

California Housing Price Prediction

```
In [56]: import pandas as pd

In [57]: import matplotlib.pyplot as plt
import seaborn as sns

In [58]: # 1. Load the data :

In [59]: # Read the "housing.csv" file from the folder into the program.

In [60]: df = pd.read_excel('1553768847_housing.xlsx')

In [61]: df

Out[61]:
(20640, 10)

In [62]: # Print first few rows of this data.

In [63]: df.head()

Out[63]:
   longitude  latitude  housing_median_age  total_rooms  total_bedrooms  population  households  median_income  ocean_proximity  median_house_value
0   -122.231  37.85    41                880          129.0      322      126      8.3252      NEAR BAY      452500
1   -122.24  37.85    21                7099          1106.0     2401     1138      8.3014      NEAR BAY      385600
2   -112.21  37.85    52                1467           190.0      496     177      7.2574      NEAR BAY      352100
3   -122.25  37.85    52                1274           235.0     558     219      5.6431      NEAR BAY      341300
4   -122.25  37.85    16                1627           280.0     565     259      3.8462      NEAR BAY      342200

In [64]: # To plot a histogram to understand the data
plt.hist(bins=50, figsize=(20,15))
plt.show()

longitude
0
500
1000
1500
2000
2500
3000
3500
4000
-134
-132
-120
-118
-114

latitude
0
500
1000
1500
2000
2500
3000
34
36
38
40
42

housing_median_age
0
400
800
1200
1600
2000
2400
2800
3200
3600
4000

total_rooms
0
1000
2000
3000
4000
5000
5000
10000
15000
20000
25000
30000
35000
40000

total_bedrooms
0
1000
2000
3000
4000
5000
5000
10000
15000
20000
25000
30000
35000
40000

population
0
2000
4000
6000
8000
10000
12000
14000
16000
18000
20000
22000
24000
26000
28000
30000
32000
34000
36000
38000
40000

households
0
1000
2000
3000
4000
5000
5000
10000
15000
20000
25000
30000
35000
40000

median_income
0
200
400
600
800
1000
1200
1400
1600
1800
2000
2200
2400
2600
2800
3000
3200
3400
3600
3800
4000
4200
4400
4600
4800
5000

median_house_value
0
200
400
600
800
1000
1200
1400
1600
1800
2000
2200
2400
2600
2800
3000
3200
3400
3600
3800
4000
4200
4400
4600
4800
5000

In [65]: import seaborn as sns
sns.jointplot(x="median_income", y="median_house_value", hue="ocean_proximity", data=df);

median_income
0
2000
4000
6000
8000
10000
12000
14000
16000
18000
20000
22000
24000
26000
28000
30000
32000
34000
36000
38000
40000
42000
44000
46000
48000
50000
median_house_value
0
100000
200000
300000
400000
500000
600000
700000
ocean_proximity
• NEAR BAY
• IN OCEAN
• NEAR OCEAN
• ISLAND
• NEAR OCEAN ISLAND

In [66]: df.tail()

Out[66]:
   longitude  latitude  housing_median_age  total_rooms  total_bedrooms  population  households  median_income  ocean_proximity  median_house_value
20635 -121.09  39.48      25                1655           374.0      845     330      1.5603      INLAND      78100
20636 -121.21  39.49      18                697           150.0      356     114      2.5568      INLAND      77100
20637 -121.22  39.43      17                2254           485.0     1007     433      1.7000      INLAND      92300
20638 -121.32  39.43      18                1860           409.0      741     349      1.8672      INLAND      84700
20639 -121.24  39.37      16                2785           616.0     1387     530      2.3886      INLAND      89400

Predict the median housing price

In [67]: df.corr()

Out[67]:
   longitude  latitude  housing_median_age  total_rooms  total_bedrooms  population  households  median_income  median_house_value
longitude  1.000000  -0.246664  -0.108187  0.045668  0.069050  0.099773  0.055310  -0.015176  -0.049957
latitude   -0.246664  1.000000  0.011173  -0.081205  -0.066983  -0.106785  -0.071035  -0.079809  -0.144180
housing_median_age -0.108187  0.011173  1.000000  -0.361262  -0.302051  -0.296044  -0.302916  -0.119034  0.105623
total_rooms  0.045668  -0.081205  -0.361262  1.000000  0.393090  0.877126  0.918484  0.188050  0.134153
total_bedrooms  0.069050  -0.066983  -0.320051  0.393090  1.000000  0.877747  0.979728  0.007722  0.046961
population    0.099773  -0.106785  -0.296044  0.877126  0.877747  1.000000  0.979728  0.004834  -0.024650
households    0.055310  -0.071035  -0.302051  0.918484  0.979728  0.977222  1.000000  0.013033  0.005843
median_income -0.015176  -0.079809  -0.119034  0.188050  -0.007723  0.004834  0.013033  1.000000  0.688075
median_house_value -0.049957  -0.144180  0.105623  0.134153  0.046961  -0.024650  0.005843  0.688075  1.000000

In [68]: df.columns

Out[68]:
Index(['longitude', 'latitude', 'housing_median_age', 'total_rooms',
      'total_bedrooms', 'population', 'households', 'median_income',
      'ocean_proximity', 'median_house_value'],
      dtype='object')

In [69]: df.isnull().sum()

Out[69]:
longitude      0
latitude       0
housing_median_age  0
total_rooms    0
total_bedrooms 267
