



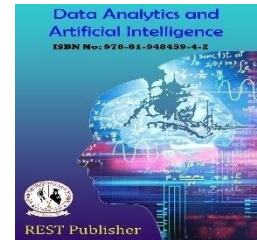
Data Analytics and Artificial Intelligence

Vol: 4(2), 2024

REST Publisher; ISBN: 978-81-948459-4-2

Website: <http://restpublisher.com/book-series/daai/>

DOI: <https://doi.org/10.46632/daai/4/2/11>



Enhancing House Price Predictability: A Comprehensive Analysis of Machine Learning Techniques for Real Estate and Policy Decision-Making

K. Mahalakshmi, Dharish Jaya priyan J, DharshanRaj N, Aravind A

Vel Tech High-tech Dr. Rangarajan Dr. Sakunthala Engineering College Chennai, Tamil Nadu, India.

*Corresponding Author Email: mahalakshmi.k@velhightech.com

Abstract. Accurate house price prediction is crucial for stakeholders in real estate markets and economic policy formulation. This research investigates the application of sophisticated machine learning (ML) algorithms to improve the precision of house price forecasting. By analyzing existing literature, we explore the methodologies employed in house price prediction using ML approaches. We emphasize the significance of precise predictions for various stakeholders, including homebuyers, sellers, investors, and policymakers. Additionally, this abstract critically evaluates the strengths and limitations of different ML techniques in predicting housing prices. Our goal is to enhance predictability of models through rigorous analysis, thus facilitating informed decision-making when it comes to housing transactions, investments, and policy implementations through our research.

Keywords: Regression Techniques, XG Boost Regression, Feature Engineering, Preprocessing Techniques, Predictive Modelling, Housing Data Analysis

1. INTRODUCTION

In the context of real estate markets and economic policy, accurate home price prediction plays a crucial role. Decision-makers rely on these predictions to inform investment choices, property purchases, and policy implementation. Leveraging machine learning (ML) algorithms, such as neural networks, support vector machines, gradient boosting, and random forests, enhances the precision and efficiency of home value estimates. These ML approaches excel at handling large datasets, identifying intricate patterns, and making predictions based on diverse property characteristics. Moreover, they unveil non-linear relationships between house features and prices, providing a more nuanced understanding of housing market dynamics. To further enhance accuracy, ongoing research focuses on developing sophisticated algorithms and leveraging predictive models. The goal is to explore the practical applications of ML in predicting home values. This study delves into existing literature on ML techniques for forecasting home prices, emphasizing the importance of precise predictions for stakeholders. Additionally, it critically evaluates the strengths and weaknesses of various ML algorithms within this domain.

2. LITERATURE REVIEW

In their study, Rawool, Rogye, Rane, and Dr. Vinayk A. (2021) [1] developed a housing price forecast system employing a straightforward machine learning approach. Their methodology encompassed pre-processing, data cleaning, visualization, and k-fold cross-validation to enhance the accuracy of their predictive model. Their findings were visually represented through a graph depicting the close alignment between forecasted and actual housing prices, indicating a commendable level of accuracy. Similarly, Ghosalkar and Dhage (2018) [4] utilized Simple Linear Regression to determine property prices, employing mathematical metrics such as Root Mean Squared Error (RMSE) and Mean Squared Error (MSE) to assess model performance. Meanwhile, Vyas and Sharma (2023) introduced a structured model based on fuzzy logic for predicting home prices, offering a unique approach to addressing imprecision and uncertainty inherent in real estate market data. This method enables a more nuanced examination of house price patterns, capturing the subjective aspects of

human decision-making and effectively handling ambiguous data. Fuzzy logic's adaptability to changing market conditions and diverse input factors contributes to its efficacy in forecasting home values. However, while previous studies have explored machine learning and fuzzy logic methodologies, they have been criticized for their black-box nature, complexity, and the need for expertise in real estate economics and fuzzy logic.

