

Multiple Linear Regression on California Housing Dataset

In this notebook, we will use the **California Housing Dataset** to predict median house values using multiple features through **Multiple Linear Regression**.

```
In [1]: 1 # Importing Libraries
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6
7 from sklearn.datasets import fetch_california_housing
8 from sklearn.model_selection import train_test_split
9 from sklearn.linear_model import LinearRegression
10 from sklearn.metrics import mean_squared_error, r2_score
```

```
c:\users\vamsi2001\appdata\local\programs\python\python39\lib\site-packages\num
py\_distributor_init.py:30: UserWarning: loaded more than 1 DLL from .libs:
c:\users\vamsi2001\appdata\local\programs\python\python39\lib\site-packages\num
py\.libs\libopenblas.EL2C6PLE4ZYW3ECEVIV30XXGRN2NRFM2.gfortran-win_amd64.dll
c:\users\vamsi2001\appdata\local\programs\python\python39\lib\site-packages\num
py\.libs\libopenblas.XWYDX2IKJW2NMTWSFYNGFUWKQU3LYTCZ.gfortran-win_amd64.dll
warnings.warn("loaded more than 1 DLL from .libs:")
```

```
In [3]: 1 # Load the California Housing Dataset
2 california = fetch_california_housing()
3 df = pd.DataFrame(california.data, columns=california.feature_names)
4 df['MedHouseVal'] = california.target
5 df.head()
6 #california
```

```
Out[3]:
```

	MedInc	HouseAge	AveRooms	AveBedrms	Population	AveOccup	Latitude	Longitude	MedHou
0	8.3252	41.0	6.984127	1.023810	322.0	2.555556	37.88	-122.23	
1	8.3014	21.0	6.238137	0.971880	2401.0	2.109842	37.86	-122.22	
2	7.2574	52.0	8.288136	1.073446	496.0	2.802260	37.85	-122.24	
3	5.6431	52.0	5.817352	1.073059	558.0	2.547945	37.85	-122.25	
4	3.8462	52.0	6.281853	1.081081	565.0	2.181467	37.85	-122.25	

```
In [9]: 1 #california
```

About this file

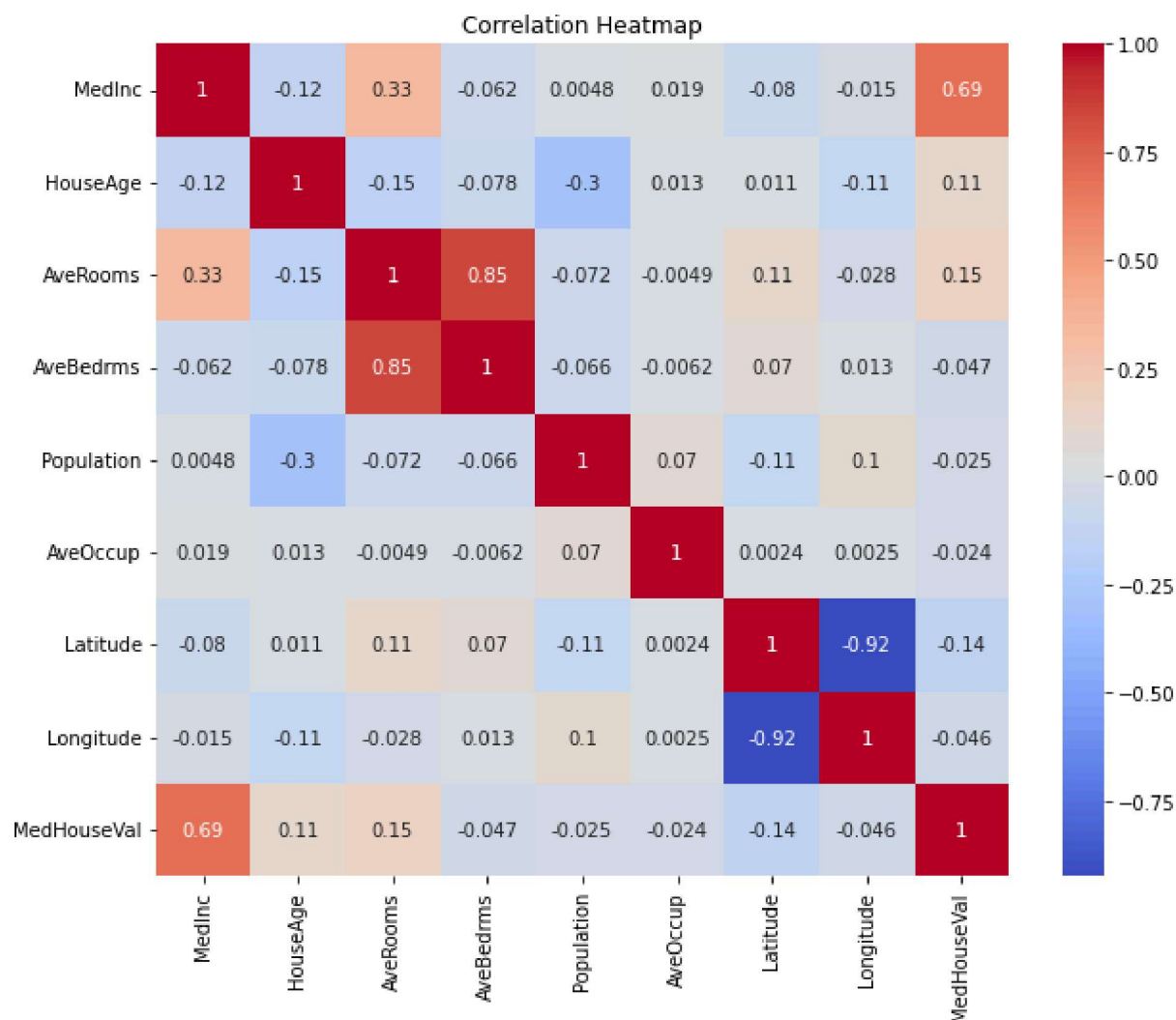
1. longitude: A measure of how far west a house is; a higher value is farther west
2. latitude: A measure of how far north a house is; a higher value is farther north
3. housingMedianAge: Median age of a house within a block; a lower number is a newer building

