# Artificial Intelligence and Machine Learning

Project Report

Semester-IV (Batch-2022)

Pridicting House Prices

**Supervised By:**

Mr. Talib

**Submitted By:**

Vishnu Bansal (2210992542), G-28
Aditya Sharma (2210992588), G-28
Swayam Sharma (2210992434),G-28
Tanishjot Kaur (2210992443), G-28

**Department of Computer Science and Engineering**
**Chitkara University Institute of Engineering & Technology,**
**Chitkara University, Punjab**

# TABLE CONTENTS

# 1. <u>INTRODUCTION</u>

The real estate market is a dynamic and vital sector of the economy, influencing various stakeholders ranging from homeowners and investors to policymakers and financial institutions. Accurately predicting house prices is crucial for making informed decisions in buying, selling, or investing in properties. In recent years, with the advent of advanced technologies and machine learning algorithms, predicting house prices has become more data-driven and precise.

This project aims to leverage machine learning techniques to predict house prices based on various features such as location, size, number of rooms, amenities, and other relevant factors. By analyzing historical data and employing predictive modeling, we seek to develop a robust and accurate pricing model that can assist homeowners, real estate agents, and investors in estimating property values.

The project will utilize a dataset containing information about past real estate transactions, including attributes like the number of bedrooms and bathrooms, square footage, neighborhood demographics, proximity to amenities, and sale prices. By preprocessing and cleaning the data, we will ensure its quality and relevance for training our predictive models.

## 1.1 Background:

The real estate industry has always been a focal point of economic activity, with property transactions representing significant investments for individuals, businesses, and institutions. Traditionally, determining the value of a property has relied heavily on the expertise and subjective judgment of real estate professionals, often supplemented by market research and comparable sales data. However, this approach may be prone to biases and inaccuracies, leading to suboptimal decision-making and potential financial losses.

In recent years, there has been a paradigm shift in how real estate valuation is approached, driven by advancements in technology, data availability, and machine learning techniques. With the proliferation of online real estate platforms and the digitization of property data, there is now a wealth of information that can be leveraged to develop more data-driven and precise pricing models.

Machine learning, in particular, has emerged as a powerful tool for predicting house prices by analyzing large volumes of historical transaction data and identifying complex patterns and relationships among various property attributes. By training predictive models on diverse datasets containing information such as property characteristics, location features, economic indicators, and market trends, it is possible to generate accurate price estimates and insights into the factors driving property values.

Moreover, the increasing availability of open data sources, such as government housing databases, real estate listings, and demographic information, has facilitated the development of more comprehensive and sophisticated pricing models. These models can incorporate a wide range of features and variables, allowing for more nuanced and granular predictions tailored to specific markets and property types.

## 1.2 Objectives:

**Data Collection and Preprocessing:** Gather a comprehensive dataset containing information on past real estate transactions, including property attributes, location features, and sale prices. Cleanse and preprocess the data to ensure its quality and relevance for training predictive models.

**Exploratory Data Analysis (EDA):** Conduct exploratory data analysis to gain insights into the distribution of key variables, identify correlations, outliers, and missing values, and understand the relationships between different features and house prices.

**Feature Engineering:** Engineer new features and transform existing variables to capture additional information and enhance the predictive power of the models. This may involve techniques such as creating interaction terms, encoding categorical variables, and deriving meaningful indicators from raw data.

**Model Selection and Training:** Evaluate a range of machine learning algorithms, including linear regression, decision trees, random forests, gradient boosting, and neural networks, to identify the most suitable model for predicting house prices. Train the selected models on the training dataset using appropriate techniques such as cross-validation and regularization.

**Model Evaluation and Validation:** Assess the performance of the trained models using evaluation metrics such as mean absolute error (MAE), mean squared error (MSE), root mean squared error (RMSE), and R-squared (R2) to measure their accuracy, precision, and generalization ability. Validate the models on a separate test dataset to ensure their reliability and robustness.

**Hyperparameter Tuning:** Optimize the hyperparameters of the selected models through techniques such as grid search, random search, or Bayesian optimization to improve their performance and fine-tune their predictive capabilities.

