



DEVELOPERS HUB CORPORATION

PROJECT TITLE: HOUSING PRICE PREDICTION PROJECT

Name: Ayesha Liaqat

Intern ID: DHC-1059

Department: AI / Machine Learning

ayeshaliaqat908@gmail.com

Machine Learning –TASK 2 -

California Housing Price Prediction Project

Objective:

The goal of this project is to walk through the complete machine learning pipeline using the California Housing Prices dataset. It follows the steps outlined in Chapter 2 of *Hands-On Machine Learning with Scikit-Learn*. By the end of the task, we apply all stages of ML including data loading, cleaning, feature engineering, model training, evaluation, and comparison.

California Housing Price Prediction Project Summary

Step 1: Data Loading

Dataset used: *California housing dataset* (CSV or via `fetch_california_housing()`).

Loaded using pandas `read_csv()` function.

Displayed first few rows using `head()` and checked data types, shape, and null values.

Step 2: Exploratory Data Analysis

Plotted histograms for numerical columns like population, rooms, and income.

Checked correlation matrix using seaborn heatmap.

Scatter plots of longitude vs latitude colored by house value to observe geographical trends.

Step 3: Data Cleaning

Handled missing values in 'total_bedrooms' using median imputation.

Converted categorical variable 'ocean_proximity' using label encoding.

Standardized the features using `StandardScaler` to bring them to a common scale.

Step 4: Feature Selection

Selected all columns except the target column `median_house_value`.

Split the data into features (X) and target (y).

Performed train-test split using 80:20 ratio.

Step 5: Model Training

Trained and evaluated multiple models:

Linear Regression

Lasso Regression

Ridge Regression

Used `r2_score`, MAE, MAPE, and MSE to evaluate each model.

Compared models based on performance metrics.

Step 6: Visualization of Predictions

Plotted actual vs predicted prices for each model.

Used line plots to compare first 20 predicted and actual values.

Displayed side-by-side graphs for all three models.

Step 7: Results Summary (Model Comparison)

Created a performance table showing:

R^2 Score

Adjusted R^2

MAE

MAPE

MSE

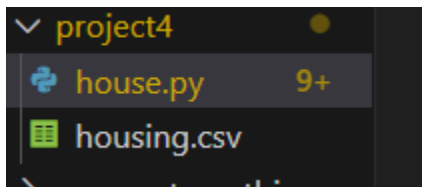
This helped identify which model performed best on unseen test data.

Conclusion:

This project allowed hands-on application of the complete machine learning pipeline. It strengthened understanding of core steps like data preprocessing, model fitting, and evaluation. The comparison between models showed how regularization (Lasso, Ridge) affects performance. This task directly applied the theoretical knowledge from Chapter 2 and provided practical experience with real-world data.

• .

Files and Folders



Libraries

```
# 📦 Import required libraries
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

# 🔧 Preprocessing & model tools
from sklearn.preprocessing import LabelEncoder, StandardScaler
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression, Lasso, Ridge
from sklearn.metrics import r2_score, mean_absolute_error, mean_absolute_percentage_error, mean_squared_error

# 📊 Statsmodels for VIF (optional advanced diagnostics)
from statsmodels.stats.outliers_influence import variance_inflation_factor
from statsmodels.tools.tools import add_constant
```

House.py

```
21 # 🏠 Read dataset
22 df = pd.read_csv(r"f:\artificial intelligence\project4\housing.csv")
23 print(df.head())
24
25 # -----
26 # 📊 Basic data exploration
27 # -----
28 print(f"Shape of dataset: {df.shape}")
29 print("\nMissing values per column:\n", df.isnull().sum())
30 print("\nDuplicate rows:", df.duplicated().sum())
31 print("\nData types:\n", df.dtypes)
32
33 # 🔥 Correlation heatmap
34 plt.figure(figsize=(16, 12))
35 sns.heatmap(df.corr(numeric_only=True), annot=True, cmap='coolwarm')
36 plt.title("Correlation Matrix")
37 plt.show()
38
39 # ✂ Fill missing values
40 df['total_bedrooms'].fillna(df['total_bedrooms'].median(), inplace=True)
41
42 # 🔍 Re-check for missing values
43 print("\nMissing values after fill:\n", df.isnull().sum())
44
45 # -----
46 # 📊 Distributions of numeric features
47 # -----
48 fig, ax = plt.subplots(4, 2, figsize=(14, 12))
49 cols = ['longitude', 'latitude', 'housing_median_age', 'total_rooms',
50         'total_bedrooms', 'population', 'households', 'median_house_value']
51 for i, col in enumerate(cols):
52     sns.histplot(df[col], kde=True, ax=ax[i//2, i%2])
53 plt.tight_layout()
54 plt.show()
55
```

```

46 # 📊 Distributions of numeric features
47 # -----
48 fig, ax = plt.subplots(4, 2, figsize=(14, 12))
49 cols = ['longitude', 'latitude', 'housing_median_age', 'total_rooms',
50         'total_bedrooms', 'population', 'households', 'median_house_value']
51 for i, col in enumerate(cols):
52     sns.histplot(df[col], kde=True, ax=ax[i//2, i%2])
53 plt.tight_layout()
54 plt.show()
55
56 # 📊 Categorical vs Target
57 sns.barplot(x='ocean_proximity', y='median_house_value', data=df, palette="Set1")
58 plt.title("House Value by Ocean Proximity")
59 plt.xticks(rotation=45)
60 plt.show()
61
62 # 📍 House value across geographic coordinates
63 plt.figure(figsize=(12, 8))
64 sns.scatterplot(data=df, x='longitude', y='latitude', hue='median_house_value', palette='coolwarm', alpha=0.5)
65 plt.title("Geographic Distribution of Median House Value")
66 plt.show()
67
68 # 🏷️ Label encode 'ocean_proximity'
69 le = LabelEncoder()
70 df['ocean_proximity'] = le.fit_transform(df['ocean_proximity'])
71
72 # 📊 Re-check Correlation heatmap
73 plt.figure(figsize=(12, 10))
74 sns.heatmap(df.corr(numeric_only=True), annot=True, cmap="Pastel1")
75 plt.title("Correlation Matrix After Encoding")
76 plt.show()
77