population     0
households     0
median_income  0
ocean_proximity 0
median_house_value 0
dtype: object

In [70]: # Fill missing values with mean

In [71]: df = df.fillna(df.mean())

In [72]: df.isnull().sum()

Out[72]:
longitude      0
latitude       0
housing_median_age  0
total_rooms    0
total_bedrooms 0
population     0
households     0
median_income  0
ocean_proximity 0
median_house_value 0
dtype: object

In [73]: df.dtypes

Out[73]:
longitude      float64
latitude       float64
housing_median_age  int64
total_rooms    int64
total_bedrooms int64
population     int64
households     int64
median_income  float64
ocean_proximity object
median_house_value  float64
dtype: object

In [74]: # select categorical data

In [75]: df_obj = df.select_dtypes(include=['object']).copy()

In [76]: df_obj

Out[76]:
   ocean_proximity
0      NEAR BAY
1      NEAR BAY
2      NEAR BAY
3      NEAR BAY
4      NEAR BAY
...
20635      INLAND
20636      INLAND
20637      INLAND
20638      INLAND
20639      INLAND
20640 rows x 1 columns

In [77]: # convert categorical data to dummy variables

In [78]: df_dummy = pd.get_dummies(df[['ocean_proximity']])

In [79]: df_dummy

Out[79]:
   ocean_proximity<IN OCEAN  ocean_proximity<INLAND  ocean_proximity<ISLAND  ocean_proximity<NEAR BAY  ocean_proximity<NEAR OCEAN
0                0                0                0                1                0
1                0                0                0                1                0
2                0                0                0                1                0
3                0                0                0                1                0
4                0                0                0                1                0
...
20635                0                1                0                0                0
20636                0                1                0                0                0
20637                0                1                0                0                0
20638                0                1                0                0                0
20639                0                1                0                0                0
20640 rows x 5 columns

In [80]: df_num = df.select_dtypes(include=['int64', 'float64']).copy()
df_num

Out[80]:
   longitude  latitude  housing_median_age  total_rooms  total_bedrooms  population  households  median_income  median_house_value  ocean_proximity<IN OCEAN  ocean_proximity<INLAND
0   -122.23  37.88      41                880          129.0      322      126      8.3252      452500                0                0
1   -122.22  37.86      21                7099          1106.0     2401     1138      8.3014      385600                0                0
2   -122.24  37.85      52                1467           190.0      496     177      7.2574      352100                0                0
3   -122.25  37.85      52                1274           235.0     558     219      5.6431      341300                0                0
4   -122.25  37.85      16                1627           280.0     565     259      3.8462      342200                0                0
...
20635 -121.09  39.48      25                1655           374.0      845     330      1.5603      78100                0                0
20636 -121.21  39.49      18                697           150.0      356     114      2.5568      77100                0                0
20637 -121.22  39.43      17                2254           485.0     1007     433      1.7000      92300                0                0
20638 -121.32  39.43      18                1860           409.0      741     349      1.8672      84700                0                0
20639 -121.24  39.37      16                2785           616.0     1387     530      2.3886      89400                0                0
20640 rows x 14 columns

In [81]: df = pd.concat([df_num, df_dummy], axis=1)

In [82]: df

Out[82]:
   longitude  latitude  housing_median_age  total_rooms  total_bedrooms  population  households  median_income  median_house_value  ocean_proximity<IN OCEAN  ocean_proximity<INLAND
0   -122.23  37.88      41                880          129.0      322      126      8.3252      452500                0                0
1   -122.22  37.86      21                7099          1106.0     2401     1138      8.3014      385600                0                0
2   -122.24  37.85      52                1467           190.0      496     177      7.2574      352100                0                0
3   -122.25  37.85      52                1274           235.0     558     219      5.6431      341300                0                0
4   -122.25  37.85      16                1627           280.0     565     259      3.8462      342200                0                0
...
