## PHASE 3 Submission

**Student Name:** S.Ledavathi

**Register Number:** 510123106025

**Institution:** ADHIPARASAKTHI
COLLAGE OF ENGINEERING

**Department:** B.E ELECTRONICS AND
COMMUNICATION ENGINEERING

Date of Submission:[08-06-2025]

**Github repository link:** https://github.com/Ledavathi-123/orcasting-house-price-accurately-using-smart-regression-techniques-in-data-science.git

# Forecasting house price accurately using smart regression techniques in data science

## Problem Statement

Real estate price prediction is a crucial task in today's property market. Stakeholders such as buyers, sellers and investors often make decisions based on their estimated market value of properties. Traditional valuation methods are often subjective and time-consuming. This project aims to build a data-driven solution using regression techniques in machine learning to predict house prices more accurately and efficiently. The objective is to develop a predictive model using historical housing data that considers features like square footage, location, number of rooms, and amenities to forecast property prices. Since the target variable (SalePrice) is continuous, this is a supervised regression problem.

## 1. Abstract

This project focuses on leveraging regression techniques to train prediction house to predict house prices of house prices based on various property features. The aim is to develop

questions to vividly provide accurate and useful recent estimates of home prices, helping users make informed buying or selling decisions. The dataset used is sourced from Kaggle and contains 80 features related to residential homes in Ames, Iowa. After preprocessing and exploratory analysis, several regression models including Linear Regression, Random Forest, and XGBoost are implemented. XGBoost was identified as the most accurate model, offering an R^score of 0.91. The model is deployed using Streamlit for easy user interaction.

## 3.System Requirements

### Hardware

- **RAM**: 8GB minimum

- **Storage**: 10GB of free disk space

- **Processor**: Intel i5 or AMD Ryzen 5 or higher

### Software

- **Programming language**: Python 3.10+

- **IDE**: Jupyter Notebook or Google Colab for development

- **Required Libraries**: pandas, numpy, seaborn, matplotlib, scikit-learn

- **Web scraping Tools**: BeautifulSoup

- **Deployment platform**: Jupyter notebook or cell

## 2.Objectives

The main objectives of the project are:

- To build a house price prediction model capable of predicting home prices based on various features.

- To explore and analyze relationships between different property characteristics and sale prices.

- To optimize model performance using robust and feature selection or tuning techniques

- To deploy the model and build a user-friendly interface, allowing users to predict prices interactively

- To assist stakeholders in making informed decisions based on data insights

## 3.Flowchart of Project Workflow

Stages of the Project Workflow.

- **Data Collection:** Downloaded from Kaggle

- **Data Preprocessing:** Cleaning missing values, encoding categorical variables

  **Exploratory Data Analysis (EDA):** Visual analysis and correlation checks

  **Feature Engineering:** New features creation, selection of relevant variables

  **Modeling:** Training with multiple regression algorithms

- **Evaluation:** Assessing model accuracy using statistical metrics

- **Deployment:** Building web interface with Streamlit

## FORECASTING HOUSE PRICES USING SMART REGRESSION TECHNIQUES IN DATA SCIENCE

↓

### DATA COLLECTION

↓

### DATA PREPROCESSING

↓

### EXPLORATORY DATA ANALYSIS (EDA)

↓

### FEATURE ENGINEERING

↓

### MODEL SELECTION

↓

### MODEL TRAINING

↓

### HYPERPARAMETER TUNING

↓

### MODEL DEPLOYMENT

↓

### MONITORING & MAINTENANCE

## 4.Dataset Description

https://www.kaggle.com/datasets/jordanin/llm-prompts

- Source:

Scikit-learn (fetch_california_housing)

(Originally from the StatLib repository, published by the UCI Machine Learning Repository)

Publicdataset (sci-wurdd)

- **Size and Structure**
  - **Rows:** 20,640
  - **Columns** 9 (8 features + 1 target)
  - **Target Variable:** Price (Median house value in $100,000s) - **Features Include:**
    - MedInc - Median income/block.
      HouseAge - Median housing age. AveRooms - Average
      number of rooms . AveBedrms -
      Average number of bedrooms
    - Population - Block population. AveOccup -
      Average occupancy - Latitude, Longitude -
      Geographic coordinates

- Sample of Dataset (first five):

  ```
  # For demonstration, we use sklearn's Boston housing
  dataset from sklearn.datasets import
  fetch_california_housing data = fetch_california_housing()
  df=pd.DataFrame(data.data,columns=data.feature_names)
  df['Price'] = data.target
  ```