## 1.3 Significance:

- Enhanced Decision-Making: Accurate house price prediction models can significantly improve decision-making for various stakeholders, including homeowners, buyers, sellers, and investors. By providing reliable estimates of property values, these models enable informed decisions regarding purchasing, selling, or investing in real estate assets, ultimately maximizing returns and minimizing risks.

- Market Transparency: Predictive models contribute to increasing market transparency by demystifying the factors influencing house prices. By analyzing historical transaction data and identifying key drivers of property values, these models help stakeholders understand market dynamics, trends, and fluctuations, fostering a more transparent and efficient real estate market.

- Risk Mitigation: Reliable house price prediction models empower individuals and organizations to assess and mitigate risks associated with real estate transactions. By accurately estimating property values and identifying potential overvaluation or undervaluation, stakeholders can make more prudent investment decisions and mitigate financial risks, such as overpaying for a property or experiencing unexpected losses.

- Empowerment of Stakeholders: Access to accurate house price predictions empowers homeowners, buyers, and sellers with valuable insights into the value of their properties. This knowledge enables them to negotiate effectively, set appropriate listing prices, and make informed decisions regarding property transactions, ultimately empowering individuals to achieve their real estate goals more effectively.

- Financial Planning and Investment Strategies: For investors and financial institutions, house price prediction models serve as essential tools for financial planning and investment strategies. By forecasting future property values and assessing market trends, investors can allocate capital more efficiently, optimize portfolio diversification, and make strategic decisions to maximize returns on real estate investments.

- Policy Formulation: House price prediction models provide policymakers and government agencies with valuable insights into housing market trends and dynamics. By analyzing large-scale data on property transactions and price movements, policymakers can formulate evidence-based policies and interventions to address housing affordability, promote sustainable development, and ensure the stability and resilience of the housing market.

- Real Estate Valuation Practices: The development of advanced machine learning models for house price prediction contributes to the evolution of real

- estate valuation practices. By integrating data-driven methodologies and predictive analytics into traditional valuation approaches, stakeholders can enhance the accuracy, objectivity, and efficiency of property appraisal processes, leading to more reliable and consistent valuation outcomes.

- Innovation in Real Estate Technology: The deployment of house price prediction models fosters innovation in real estate technology and digital platforms. By integrating predictive analytics capabilities into online real estate marketplaces and property valuation tools, stakeholders can access real-time pricing information, personalized recommendations, and interactive visualization tools, enhancing user experience and facilitating more efficient property transactions.

- Economic Impact: Accurate house price prediction models have broader economic implications, influencing consumer spending, investment behavior, and overall economic stability. By providing stakeholders with confidence in property valuations and market conditions, these models contribute to a more resilient and vibrant real estate sector, which, in turn, supports economic growth, job creation, and wealth accumulation.

- Long-term Socioeconomic Benefits: The development and adoption of house price prediction models have the potential to generate long-term socioeconomic benefits, including improved housing affordability, increased access to homeownership, and enhanced wealth accumulation for individuals and communities. By promoting transparency, efficiency, and fairness in the real estate market, these models contribute to building more inclusive and sustainable societies.

# 2.        PROBLEM DEFINITION AND REQUIREMENT

## 2.1        PROBLEM STATEMENT

In the ever-changing landscape of real estate, the need for accurate predictions of future property prices is paramount. Leveraging a comprehensive housing dataset, we aim to develop a predictive model that takes into account crucial factors such as location, size, bedrooms, and amenities. This model will empower real estate stakeholders and potential buyers with valuable insights, facilitating well-informed decisions in a dynamic and competitive housing market.

The problem statement for this project revolves around developing a machine learning-based solution to predict house prices accurately. Specifically, the project aims to address the following challenges:

1. **Model Accuracy**: Existing house price prediction models may suffer from limitations in accuracy due to factors such as incomplete data, model complexity, or inadequate feature representation. Improving model accuracy is essential for building stakeholders' confidence in the predictive capabilities of the system.

