

Predicting House Prices in Emerging Markets: A Data-Driven Approach to Urban Growth in India

¹Lilly Florence, ²Charu Prabha P, ³Arul Maria Agnes, ⁴Anandha Bhairavi, ⁵Geo

¹*Professor, Adhiyamaan College of Engineering*

^{2,3,4,5}*UG Student, Adhiyamaan College of Engineering*

Abstract-The paper presents a data-driven approach to predicting house prices in emerging markets, with a focus on urban growth in India. Using a blend of historical housing data, socioeconomic factors, and machine learning models, we aim to identify key predictors of housing prices. The study explores various algorithms and evaluates their accuracy in forecasting price trends. Our findings offer insights into market dynamics and urbanization's role in price fluctuations. The results highlight the potential of predictive models in assisting stakeholders, from policymakers to investors, in making informed decisions.

I. INTRODUCTION

The real estate sector, particularly housing markets, plays a crucial role in the economic development of nations. Accurate forecasting of house prices is essential for a variety of stakeholders, including homeowners, real estate investors, financial institutions, and policymakers. In recent years, the rapid urbanization of emerging markets, such as India, has led to significant fluctuations in property values, making it increasingly challenging to predict housing prices with traditional methods. A deeper understanding of the factors driving these changes, along with advancements in predictive modeling techniques, offers new avenues for addressing this issue.

In India, the housing market is influenced by a multitude of factors, including population growth, migration to urban areas, changing income levels, infrastructure development, and government policies. These dynamics are further complicated by regional disparities, as property values can vary significantly between major metropolitan areas and smaller cities. Additionally, the availability of housing data, both structured and unstructured, presents an opportunity to leverage data science techniques to model and predict future trends.

This research aims to explore the use of various machine learning models to predict house prices in India's urban centers. By utilizing a combination of historical sales data, socioeconomic indicators, and

geographical features, we seek to identify the key drivers of price changes and build models that can provide accurate predictions. The outcomes of this study can contribute to more informed decision-making for buyers, sellers, investors, and policymakers alike, offering insights into future market trends and price movements.

II. LITERATURE SURVEY

[1] Rosen, S, The hedonic pricing model introduced by Rosen laid the foundation for understanding how various attributes of a property, such as location, size, and amenities, contribute to its overall price. This model continues to be a fundamental tool in real estate price estimation and is widely applied in housing market analysis. [2] Mayo, S. This study focuses on the determinants of housing prices in developing nations, particularly examining income levels, supply constraints, and urban infrastructure as key influencers of price variations. [3] Worzala, E., Lenk, M., Silva, A. One of the pioneering works exploring the application of artificial neural networks (ANNs) for predicting house prices, this study demonstrated how machine learning techniques could outperform traditional regression methods in real estate forecasting. [4] Antipov, E., Pokryshevskaya, E. This paper compares various machine learning techniques, such as decision trees and support vector machines, for their effectiveness in predicting real estate prices. It concludes that ensemble methods, such as random forests, generally offer higher predictive accuracy. [5] Goodman, A. Goodman examines the influence of income levels, employment rates, and demographic changes on housing prices, providing a socioeconomic framework for understanding market fluctuations.

[6] Bhan, G. This paper investigates the role of urban infrastructure, such as transportation networks, schools, and healthcare facilities, in determining housing prices across Indian cities. It highlights the premium placed on properties located near key

infrastructure. [7] Glaeser, E. Glaeser's work focuses on how regional disparities in income, infrastructure, and urban growth contribute to significant price differences in global housing markets. The paper's insights are applicable to emerging markets like India, where regional disparities are stark. [8] Bourassa, S., Hoesli, M. This study evaluates the effectiveness of traditional regression models for predicting house prices and contrasts them with newer machine learning techniques, noting the limitations of linear models in capturing nonlinear market trends. [9] Zhang, Y., Liu, P. This paper explores the impact of rapid urbanization in emerging markets, particularly in Asia, on real estate prices. It highlights how migration patterns and government policies significantly influence price trends. [10] Khashman, A. This study compares multiple machine learning approaches, such as support vector machines (SVMs), decision trees, and neural networks, in predicting house prices, and concludes that ensemble learning techniques offer superior performance. [11] Ball, M. Ball discusses how government regulations, tax incentives, and housing policies influence real estate prices, particularly in emerging economies. This is particularly relevant to India's housing market, where policy shifts significantly impact prices.