4. totalRooms: Total number of rooms within a block
5. totalBedrooms: Total number of bedrooms within a block
6. population: Total number of people residing within a block
7. households: Total number of households, a group of people residing within a home unit, for a block
8. medianIncome: Median income for households within a block of houses (measured in tens of thousands of US Dollars)
9. medianHouseValue: Median house value for households within a block (measured in US Dollars)

```
In [4]: 1 #Check for missing values  
        2 df.isnull().sum()
```

```
Out[4]: MedInc          0  
        HouseAge       0  
        AveRooms       0  
        AveBedrms      0  
        Population     0  
        AveOccup       0  
        Latitude       0  
        Longitude      0  
        MedHouseVal     0  
        dtype: int64
```

```
In [11]: 1 # Correlation heatmap
2 plt.figure(figsize=(10, 8))
3 sns.heatmap(df.corr(), annot=True, cmap='coolwarm')
4 plt.title("Correlation Heatmap")
5 plt.show()
```



```
In [5]: 1 # Train-Test Split
2 X = df.drop("MedHouseVal", axis=1)
3 y = df["MedHouseVal"]
4
5 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, ran
```

```
In [6]: 1 # Train the Model
2 model = LinearRegression()
3 model.fit(X_train, y_train)
```

Out[6]: LinearRegression()

```
In [7]: 1 # Predictions
        2 y_pred = model.predict(X_test)
        3 y_pred
```

```
Out[7]: array([0.71912284, 1.76401657, 2.70965883, ..., 4.46877017, 1.18751119,
                2.00940251])
```

```
In [8]: 1 # Model Evaluation
        2 mse = mean_squared_error(y_test, y_pred)
        3 rmse = np.sqrt(mse)
        4 r2 = r2_score(y_test, y_pred)
        5
        6 print(f"RMSE: {rmse}")
        7 print(f"R-squared: {r2}")
```

```
RMSE: 0.7455813830127761
```

```
R-squared: 0.5757877060324511
```

```
In [11]: 1 # Coefficients of the Model
        2 coefficients = pd.DataFrame({
        3     "Feature": X.columns,
        4     "Coefficient": model.coef_
        5 })
        6 coefficients
```

```
Out[11]:
```

	Feature	Coefficient
0	MedInc	0.448675
1	HouseAge	0.009724
2	AveRooms	-0.123323
3	AveBedrms	0.783145
4	Population	-0.000002
5	AveOccup	-0.003526
6	Latitude	-0.419792
7	Longitude	-0.433708

```
In [10]: 1 model.intercept_
```

```
Out[10]: -37.02327770606416
```

```
In [ ]: 1
```