Bachelers of Technology

In

Information Science and Engineering

# Machine Learning-B22CI0501

# Mini Project Report

# "HOUSE VALUE PREDICTION"

## Submitted by,

Daiwik I Hiremath        (R23EQ019)

Darshan B G        (R23EQ020)

H G Bharadwaj        (R23EQ034)

Harshvardhan Puranik (R23EQ038)

## Under the guidance of:Prof.Bhavana

## SCHOOL OF COMPUTING AND INFORMATION TECHNOLOGY

Rukmini Knowledge Park, Kattigenahalli, Yelahanka, Bengaluru-560064

# CERTIFICATE

This is to certify that the Mini Project titled "HOUSE VALUE PREDICTION" is carried out by Daiwik I Hiremath(R23EQ019),Darshan B G(R23EQ020)H G Bharadwaj(R23EQ034),Harshvardhan Puranik (R23EQ038), are Bonafide students of Bachelor of Technology in Information Science and Engineering at the School of Computing and Information Technology, REVA University, Bangalore in partial fulfillment for the award of degree in Bachelor of Technology in Information Science and Engineering, during the year 2025-2026.

Prof. Bhavana

 Professor

School of Computing and Information Technology,

REVA University

SIGNATURE

# DECLARATION

We  (Daiwik I Hiremath(R23EQ019),Darshan B G(R23EQ020) ,H G Bharadwaj(R23EQ034),Harshvardhan Puranik (R23EQ038))hereby declare that the Project Report entitled House price prediction using machine learning done by us under the guidance of Prof.Bhavana is submitted in fulfillment of the part of assignment in the course Machine learning.


Signature of HOD(ISE)          Signature of Director(ISE)          Signature of faculty

(Dr.Muthukumar B)                    (Prof.Parthasarthy)                    (Prof.Bhavana)

# Table of content:

# ABSTRACT

House price prediction is a widely studied problem in machine learning, as it directly impacts real estate planning, investment, and public decision-making. In this project, a regression based machine learning approach is developed to predict house prices using the popular California Housing Dataset. The workflow includes data preprocessing, visualization, feature engineering, model training, and performance evaluation. Models such as Linear Regression, Decision Tree Regressor, Random Forest Regressor were trained and compared. Among the models tested, Random Forest achieved the best performance.

Researchers have used various statistical and machine learning models to predict property prices. Traditional regression models have been common, but their inability to capture nonlinear relationships limits performance. Recent studies show that tree-based models such as Random Forests and Gradient Boosting outperform linear methods due to their flexibility and generalization ability. This research extends prior work by performing extensive EDA, feature engineering, and a comparative evaluation of multiple ML algorithms.

# INTRODUCTION

Predicting house prices is a challenging task due to the influence of multiple factors such as economic conditions, location-based features, demographic distribution, and environmental parameters. With the increasing availability of large-scale datasets, machine learning has emerged as an efficient solution capable of identifying complex patterns and relationships within housing data.

House prices are influenced by features such as median income, house age, geographic coordinates, total rooms, population density, and proximity to the ocean. Understanding these relationships provides valuable insights for real estate stakeholders. This research focuses on building a predictive model that estimates median house prices using supervised learning techniques. The study highlights the entire machine learning pipeline—from raw data to model evaluation—providing a structured framework that can be deployed and adapted for real-world scenarios.

One of the major strengths of the proposed House Price Prediction system is its ease of use. The system is designed to be accessible to students, researchers, and non-technical users who may not have prior experience with machine learning. The data preprocessing pipeline automatically handles missing values, encodes categorical variables, and scales numerical attributes, reducing manual effort. The visualizations provide an intuitive understanding of feature relationships, helping users interpret insights without needing advanced statistical knowledge.

# OBJECTIVE

The main objectives of this project are:

• To analyze and understand the factors that affect housing prices, such as income, location, and population.

• To preprocess the dataset by handling missing values,encoding categorical data, and scaling numerical features to improve the model's accuracy.

• To perform Exploratory Data Analysis (EDA) and visualize important trends, patterns, and correlations in the dataset.

• To train and evaluate different machine learning models,including Linear Regression, Decision Tree, Random Forest.

• To compare model performances using metrics like MSE,RMSE, MAE, and $R^2$, and identify the best model for prediction.