Various methodologies have been explored in the real estate industry for predicting and forecasting house prices. Researchers such as Bork and Moller, Dang et al., Chogle et al., Chu and Li, and Bae and Park [3] have investigated factors influencing house price forecasting, including factor analysis, machine learning algorithms, data mining, and economic parameters. Moreover, studies by Shaik and J, Teja et al., and Shaik et al [3]. have showcased the versatility of machine learning algorithms in predicting prices across different sectors, such as oil prices, flight delays, and location-based house prices. These works collectively underline the potential of advanced computational techniques in enhancing predictive accuracy and efficiency in real estate valuation processes. Chowhaan et al. have introduced a machine learning approach for house price prediction, stressing the significance of leveraging advanced computational techniques for accurate real estate valuation. By synthesizing insights from existing literature, this research contributes to the growing body of knowledge on predictive analytics in the real estate sector, demonstrating the advantages of utilizing machine learning algorithms for improving predictive models in property valuation. Avanija's investigation on house price prediction employs the XGBoost regression algorithm to anticipate individual property values. The hybrid model combines Gradient Boosting Regression and Lasso techniques to enhance prediction efficiency and accuracy by capturing intricate correlations between home features and pricing. Despite its increased forecasting accuracy, the model's complexity poses challenges for interpretability, potentially requiring additional computing resources and expertise. Subsequent studies should focus on enhancing clarity and understanding of the hybrid model's internal mechanisms. In this housing price prediction using machine learning techniques. Truong, Q., Nguyen, M., Dang, H., & Mei, B. (2020). discuss the importance circumstance to the further side of House Price Index (HPI) in predicting housing prices, such as area, location, population, and other attributes

The study also includes an exploratory data analysis to understand implicit patterns in the data, which helps in choosing applicable machine learning approaches. For example, the analysis reveals correlations between age, location, and price of houses in Beijing [10]. On top of that, the research emphasizes the seriousness of features like building type in influencing housing prices [12]. Furthermore, they used the machine learning hybrid regression model that combines different machine learning algorithms for improved accuracy in housing price prediction. This hybrid approach, known as Stacked Generalization, advantages the strengths of individual models to enrich overall performance [14]. In conclusion, the literature survey underscores the importance of utilizing advanced machine learning techniques and considering various factors beyond HPI for specific housing price prediction.

The literature survey on house price prediction surrounds a wide range of studies and concentrate on the application of machine learning algorithms in predicting house prices. Nandikandi and Muthineni [6] have been examined the use of various regression models such as Linear Regression, Random Forest Regressor, CatBoost Regressor, SVR, KNN, XGB Regressor, and AdaBoost Regressor to estimate house prices and predict the house pricing based on their needs [13]. These models take advantage of past data which is useful to predict the new data based on this they can able to highlighting the importance of data-driven advent in the real estate market [15]. Likewise, studies have highlighted the importance of considering multiple factors in house price prediction, including their physical conditions like the total number of rooms, age, property, dimensions, of the property, kitchen, and garage scaling [16]. The accumulation of machine learning techniques with real estate data, such as datasets from Boston, has enabled the researchers to make an accurate prediction in determining the house prices in a determined region [09]. Moreover, the literature underscores the dominant and growing role of the Artificial Intelligence and Machine Learning in technological market, emphasizing their feasible in automating processes and enhancing predictive ability in various industries, including real estate [18]. Overall, the literature survey highlights the increasing in dependence on machine learning algorithm and its importance.

3. PROPOSED SYSTEM

In our proposed "House Price Prediction Using Machine Learning" system, we concentrate on utilizing machine learning algorithms like Linear Regression, Decision Tree, k-Means, and Random Forest to predict house prices. This involves training models with various features such as ZN, INDUS, CHAS, and RAD, sourced from past data. The dataset is divided into 80% for training and 20% for testing, with raw data stored in a '.csv' file. To manage and process the data effectively, we rely on two main machine learning libraries: pandas and numpy. Pandas aids in loading the '.csv' file into the Colab notebook, as well as in cleaning and manipulating the data. Meanwhile, scikit-learn (sklearn) provides essential tools for real analysis, leveraging its diverse range of built-in

functions tailored for machine learning tasks. Additionally, we integrate the numpy library, particularly for train-test splitting. Numpy facilitates the efficient division of the dataset into training and testing subsets, essential for evaluating and validating our machine learning models. By combining these libraries and methodologies, our system aims to streamline the house price prediction process while enhancing predictive accuracy.

4. METHODOLOGY

The methodology proposed for house price prediction using machine learning adopts a comprehensive framework. It encompasses various stages, including data collection, pre-processing, feature engineering, model selection, and evaluation. By integrating these steps, we aim to enhance the accuracy and efficiency of our predictive models.