[12] Nath, S., Ray, R. This paper offers a comprehensive analysis of the Indian housing market, focusing on the role of economic reforms, financialization of real estate, and foreign direct investment in shaping housing demand and prices. [13] Fuerst, F., McAllister, P. This study investigates the use of big data and predictive analytics in real estate price forecasting, showcasing how large datasets on property transactions, economic indicators, and social trends can improve predictive accuracy. [14] Kim, S., Lee, Y. This paper explores the application of deep learning techniques, such as recurrent neural networks (RNNs) and convolutional neural networks (CNNs), for predicting house prices. The study finds that deep learning models outperform traditional models, particularly when dealing with large datasets. [15] Mohanty, R., Chakravarty, S. This paper examines housing affordability in India's urban centers, focusing on how income inequality, inflation, and government policies drive housing price dynamics. The study also discusses the social implications of unaffordable housing in rapidly urbanizing cities.

III. EXISTING SYSTEM

In the real estate market, predicting house prices has traditionally been done using manual methods such as comparative market analysis, where professionals analyze the prices of similar properties in a neighborhood. This approach, while still prevalent, has significant limitations due to its reliance on human judgment and the subjective selection of comparable properties. Additionally, this manual system is time-consuming and cannot efficiently handle large datasets or rapidly changing market conditions.

With the advancement of technology, more automated systems have emerged, leveraging statistical models and machine learning techniques to predict house prices. Regression models, particularly linear regression, have been commonly used to model the relationship between house prices and features such as square footage, number of bedrooms, bathrooms, and location. However, traditional regression models often struggle with capturing complex relationships between features, such as nonlinear interactions and dependencies, limiting their accuracy.

In recent years, machine learning systems have started replacing traditional methods. Artificial Neural Networks (ANNs), Random Forests, and Support Vector Machines (SVMs) are frequently employed to model the nonlinear nature of housing markets. These models can process a large volume of data and handle the complex interdependencies between features like location, amenities, and market conditions. For example, Random Forests and Gradient Boosting models offer greater predictive power by learning patterns from the data through decision trees, while neural networks capture even more complex relationships using layers of interconnected nodes.

Additionally, ensemble methods that combine multiple models are increasingly used to improve prediction accuracy. These systems leverage the strengths of various algorithms, improving the robustness of price predictions by reducing the risk of overfitting and underfitting. Furthermore, some systems integrate big data analytics, using real-time market trends, economic indicators, and property-specific data (e.g., satellite imagery or urban development plans) to improve predictive accuracy.

Despite these advances, current systems still face challenges. While machine learning models have improved accuracy, they often require large datasets and extensive feature engineering. Moreover, many of

these systems are unable to account for sudden market shocks or non-quantitative factors such as buyer sentiment or government policy changes. Existing systems are still evolving to integrate dynamic variables, such as market volatility and changing consumer preferences, which are difficult to capture using purely statistical methods.

The existing systems for house price prediction have evolved from manual, comparative methods to sophisticated, data-driven models powered by machine learning algorithms.

IV. PROPOSED SYSTEM

The proposed system aims to address the limitations of existing house price prediction models by combining advanced machine learning algorithms with enhanced data processing techniques. This system will utilize a hybrid approach that incorporates both structured and unstructured data, allowing for more accurate and dynamic price predictions. By leveraging the power of machine learning, the proposed system will be able to model the complex relationships between housing features and prices, while also accounting for real-time market fluctuations, economic indicators, and evolving buyer preferences.

The system will start by collecting data from multiple sources, including real estate listings, historical transaction data, economic indicators (such as interest rates and inflation), and neighborhood-specific data (such as crime rates, school quality, and proximity to amenities). Additionally, unstructured data such as social media trends, reviews, and buyer sentiment will be incorporated using natural language processing (NLP) techniques.

