

✓ Building and Evaluating a Linear Regression Model on Real-World Data

California Housing Price Prediction

Objective: Estimate house prices using linear regression, evaluate using MSE, RMSE, and R^2 , and analyze model behavior.

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.datasets import fetch_california_housing
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error, r2_score
```

✓ 1. Data Loading and Exploration

```
housing = fetch_california_housing(as_frame=True)
df = housing.frame
```

```
print(df.head())
print(df.info())
print(df.describe())
```

```

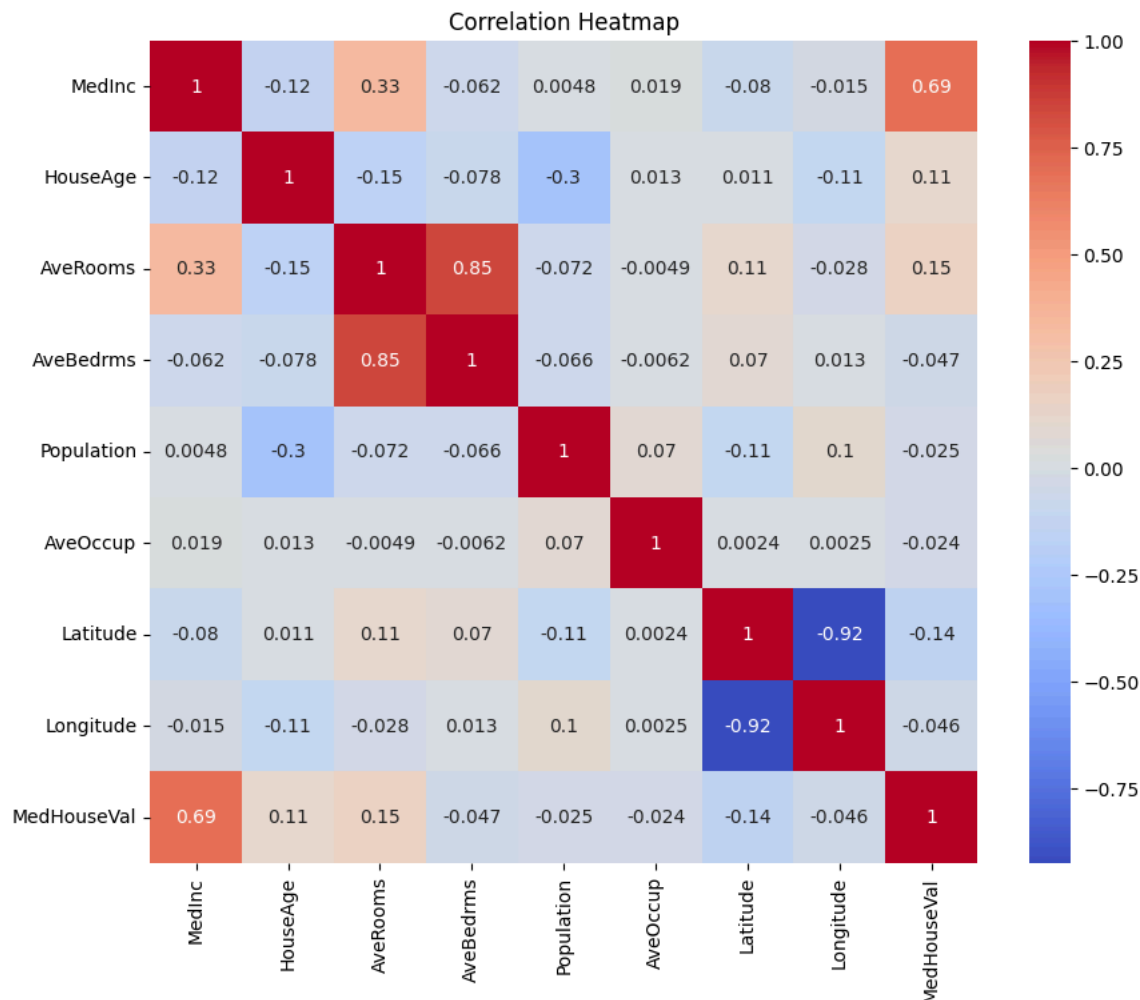
➡      MedInc  HouseAge  AveRooms  AveBedrms  Population  AveOccup  Latitude  \
0    8.3252      41.0    6.984127   1.023810      322.0    2.555556     37.88
1    8.3014      21.0    6.238137   0.971880     2401.0    2.109842     37.86
2    7.2574      52.0    8.288136   1.073446      496.0    2.802260     37.85
3    5.6431      52.0    5.817352   1.073059      558.0    2.547945     37.85
4    3.8462      52.0    6.281853   1.081081      565.0    2.181467     37.85

