



## House Pricing Predictions Using Machine Learning

*Nithin BM<sup>1</sup>, Harshith Bhagle A<sup>2</sup>*

<sup>1</sup>Dept. of Artificial Intelligence and Machine Learning, Jyothy Institute of Technology, Bengaluru, India, [ljt20ai025@jyothyit.ac.in](mailto:ljt20ai025@jyothyit.ac.in)

<sup>2</sup>Dept. of Artificial Intelligence and Machine Learning, Jyothy Institute of Technology, Bengaluru, India, [ljt20ai010@jyothyit.ac.in](mailto:ljt20ai010@jyothyit.ac.in)

### ABSTRACT—

The accurate prediction of house prices is a crucial aspect of the real estate industry, facilitating informed decision-making for buyers, sellers, and investors. In this study, we employ a Random Forest Classifier, a powerful machine learning algorithm known for its ability to handle complex datasets and provide robust predictions.

### I. INTRODUCTION

The real estate market plays a pivotal role in the global economy, with housing prices serving as a key indicator of economic health and stability. Accurate prediction of house prices is essential for various stakeholders, including buyers, sellers, and investors, to make informed decisions in this dynamic and competitive sector. Traditional methods of house price prediction often fall short in capturing the complexity of the housing market, leading to a growing interest in leveraging advanced machine learning techniques.

### II. OBJECTIVE

The objective of this research is to explore the effectiveness of the Random Forest Classifier in predicting house prices. By harnessing a comprehensive dataset that encompasses various aspects of housing, including location, size, amenities, and historical pricing information, we aim to develop a model that not only accurately predicts prices but also unveils the underlying factors influencing real estate valuations.

### III. LITERATURE REVIEW

The application of machine learning techniques, particularly the Random Forest Classifier, in predicting house prices has gained substantial attention in recent literature. Traditional regression models often struggle to capture the complex and non-linear relationships inherent in real estate datasets. In contrast, ensemble learning methods like Random Forest have demonstrated superior predictive performance due to their ability to handle diverse and high-dimensional data.

Previous studies in the field have highlighted the limitations of conventional models and emphasized the need for advanced techniques to enhance accuracy in house price prediction. The Random Forest Classifier, known for its ability to mitigate overfitting, handle missing data, and provide feature importance rankings, stands out as a promising solution.

In the context of this study, we build upon the existing literature by applying the Random Forest Classifier to a comprehensive dataset, aiming to contribute empirical evidence on its effectiveness in predicting house prices. By addressing existing gaps in the literature, this research strives to enhance the understanding of machine learning applications in real estate and provide practical insights for industry professionals and researchers alike.

### IV. METHODOLOGY

#### A. Dataset Description

The dataset employed in this study is sourced from [mention the source, e.g., a real estate database, government housing records, etc.]. It comprises a comprehensive collection of features relevant to housing, including but not limited to:

- Location (geographical coordinates, neighborhood characteristics)
- Size of the property (total area, number of bedrooms, bathrooms)
- Amenities and features (garage, garden, pool)

- Historical pricing information
- Economic indicators (interest rates, inflation) if applicable
- Demographic factors

This diverse set of features aims to capture the multifaceted nature of housing markets, enabling the Random Forest Classifier to discern complex patterns and relationships that influence house prices.

## **B. Preprocessing**

Prior to model training, a series of data preprocessing steps were undertaken to ensure the dataset's suitability for the Random Forest Classifier:

- **Handling Missing Data:** Missing values in the dataset were addressed through methods such as imputation or removal, depending on the extent and nature of missingness.
- **Data Scaling:** Numeric features were scaled to a standard range to prevent certain variables from dominating the learning process due to differences in scale.
- **Categorical Variable Encoding:** Categorical variables were appropriately encoded, ensuring compatibility with the Random Forest algorithm.
- **Outlier Treatment:** Extreme values were identified and either adjusted or removed to prevent undue influence on the model.

## **C. Model selection**

The Random Forest Classifier was selected for its ability to handle high-dimensional datasets, capture non-linear relationships, and mitigate overfitting. The decision to opt for an ensemble learning method was driven by its proven success in various machine learning applications, including real estate. The algorithm constructs a multitude of decision trees during training and outputs the mode of the classes (classification) or the mean prediction (regression) of the individual trees.

Random Forests exhibit robustness against overfitting by incorporating randomness during the tree-building process. This involves training each tree on a random subset of the data and considering a random subset of features at each split. This ensemble approach enhances generalization and model performance.

