# 🏠 House Price Prediction: Step-by-Step Training Guide

This guide provides a clear, practical walkthrough for building a complete house price prediction model — from loading data to saving your trained model as a pickle file. Each step includes what you should do, why it's important, and how it fits into the full workflow.

## 1. Data Loading (Pandas, Numpy)

• Start by importing the necessary Python libraries (pandas, numpy, sklearn, matplotlib, seaborn).
• Load the dataset (e.g., house_prices_dataset.csv) using pandas' `read_csv()`.
• Check the first few rows with `.head()` and inspect summary statistics with `.describe()`.
• Confirm data types and missing values using `.info()` and `.isnull().sum()`.

## 2. Data Cleaning (Pandas, Numpy)

• Handle missing values appropriately — use median or mean for numerical features and mode for categorical ones.
• Drop irrelevant columns such as 'id' or duplicates if any.
• Convert columns to proper data types (e.g., ensure date columns are datetime, categories are categorical).

## 3. Exploratory Data Analysis (EDA)-(Pandas, Numpy, Matplotlib, Seaborn)

• Visualize data distributions (e.g., histograms of prices, scatter plots of price vs. square footage).
• Check correlations between features and target variable using `.corr()`.
• Identify outliers using boxplots or z-score thresholds.
• Summarize key patterns and insights (e.g., which features most affect price).

## 4. Feature Engineering (Scipy-stats, Pandas, Scikit-learn, Statsmodels)

• Create new features such as 'price_per_sqft', 'house_age', or interaction terms.
• Encode categorical variables using One-Hot Encoding or Label Encoding depending on the model.
• Remove multicollinear features if needed. - Use variance inflation factor (VIF >5; drop it)

## 5. Data Splitting (Scikit-learn)

• Split your dataset into training and testing sets using `train_test_split()` from sklearn.
• Typical splits are 80/20 or 70/30.
• Always use a random seed for reproducibility (e.g., `random_state=42`).

## 6. Feature Scaling (Scikit-learn)

• Apply scaling to normalize numerical variables — this ensures all features contribute equally.
• Common scalers include:
  - StandardScaler (mean = 0, std = 1) - standard normal distribution
  - MinMaxScaler (range between 0 and 1)
• Scaling is especially important for algorithms like Linear Regression, KNN, and SVM.

## 7. Model Selection and Training(Scikit-learn, Statsmodels)

• Choose a regression model that suits your data:
  - Linear Regression (for linear relationships)
  - Random Forest or Gradient Boosting (for complex nonlinear patterns)
  - XGBoost or LightGBM for performance-oriented modeling
• Train your model using the `.fit()` method on the training data.

## 8. Model Evaluation (Scikit-learn)

• Use the test dataset to assess model performance.
• Evaluate using metrics like:
  - Mean Absolute Error (MAE)
  - Mean Squared Error (MSE)
  - Root Mean Squared Error (RMSE)
  - $R^2$ score (explained variance)

Hyperparameter tuning - getting the best model parameters
• Visualize actual vs predicted values using scatter plots.

## 9. Model Serialization (Pickle, Joblib)

• Once your model performs well, save it using Pickle so it can be reloaded without retraining.
• Example:
```python
import pickle
pickle.dump(model, open('house_price_model.pkl', 'wb'))
```

```
```
- To load the model again:
  ```python
  model = pickle.load(open('house_price_model.pkl', 'rb'))
  ```

## 10. Deployment and Maintenance
- Integrate the model into an application or API for prediction.
- Ensure consistent preprocessing between training and inference.
- Monitor model performance over time and retrain when necessary.
- Document your pipeline for transparency and reproducibility.

## Summary
This workflow represents an end-to-end pipeline for house price prediction. It emphasizes clean data, meaningful features, proper scaling, and robust evaluation. Saving the model ensures it can be reused in production or integrated into web applications easily.