Preprocessing will involve cleaning and transforming the dataset, handling missing or inconsistent values, encoding categorical variables, and scaling numerical features. Feature engineering will also be employed to create new attributes, such as proximity to city centers, public transportation, and environmental factors like air quality, which can have significant impacts on house prices.

The system will use advanced feature selection techniques like recursive feature elimination (RFE) and mutual information to identify the most relevant variables influencing house prices. These techniques will help filter out irrelevant or redundant features, improving the efficiency and accuracy of the model.

Furthermore, the proposed system will employ domain-specific feature engineering to capture intricate aspects of the housing market. For instance, new features like the rate of property appreciation in a particular neighborhood or the frequency of nearby development projects will be included to reflect long-term price trends.

The core of the proposed system will be built using a combination of ensemble machine learning models, such as Gradient Boosting Machines (GBM), XGBoost, and Random Forests, along with Deep Learning models like Artificial Neural Networks (ANNs). The ensemble models will improve prediction accuracy by aggregating the predictions of multiple weak learners, effectively capturing both linear and nonlinear patterns in the data.

In addition to these models, the system will incorporate time series forecasting methods, such as Long Short-Term Memory (LSTM) networks, to predict future trends in house prices based on historical data and market fluctuations. LSTM models are particularly suited for capturing time-dependent trends and predicting long-term price movements based on past patterns.

A key innovation of the proposed system is its ability to dynamically integrate real-time market data. Economic indicators like interest rates, inflation, and government policies will be regularly updated and fed into the system to ensure that the model reflects the latest market conditions. Additionally, data on real estate trends, buyer sentiment, and social factors will be updated in real-time using web scraping and API integration from real estate platforms and social media.

This dynamic integration ensures that the model can adapt to sudden market changes, such as economic crises, policy changes, or shifts in buyer behavior, providing more accurate and up-to-date price predictions.

The proposed system will employ multiple evaluation metrics to assess model performance, including Mean Absolute Error (MAE), Root Mean Squared Error (RMSE), and R-squared (R^2). Cross-validation techniques like K-fold validation will be used to avoid overfitting and ensure that the model generalizes well to unseen data.

Hyperparameter optimization will be performed using techniques like Grid Search and Random Search to

fine-tune model parameters for the best performance. In addition, the system will use model interpretability techniques, such as SHAP (Shapley Additive Explanations), to provide insights into how individual features impact house price predictions.

The system will feature a user-friendly interface that allows real estate professionals, buyers, and sellers to interact with the model. Users will be able to input house characteristics (e.g., area, number of bedrooms, location, amenities) and receive instant price predictions. The system will also offer a comparative analysis feature, allowing users to compare their property with similar listings in the market.

The proposed system will be deployed as a cloud-based platform, ensuring scalability and accessibility for users across different regions. The system's backend will be designed to handle large datasets and multiple concurrent requests, making it suitable for use by real estate agencies, investors, and individual buyers.

By using a combination of ensemble learning, deep learning, and time series forecasting, the proposed system will provide more accurate predictions than traditional regression models. The integration of real-time market data ensures that the model remains up-to-date and responsive to market changes. The system combines both structured and unstructured data (such as buyer sentiment and market trends), offering a more holistic view of the factors influencing house prices. The use of SHAP values will provide transparency into how different features influence the predicted price, offering valuable insights for users. In summary, the proposed system offers a powerful, data-driven approach to predicting house prices, combining advanced machine learning techniques with real-time data integration to provide accurate, dynamic, and interpretable predictions.