```

```

house.py 9+ ×  california-housing-prices-prediction.ipynb
artificial intelligence > project4 > house.py > ...

84 # 🛠️ Train-test split
85 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=123)
86
87 # 📏 Standardize the features
88 scaler = StandardScaler()
89 X_train = scaler.fit_transform(X_train)
90 X_test = scaler.transform(X_test)
91
92 # -----
93 # 📊 Fit Linear, Lasso & Ridge Models
94 # -----
95 # Linear Regression
96 lr = LinearRegression()
97 lr.fit(X_train, y_train)
98 y_pred = lr.predict(X_test)
99
100 # Lasso Regression
101 lasso = Lasso()
102 lasso.fit(X_train, y_train)
103 y_pred_lasso = lasso.predict(X_test)
104
105 # Ridge Regression
106 ridge = Ridge()
107 ridge.fit(X_train, y_train)
108 y_pred_ridge = ridge.predict(X_test)
109
110 # -----
111 # 📏 Metrics Function
112 # -----
113 def get_metrics(p, y, y_pred):
114     n = len(y)
115     r2 = r2_score(y, y_pred)
116     adjusted_r2 = 1 - (((1 - r2) * (n - 1)) / (n - p - 1))
117     mae = mean_absolute_error(y, y_pred)
118     mape = mean_absolute_percentage_error(y, y_pred)
119     mse = mean_squared_error(y, y_pred)
120     return r2, adjusted_r2, mae, mape, mse

```

```

project3      127 # Linear Regression
templates     128 df_lr = pd.DataFrame({'Actual': y_test, 'Predicted': y_pred}).reset_index(drop=True)
              129 axes[0].plot(df_lr[:20])
venv          130 axes[0].set_title("Linear Regression")
              131 axes[0].legend(["Actual", "Predicted"], loc="upper left")
              132
              133 # Lasso Regression
              134 df_lasso = pd.DataFrame({'Actual': y_test, 'Predicted': y_pred_lasso}).reset_index(drop=True)
              135 axes[1].plot(df_lasso[:20])
              136 axes[1].set_title("Lasso Regression")
              137 axes[1].legend(["Actual", "Predicted"], loc="upper left")
              138
segment.py    139 # Ridge Regression
              140 df_ridge = pd.DataFrame({'Actual': y_test, 'Predicted': y_pred_ridge}).reset_index(drop=True)
project4      141 axes[2].plot(df_ridge[:20])
house.py      142 axes[2].set_title("Ridge Regression")
              143 axes[2].legend(["Actual", "Predicted"], loc="upper left")
              144
              145 plt.tight_layout()
              146 plt.show()
              147
              148 # -----
              149 # 📊 Compare Model Performance Metrics
              150 # -----
              151 p = x_train.shape[1] # Number of features
              152 performance_df = pd.DataFrame([
              153     get_metrics(p, y_test, y_pred),
              154     get_metrics(p, y_test, y_pred_lasso),
              155     get_metrics(p, y_test, y_pred_ridge)
              156 ],
              157     columns=['R2', 'Adjusted R2', 'MAE', 'MAPE', 'MSE'],
              158     index=['Linear Regression', 'Lasso Regression', 'Ridge Regression'])
              159
              160 print("\n📊 Model Performance Comparison:")
              161 print(performance_df)

```

Outputs:

```

Arsalan@Aayesha MINGW64 /f/artificial intelligence (main)
$ python -u "f:\artificial intelligence\project4\house.py"
  longitude  latitude  housing_median_age  total_rooms  total_bedrooms  population  households  median_income  median_house_value  ocean_proximity
0   -122.23    37.88           41.0         880.0           129.0        322.0         126.0         8.3252         452600.0        NEAR BAY
1   -122.22    37.86           21.0       7099.0          1106.0       2401.0        1138.0         8.3014         358500.0        NEAR BAY
2   -122.24    37.85           52.0       1467.0           190.0         496.0         177.0         7.2574         352100.0        NEAR BAY
3   -122.25    37.85           52.0       1274.0           235.0         558.0         219.0         5.6431         341300.0        NEAR BAY
4   -122.25    37.85           52.0       1627.0           280.0         565.0         259.0         3.8462         342200.0        NEAR BAY
Shape of dataset: (20640, 10)

Missing values per column:
  longitude      0
  latitude       0
housing_median_age  0
total_rooms      0
total_bedrooms   207
population       0
households       0
median_income    0
median_house_value  0
ocean_proximity  0
dtype: int64

Duplicate rows: 0

Data types:
  longitude      float64
  latitude       float64
housing_median_age  float64
total_rooms       float64
total_bedrooms    float64
population        float64
households        float64
median_income     float64
median_house_value float64
ocean_proximity   object
dtype: object

```

