

Instructions: Read the comments in code cells as you complete the problems. Submit your completed notebook and PDF files to Blackboard.

```
In [13]: 1 #0 Run this cell
2 import numpy as np
3 import pandas as pd
4 import seaborn as sns
5 import matplotlib.pyplot as plt
6 from sklearn.model_selection import train_test_split
7 from sklearn.linear_model import LinearRegression
8 from sklearn.model_selection import cross_val_score
```

```
In [14]: 1 #0 Run this cell and enter your first name when prompted.
2 # Do not modify the code here.
3 name = input("Enter your first name: ")
4 url = "https://raw.githubusercontent.com/babdelfa/ML/main/california_housing.csv"
5 name = pd.read_csv(url)
6 print("\nSample of the data: \n")
7 print(name.head())
8 print('The data has been loaded. Refer to the dataframe using your name.')
9 print('\nBegin the analysis below.')
```

Enter your first name: jij

Sample of the data:

	MedInc	HouseAge	AveRooms	AveBedrms	Population	AveOccup	Latitude	\
0	3.2083	34.0	5.471264	1.063218	1627.0	3.116858	32.78	
1	5.7049	9.0	6.699571	1.085837	731.0	3.137339	34.66	
2	1.6125	52.0	3.135135	1.364865	286.0	1.932432	37.78	
3	5.2586	17.0	6.945035	1.131206	1809.0	3.207447	38.69	
4	5.0380	9.0	5.428415	0.967213	2581.0	2.820765	38.42	

	Longitude	MedHouseValue
0	-115.56	0.762
1	-118.17	1.732
2	-122.40	1.125
3	-121.26	1.370
4	-122.79	1.856

The data has been loaded. Refer to the dataframe using your name.

Begin the analysis below.

```
In [15]: 1 #1 Provide summary statistics about the data using the describe method.
2 name.describe()
```

Out[15]:

	MedInc	HouseAge	AveRooms	AveBedrms	Population	AveOccup	Latitude	Long
count	4128.000000	4128.000000	4128.000000	4128.000000	4128.000000	4128.000000	4128.000000	4128.00
mean	3.928064	28.267684	5.436901	1.091814	1421.930717	2.988522	35.644944	-119.58
std	1.904240	12.485409	1.903414	0.351354	1051.047314	3.698413	2.150855	1.99
min	0.499900	1.000000	1.885057	0.444444	8.000000	0.970588	32.560000	-124.23
25%	2.612150	18.000000	4.488618	1.004777	795.000000	2.432214	33.930000	-121.79
50%	3.595100	28.000000	5.300644	1.048780	1173.000000	2.815526	34.270000	-118.51
75%	4.813000	37.000000	6.083468	1.101280	1752.000000	3.265220	37.700000	-118.00
max	15.000100	52.000000	47.515152	11.181818	11935.000000	230.172414	41.860000	-114.56

```
In [16]: 1 #2 Provide a correlation matrix on the dataframe using the corr method.
        2 name.corr()
```

Out[16]:

	MedInc	HouseAge	AveRooms	AveBedrms	Population	AveOccup	Latitude	Longitude
MedInc	1.000000	-0.098871	0.430025	-0.079185	-0.009787	0.016813	-0.054324	-0.044666
HouseAge	-0.098871	1.000000	-0.179765	-0.097652	-0.311671	0.031231	0.004006	-0.111073
AveRooms	0.430025	-0.179765	1.000000	0.742803	-0.094080	0.021448	0.137633	-0.059797
AveBedrms	-0.079185	-0.097652	0.742803	1.000000	-0.077782	0.005210	0.078649	0.010896
Population	-0.009787	-0.311671	-0.094080	-0.077782	1.000000	0.109704	-0.131691	0.132512
AveOccup	0.016813	0.031231	0.021448	0.005210	0.109704	1.000000	-0.015511	0.027616
Latitude	-0.054324	0.004006	0.137633	0.078649	-0.131691	-0.015511	1.000000	-0.921209
Longitude	-0.044666	-0.111073	-0.059797	0.010896	0.132512	0.027616	-0.921209	1.000000
MedHouseValue	0.683157	0.133851	0.186472	-0.066514	-0.034604	-0.049769	-0.148020	-0.048220

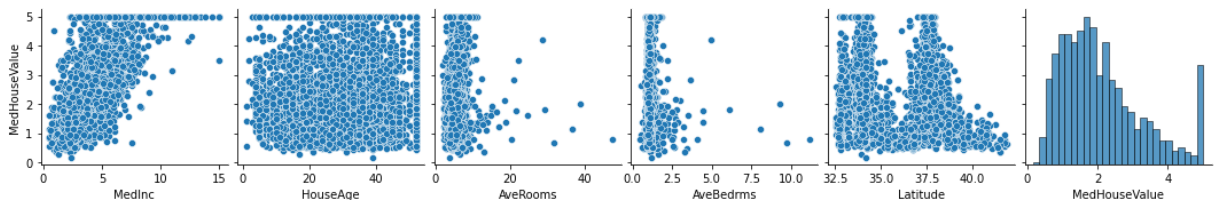
```
In [17]: 1 #3 Review median house value correlaton coefficients with the other variables
        2 # Using the drop method, remove the three variables with the weakest correlati
        3
        4
        5 drop = ["Population", "AveOccup", "Longitude"]
        6 name.drop(columns=drop, inplace=True)
        7 print(name)
```