2. **Feature Selection and Engineering**: Identifying the most relevant features and engineering informative variables are critical for enhancing the predictive power of the models. This involves selecting features that capture essential aspects of the property, neighborhood, and market dynamics while minimizing noise and redundancy.

3. **Model Interpretability**: While complex machine learning algorithms may achieve high prediction accuracy, they often lack interpretability, making it challenging to understand the underlying factors driving the predictions. Ensuring model interpretability is essential for building trust and facilitating informed decision-making by stakeholders.

4. **Generalization and Robustness**: House price prediction models must generalize well to unseen data and be robust to variations in market conditions, property types, and geographic regions. Overfitting to specific datasets or failing to capture the heterogeneity of real estate markets can compromise the reliability and applicability of the models.

5. **Scalability and Deployment**: Deploying house price prediction models at scale requires considerations such as computational efficiency, scalability, and integration with existing real estate platforms and workflows. Developing a deployable solution that meets the performance and usability requirements of stakeholders is crucial for real-world adoption.

## 2.2  SOFTWARE REQUIREMENT

Programming Language: Choose a programming language suitable for data analysis, machine learning, and model development. Common choices include Python and R. Python is widely preferred for its extensive libraries and frameworks, such as scikit-learn, pandas, and NumPy, which are well-suited for building machine learning models.

Integrated Development Environment (IDE): Select an IDE for coding, debugging, and running scripts efficiently. Popular options include PyCharm, Jupyter Notebook, and VS Code. Jupyter Notebook is particularly useful for exploratory data analysis and interactive development.

Machine Learning Libraries: Install machine learning libraries and frameworks to implement predictive models. Key libraries include scikit-learn, TensorFlow, Keras, and PyTorch. Scikit-learn offers a wide range of algorithms and tools for model training, evaluation, and deployment, making it suitable for this project.

Data Manipulation and Analysis Tools: Utilize libraries for data manipulation and analysis to preprocess and explore the dataset. Pandas is a powerful library for working with structured data, providing functionalities for data cleaning, transformation, and aggregation.

Data Visualization Tools: Employ data visualization libraries to create insightful visualizations of the dataset and model outputs. Matplotlib and Seaborn are popular libraries for creating static plots, while Plotly and Bokeh offer interactive visualization capabilities.

Database Management System (DBMS): Depending on the size and nature of the dataset, consider using a DBMS for data storage and retrieval. Common choices include SQLite for small-scale projects and PostgreSQL or MySQL for larger datasets requiring more advanced features and scalability.

# 3. PROPOSED DESIGN / METHODOLOGY

Proposed Solution:

The solution involves building a machine learning model that can learn from the dataset to predict the selling price of houses. The model will be trained on historical data, validated using a portion of the dataset, and evaluated based on its predictive performance.

Methodology:

- Data Preprocessing: Cleanse and preprocess the dataset, handling missing values, outliers, and encoding categorical variables as necessary.
- Feature Engineering: Create additional features or transform existing ones to enhance the predictive power of the model.
- Model Selection: Experiment with different regression algorithms (e.g., linear regression, decision trees, ensemble methods) and select the best-performing model based on evaluation metrics.
- Model Training: Train the selected model using the training data, tuning hyperparameters as needed to optimize performance.
- Model Evaluation: Evaluate the trained model's performance using validation data, assessing metrics such as Mean Absolute Error (MAE), Mean Squared Error (MSE), and Root Mean Squared Error (RMSE).
- Deployment: Integrate the trained model into an AIML framework to create a conversational interface where users can input property features and receive predicted prices.
- Continuous Improvement: Monitor the model's performance over time, retraining with updated data as necessary to maintain accuracy and relevance

### 3.1 Technical details:

NumPy:
- Perform mathematical operations on data arrays, such as calculating means, medians, and standard deviations.
- Use NumPy arrays to represent and manipulate numerical data efficiently.

Pandas:
- Use pandas to load, clean, and preprocess your dataset.
- Perform data manipulation tasks such as filtering, grouping, and aggregating data.
- Handle missing values, outliers, and data formatting issues.