## OUTPUT:

| MedInc | HouseAge | AveRooms | AveBedrms | Population | AveOccup | Latitude | Longitude | Price |
|--------|----------|----------|-----------|------------|----------|----------|-----------|-------|
| 8.3252 | 41.0 | 6.984127 | 1.023810 | 322.0 | 2.555556 | 37.88 | -122.23 | 4.526 |
| 8.3014 | 21.0 | 6.238137 | 0.971880 | 2401.0 | 2.109842 | 37.86 | -122.22 | 3.585 |
| 7.2574 | 52.0 | 8.288136 | 1.073446 | 496.0 | 2.802260 | 37.85 | -122.24 | 3.521 |
| 5.6431 | 52.0 | 5.817352 | 1.073059 | 558.0 | 2.547945 | 37.85 | -122.25 | 3.413 |
| 3.8462 | 52.0 | 6.281853 | 1.081081 | 565.0 | 2.181467 | 37.85 | -122.25 | 3.422 |

## 3. Data Preprocessing

1. Handle Missing Values

**Numerical Features:** Missing values filled using mean/median imputation.
**Categorical Features:** Filled with the most frequent value or 'None' if applicable.

```
#Fill numeric NaNs with median for col in
df.select_dtypes(include=['number']).columns
    df[col].fillna(df[col].median(),inplace=True)

#Fill categorical NaN with mode or 'None' for col in
df.select_dtypes(include='object').columns
    df[col].fillna(df[col].mode()[0],inplace=True)
```

## 2. Handle Duplicates: df.drop_duplicates(inplace=True)

## 3. Handle Outliers

Used Z-score/IQR method to remove outliers in key numeric columns like GrLivArea, TotalBsmtSF, etc.

```
from scipy import stats
df = df[(np.abs(stats.zscore(df.select_dtypes(include='number')))
    <3).all(axis=1)]
```

## 4. Feature Encoding and Scaling

**Encoding:** Used One-Hot Encoding for categorical variables.
**Scaling:** Used StandardScaler to normalize numerical features.

```
from sklearn.preprocessing import OneHotEncoder,StandardScaler

#One-hot encode categorical features
df_encoded=pd.get_dummies(df,drop_first=True)

# Standard scaling for numerical features scaler = StandardScaler()
num_cols=df_encoded.select_dtypes(include='number').columns
df_encoded[num_cols]=scaler.fit_transform(df_encoded[num_cols])
```

## 5. Before/After Transformation Screenshots Before Cleaning: df.info()
df.describe() df.isnull().sum()

Shows row data, missing values, standard types.

AfterCleaningand Encoding

df.encodeddata(df.encodedhead())

Shows no more null valueDNaNs- goesall currectype, and no-boxes-dataofcolumns.

## 6.Exploratory Data Analysis (EDA)

```
import pandas as pd import seaborn
as sns import matplotlib.pyplot as
plt
from sklearn.datasets import fetch_california_housing

#Loaddata
df = fetch_california_housing(as_frame=True).frame
df.rename(columns={"MedHouseVal": "Price"}, inplace=True)

#histogram df.hist(figsize=(12,8),bins=30,
edgecolor='black') plt.suptitle("Feature
Distributions", y=1.02) plt.tight_layout()
plt.show()

#Boxplot sns.boxplot(data=df)
orient='h') plt.title("Boxplot of
Features") plt.tight_layout()
plt.show()

#Correlation Heatmap
sns.heatmap(df.corr(),annot=True,cmap="coolwarm",fmt=".2f") plt.title("Correlation
Heatmap")
plt.tight_layout()
plt.show()
```

OUTPUT:

## 7. Feature Engineering

1. **New Feature Creation:** Adding meaningful features to improve model performance.
   - Example: Price-per-squarefoot as an average.

2. **Feature Selection:** Removes irrelevant features to reduce overfitting and improve efficiency.
   - Method: Filter, Wrapper, and Embedded (e.g., Lasso, Random Forest).

3. **Transformation Techniques:** Adjust features to improve model fit (scaling, skewness reduction).
   - Techniques: Standardization, Log Transformation, Polynomial features.