V. DATASET DESCRIPTION

id	price	area	bedrooms	bathrooms	stories	material	guestroom	basement	hotwaterhe	aircondition	parking	prefere	furnishingstatus	
1	13300000	7420	4	2	3	yes	no	no	no	yes	no	2	yes	furnished
2	12200000	8960	4	4	4	yes	no	no	no	yes	no	2	no	furnished
3	12200000	9960	3	2	2	yes	no	no	no	no	no	2	yes	semi furnished
4	11100000	7960	4	2	2	yes	no	no	no	yes	no	2	no	furnished
5	11450000	7420	4	1	2	yes	yes	no	no	yes	no	2	no	furnished
6	10800000	7560	3	3	1	yes	no	no	no	no	no	2	yes	semi furnished
7	10100000	8580	4	3	4	yes	no	no	no	yes	no	2	yes	semi furnished
8	10100000	16200	5	3	2	yes	no	no	no	no	no	0	no	unfurnished
9	9870000	8160	4	1	2	yes	yes	no	no	yes	no	2	yes	furnished
10	9800000	5750	3	2	4	yes	yes	no	no	yes	no	1	yes	unfurnished
11	9800000	13260	3	1	2	yes	no	no	no	no	no	2	no	furnished
12	9681000	6500	4	3	2	yes	yes	yes	no	no	no	2	no	semi furnished
13	9120000	6150	4	1	2	yes	no	no	no	yes	no	2	no	semi furnished
14	9040000	3500	4	2	2	yes	no	no	no	yes	no	2	no	furnished
15	9240000	7800	3	2	2	yes	no	no	no	no	no	0	yes	semi furnished
16	9100000	6050	4	1	2	yes	no	yes	no	no	no	2	no	semi furnished
17	9100000	4600	4	2	2	yes	yes	yes	no	yes	no	1	yes	unfurnished
18	9100000	8600	3	2	4	yes	no	no	no	no	no	2	no	furnished
19	8900000	8600	3	2	2	yes	yes	no	no	yes	no	2	no	furnished
20	8890000	4600	3	2	2	yes	yes	no	no	yes	no	2	no	furnished
21	8850000	6420	3	2	2	yes	no	no	no	no	no	1	no	semi furnished
22	8750000	4120	3	1	2	yes	no	yes	yes	no	no	2	no	semi furnished
23	8660000	7151	3	2	1	yes	yes	yes	no	yes	no	2	no	unfurnished
24	8640000	8050	3	1	1	yes	yes	yes	no	no	no	1	no	furnished
25	8640000	4540	3	2	2	yes	yes	yes	no	yes	no	1	no	furnished
26	8570000	8860	3	2	2	yes	no	no	no	no	no	2	no	furnished
27	8540000	6540	4	2	2	yes	yes	yes	no	yes	no	2	yes	furnished
28	8460000	6050	3	2	4	yes	yes	no	no	yes	no	2	no	semi furnished
29	8450000	8875	3	1	1	yes	no	no	no	no	no	1	no	semi furnished

Figure 1: Datasets

The dataset provided consists of multiple features that are significant determinants of house prices, reflecting various characteristics of the properties. The primary

variable of interest is price, which serves as the target for prediction. The dataset includes area, measured in square feet, a key numerical factor as larger houses tend to command higher prices. Bedrooms, bathrooms, and stories represent the internal structure of the house, with more of these features often correlating with higher value. The dataset also includes parking spaces, where additional spaces likely add to the convenience and thus the price of the property.

Several categorical or binary features describe the amenities and location advantages of the houses. Main road access is a crucial factor, as homes with direct access to main roads are generally more desirable due to better connectivity. Features like guest room, basement, hot water availability, and air conditioning capture the presence of specific amenities that can make a house more comfortable and, as a result, more expensive. Lastly, the dataset accounts for the furnishing status, where properties are labeled as furnished, semi-furnished, or unfurnished. Furnished houses tend to have higher prices, appealing to buyers who prefer ready-to-move-in options.

Overall, this dataset includes a mix of numerical and categorical attributes, all of which contribute to the final price of the house. Such a rich feature set makes it well-suited for building machine learning models to predict housing prices, as it captures both the physical dimensions of the house and additional factors that can influence buyer demand.

VI. METHODOLOGY

The methodology of this project follows a structured, multi-phase approach, leveraging advanced data processing and machine learning techniques to predict house prices accurately. Initially, the data collection phase involves gathering a comprehensive dataset containing housing attributes such as area, number of bedrooms, bathrooms, amenities (like air conditioning, parking spaces, etc.), and location-specific features. This data is preprocessed to handle missing values, standardize formats, and encode categorical variables. Feature engineering is then applied to create new, relevant features that may influence prices, such as proximity to city centers, transport links, and neighborhood demographics.