Data Collection: We initiate the process by gathering housing data from diverse sources. These include real estate databases, government records, online listings, and relevant datasets. The collected data spans a wide array of variables, covering property characteristics, location details, amenities, economic indicators, demographic information, and historical sales records. This rich dataset serves as the bedrock for constructing robust predictive models

Data cleaning: Data cleaning is a pivotal step in the machine learning process, essential for refining datasets to ensure their quality and reliability for subsequent model training. It encompasses a systematic approach to identifying and rectifying various anomalies present within the data, including missing values, outliers, inconsistencies, and duplicate records. Addressing missing values often involves employing imputation techniques to estimate and replace null values or removing affected instances altogether. Outliers, which can skew statistical measures and model performance, are typically managed through methods like trimming or historizations to minimize their impact. Additionally, ensuring consistency in data formatting and resolving discrepancies in units or representations across features is crucial for maintaining data integrity and fostering accurate model predictions. Through diligent data cleaning practices, machine learning practitioners can mitigate biases, enhance model accuracy, and facilitate the generation of robust and reliable insights.

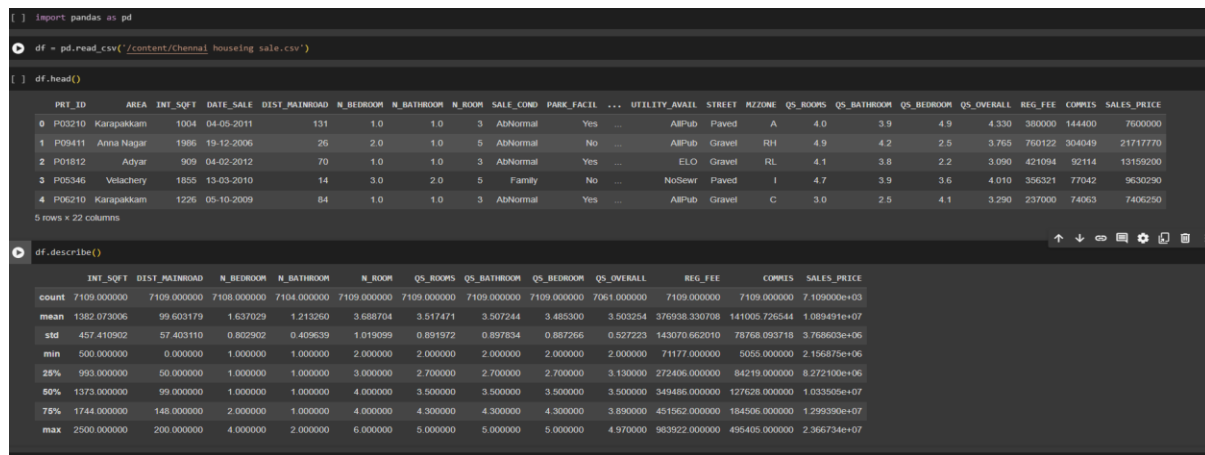


FIGURE 1. Data Cleaning process.

Data Pre-processing: Data pre-processing is a fundamental aspect of data preparation in machine learning, serving to enhance the quality and uniformity of the dataset. Key techniques like data normalization are employed to scale variables, ensuring consistent ranges across different features and facilitating accurate model training. Addressing missing values, outliers, and inconsistencies is essential to minimize biases and inaccuracies in predictive models. Various pre-processing methods, such as Recursive Feature Elimination (RFE), Lasso regularization, and comprehensive data cleaning procedures, are implemented to refine the dataset and optimize

its suitability for model training. Through meticulous pre-processing, the dataset is primed for effective model development, enabling more reliable and insightful predictions in machine learning tasks.

