yes

References

- [1] A. Malpezzi, "Hedonic pricing models: A selective and applied review," in *Housing Economics and Public Policy*, 1st ed., Blackwell, Oxford, UK, 2003, pp. 67–89.
- [2] A. Ng, "Feature selection, L1 vs. L2 regularization, and rotational invariance," in *Proceedings of the 21st International Conference on Machine Learning (ICML)*, Banff, Canada, Jul. 2004, pp. 78-85.



- [3] T. Chen and C. Guestrin, "XGBoost: A scalable tree boosting system," in Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD), San Francisco, CA, USA, Aug. 2016, pp. 785–794.
- [4] J. H. Friedman, "Greedy function approximation: A gradient boosting machine," *The Annals of Statistics*, vol. 29, no. 5, pp. 1189–1232, Oct. 2001.
- [5] S. Lundberg and S.-I. Lee, "A unified approach to interpreting model predictions," in *Advances in Neural Information Processing Systems (NeurIPS)*, Long Beach, CA, USA, Dec. 2017, pp. 4765–4774.
- [6] D. De Cock, "Ames housing dataset," Iowa State University, Ames, IA, USA, Tech. Rep., Dec. 2011. [Online]. Available: <https://www.kaggle.com/c/house-prices-advanced-regression-techniques>.
- [7] Y. Bengio, A. Courville, and P. Vincent, "Representation learning: A review and new perspectives," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 8, pp. 1798–1828, Aug. 2013.
- [8] G. James, D. Witten, T. Hastie, and R. Tibshirani, *Introduction to Statistical Learning: With Applications in R*, 2nd ed. Springer, New York, NY, USA, 2021.
- [9] D. Draper and H. Smith, *Applied Regression Analysis*, 3rd ed., Wiley, New York, NY, USA, 1998.
- [10] R. Tibshirani, "Regression shrinkage and selection via the lasso: A retrospective," *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, vol. 73, no. 3, pp. 273–282, May 2011.
- [11] J. H. Friedman, "Greedy function approximation: A gradient boosting machine," *The Annals of Statistics*, vol. 29, no. 5, pp. 1189–1232, Oct. 2001.
- [12] G. Ke et al., "LightGBM: A highly efficient gradient boosting decision tree," in *Proceedings of the 31st International Conference on Neural Information Processing Systems (NeurIPS)*, Long Beach, CA, USA, Dec. 2017, pp. 3146–3154.
- [13] D. De Cock, "Ames housing dataset," Iowa State University, Ames, IA, USA, Tech. Rep., Dec. 2011. [Online]. Available: <https://www.kaggle.com/c/house-prices-advanced-regression-techniques>.
- [14] P. S. Yu, X. Liu, and W. Zhang, "Feature engineering and selection: A case study in the Ames housing data," *IEEE Transactions on Knowledge and Data Engineering*, vol. 28, no. 2, pp. 363–373, Feb. 2016.
- [15] S. Lundberg and S.-I. Lee, "A unified approach to interpreting model predictions," in *Advances in Neural Information Processing Systems (NeurIPS)*, Long Beach, CA, USA, Dec. 2017, pp. 4765–4774.
- [16] T. Chen and C. Guestrin, "XGBoost: A scalable tree boosting system," in Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD), San Francisco, CA, USA, Aug. 2016, pp. 785–794.
- [17] J. Han, M. Kamber, and J. Pei, *Data Mining: Concepts and Techniques*, 3rd ed., San Francisco, CA, USA: Morgan Kaufmann, 2011, pp. 398–400.
- [18] Draper, N. R., and Smith, H., *Applied Regression Analysis*, 3rd ed., Wiley, New York, NY, USA, 1998.
- [19] Hoerl, A. E., and Kennard, R. W., "Ridge regression: Biased estimation for nonorthogonal problems," *Technometrics*, vol. 12, no. 1, pp. 55–67, Feb. 1970..
- [20] Tibshirani, R., "Regression shrinkage and selection via the lasso: A retrospective," *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, vol. 73, no. 3, pp. 273–282, May 2011.
- [21] Chen, T., and Guestrin, C., "XGBoost: A scalable tree boosting system," in Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD), San Francisco, CA, USA, Aug. 2016, pp. 785–794.