Following preprocessing, the next phase involves selecting appropriate machine learning models. The methodology employs both ensemble learning methods, such as Random Forests and XGBoost, and

deep learning techniques like Artificial Neural Networks (ANNs). These models are chosen for their ability to capture complex, nonlinear relationships between the features and the target variable, house price. Time series analysis using Long Short-Term Memory (LSTM) networks is also incorporated to forecast long-term price trends based on historical data.

To improve the accuracy and robustness of predictions, hyperparameter tuning is conducted using techniques such as Grid Search and Random Search. Cross-validation is applied to ensure that the models generalize well across different datasets, reducing the risk of overfitting. The system is dynamically integrated with real-time market data, including economic indicators and social sentiment analysis, to enhance the predictive capacity and adapt to market fluctuations.

Finally, the model is evaluated using standard performance metrics such as Mean Absolute Error (MAE), Root Mean Squared Error (RMSE), and R-squared (R^2). SHAP (Shapley Additive Explanations) values are used to interpret the impact of individual features on price predictions, offering insights into the factors driving house prices. The methodology ensures that the system is robust, accurate, and adaptable, providing a reliable tool for real estate price prediction.

VII. ALGORITHMS USED

Linear Regression serves as the baseline model for this study. It assumes a linear relationship between the independent variables (house features) and the dependent variable (price). Despite its simplicity, Linear Regression helps provide an initial understanding of the data and the direct impact of individual features on house prices. Random Forest is an ensemble learning algorithm that operates by constructing a multitude of decision trees during training. It reduces overfitting by averaging the results of multiple trees, making it robust and effective in capturing complex interactions between features. In this project, Random Forest is used to predict house prices by considering various attributes like area, number of bedrooms, bathrooms, and other categorical features (e.g., furnishing status, access to the main road). Gradient Boosting is another ensemble method that builds models sequentially, where each subsequent model corrects the errors of its predecessor. GBM is known for its high accuracy and ability to handle imbalanced data. It works well with

both numerical and categorical data, making it suitable for predicting house prices by improving the prediction gradually with each iteration. XGBoost (Extreme Gradient Boosting) is an optimized version of Gradient Boosting, designed for speed and performance. It uses regularization techniques to avoid overfitting, making it highly efficient for large datasets. XGBoost is used in this project to improve prediction accuracy by fine-tuning model parameters and ensuring a high-performance outcome when predicting house prices.

ANNs are powerful deep learning algorithms inspired by the human brain. They consist of interconnected layers of neurons that process input data and learn complex patterns. In this project, ANNs are used to capture the nonlinear relationships between house features and prices. The network learns from the data through multiple hidden layers, making it capable of identifying intricate patterns that other algorithms may miss. LSTM networks are a type of recurrent neural network (RNN) specifically designed for time series forecasting. In the context of this project, LSTMs are used to capture the time-dependent trends in house prices, based on historical data. They are effective in predicting long-term changes in the housing market by considering past price movements, making them valuable for forecasting future prices. SVM is a supervised machine learning algorithm that is used for both classification and regression tasks. In this project, SVM is applied to predict house prices by constructing hyperplanes in a multidimensional space that separate different price categories. It is particularly useful in handling high-dimensional data and works well with non-linear datasets by employing the kernel trick to transform input data into higher dimensions.

These optimization techniques are used for hyperparameter tuning to enhance the performance of the machine learning models. Grid Search systematically searches through predefined hyperparameter values, while Random Search randomly samples combinations of hyperparameters to identify the optimal set that improves the model's performance. To assess the performance of these algorithms, evaluation metrics like Mean Absolute Error (MAE), Root Mean Squared Error (RMSE), and R-squared (R^2) are used. These metrics help quantify the prediction accuracy and ensure the reliability of the models. By combining these algorithms, the system benefits from both interpretability and high predictive power, ensuring that the house price

predictions are accurate, dynamic, and reflective of real-world trends.