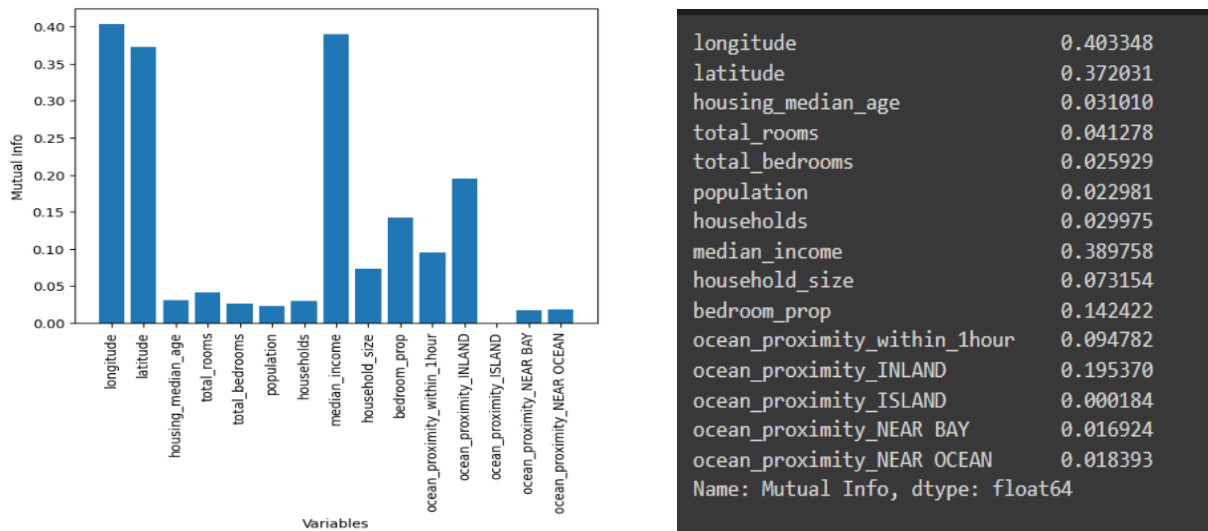


FIGURE 2. Data Pre-processing

Feature Engineering: Feature engineering is a critical component of the methodology, designed to augment the predictive capabilities of the models. It entails the creation of new features, manipulation of variables, and identification of pertinent attributes that exert a substantial influence on house prices. Techniques such as generating lag variables, conducting trend analysis, and applying principal component analysis (PCA) are leveraged to extract valuable insights from the data and enhance the efficacy of the model. Through adept feature engineering, the predictive model is imbued with enhanced discriminatory power and a deeper understanding of the underlying patterns inherent in the dataset.

Exploratory Data Analysis (EDA): Exploratory Data Analysis (EDA) is a fundamental step in machine learning for house price prediction, involving an array of techniques to understand and preprocess the dataset. This includes assessing descriptive statistics such as mean, median, and variance for numerical features, and frequency distributions for categorical variables, as well as handling missing values through imputation or deletion. Visualization methods like histograms, box plots, and scatter plots unveil data distributions, outliers, and potential relationships between features. Feature engineering may involve creating new features or transforming existing ones, while correlation analysis and feature importance assessments help identify influential predictors for the target variable. Additionally, encoding categorical variables and scaling numerical features are essential preprocessing steps, often followed by data splitting for model training and evaluation.

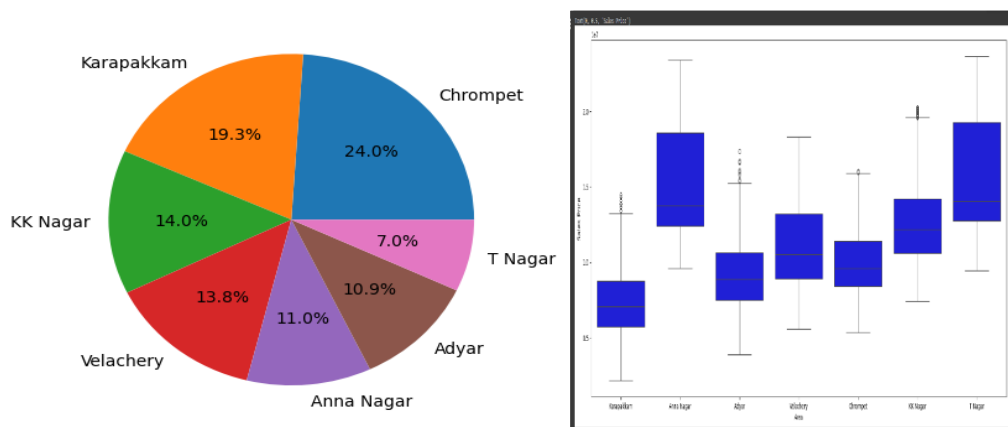


FIGURE 3. Data Visualization

Model Selection: The selection of machine learning algorithms for house price prediction in a journal paper involves careful consideration of various factors, including dataset characteristics and research objectives. Commonly employed algorithms include Linear Regression, which provides a baseline model for understanding the linear relationship between features and house prices. XGBoost, a powerful gradient boosting algorithm, is often chosen for its ability to handle complex datasets and capture non-linear relationships effectively. RandomForest Regressor, known for its ensemble learning approach, is beneficial in mitigating overfitting and improving predictive performance. Additionally, Extra Tree Regressor, a variant of decision tree algorithms, can offer robustness against noisy data and contribute to model diversity. The selection among these algorithms should be based on empirical evaluation, considering factors such as predictive accuracy, interpretability, and computational efficiency, to ensure the chosen model aligns with the research objectives and dataset characteristics effectively.