Matplotlib:

- Create static, publication-quality visualizations using Matplotlib.

- Plot various types of charts, including line plots, scatter plots, histograms, and bar charts

Seaborn:

- Use Seaborn to create more visually appealing and informative statistical visualizations.
- Generate complex plots such as scatter plots with regression lines, box plots, violin plots, and pair plots
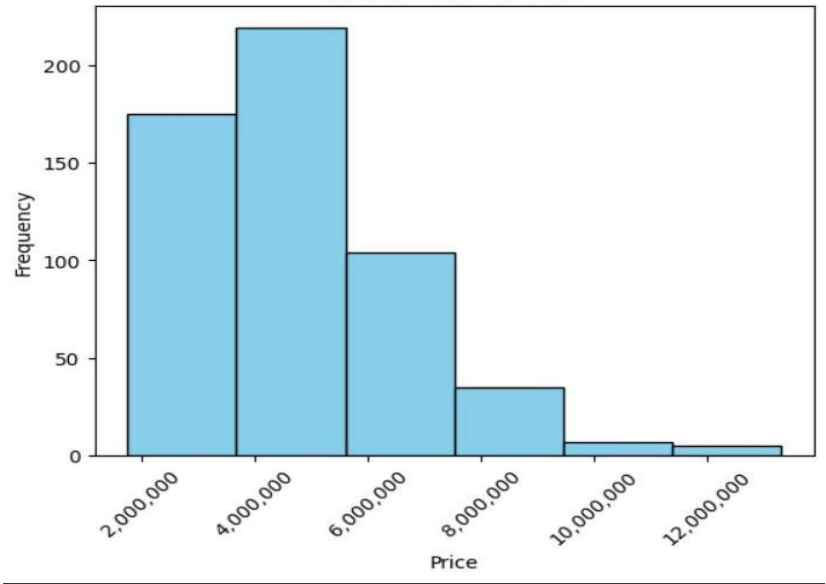
Scikit Learn.

- Scikit-learn, often referred to as sklearn, is a popular machine learning library for Python