VIII. PERFORMANCE AND EFFICIENCY

The performance and efficiency of the algorithms used in this research are evaluated based on their ability to predict house prices accurately, generalize well to unseen data, and process large datasets efficiently. Each algorithm offers distinct advantages in terms of performance, with a combination of traditional machine learning models and advanced deep learning techniques providing a comprehensive solution for house price prediction.

Linear regression serves as a benchmark model. While it performs well when the relationship between features and house prices is roughly linear, its predictive power is limited when the dataset contains nonlinear patterns, interactions, or complex relationships. In terms of performance, it is the least accurate of the models used, often underfitting the data. Linear regression is computationally efficient and fast to train, making it suitable for small datasets or when interpretability is more important than prediction accuracy. However, its simplicity results in limited usefulness for complex datasets like housing data, where interactions between multiple features significantly influence price. Random Forest delivers strong predictive performance by aggregating the results of multiple decision trees, reducing variance and improving generalization. Its ability to handle both categorical and numerical features, along with capturing non-linear relationships, makes it effective for house price prediction. It typically performs better than linear regression, particularly in terms of handling overfitting and capturing complex interactions between features such as area, furnishing status, and number of bedrooms.

Random Forest can be computationally expensive, especially when the number of trees and depth of each tree increases. However, it parallelizes well, meaning that it can scale effectively on modern hardware, which compensates for its relatively slower training times. Its performance-accuracy trade-off makes it one of the more balanced models for this type of problem. GBM provides superior accuracy compared to Random Forest by building models sequentially, with each model correcting the errors of the previous one. This makes GBM particularly adept at handling complex datasets like house prices, where multiple variables interact in intricate ways. GBM is less prone to overfitting due to its regularization capabilities and

typically achieves high performance, especially when predicting non-linear relationships in the housing market. While GBM is highly accurate, it is computationally more expensive and slower to train than Random Forest due to its sequential nature. Each tree must be built one after the other, which makes it less efficient for large datasets unless optimized through parallelization techniques or early stopping criteria. Nonetheless, its performance justifies the extra computation time.

XGBoost improves upon GBM by adding further optimizations such as handling missing data and utilizing both L1 and L2 regularization to reduce overfitting. XGBoost consistently delivers high accuracy and is often considered one of the best-performing algorithms for structured data like housing datasets. It captures complex, non-linear relationships between features and the target variable, making it highly suited for predictive tasks in real estate. XGBoost is significantly more efficient than traditional GBM, offering faster training times due to its optimized memory usage and parallelized tree-building process. It is designed for high performance and scalability, allowing it to handle very large datasets effectively. The efficiency of XGBoost makes it a popular choice for competitive machine learning tasks, as it balances performance with computational cost. ANNs excel at learning complex, non-linear patterns, making them particularly powerful for house price prediction, where interactions between features are not always obvious. They can capture subtle relationships between variables, such as how area, amenities, and location interact to influence price. ANNs generally outperform traditional models like Random Forest and GBM on large datasets due to their deep learning architecture, which enables them to model more intricate patterns.

While ANNs offer superior performance on complex tasks, their training process can be computationally intensive, especially as the number of hidden layers and neurons increases. They require significant computational resources and time to train, particularly for large datasets. Despite this, ANNs are capable of achieving higher performance than traditional machine learning models when provided with sufficient data and computation power. LSTM networks are particularly effective for time-series forecasting tasks, capturing long-term dependencies in data. In the context of house price prediction, LSTMs perform well by analyzing historical trends

and market fluctuations to predict future price movements. Their ability to remember past states makes them useful for predicting how house prices will evolve based on past trends, and they generally outperform other models in this specific domain.

LSTMs, while highly accurate for time-dependent data, can be slower to train compared to traditional models due to their recurrent nature. They require more computational power and training time because each time step depends on the previous one. However, the efficiency of LSTMs can be improved by using techniques such as GPU acceleration and optimizing the network architecture. SVMs perform well for regression tasks, particularly when the data has clear, separable patterns. They are effective in capturing non-linear relationships by transforming the input data into higher-dimensional spaces. In house price prediction, SVMs can handle smaller datasets effectively but might not scale as well for larger, more complex datasets like housing data. SVMs are computationally efficient for small to medium-sized datasets but can become slow and memory-intensive when dealing with larger datasets or a high number of features. The kernel trick helps SVMs achieve better performance in non-linear tasks, but the computation of the kernel can be costly when the dataset size increases.