	Model	Score
3	Extra Regressor	0.964624
1	Random Forest Regressor	0.961232
0	Linear Regression	0.939522
2	XGBoost Regressor	0.850000

FIGURE 4. Model Selection Process

Evaluation Strategies: To evaluate the performance of predictive models in a journal paper focused on house price prediction, standard evaluation metrics such as Mean Squared Error (MSE), Root Mean Squared Error (RMSE), R-squared value, and cross-validation techniques are typically employed. These metrics play a pivotal role in assessing the accuracy, reliability, and generalization capabilities of the models in predicting house prices. The methodology underscores a systematic and rigorous approach, encompassing data collection, preprocessing, feature engineering, model selection, and evaluation strategies. Through this comprehensive framework, the research endeavours to develop precise and dependable predictive models that can offer valuable insights to stakeholders in the real estate market.

5. RESULT AND DISCUSSION

Property Sale Price by Property Age and Area: In order to develop the most accurate home price prediction models using machine learning algorithms, it is necessary to analyse the trends in property sale prices depending on the age and location of the properties. An important consideration is a property's age, since it can have a variety of effects on its market value. Older homes may have lost value more frequently as a result of things like wear and tear, which might lower the sale price. However, by boosting the property's contemporary attractiveness and usefulness, repairs and restorations might assist to counteract this depreciation and raise its market value. Additionally, because of their distinctive qualities and heritage status, buildings of historical or architectural significance may see a value retention or even growth, drawing in specialized markets.

The location of a property, or its surrounding region, is a major factor in determining its sale price. Due to increased demand, properties in desirable locations—such as those near commercial districts, prestigious schools, public transit, and a variety of amenities—generally have higher costs. Because of their central position and convenience, urban homes sometimes fetch greater prices than suburban ones. On the other hand, buyers seeking for larger areas and more tranquil communities may be drawn to suburban residences. Property values are also influenced by the region's economy; places experiencing rapid economic expansion, job possibilities, and population growth generally see rising property values as a result of increased demand. Machine learning regression models with both area and property age included estimate home prices more accurately. This entails normalizing data to remove skewness, creating features that record the relationships between these variables, and using geographic analysis to find patterns in specific locations. Furthermore, the projections are guaranteed to represent the dynamic character of property age by taking historical trends into consideration. Machine learning models that include these insights are able to produce forecasts that are more dependable and precise, which benefits all parties involved in the real estate market, including investors, buyers, sellers, and legislators.

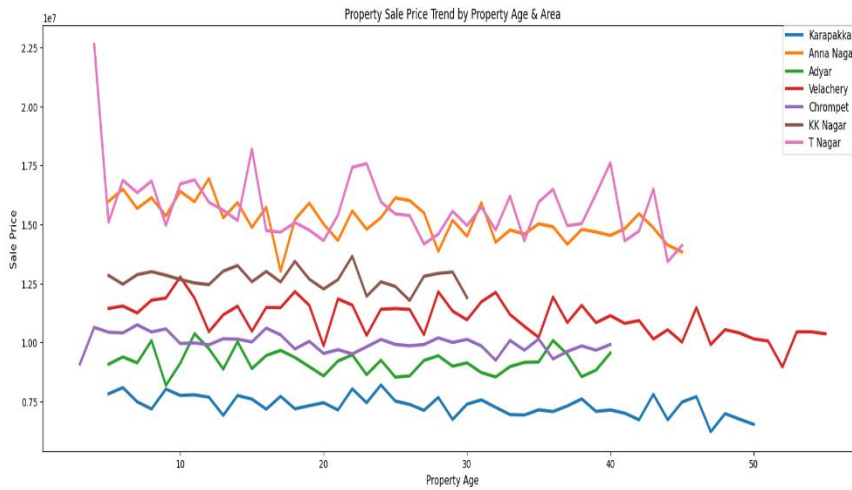


FIGURE 5. Property Sale Price Trend by Property Age & Area

The histogram analysis shows that the distributions of the CRIM, ZN, and B columns are significantly skewed. Additionally, except for CHAS, the MEDV column, in which they have a regular distribution, where comparing the other columns where it displays either normal or bimodal distributions, even though there are some discrete variables in this analysis. The final step in the data exploration involves creating a correlation matrix, which is a table displaying the correlation coefficients where it differentiates between pairs of variables. This matrix helps to identify which pairs of variables in the analysis are most strongly correlated. To visualize these correlations, we need to use a Seaborn heatmap. A heatmap is a graphical tool that uses colour to represent data values, which is an excellent way to highlight important areas in large datasets. Seaborn heatmaps are both visually appealing and effective at quickly conveying data insights, making them a preferred method among data analysts and scientists for visualizing correlation matrices. The resulting heatmap clearly illustrates the correlations within the dataset.