	MedInc	HouseAge	AveRooms	AveBedrms	Latitude	MedHouseValue
0	3.2083	34.0	5.471264	1.063218	32.78	0.762
1	5.7049	9.0	6.699571	1.085837	34.66	1.732
2	1.6125	52.0	3.135135	1.364865	37.78	1.125
3	5.2586	17.0	6.945035	1.131206	38.69	1.370
4	5.0380	9.0	5.428415	0.967213	38.42	1.856
...
4123	5.7480	28.0	6.002427	1.029126	37.63	3.021
4124	2.0625	52.0	4.323353	1.026946	34.01	1.344
4125	2.1944	33.0	4.131313	0.984848	34.04	2.033
4126	5.6360	49.0	5.117647	0.894118	34.16	2.679
4127	3.3438	25.0	4.297674	1.108140	33.84	1.533

[4128 rows x 6 columns]

```
In [21]: 1 #4 Use seaborn's pairplot function on the data showing median house value on t
        2 sns.pairplot(name, y_vars = "MedHouseValue")
```

Out[21]: <seaborn.axisgrid.PairGrid at 0x7fb3078a6fd0>



```
In [27]: 1 #5 Use a for statement to run a simple learning regression model for each of the
2 # features predicting the target (median house value)
3
4 for x in name.columns[:-1]:
5     X = name[x] #feature variable
6     y = name.MedHouseValue #target variable
7     reg1 = LinearRegression() #instantiate the model
8     reg1.fit(X, y) #Train the model using X and y
9     print(x, reg1.score(X, y))
```

```
MedInc 0.46670303605951935
HouseAge 0.017916166265603772
AveRooms 0.03477164709990799
AveBedrms 0.004424072299616277
Latitude 0.021909951040249176
```

```
In [44]: 1 #6 Use sklearn to do a train_test_split on the data
2 # random_state to 433.
3 # model using the LinearRegression estimator.
4 # using the train data.
5 # 50% train data.
6 # 50% test data.
7
8 X_train, y_train = train_test_split(name[["MedInc", "HouseAge", "AveRooms", "AveBedrms",
9 "Latitude"], name.MedHouseValue, random_state=433)
10 reg2 = LinearRegression() # instantiate the model using the estimator
11 reg2.fit(X_train, y_train) # train the model
12 score_train = reg2.score(X_train, y_train)
13 score_test = reg2.score(X_test, y_test)
```

```
Train: 0.5572196613842659
Test: 0.5229219735139651
```

```
In [45]: 1 X_test[:1]
```

Out[45]:

	MedInc	HouseAge	AveRooms	AveBedrms	Latitude
403	3.75	29.0	5.318841	1.123188	34.04

```
In [47]: 1 #7 Use the predict method to find the target value (i.e., median house value)
2 # the features' values in the first row of X_test
3 reg2.predict([[3.75, 29.0, 5.318841, 1.123188, 34.04]])
```

Out[47]: array([2.14247342])

```
In [39]: 1 y_test[:1]
```

Out[39]: 403 1.906
Name: MedHouseValue, dtype: float64

```
In [41]: 1 #8 What is the amount difference between the predicted median house value from
2 # and the observed median house value? (found in y_test)
3 difference = 1.906 - 2.14247342
4 print(difference)
```

-0.23647342000000005

```
In [48]: sume a model where the target variable is the MedHouseValue and the remaining columns
a 20-Fold cross-validation showing the ten scores. Also show the average of the ten scores.
reg3 = LinearRegression()
results = cross_val_score(reg3, name[["MedInc", "HouseAge", "AveRooms", "AveBedrms", "Latitude", "Longitude"], y, cv=5)
results.mean()
```

```
[0.57357707 0.56383348 0.54648759 0.55034607 0.59756961 0.54713964
 0.51005965 0.54130444 0.52998809 0.49594595]
```

```
Out[48]: 0.5456251600969492
```

```
In [50]: 1 #10 What is the amount difference between the model score using the train-test
2 # the cross-validation approach?
3 difference1 = 0.57357707 - 0.5456251600969492
4 print(difference1)
```

```
0.027951909903050853
```

```
1 ##### 11 Between reg2 and reg3, which model seems to be better in predicting median house
2 values? Justify your answer.
```

```
2
3 reg2>reg3 is better because it is closer to 1
```

```
1 ##### 12 Based on your analysis from # 5, which feature seems to have less weight in
2 predict the target? Justify your answer.
```

```
2
3 Avebedroom. because it is closet to 0
```