4. **Feature Impact:** Understand how features influence predictions.
   - Linear models: Coefficients show feature impact.

```
import pandas as pd, numpy as np from
sklearn.preprocessing import StandardScaler from
sklearn.ensemble import RandomForestRegressor
from sklearn.linear_model import LinearRegression

# Create a dataframe from the features
```

```
data['price_per_sqft']=data['price']/data['square_footage']data['house_age']
=2025-data['year_built']
data['square_footage']=pd.to_numeric(df_to_numeric(data['square_footage']))
data['log_price']=np.log1p(data['price']) + G
```

```
#Feature selection and regression:
rf_model = RandomForestRegressor().fit(data.drop(columns=['price']), data['price'])
importances_features= data.drop(columns=['price']).columns[rf_model.feature_importances_ >
0.05]coefficients=LinearRegressor().fit(data.drop(columns=['price']),
data['price']).coef_
```

```
#Display principal importance, features,
    coefficient
```

## 8. Model Building

- **LinearRegression:** Simple baseline for comparison.
- **RandomForest:** Handles non-linear feature/interactions.
- **GradientBoosting:** Minimizes number/models for super-complexmodels.
- **XGBoost:** Optimized, faster, and more efficient than GradientBoosting.
- **EvaluationMSE:** Compares Mean Squared Error for model performance.
- **ScreenshotOutputs:** Captured MSE for types/each model.
- **Code:**

```
X_train,X_test,y_train,y_test=train_test_split(X,y,test_size=0.2,
random_state=42)
model=RandomForestRegressor(random_state=42)model.fit(X_train,
y_train)
```

## 9. Model Evaluation

1. Metrics for Regression:
   - MSE: Mean Squared Error.
   - RMSE: Root Mean Squared Error.
   - R Score: Coefficient of Determination.
2. Error Analysis: Plot Actual vs Predicted values (scatter plot).
3. ROC Curve: For classification, plot ROC Curve to evaluate performance.
4. ConfusionMatrix: Showing true/false predictions ConfusionMatrix for classification.
5. ModelComparison: Compares models using MSE, RMSE, and R² in a table.
6. Visuals: Create subplots to show how model explores trends and innovativeness.

7. Code:

```
y_pred=model.predict(X_test)print("MeanSquaredError:",
mean_squared_error(y_test, y_pred))
print("R2Score",r2_score(y_test,y_pred))
```

## 10. Deployment

- **Platform:** Deployed using Streamlit Cloud
- **Method:** GitHub repo linked to Streamlit
- **UI Screenshot:** Attach screenshot of app
- **Prediction:** UI accepts features → model predicts price
- **Output Example:** Predicted Price: ₹45,00,000
- **Alternate Options:** Gradio + Hugging Face or Flask + Render
- **Tip:** Use Streamlit/Gradio for quick, free deployment with UI