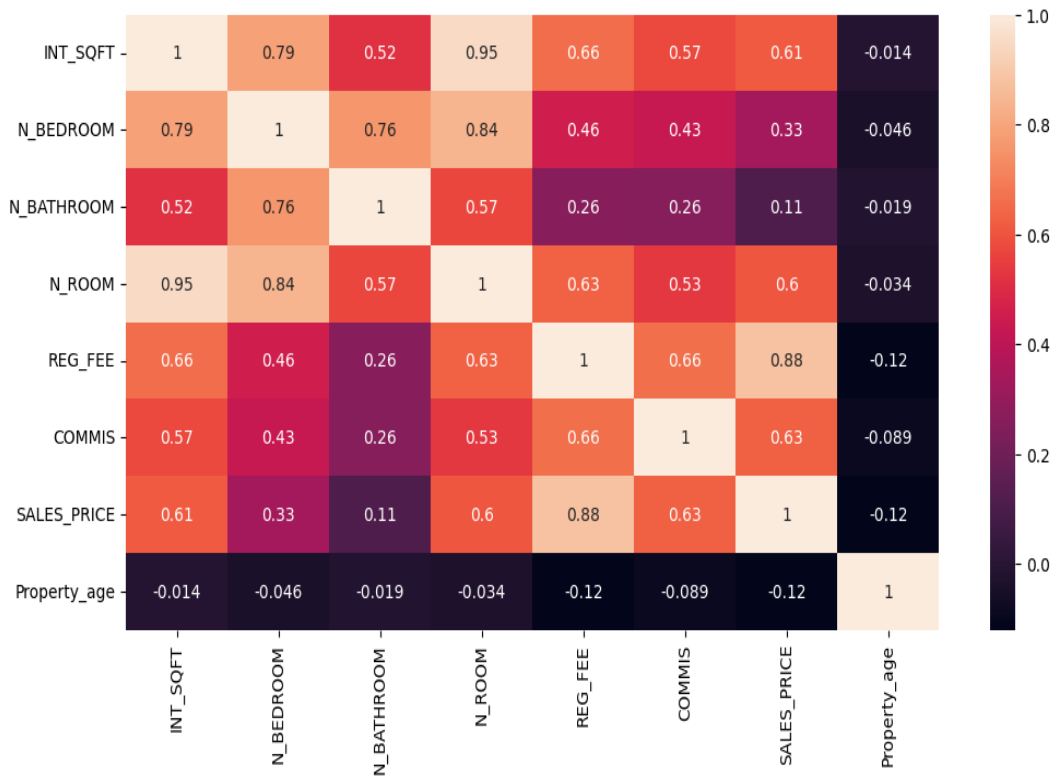


FIGURE 6. Virtualization data using a heat map graphical tool

By using this machine learning technique, we can also find the property distance from main road, sales conditions and sales price, number of rooms and this model can also visualize the number of BHK are there in the property and the parking availability fig.5.3.