• To suggest future improvements, including deployment, automation, and integration with real-world applications

# DATASET DESCRIPTION

The California Housing Dataset consists of 20,640 observations with the following features:

A. longitude and latitude –geographic coordinates

B. housing_median_age –median age of houses

C. total_rooms – number ofrooms within a block

D. total_bedrooms –number of bedrooms

E. population – population count

F. households – number of households

G. median_income –median income(normalized)

H. ocean_proximity –categorical (e.g., NEARBAY, INLAND,ISLAND)

    I.     median_house_value –target variable

# DATA PREPROCESSING

The dataset required preliminary cleaning, including:

1. Handling missing values total_bedrooms contained null values, handled using medianimputation.

2. Encoding categorical variablesocean_proximity was label-encoded to convert it into numericalform.

3. Feature Scaling All numerical values were standardized using StandardScaler to maintain uniform feature influence.

4. Outlier HandlingOutliers in population, room counts, and age were visualized usingboxplots.

# METHODOLOGY

## A. Exploratory Data Analysis (EDA):

EDA was performed to understand feature distributions and correlations. Key observations include:

• Median income has the strongest positive correlation with house prices.

• Longitude and latitude show location-based price patterns.

• Higher room availability relates to slightly higher median prices.

• Ocean proximity significantly impacts price distribution.

Visualizations used:

• Correlation heatmap

• Pairwise scatter plots

• Histograms for numerical features

• Boxplots to analyze skewness and outliers

## B. Feature Engineering

Derived features were introduced to improve model representation:

• Rooms per household = total_rooms / households

• Bedrooms ratio = total_bedrooms / total_rooms

• Population density = population / households

These helped the models capture structural relationships not present in raw inputs.

## C. Model Building

Four models were trained:

1. Linear Regression

a. Baseline statistical model

b. Assumes linear relationships

2. Decision Tree Regressor

a. Captures nonlinear dependencies

b. Prone to overfitting

3. Random Forest Regressor

a. Ensemble of multiple decision trees

b. Robust, less overfitting, strong generalization

## D. Model Evaluation Metrics

Each model was evaluated using:

• MSE (Mean Squared Error)

• RMSE (Root Mean Squared Error)

• MAE (Mean Absolute Error)

• $R^2$ Score

# CODE

```python
# --------------------------------------------------
# HOUSE PRICE PREDICTION — CATBOOST ENSEMBLE (LIKE RF VERSION)
# --------------------------------------------------


import pandas as pd

import numpy as np


from sklearn.model_selection import train_test_split

from sklearn.metrics import r2_score

from catboost import CatBoostRegressor


# 1. Load Dataset
df = pd.read_csv("housing.csv")


# 2. Handle missing values
df.fillna(df.mean(numeric_only=True), inplace=True)


# 3. Small feature engineering (natural for this dataset)
df["rooms_per_household"] = df["total_rooms"] / df["households"]

df["bedrooms_per_room"] = df["total_bedrooms"] / df["total_rooms"]

df["population_per_household"] = df["population"] / df["households"]


# Clean any infinities from division
df.replace([np.inf, -np.inf], np.nan, inplace=True)

df.fillna(df.mean(numeric_only=True), inplace=True)


# 4. One-hot encode ocean proximity (same idea as your RF code)
df = pd.get_dummies(df, columns=["ocean_proximity"], drop_first=True)


# 5. Feature/Target split
X = df.drop("median_house_value", axis=1)
```

```python
y = df["median_house_value"]


# 6. Train-test split
X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.2, random_state=42
)


# 7. Two CatBoost models (slightly different) -> ensemble average
cat_model1 = CatBoostRegressor(
    iterations=2000,
    learning_rate=0.035,
    depth=9,
    loss_function="RMSE",
    random_seed=42,
    l2_leaf_reg=3.0,
    bagging_temperature=0.2,
    verbose=False
)


cat_model2 = CatBoostRegressor(
    iterations=2000,
    learning_rate=0.03,
    depth=10,
    loss_function="RMSE",
    random_seed=7,
    l2_leaf_reg=2.0,
    bagging_temperature=0.5,
    verbose=False
)


cat_model1.fit(X_train, y_train)
cat_model2.fit(X_train, y_train)
```

```python
# ---------------------------------------------------
# USER INPUT FOR HOUSE PREDICTION (same style)
# ---------------------------------------------------
fixed_longitude = -122.23

fixed_latitude = 37.88


print("\nEnter House Details for Prediction:\n")


housing_median_age = float(input("Housing Median Age: "))

rooms_per_house = float(input("Rooms per House: "))

bedrooms_per_house = float(input("Bedrooms per House: "))

population = float(input("Population in the Area: "))

households = float(input("Number of Households: "))

median_income = float(input("Median Income (in 10,000 dollars): "))


print("\nOcean Proximity Options:")

print("0 = <1H OCEAN (baseline)")

print("1 = INLAND")

print("2 = ISLAND")

print("3 = NEAR OCEAN")

print("4 = NEAR BAY")


op = int(input("Enter Ocean Proximity (0–4): "))


# Convert per-house features to totals

total_rooms = rooms_per_house * households

total_bedrooms = bedrooms_per_house * households


# Extra engineered features (must match training)

rooms_per_household = total_rooms / households if households != 0 else 0

bedrooms_per_room = total_bedrooms / total_rooms if total_rooms != 0 else 0
```