      Longitude  MedHouseVal
0    -122.23      4.526
1    -122.22      3.585
2    -122.24      3.521
3    -122.25      3.413
4    -122.25      3.422
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 20640 entries, 0 to 20639
Data columns (total 9 columns):
#   Column          Non-Null Count  Dtype
---  -
0   MedInc          20640 non-null   float64
1   HouseAge        20640 non-null   float64
2   AveRooms        20640 non-null   float64
3   AveBedrms       20640 non-null   float64
4   Population      20640 non-null   float64
5   AveOccup        20640 non-null   float64
6   Latitude        20640 non-null   float64
7   Longitude       20640 non-null   float64
8   MedHouseVal     20640 non-null   float64
dtypes: float64(9)
memory usage: 1.4 MB
None

      MedInc      HouseAge      AveRooms      AveBedrms      Population  \
count  20640.000000  20640.000000  20640.000000  20640.000000  20640.000000
mean     3.870671    28.639486     5.429000     1.096675    1425.476744
std     1.899822    12.585558     2.474173     0.473911    1132.462122
min      0.499900     1.000000     0.846154     0.333333     3.000000
25%     2.563400    18.000000     4.440716     1.006079     787.000000
50%     3.534800    29.000000     5.229129     1.048780    1166.000000
75%     4.743250    37.000000     6.052381     1.099526    1725.000000
max    15.000100    52.000000    141.909091    34.066667   35682.000000

      AveOccup      Latitude      Longitude      MedHouseVal
count  20640.000000  20640.000000  20640.000000  20640.000000
mean     3.070655    35.631861    -119.569704     2.068558
std      10.386050     2.135952     2.003532     1.153956
min      0.692308    32.540000    -124.350000     0.149990
25%     2.429741    33.930000    -121.800000     1.196000
50%     2.818116    34.260000    -118.490000     1.797000
75%     3.282261    37.710000    -118.010000     2.647250
max    1243.333333    41.950000    -114.310000     5.000010
```

```
plt.figure(figsize=(10, 8))
sns.heatmap(df.corr(), annot=True, cmap='coolwarm')
plt.title("Correlation Heatmap")
plt.show()
```



2. Data Preparation: Train-Test Split

```
X = df.drop('MedHouseVal', axis=1)
y = df['MedHouseVal']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

3. Model Development: Linear Regression

```
lr = LinearRegression()
lr.fit(X_train, y_train)
print("Coefficients:", lr.coef_)
print("Intercept:", lr.intercept_)

Coefficients: [ 4.48674910e-01  9.72425752e-03 -1.23323343e-01  7.83144907e-01
 -2.02962058e-06 -3.52631849e-03 -4.19792487e-01 -4.33708065e-01]
Intercept: -37.02327770606409
```

4. Prediction and Evaluation

```
y_pred = lr.predict(X_test)
mse = mean_squared_error(y_test, y_pred)
rmse = np.sqrt(mse)
r2 = r2_score(y_test, y_pred)

print(f"MSE: {mse:.2f}")
print(f"RMSE: {rmse:.2f}")
print(f"R² Score: {r2:.2f}")
```

```
MSE: 0.56
RMSE: 0.75
R² Score: 0.58
```

```
residuals = y_test - y_pred
```

```
plt.figure(figsize=(8, 6))  
sns.scatterplot(x=y_test, y=y_pred)  
plt.xlabel("Actual Median House Value")  
plt.ylabel("Predicted Median House Value")  
plt.title("Actual vs Predicted Median House Values")  
plt.show()
```

```
plt.figure(figsize=(8, 6))  
sns.histplot(residuals, kde=True)  
plt.title("Residuals Distribution")  
plt.xlabel("Residuals")  
plt.show()
```