- Scikit-learn, often referred to as sklearn, is a popular machine learning library for Python.
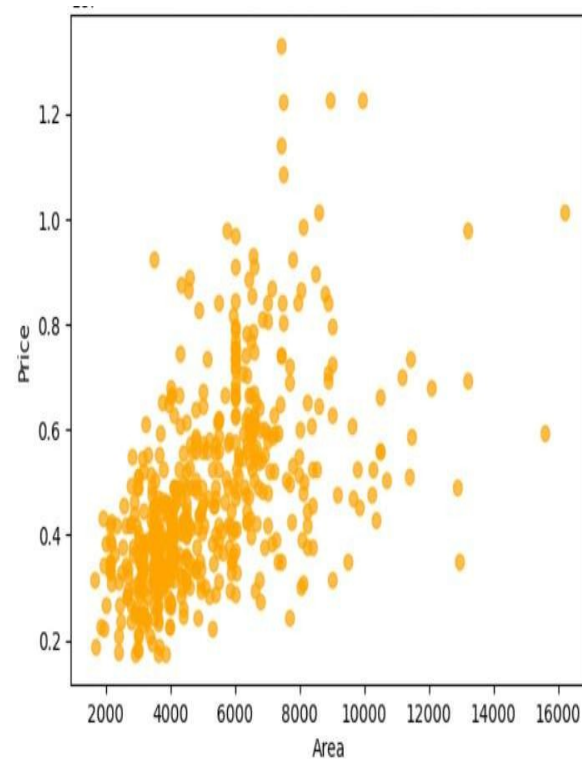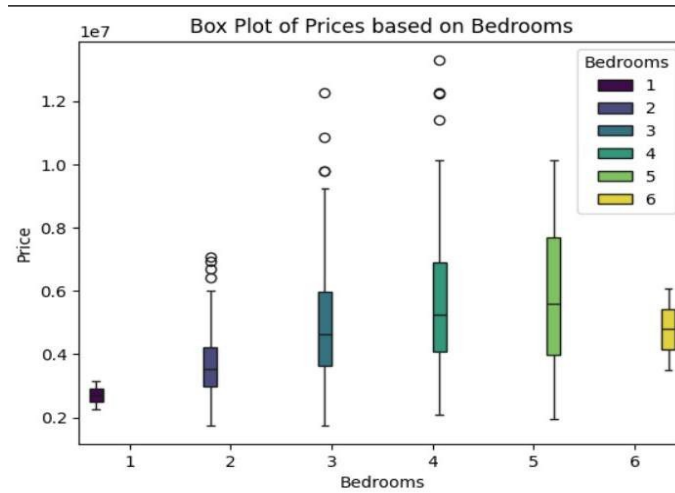
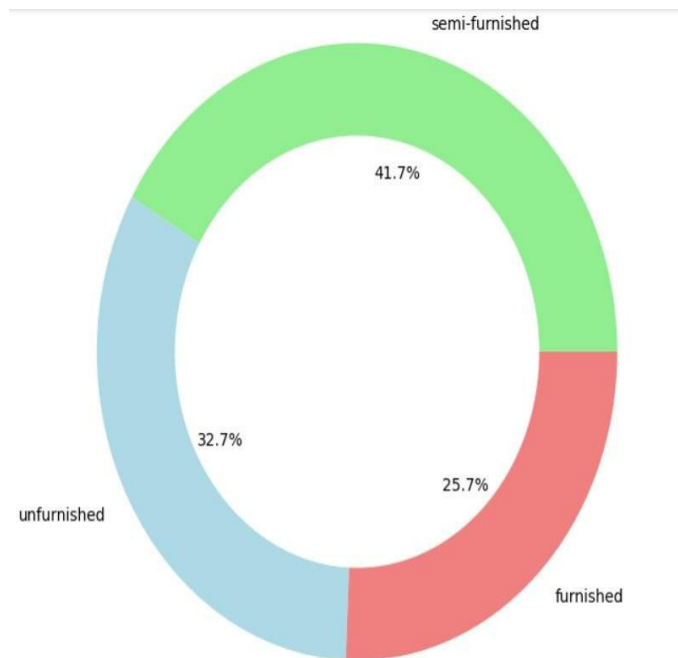## 3.2 PLOTS USED :