20635 -121.09  39.48      25                1655           374.0      845     330      1.5603      78100                0                0
20636 -121.21  39.49      18                697           150.0      356     114      2.5568      77100                0                0
20637 -121.22  39.43      17                2254           485.0     1007     433      1.7000      92300                0                0
20638 -121.32  39.43      18                1860           409.0      741     349      1.8672      84700                0                0
20639 -121.24  39.37      16                2785           616.0     1387     530      2.3886      89400                0                0
20640 rows x 14 columns

In [83]: df.corr()

Out[83]:
   longitude  latitude  housing_median_age  total_rooms  total_bedrooms  population  households  median_income  median_house_value  ocean_proximity<IN OCEAN  ocean_proximity<INLAND
longitude  1.000000  -0.246664  -0.108187  0.045668  0.069050  0.099773  0.055310  -0.015176  -0.049957  -0.321121  -0.446969
latitude   -0.246664  1.000000  0.011173  -0.081205  -0.066983  -0.106785  -0.071035  -0.079809  -0.144180  -0.446969  -0.105623
housing_median_age -0.108187  0.011173  1.000000  -0.361262  -0.302051  -0.296044  -0.302916  -0.119034  0.105623  -0.446969  -0.000301
total_rooms  0.045668  -0.081205  -0.361262  1.000000  0.393090  0.877126  0.918484  0.188050  0.134153  -0.000301  -0.000301
total_bedrooms 0.069050  -0.066983  -0.320051  0.393090  1.000000  0.879120  0.974725  0.007682  0.046961  -0.000301  -0.018220
population    0.099773  -0.106785  -0.296044  0.879120  0.879120  1.000000  0.974725  0.004834  -0.024650  0.007682  -0.018220
households    0.055310  -0.071035  -0.302051  0.918484  0.974725  0.977222  1.000000  0.013033  0.005843  0.007682  -0.018220
median_income -0.015176  -0.079809  -0.119034  0.188050  0.074725  0.907222  0.004834  1.000000  0.688075  0.007682  -0.018220
median_house_value -0.049957  -0.144180  0.105623  0.134153  0.046961  -0.024650  0.005843  0.688075  1.000000  0.688075  -0.018220
ocean_proximity<INLAND 0.321121  -0.446969  0.045306  -0.000301  -0.000301  0.018220  0.074613  0.024650  0.168876  0.256617  1.000000
ocean_proximity<NEAR OCEAN 0.055675  0.351166  -0.236645  0.025624  -0.005433  -0.020732  -0.029402  -0.237496  -0.446969  -0.607669  -0.013872
ocean_proximity<NEAR BAY 0.047489  0.258771  0.255172  -0.023022  -0.019785  -0.000880  -0.010093  0.056197  0.160238  -0.314813  -0.013872
ocean_proximity<NEAR OCEAN 0.045509  -0.160818  0.021622  -0.009175  0.000676  -0.024264  0.001714  0.027344  0.141862  -0.342620  -0.013872

In [84]: # Finding the correlation
corr_matrix = df.corr()
corr_matrix["median_house_value"].sort_values(ascending=False)

Out[84]:
median_house_value      1.000000
ocean_proximity<IN OCEAN  0.258617
ocean_proximity<NEAR BAY  0.168284
ocean_proximity<NEAR OCEAN  0.141862
total_rooms              0.134153
housing_median_age       0.105623
households               0.068843
median_income            0.023415
ocean_proximity<ISLAND   0.023415
population               0.024650
longitude                -0.045967
latitude                 -0.144180
ocean_proximity<INLAND  -0.484859
name: median_house_value, dtype: float64

In [85]: df.columns

Out[85]:
Index(['longitude', 'latitude', 'housing_median_age', 'total_rooms',
      'total_bedrooms', 'population', 'households', 'median_income',
      'median_house_value', 'ocean_proximity<IN OCEAN', 'ocean_proximity<INLAND',
      'ocean_proximity<NEAR BAY', 'ocean_proximity<NEAR OCEAN'],
      dtype='object')

In [86]: sns.jointplot(df["median_income"], df["median_house_value"], kind='reg')

C:\Users\srivivas\anaconda3\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following variables as keyword args: x, y. From version 0.8.12, the only valid positional argument will be 'data', and passing other arguments without an explicit keyword will result in an error or misinterpretation.
warnings.warn()

Out[86]:
<seaborn.axisgrid.JointGrid at 0x27dc17fb138>

median_house_value
700000
600000
500000
400000
300000
200000
100000
0
2 4 6 8 10 12 14
longitude

In [87]: df.plot(kind='scatter', x='longitude', y='latitude', figsize=(8,6), alpha=0.1)

Out[87]:
<AxesSubplot: xlabel='longitude', ylabel='latitude'>

longitude
-134
-132
-120
-118
-116
latitude
34
36
38
40
42

In [88]: df.plot(kind='scatter', x='median_income', y='median_house_value', alpha=0.1)

Out[88]:
<AxesSubplot: xlabel='median_income', ylabel='median_house_value'>

median_house_value
500000
400000
300000
200000
100000
0
0 2 4 6 8 10 12 14
median_income

In [89]: # Extract input (X) and output (Y) data from the dataset.