Ensemble learning models like Random Forest, GBM, and XGBoost offer a good balance between performance and efficiency for house price prediction, particularly for datasets with complex, non-linear relationships. XGBoost stands out as one of the most efficient and high-performing models due to its optimization techniques. ANNs and LSTMs provide superior performance in capturing deep, complex patterns, but at the cost of higher computational requirements. Each model's efficiency can be further enhanced with techniques like hyperparameter tuning, parallelization, and hardware acceleration (e.g., GPU usage), making them suitable for large-scale house price prediction systems.

IX. CONCLUSION

In the research based project, we explored various machine learning algorithms for predicting house prices based on a diverse set of features. Ensemble methods such as Random Forest, Gradient Boosting, and XGBoost provided accurate and reliable results by capturing complex, non-linear relationships between features. Deep learning models like Artificial Neural Networks (ANNs) and Long Short-Term

Memory (LSTM) networks further enhanced prediction accuracy, particularly in scenarios involving historical price trends. XGBoost stood out as the most efficient and high-performing algorithm, balancing computational cost with prediction accuracy. Overall, the study demonstrates that machine learning, combined with feature engineering and optimization techniques, can significantly improve the accuracy of house price predictions, making it a valuable tool for the real estate market.

X. REFERENCES

- [1] Rosen, S. (1974). Hedonic prices and implicit markets: Product differentiation in pure competition. *Journal of Political Economy**, 82(1), 34-55.
- [2] Malpezzi, S., & Mayo, S. K. (1997). Getting housing incentives right: A case study of the effects of regulation, taxes, and subsidies on housing supply in Malaysia. *Land Economics**, 73(3), 372-391.
- [3] Worzala, E., Lenk, M., & Silva, A. (1995). An exploration of neural networks and its application to real estate valuation. *Journal of Real Estate Research**, 10(2), 185-201.
- [4] Antipov, E. A., & Pokryshevskaya, E. B. (2012). Mass appraisal of residential apartments: An application of random forest for valuation and a CART-based approach for model diagnostics. *Expert Systems with Applications**, 39(2), 1772-1778.
- [5] Goodman, A. C. (1989). Topics in empirical urban housing research. *Journal of Urban Economics**, 26(3), 291-310.
- [6] Bhan, G. (2009). "This is no longer the city I once knew": Evictions, the urban poor, and the right to the city in millennial Delhi. *Environment and Urbanization**, 21(1), 127-142.
- [7] Glaeser, E. L. (2011). *Triumph of the city: How our greatest invention makes us richer, smarter, greener, healthier, and happier**. Penguin Press.
- [8] Bourassa, S. C., & Hoesli, M. (2006). The price of aesthetic externalities. *Journal of Real Estate Literature**, 14(1), 47-65.
- [9] Zhang, Y., & Liu, P. (2013). The impact of urbanization on housing prices: A panel data analysis for Chinese cities. *Journal of Cleaner Production**, 42(1), 75-82.
- [10] Khashman, A. (2010). Neural networks for credit risk evaluation: Investigation of different

- neural models and learning schemes. **Expert Systems with Applications**, 37(9), 6233-6239.
- [11] Ball, M. (2012). The impact of regulation on housing affordability. **Economic Affairs**, 32(3), 28-33.
- [12] Nath, S., & Ray, R. (2017). An empirical study of housing price determinants in the Indian housing market. **Real Estate Management and Valuation**, 25(3), 55-64.
- [13] Fuerst, F., & McAllister, P. (2016). The role of big data in real estate research. **Journal of Property Investment & Finance**, 34(2), 110-121.
- [14] Kim, S., & Lee, Y. (2018). House price prediction with deep learning algorithms. **International Journal of Housing Markets and Analysis**, 11(4), 494-507.
- [15] Mohanty, R., & Chakravarty, S. (2020). Housing affordability and price dynamics in urban India. **Urban Studies**, 57(14), 2825-2842.