



Predicting House Prices

USING MACHINE LEARNING

DEFINITION:

House pricing prediction, also known as real estate price prediction or property valuation, refers to the process of using data and machine learning techniques to estimate the market value of a residential property, such as a house or apartment. The goal of house pricing prediction is to provide an accurate and data-driven estimate of how much a property should be worth based on various factors and features associated with the property.

HOUSE PRICE PREDICTION

USING MACHINE LEARNING TECHNIQUES



DESIGN THINKING:

- ❖ Data Sources
- ❖ Data Preprocessing
- ❖ Features Selection
- ❖ Model Selection
- ❖ Model Training
- ❖ Evaluation



DATA SOURCES:

- **Real Estate Listings Websites:** Websites like Zillow, Realtor.com, and Redfin provide a wealth of data on property listings, including price, location, square footage, number of bedrooms and bathrooms, and more. You can scrape or obtain data through APIs from these websites.
- **Historical Sales Data:** Historical sales data from local government agencies or real estate associations can provide valuable information about past property transactions in a specific area.
- **Economic Indicators:** Economic indicators such as local employment rates, income levels, and GDP growth can influence property prices. This data can be obtained from government reports, statistical agencies, or economic databases.
- **Environmental Data:** Environmental factors like air quality, flood risk, and natural disasters can affect property prices. Sources like the Environmental Protection Agency (EPA) provide relevant data.
- **Weather Data:** In some regions, weather conditions like temperature, precipitation, and climate can impact property prices. Weather data can be obtained from meteorological agencies.
- **Public Transportation Data:** Data on public transportation routes, stations, and schedules can influence property prices, especially in urban areas. Transit agencies often provide this information.

DATA PREPROCESSING:

- **Data Collection:** Collect data from various sources as mentioned in the previous response. Ensure that the data is comprehensive and covers relevant features such as property details, location, and economic indicators.
- **Data Cleaning:** Handle Missing Values: Identify and handle missing data. You can either remove rows or columns with missing values or impute them using techniques like mean, median, or more advanced methods.
- **Handling Skewed Data:** If your target variable (house prices) is skewed, consider applying transformations like logarithmic transformation to make it more normally distributed. This can improve model performance.
- **Data Scaling:** Scale numerical features to ensure they have similar scales. Common techniques include Standardization (Z-score normalization) or Min-Max scaling.

FEATURES SELECTION:

- **Location:** The location of a property is one of the most critical factors influencing its price. Factors like proximity to schools, parks, public transportation, shopping centers, safety, and the overall desirability of the neighborhood play a significant role.
- **Property Size and Layout:** The size of the property, including the square footage of the house and lot, as well as the number of bedrooms, bathrooms, and floors, is highly influential. Larger and more functional properties tend to have higher prices.
- **Property Condition and Age:** The condition of the property and its age can significantly impact the price. Newly constructed or recently renovated homes often command higher prices.
- **Environmental Factors:** Factors such as air quality, flood risk, and susceptibility to natural disasters can influence property prices.
- **Accessibility:** Access to highways, public transportation, and major airports can affect property prices, especially in urban areas.

MODEL SELECTION:

➤ Linear Regression:

Linear Regression is a simple and interpretable regression algorithm that assumes a linear relationship between the independent variables (features) and the target variable (house prices).

Pros:

- Fast training and prediction: Linear regression models are computationally efficient.

Cons:

- Sensitive to outliers: Outliers can disproportionately affect the model's performance.

➤ Random Forest Regressor:

Random Forest Regressor is an ensemble learning method based on decision trees. It combines the predictions of multiple decision trees to provide more accurate and robust predictions.

Pros:

- Nonlinearity: Random Forest can capture complex, nonlinear relationships in the data.

Cons:

- Complexity: Random Forest models can be more complex and harder to interpret compared to linear regression.
- Training time: Random Forests can be slower to train, especially with a large number of trees.

MODEL TRAINING:

- **Load Preprocessed Data:** Load the preprocessed dataset that you've prepared for house price prediction. Make sure it includes both the features and the target variable (house prices).
- **Split the Data:** Split the dataset into three subsets: a training set, a validation set, and a test set. A common split ratio is 70% for training, 15% for validation, and 15% for testing.
- **Select the Model:** Choose the machine learning model you want to train for house price prediction. You mentioned considering Linear Regression and Random Forest Regressor.
- **Train the Model:** Train the selected model using the training data.
- **Final Testing:** After you've trained and fine-tuned the model, evaluate its performance on the test set to assess how well it generalizes to unseen data.
- **Save the Trained Model:** If you're satisfied with the model's performance, you can save it for future use without having to retrain it.

EVALUATION:

- **Mean Absolute Error (MAE)** represents the average absolute difference between the actual and predicted values. Lower MAE indicates better model performance, and it is easy to interpret because it's in the same unit as the target variable.
- **Root Mean Squared Error (RMSE)** is similar to MAE but penalizes larger errors more heavily because it involves the square of the differences. It is also in the same unit as the target variable. Lower RMSE values are better.
- **R-squared (R^2)** measures the proportion of the variance in the target variable that is predictable from the independent variables (features). R^2 values range from 0 to 1, with higher values indicating a better fit. An R^2 of 1 indicates a perfect fit, while an R^2 of 0 suggests that the model does not explain any of the variance in the target variable.

CONCLUSION:

In conclusion, this house price prediction project has provided valuable insights into the factors that influence property prices. Through rigorous data analysis and the application of machine learning models, we have uncovered significant features such as square footage, the number of bedrooms and bathrooms, and location, all of which play crucial roles in determining house prices. Our model(s) demonstrated strong predictive capabilities, as evidenced by [mention relevant evaluation metric scores].

However, it's important to acknowledge that there are inherent challenges in predicting house prices accurately, including the complexity of the real estate market and the influence of external economic factors. Furthermore, while our model(s) provide a robust foundation, there is room for improvement, particularly in terms of fine-tuning and incorporating additional data sources.





THANK YOU シ