In [90]: x = df[['median_income', 'total_rooms', 'housing_median_age', 'total_bedrooms', 'households', 'population', 'ocean_proximity<IN OCEAN', 'ocean_proximity<INLAND', 'ocean_proximity<ISLAND', 'ocean_proximity<NEAR BAY', 'ocean_proximity<NEAR OCEAN']]
y = df["median_house_value"]

In [91]: X.head()

Out[91]:
   median_income  total_rooms  housing_median_age  total_bedrooms  households  population  ocean_proximity<IN OCEAN  ocean_proximity<INLAND  ocean_proximity<ISLAND  ocean_proximity<NEAR BAY  ocean_proximity<NEAR OCEAN
0      8.3252      880          41                129.0      322      126      0                0                0                1                0
1      8.3014      7099          21                1106.0     2401     1138      0                0                0                1                0
2      7.2574      1467           52                190.0      496      177      0                0                0                1                0
3      5.6431      1274           52                235.0     558      219      0                0                0                1                0
4      3.8462      1627           16                280.0     565      259      0                0                0                1                0

In [92]: y.head()

Out[92]:
0      452500
1      385600
2      352100
3      341300
4      342200
Name: median_house_value, dtype: int64

In [93]: # pip install sklearn

In [94]: from sklearn.model_selection import train_test_split

In [95]: # need to 'random' split data in train and test set
X_train, X_test, y_train, y_test = train_test_split(X, y, train_size = 0.8, random_state= 100)

# 80% data in train
# 20% data in test

In [96]: print(X_train.shape)
print(X_test.shape)
print(y_train.shape)
print(y_test.shape)

(16512, 13)
(4128, 13)
(16512,)
(4128,)

In [97]: #standardize training and test datasets.
=====
# feature scaling is to bring all the independent variables in a dataset into
# same scale, to avoid any variable dominating the model. Here we will not
# transform the dependent variables.
=====
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.fit_transform(X_test)
print(X_train[0:5,:])
print(X_test[0:5,:])

[[[-0.24995448 -0.67047275 -1.24008298 -1.02650858 -1.15100777 -1.06404024
  -0.88684891 -1.47409087 -0.01556621 -0.35650943 -0.39080992 -1.55444193
  2.2677527]
 [-0.58988168 -0.63398882 -0.84375845 -0.42382888 -0.5311873  0.26233518
  -0.89648491 -1.47409087 -0.01556621 -0.35650943 -0.39080992 -1.55444193
  2.2677527]
 [-0.58988168 -0.63398882 -0.84375845 -0.42382888 -0.5311873  0.26233518
  -0.89648491 -1.47409087 -0.01556621 -0.35650943 -0.39080992 -1.55444193
  2.2677527]
 [-0.58988168 -0.63398882 -0.84375845 -0.42382888 -0.5311873  0.26233518
  -0.89648491 -1.47409087 -0.01556621 -0.35650943 -0.39080992 -1.55444193
  2.2677527]
 [-0.58988168 -0.63398882 -0.84375845 -0.42382888 -0.5311873  0.26233518
  -0.89648491 -1.47409087 -0.01556621 -0.35650943 -0.39080992 -1.55444193
  2.2677527]]

test data
[[[-0.98451931 -0.78469094 -0.58304336 -0.85298494 -0.85014856 -0.86294698
  1.12860781 -0.6738423 -0.01556621 -0.35650943 -0.39080992 0.7366906
  -0.85400486]
 [-1.38999368 -0.46645129 1.85130452 -0.29218887 -0.26674345 -0.23516715
  -0.88648491 -0.6738423 -0.01556621 2.80497487 -0.39080992 -1.34627085
  0.01802521]
 [-0.58988168 -0.63398882 -0.84375845 -0.42382888 -0.5311873  0.26233518
  -0.89648491 -1.47409087 -0.01556621 -0.35650943 -0.39080992 0.7366906
  -0.85400486]
 [-0.58988168 -0.63398882 -0.84375845 -0.42382888 -0.5311873  0.26233518
  -0.89648491 -1.47409087 -0.01556621 -0.35650943 -0.39080992 0.7366906
  -0.85400486]
 [-0.58988168 -0.63398882 -0.84375845 -0.42382888 -0.5311873  0.26233518
  -0.89648491 -1.47409087 -0.01556621 -0.35650943 -0.39080992 0.7366906
  -0.85400486]]

from sklearn.linear_model import LinearRegression