### 1. Bar plot :



### 2. HISTOGRAM:

# 3. SCATTER PLOT:



# 4. BOX PLOT:

**5.PIE CHART:**

semi-furnished

41.7%

32.7%

25.7%

unfurnished

furnished

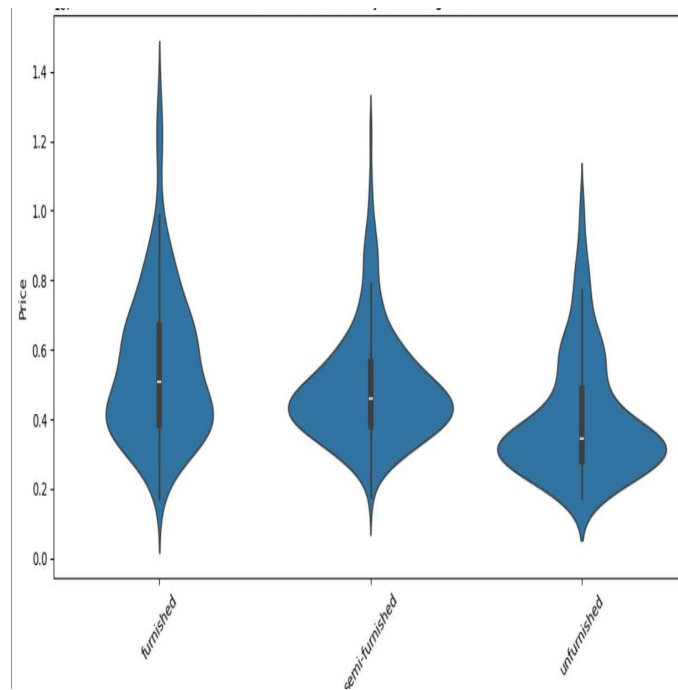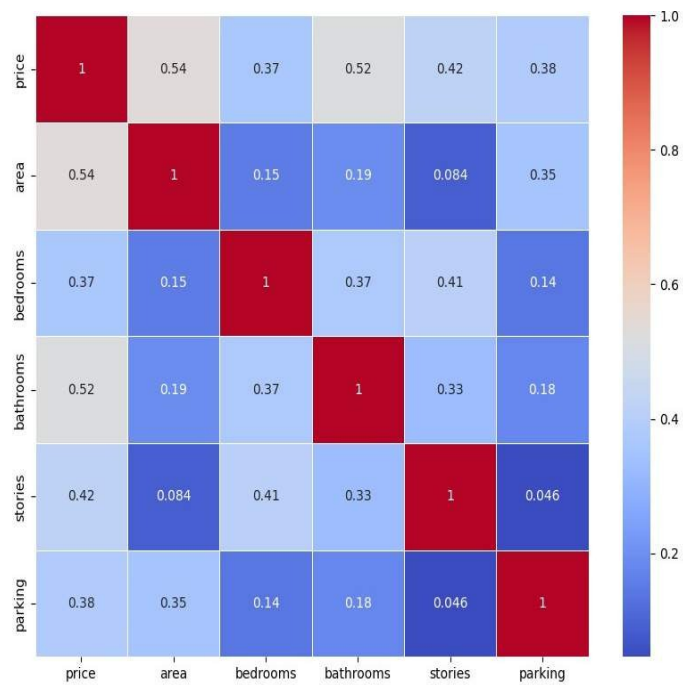# 6.KDE PLOT:

KDE Plot of Price

# 5. VIOLIN PLOT :



# 6. HEATMAP:



Page no: 14

## 4. KEY FEATURES

1.  Scatter plot: Scatter plots are effective for visualizing the relationship between two continuous variables, such as calories burnt vs duration . Key features include:

    *   Trend lines or regression lines to highlight the overall trend in the data.

    *   Different colors or markers represent categorical variables (e.g.,Calories burnt vs Duration).

    *   Transparency or alpha blending to visualize overlapping data points.

2.  Histogram plots: Histograms are useful for visualizing the distribution of numerical variables, such as calories burnt over frequency. Key features include:

    *   Customizable bin sizes to control the granularity of the distribution.
    *   Multiple histograms or density plots for different categories or subsets of data.

3.  Box plot or violin plot : Box plots and violin plots provide insights into the distribution of numerical variables across different categories or groups, such as gender vs duration Key features include:

    *   Median lines and interquartile ranges to summarize the central tendency and spread of the data.
    *   Outlier detection to identify unusual or extreme data points.

    *   Grouping and stacking to compare distributions side by side or within the same plot.

4.  Line plot: Time series plots are useful for visualizing changes in calorie burnt over time. Key features include:

    *   Line plots with time on the x-axis and the variable of interest on the y-axis.

    *   Smoothed or aggregated time series to highlight long-term trends while reducing noise.

    *   Annotations or vertical lines to mark significant events or milestones.

5.  Heatmaps and Correlation Matrices: Heatmaps and correlation matrices are effective for visualizing relationships between multiple variables in the dataset. Key features include:

    *   Color-coded cells to represent the strength and direction of correlations between variables.
    *   Hierarchical clustering to group similar variables based on their correlation patterns.
    *   Annotations or labels to identify variables and highlight interesting patterns.

6.  Pair Plot: Pair plots and joint plots allow for the visualization of pairwise relationships between multiple variables in the dataset. Key features include:

    *   ✓ Scatter plots with marginal histograms or density plots to visualize individual distributions.
    *   ✓ Pearson correlation coefficients or other measures of association between variables.
        Non-linear fits or regression lines to capture complex relationships.

# 5.REFERENCES

- https://colab.research.google.com/
- Talib sir