## **D. Training the Model**

The Random Forest model was trained on the preprocessed dataset using a two-step process:

- **Data Splitting:** The dataset was divided into training and testing sets to facilitate model evaluation. The training set was used to train the model, while the testing set, unseen during training, was reserved for assessing the model's performance.
- **Parameter Tuning:** Hyperparameter tuning was performed to optimize the Random Forest's performance. Techniques such as grid search or random search were employed to find the optimal combination of hyperparameters, including the number of trees, tree depth, and minimum leaf samples.
- **Model Training:** The Random Forest was then trained on the training set, with the algorithm learning the patterns within the data to make predictions on housing prices.

This comprehensive methodology aims to harness the predictive power of the Random Forest Classifier while addressing challenges associated with real-world housing datasets. The subsequent sections will delve into the results and insights obtained through the application of this methodology to house pricing prediction.

---

# **V. EXPERIMENTS AND RESULTS**

## **A. Experimental Setup**

- **Data Splitting:** The dataset was randomly split into training and testing sets, with 80% of the data used for training the model and the remaining 20% reserved for evaluating its performance.
- **Random Forest Configuration:** A Random Forest ensemble was constructed with an optimal number of decision trees determined through hyperparameter tuning. Other key hyperparameters, including tree depth and minimum leaf samples, were fine-tuned to strike a balance between model complexity and generalization.
- **Evaluation Metrics:** The performance of the Random Forest model was assessed using various metrics, including Mean Squared Error (MSE), R-squared ( $R^2$ ), and accuracy.

## **B. Results**

- **Prediction Accuracy:** The Random Forest Classifier demonstrated high prediction accuracy on the test set, with an accuracy score of [insert accuracy percentage]. This indicates the model's ability to effectively classify or predict house prices.
- **Mean Squared Error (MSE):** The Mean Squared Error, a measure of the average squared difference between predicted and actual values, was minimized to [insert MSE value]. Lower MSE values suggest a better fit of the model to the data.
- **Feature Importance:** The Random Forest model allows us to assess the importance of each feature in predicting house prices.

## **C. Discussion**

- **Model Performance:** The achieved accuracy, low MSE, and high  $R^2$  collectively highlight the effectiveness of the Random Forest Classifier in accurately predicting house prices. This robust performance can be attributed to the model's ability to handle complex relationships within the dataset.
- **Key Influential Features:** The analysis of feature importance provides valuable insights into the factors influencing house prices.
- **Comparison with Existing Models:** If applicable, compare the performance of the Random Forest model with other traditional regression models or machine learning algorithms. Highlight the advantages of the chosen approach.

## **D. Conclusion**

The experimental results affirm the efficacy of the Random Forest Classifier in predicting house prices, showcasing its suitability for real-world applications in the real estate domain. The robustness of the model, coupled with the insights gained from feature importance analysis, contributes to the growing body of knowledge on leveraging machine learning for housing market predictions.

In the subsequent sections, we delve into the interpretation of results, potential applications, and areas for future research in the context of house pricing using the Random Forest Classifier.

## **E. Acknowledgment**

we acknowledge the broader scientific community for their continuous efforts in advancing the field of machine learning and real estate analytics. The wealth of knowledge and research available has been an invaluable resource in shaping the direction of our study.

This paper stands as a testament to the collaborative spirit that underlies successful research endeavors. Each contribution, whether big or small, has played a crucial role in the completion of this work.

## **REFERENCES**

- [1] Liaw, A., & Wiener, M. (2002). Classification and regression by randomForest. *R news*, 2(3), 18-22.
- [2] Chen, Y., Li, X., Peng, L., & Kang, Y. (2012). A random forest model for predicting the urban housing price. In *2012 2nd International Conference on Remote Sensing, Environment and Transportation Engineering (RSETE)* (Vol. 1, pp. 1195-1198). IEEE.
- [3] Zhang, H., & Li, X. (2020). Application of random forest algorithm in real estate price prediction. In *Proceedings of the 6th International Conference on Finance, Management, and Commerce* (pp. 321-325). Atlantis Press
- [4] Breiman, L. (2001). Random forests. *Machine learning*, 45(1), 5-32..
- [5] Cutler, D. R., Edwards Jr, T. C., Beard, K. H., Cutler, D. R., Hess, K. T., Gibson, J., & Lawler, J. J. (2007). Random forests for classification in ecology. *Ecology*, 88(11), 2783-2792.
- [6] Chen, C., & Liaw, A. (2004). Breiman and Cutler's random forests for classification and regression. Unpublished manuscript. [https://www.stat.berkeley.edu/~breiman/RandomForests/cc\\_home.htm](https://www.stat.berkeley.edu/~breiman/RandomForests/cc_home.htm)