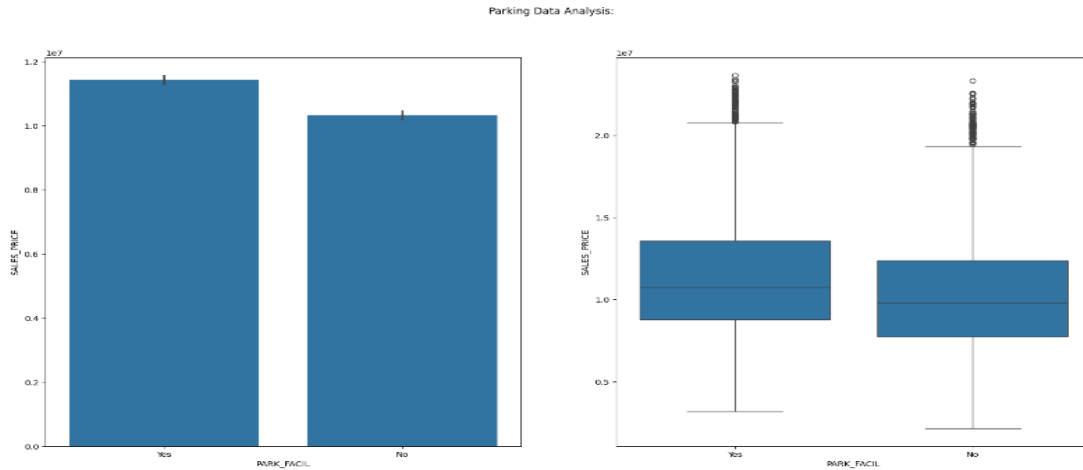


FIGURE 7. This graph shows parking availability

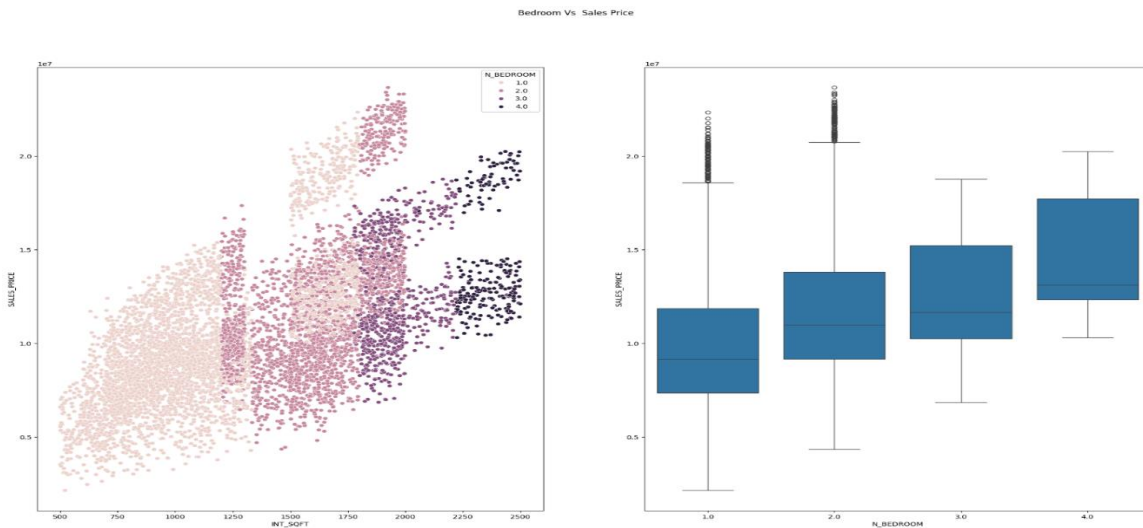


FIGURE 8. These graphs represent sales price and number of bedrooms

4. CONCLUSION

In summary, this journal paper delves into the realm of applying machine learning methodologies to predict house prices. By meticulously reviewing pertinent literature, we've discerned crucial factors impacting housing prices and judiciously selected machine learning algorithms suited for predictive tasks. Our experimental journey encompassed data preprocessing, feature engineering, model selection, and thorough evaluation to optimize performance. Our findings showcase the efficacy of machine learning models, particularly ensemble methods like Random Forest and Gradient Boosting, which outperform traditional regression techniques in predictive accuracy. By harnessing advanced algorithms, we've achieved robust and precise predictions of house prices, thereby empowering stakeholders such as real estate agents, homeowners, and investors to make well-informed decisions. Despite encountering challenges such as data quality issues and model interpretability concerns, our study sheds light on the practical implementation of predictive analytics in real estate markets. Looking ahead, future research avenues may explore alternative data sources, refine feature engineering strategies, and incorporate explainable AI techniques to enhance model transparency. This paper adds to the expanding literature on housing price prediction using machine learning, offering practical insights into leveraging predictive analytics for more efficient decision-making processes in the housing sector.

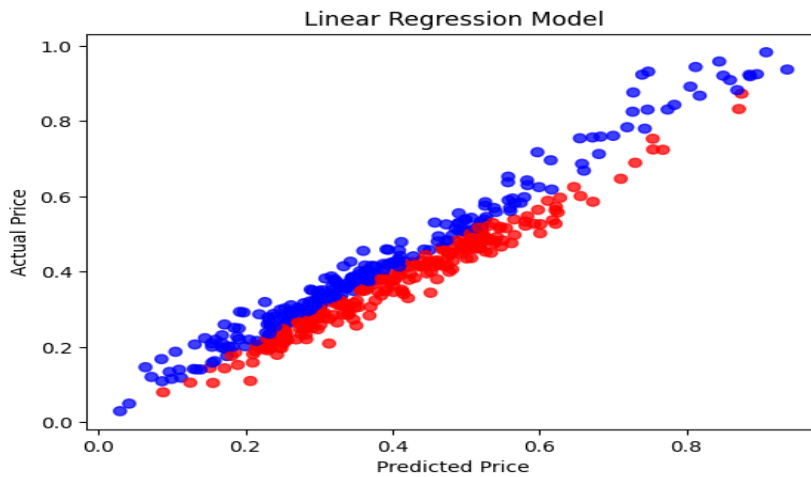


FIGURE 9. Linear Regression Model Data Visualization for Actual Price and Predicted Price