## 11. Source code

```
import pandas as pd import numpy as np import matplotlib.pyplot as plt import
sklearn.model.sklearn.model_selection import train_test_split GridSearchCV from
sklearn.preprocessing import StandardScaler, PolynomialFeatures from
sklearn.linear_model import LinearRegression, Ridge, Lasso from sklearn.tree
import DecisionTreeRegressor from sklearn.ensemble import
RandomForestRegressor, GradientBoostingRegressor import
XGBoostRegressor from sklearn.svm import SVR from sklearn.metrics import
mean_absolute_error, mean_squared_error, r2_score from sklearn.pipeline import
make_pipeline print("==== Loading Data =====") # load and explore dataset
url =
"https://raw.githubusercontent.com/.../data.csv"
import seaborn as seaborn.csv" data =
pd.read_csv(url) print("Data Shape:", data.shape)"
```

```python
print("39 bedia over ?.print data.head().#Lines 12 4.

Visualization

plt.hist(house[0]_10) # Distribution of house price
plt.x(price) 1, D on boxplot for(? radius_house_value1
sds=True Data.corr plt.title(house Price Distribution)

#Wordshole clear map

plt.subplot(1,1,1)

#plot correlation matrix among_data-
(data.dist_dtypes(include=[number])

#Corps correlation matrix or numeric_data.corr(#t the)
(untmap(nul/extmap(corr), annot=True, cmap= coolwarm),
(vil='.1f')
plt.title('Correlation matrix)#'

Price vs. median income

plt.subplot(1,1,3)
sns.scatterplot(x=median_income,y=radius_house_value data=data,alpha=0.3)
plt.title('Price vs. Income') # Price by income proximity

plt.subplot(1,2,4)
sns.boxplot(x='ocean_proximity',y=median_house_value,data=data)
plt.xticks(rotation=45) plt.title('Price by Location) plt.tight_layout()
plt.show(# print("\n=== Preprocessing Data ===")
#.unif.onscategorical variable.like ocean select_dtypes(which= numeric)onehot),
(nplace=True)

# Feature engineering (half more (_per_household)) =
data[total_rooms]/data['households'] data['bedrooms_per_room'] =
data[total_bedrooms]/data['total_rooms'] # Convert categorical to
numerical using pd.get_dummies(data) (data=data.drop['rooms_proximity'])
```

```python
# Select feature and target
X = data[['feature1', 'feature2', 'feature3']]
data[value_house_value]

# train test split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# feature scaling scaler
scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)
print("\n--- Training Models ---")

{
    "LinearRegression": LinearRegression(),

    "Ridge Regression": Ridge(alpha=1.0),

    "Lasso Regression": Lasso(alpha=0.1),

    "DecisionTree": DecisionTreeRegressor(max_depth=5),

    "RandomForest": RandomForestRegressor(n_estimators=100, random_state=42),

    "GradientBoosting": GradientBoostingRegressor(n_estimators=100, random_state=42),

    "SVR": SVR(kernel='rbf')
}
results = {}
for name, model in models.items():
    print(f"Training {name}...")
    model.fit(X_train_scaled, y_train)  y_pred
    = model.predict(X_test_scaled)   results[name] = {

        "MAE": mean_absolute_error(y_test, y_pred),

        "RMSE": np.sqrt(mean_squared_error(y_test, y_pred)),
```

```python
            'XGB':xgb_accuracy, 'RF':rfc_pred}

}

# Display evaluation results, if …
print("\n===================================== print("\n====
Model Performance ====")
print(results_df.sort_values(by='R2',ascending))

print("\n===== Optimizing Best Model
=====")

# Let's optimize based on our analysis typically performs well
from sklearn.model_selection import RandomizedSearchCV

# Define parameter grid or use RandomizedSearchCV
param_dist = {
    'n_estimators': [50,100,200],
    'max_depth': [None,10,20],
    'min_samples_split': [2,5,10]
}  df =

# rf model = RandomForestRegressor(random_state=42)
random_search=RandomizedSearchCV(param_distributions=param_dist,n_iter=5, cv=3,
scoring='neg_mean_squared_error',n_jobs=-1)
random_search.fit(X_train,y_train)
best_model = random_search.best_estimator_

# Evaluate optimized model y_pred =
best_model.predict(X_test,y_valid)print("Optimized
Model Performance") print("MAE)

{:.2f}".format(mean_absolute_error(y_test,y_pred,0.2f)")
print("RMSE: {}.format(mean_squared_error(y_test,
y_pred,0.2f") print("R2 Score: {r2_score(y_test, …
```

```
y_pred.txt","print">---Generating Visualizations

-->)

#Feature Importance plt.figure(figsize=(12, 8))
importances = best_model.feature_importances_ feature = X.columns
index = np.argsort(importances)[-10:] # Top 15 features plt.title('Feature
Importance') plt.barh(range(len(index)), importances[index], color='b',
align='center') plt.yticks(range(len(index)), [feature[i] for i in index])
plt.xlabel('Relative Importance') plt.show() # Actual vs Predicted
plt.figure(figsize=(12, 8)) plt.scatter(y_test, y_pred, alpha=0.5)
plt.plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], 'k--', lw=2)
plt.xlabel('Actual Prices') plt.ylabel('Predicted Prices') plt.title('Actual vs.
Predicted House Prices') plt.show() # Residual plot residuals = y_test -
y_pred plt.figure(figsize=(10, 6)) plt.scatter(y_pred, residuals, alpha=0.5)
plt.axhline(y=0, color='r', linestyle='--') plt.xlabel('Predicted Prices')
plt.ylabel('Residuals') plt.title('Residual Plot')
plt.show() print(")---Implementation Complete

-->)
```

**OUTPUT:**

## 14. Future scope

- **Geospatial Integration** Enhance location record rate with geo spatial and price to better capture regional price differences.

- **Time-Series Forecasting**
  Additional with housing data to forecast future price values over market trends.

- **Automated Feature Selection** Implement advanced techniques like Recursive Feature Elimination or SHAP for smarter feature optimization.

- **Real-Time Prediction API** Deploy as a REST API connected to live real estate data sources for real time usage.

- **User-Friendly Web Interface** Enhance the model with an interactive UI using Streamlit or Gradio for public use.

## 15. Team Members and Roles

*J. S. AFRIM – Data Collection and Integration: Responsible for sourcing datasets, connecting APIs, and preparing the raw dataset for analysis.*

2. *B. RAMYA – Data Cleaning and EDA:* Cleans and preprocesses data, performs exploratory analysis, and generates initial insights.

3. *T. VAISHNAVI – Feature Engineering and Modeling:* Works on feature extraction and selection; develops and trains machine learning models.

4. *S. LEELAVATHI – Evaluation and Optimization:* Tunes hyperparameters, validates models, and measures performance metrics.

5. *B. HEMALATHA – Documentation and Presentation:* Compiles reports, prepares visualizations, and handles presentation and optional deployment.