lm = LinearRegression()
lm.fit(X_train, y_train) # training

In [98]: lm.predict(X_test)

Out[99]:
print("Intercept is "+str(lm.intercept_))
print("coefficients is "+str(lm.coef_))

Intercept is 207155.1584226552
coefficients is [ 73137.11024933 -9970.89977627 13187.82226366 30756.38131542
 26782.72518277 43817.43879541 6862.65805594 -12587.73930296
 3077.35953444 2970.13993194 5854.39548256 -53228.96689143
-532975.09162748]

In [100]: # predict the price of house on test data
prediction = lm.predict(X_test)

In [101]: print(prediction)

[225683.93415973 157719.56928823 182780.86329646 ... 389946.73321349
180793.1203202 178805.9889341]

In [102]: # Actual value
y_test

Out[102]:
8151    245380
53      184280
3039    146280
9484    154280
9307    324280
12743    117280
5264      50800
13774    238300
13662     85780
11942    130800
Name: median_house_value, Length: 4128, dtype: int64

In [103]: ...
error = Actual - prediction
rmse = sqrt(mean_squared_error(y_test, prediction))
mape_test = sqrt(MSE_test)

# error metrics for test data
MSE_test = mean_squared_error(y_test, prediction)
RMSE_test = sqrt(MSE_test)
MAPE_test = mean_absolute_percentage_error(y_test, prediction)

print("Root Mean square error for test data is :", MSE_test)
print("Root Mean square error for train data is :", RMSE_train)
print("Mean Absolute error for test data is :", MAE_test)
print("Mean Absolute percentage error for test data is :", MAPE_test)

Mean square error for test data is : 458065795.673236
Root Mean square error for test data is : 67678.83316687138
Mean Absolute error for test data is : 49052.19405515740
Mean Absolute percentage error for test data is : 0.28348997603528245

In [104]: # predict the price of house on train data
prediction_train = lm.predict(X_train)

In [105]: prediction_train

Out[105]:
array([[132441.18647085, 277074.8728664, 188272.5795254, ...,
       1404641.85311888, 171218.58837806, 278179.79628821]])

In [106]: # error metric for train data
error = Actual - prediction
Actual = y_train
Prediction = prediction_train
rmse_train = sqrt(MSE_train)
mape_train = sqrt(MAPE_train)

# error metrics for train data
MSE_train = mean_squared_error(y_train, prediction_train)
RMSE_train = sqrt(MSE_train)
MAPE_train = mean_absolute_percentage_error(y_train, prediction_train)

print("Mean square error for train data is :", MSE_train)
print("Root Mean square error for train data is :", RMSE_train)
print("Mean Absolute error for train data is :", MAE_train)
print("Mean Absolute percentage error for train data is :", MAPE_train)

=====
Mean square error for train data is : 477492941.829318
Root Mean square error for train data is : 69100.8641781923
Mean Absolute error for train data is : 50929.20647093294
Mean Absolute percentage error for train data is : 0.28841643462182376
R square for train data is : 0.6403184439877669

In [107]: print(prediction[0:5])
print(y_test[0:5])

8151    245380
53      184280
3039    146280
9484    154280
9307    324280
Name: median_house_value, dtype: int64

In [108]: test = pd.DataFrame({"Predicted": prediction, 'Actual': y_test})
print(test)

Out[108]:
   Predicted  Actual
0      225683  452500
53    157719.6  184280
3039   182780.9  146280
9484   120793.1  154280
9307   178806.0  324280
...
16733  120320.2  117280
13774   238300.0  238300
13662    85780.0  85780
11942   130800.0  130800
Name: median_house_value, Length: 4128, dtype: int64

In [109]: fig=plt.figure(figsize=(16,8))
test = test.reset_index()
test = test.drop(['index'],axis=1)
plt.legend(['Actual','Predicted'])
sns.jointplot(x=Actual, y=Predicted, data=test, kind='reg',);

median_house_value
600000
500000
400000
300000
200000
100000
0
0 10 20 30 40 50
Actual

Actual
0
10000
20000
30000
40000
50000
60000
70000
median_house_value
0
100000
200000
300000
400000
500000
600000
700000
Actual
0
10000
20000
30000
40000
50000
```