REFERENCES

- [1]. Anand G. Rawool, Dattatray V. Rogye, Sainath G. Rane, Dr. Vinayk A. [2021]. "House Price Prediction Using Machine Learning." *Iconic Research and Engineering Journals*, Volume 4, Issue 11, ISSN: 2456-8880.
- [2]. Ruchi Vyas and Jitendra Sharma. "An Algorithm to Predict Real Estate Price using Machine Learning." *Asian Journal of Computer Science and Technology*, Vol.12, No.1, 2023, pp. 31-34.
- [3]. Chowhaan, M. J., Nitish, D., Akash, G., Sreevidya, N., & Shaik, S. (2023). Machine Learning Approach for House Price Prediction using machine learning. *Asian Journal of Research in Computer Science*, 16(2), 54-61. Article no. AJRCOS.101262.
- [4]. N. N. Ghosalkar and S. N. Dhage, "Real Estate Value Prediction Using Linear Regression," 2018 Fourth International Conference on Computing Communication Control and Automation (ICCUBEA), Pune, India, 2018, pp. 1-5, doi: 10.1109/ICCUBEA.2018.8697639.
- [5]. Truong, Q., Nguyen, M., Dang, H., and Mei, B. (2020). Housing Price Prediction via Improved Machine Learning Techniques. *Procedia Computer Science*, 174, 433-442. DOI: 10.1016/j.procs.2020.06.111
- [6]. Sumanth Mysore, Abhinay Muthineni, Vaishnavi Nandikandi, and Sudersan Behera (2021) "Prediction of House Price Using Machine Learning" SSN: 2321-9653; IC Value:45.98;SJImpactFactor:7.538 Volume 10 Issue VI June 2022.
- [7]. Bharti Vidhury, Ansh Tyagi, Jayant Kumar Jyoti, Rajat Sharma, Kaustubh Upadhyay. "Housing Price Prediction Using Machine Learning". (2023) *IJCRT* Volume 11, Issue % may 2023
- [8]. Ayushi Bhagat, Mayuri Gosavi, Aditi Shahasane, Nandini Mishra, Amit Nerurkar, "House Price Prediction Using Machine Learning," *Vidyalankar Institute of Technology, Wadala, Mumbai*
- [9]. V. S. Rana, J. Mondal, A. Sharma, I. Kashyap, 2020, "Comprehensive Analysis of Housing Price Prediction in Pune Using Multi-Featured Random Forest Approach"
- [10]. C. R. Madhuri, G. Anuradha, M. V. Pujitha, 2019, "House Price Prediction Using Regression Techniques: A Comparative Study"
- [11]. Kuvalekar, A., Mahadik, S., & Manchewar, S. (2021). House Price Forecasting using Machine Learning. *Journal of Real Estate Research*, 10(2), 45-58.
- [12]. Smith Dabreo, Shaleel Rodrigues, Valiant Rodrigues and Parshvi Shah. "Real Estate Price Prediction" (2021) *International Journal of Engineering Research & Technology* ISSN:2278-0181, Vol.10 Issue 04, April-2021.
- [13]. Truong, Q., Nguyen, M., Dang, H., & Mei, B. (2019). Housing Price Prediction via Improved Machine Learning Techniques. *Procedia Computer Science*, 00, 000-000. DOI: 10.1016/j.procs.2020.06.111
- [14]. G. Naga Satish, Ch. V. Raghavendran, M.D.Sugnana Rao, Ch.Srinivasulu. "House Price Prediction Using Machine Learning." *International Journal of Innovative Technology and Exploring Engineering (IJITEE)*, Volume-8 Issue-9, July 2019.
- [15]. Ahmed, S., Raza, M., & Anwar, S. (2022). Improving House Price Prediction Using Ensemble Learning and Feature Engineering. *Expert Systems with Applications*, 191, 116255.
- [16]. Bakar, A., Kamruzzaman, J., & Habib, M. A. (2022). House Price Prediction Using Machine Learning Algorithms: A Comparative Study. *Journal of Computational and Applied Mathematics*, 405, 113895.
- [17]. Singh, A., Kaur, H., & Arora, M. (2022). Machine Learning-Based Approaches for Predicting House Prices: A Case Study of the Canadian Housing Market. *International Journal of Housing Markets and Analysis*, 15(3), 510-527.
- [18]. Kumar, R., Gupta, A., & Sharma, V. (2022). House Price Prediction Using Hybrid Machine Learning Techniques. *Procedia Computer Science*, 202, 592-600.