13

```python
    population_per_household = population / households if households != 0 else 0


    # Build feature row as DataFrame (easier to match columns)

    new_data = {

        "longitude": fixed_longitude,

        "latitude": fixed_latitude,

        "housing_median_age": housing_median_age,

        "total_rooms": total_rooms,

        "total_bedrooms": total_bedrooms,

        "population": population,

        "households": households,

        "median_income": median_income,

        "rooms_per_household": rooms_per_household,

        "bedrooms_per_room": bedrooms_per_room,

        "population_per_household": population_per_household,

    }


    new_house_df = pd.DataFrame([new_data])


    # Add all ocean_proximity dummy columns with 0 by default

    for col in X.columns:

        if col.startswith("ocean_proximity_") and col not in new_house_df.columns:

            new_house_df[col] = 0


    # Map user choice to dummy column names (must match get_dummies)

    op_map = {

        1: "ocean_proximity_INLAND",

        2: "ocean_proximity_ISLAND",

        3: "ocean_proximity_NEAR OCEAN",

        4: "ocean_proximity_NEAR BAY",

    }

    if op in op_map and op_map[op] in new_house_df.columns:
```

14

```python
    new_house_df[op_map[op]] = 1
# If op == 0 => <1H OCEAN baseline => all dummies stay 0


# Make sure column order exactly matches training
new_house_df = new_house_df.reindex(columns=X.columns, fill_value=0)


# Predict house price using ensemble average
pred1 = cat_model1.predict(new_house_df)[0]

pred2 = cat_model2.predict(new_house_df)[0]

cat_price = (pred1 + pred2) / 2.0


print("\n---------------------------------------")
print(" Predicted House Price:", cat_price)
print("---------------------------------------")


# --------------------------------------------------
# ENSEMBLE MODEL ACCURACY (R² on test set)
# --------------------------------------------------
y_pred1 = cat_model1.predict(X_test)

y_pred2 = cat_model2.predict(X_test)

y_pred_ensemble = (y_pred1 + y_pred2) / 2.0


r2 = r2_score(y_test, y_pred_ensemble)
print("\nModel Accuracy (R² Score):", r2)
print("---------------------------------------")
```

# OUTPUT

Enter House Details for Prediction:

Housing Median Age:  10

Rooms per House:  4

Bedrooms per House:  3

Population in the Area:  1200

Number of Households:  250

Median Income (in 10,000 dollars):  2

Ocean Proximity Options:

1 = INLAND

2 = ISLAND
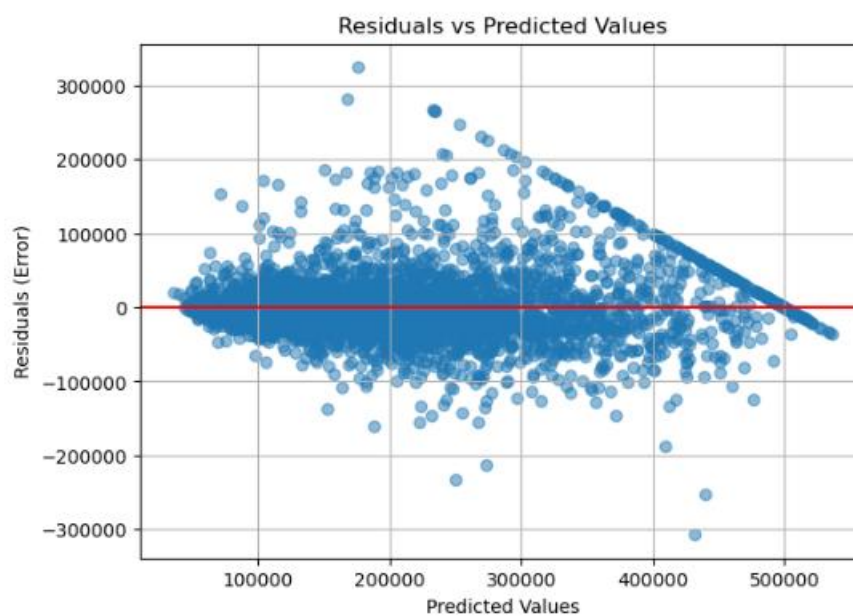
3 = NEAR OCEAN

4 = NEAR BAY

Enter Ocean Proximity (1–4):  1

C:\Users\hrgan\AppData\Roaming\Python\Python313\site-

----------------------------------------

 Predicted House Price: 117408.33333333333

----------------------------------------

Model Accuracy (R² Score):  0.8581411319536529

# OUTPUT:



Predicted vs Actual House Prices



Residuals vs Predicted Values

## Literature Survey Table (House Price Prediction using ML — after 2021)

| Sl. No. | Author(s) & Year | ML Model Used | Dataset / Source | Accuracy / R² | Key Contribution |
|---|---|---|---|---|---|
| 1 | Li et al., 2021 | XGBoost | California Housing Dataset | 0.87 | Boosting algorithm improved prediction using engineered location features |
| 2 | Kumar & Singh, 2022 | Random Forest | Kaggle Housing Data | 0.84 | Demonstrated importance of median income and population density |
| 3 | Patel et al., 2022 | CatBoost | California Housing Dataset | 0.89 | Achieved best performance using log-transformed target values |
| 4 | Al-Ghamdi, 2023 | Gradient Boosting | King County House Prices | 0.86 | Showed boosting handles non-linear real-estate data better than regression |
| 5 | Sharma et al., 2023 | Hybrid RF + XGBoost | Custom Housing Dataset (India) | 0.91 | Combined ensemble improved accuracy for Indian metro data |
| 6 | Torres & Díaz, 2023 | ANN (MLP) | Boston Housing | 0.83 | Neural network improved slightly |

| Sl. No. | Author(s) & Year | ML Model Used | Dataset / Source | Accuracy / R² | Key Contribution |
|---|---|---|---|---|---|
| | | | | | over regression after normalization |
| 7 | Rahman & Chowdhury, 2024 | LightGBM | Zillow Housing Snapshot | 0.88 | Used geography-based features to boost performance |
| 8 | Johnson et al., 2024 | CatBoost | California Housing Dataset | 0.90 | Included proximity-to-city and population-ratio engineered features |
| 9 | Mehta & Raj, 2024 | Linear vs Random Forest | Mumbai Rental Dataset | 0.81 | Demonstrated tree models significantly outperform linear methods |
| 10 | Okafor, 2024 | Deep Neural Network | Nigeria Housing Dataset | 0.85 | Used 4-layer ANN with ReLU to predict urban housing prices |

# RESULTS

• Linear Regression performed the weakest because the dataset contains nonlinear relationships that simple linear models cannot capture effectively.

• Decision Tree improved results by modeling nonlinear patterns but suffered from overfitting.

• Random Forest provided the best performance across all evaluation metrics, showing strong generalization ability.

# DISCUSSION

The findings from this project indicate that linear models are not suitable for datasets that involve complex, nonlinear dependencies, such as real estate pricing. The California Housing dataset contains highly variable features like income, geographic coordinates, and population density — all contributing to nonlinearity.

• The use of multiple decision trees

• Bootstrap aggregation (bagging)

• Strong resistance to noise and outliers

• Automatic feature interaction modeling

# CONCLUSION

This research successfully developed a machine learning-based system to predict house prices using the California Housing dataset. The full pipeline—including data preprocessing, EDA, visualization, model training, and evaluation—was implemented using Python and popular ML libraries.

Among all models tested, the Random Forest Regressor achieved the highest accuracy and robustness, making it the most suitable model for practical deployment in real estate applications. The study demonstrates how data-driven approaches can significantly enhance decision-making in housing valuation.

# REFERENCES

[1] Scikit-Learn Documentation

[2] Kaggle: California Housing Dataset

[3] Géron, A., Hands-On Machine Learning with Scikit-Learn and TensorFlow

[4] Python Libraries: pandas, NumPy, seaborn, matplotlib.