

Question 1

Display the data types of each column using the function dtypes. Take a screenshot of your code and output. You will need to submit the screenshot for the final project.

[14]:

```
print(df.dtypes)

Unnamed: 0      int64
id              int64
date           object
price          float64
bedrooms       float64
bathrooms      float64
sqft_living    int64
sqft_lot       int64
floors         float64
waterfront     int64
view           int64
condition      int64
grade          int64
sqft_above     int64
sqft_basement  int64
yr_built       int64
yr_renovated   int64
zipcode        int64
lat            float64
long           float64
sqft_living15  int64
sqft_lot15     int64
dtype: object
```

We use the method describe to obtain a statistical summary of the dataframe

[15]:

```
df.describe()
```

[15]:

	Unnamed: 0	id	price	bedrooms	bathrooms
count	21613.00000	2.161300e+04	2.161300e+04	21600.000000	21603.00000
mean	10806.00000	4.580302e+09	5.400881e+05	3.372870	2.161300
std	6239.28002	2.876566e+09	3.671272e+05	0.926657	0.707107
min	0.00000	1.000102e+06	7.500000e+04	1.000000	0.500000
25%	5403.00000	2.123049e+09	3.219500e+05	3.000000	1.750000
50%	10806.00000	3.904930e+09	4.500000e+05	3.000000	2.250000
75%	16209.00000	7.308900e+09	6.450000e+05	4.000000	2.500000
max	21612.00000	9.900000e+09	7.700000e+06	33.000000	8.000000

8 rows x 21 columns

Question 2

Drop the columns "id" and "Unnamed: 0" from axis 1 using the method `drop()`, then use the method `describe()` to obtain a statistical summary of the data. Make sure the `inplace` parameter is set to `True`. Take a screenshot of your code and output. You will need to submit the screenshot for the final project.

[16]:

Question 3

Question 4

[21]:

Question 5

[22]:

Question 6

[23]:

Question 7

Question 8

[24]:

Question 9

[25]:

[32]

Question 10

Perform a second order polynomial transform on both the training data and testing data. Create and fit a Ridge regression object using the training data, set the regularization parameter to 0.1, and calculate the  $R^2$  utilizing the test data provided. Take a screenshot of your code and the  $R^2$ . You will need to submit it for the final project.

[35]:

[36]:

```
poly = PolynomialFeatures(degree=2, include_bias=False)
x_train_poly = poly.fit_transform(x_train)
x_test_poly = poly.transform(x_test)

ridge_reg = Ridge(alpha=0.1)
ridge_reg.fit(x_train_poly, y_train)

r2_score = ridge_reg.score(x_test_poly, y_test)
print(f'R^2: {r2_score}')
```

[37]:

[38]

[39]

[40]

[41]

[42]

[43]

[44]

[45]

[46]

[47]

[48]

[49]

[50]

[51]

[52]

[53]

[54]

[55]

[56]

[57]

[58]

[59]

[60]

[61]

[62]

[63]

[64]

[65]

[66]

[67]

[68]

[69]

[70]

[71]

[72]

[73]

[74]

[75]

[76]

[77]

[78]

[79]

[80]

[81]

[82]

[83]

[84]

[85]

[86]

[87]

[88]

[89]

[90]

[91]

[92]

[93]

[94]

[95]

[96]

[97]

[98]

[99]

[100]

[101]

[102]

[103]

[104]

[105]

[106]

[107]

[108]

[109]

[110]

[111]

[112]

[113]

[114]

[115]

[116]

[117]

[118]

[119]

[120]

[121]

[122]

[123]

[124]

[125]

[126]

[127]

[128]

[129]

[130]

[131]

[132]

[133]

[134]

[135]

[136]

[137]

[138]

[139]

[140]

[141]

[142]

[143]

[144]

[145]

[146]

[147]

[148]

[149]

[150]

[151]

[152]

[153]

[154]

[155]

[156]

[157]

[158]

[159]

[160]

[161]

[162]

[163]

[164]

[165]

[166]

[167]

[168]

[169]

[170]

[171]

[172]

[173]

[174]

[175]

[176]

[177]

[178]

[179]

[180]

[181]

[182]

[183]

[184]

[185]

[186]

[187]

[188]

[189]

[190]

[191]

[192]

[193]

[194]

[195]

[196]

[197]

[198]

[199]

[200]

[201]

[202]

[203]

[204]

[205]

[206]

[207]

[208]

[209]

[210]

[211]

[212]

[213]

[214]

[215]

[216]

[217]

[218]

[219]

[220]

[221]

[222]

[223]

[224]

[225]

[226]

[227]

[228]

[229]

[230]

[231]

[232]

[233]

[234]

[235]

[236]

[237]

[238]

[239]

[240]

[241]

[242]

[243]

[244]

[245]

[246]

[247]

[248]

[249]

[250]

[251]

[252]

[253]

[254]

[255]

[256]

[257]

[258]

[259]

[260]

[261]

[262]

[263]

[264]

[265]

[266]

[267]

[268]

[269]

[270]

[271]

[272]

[273]

[274]

[275]

[276]

[277]

[278]

[279]

[280]

[281]

[282]

[283]

[284]

[285]

[286]

[287]

[288]

[289]

[290]

[291]

[292]

[293]

[294]

[295]

[296]

[297]

[298]

[299]

[300]

[301]

[302]

[303]

[304]

[305]

[306]

[307]

[308]

[309]

[310]

[311]

[312]

[313]

[314]

[315]

[316]

[317]

[

## Question 2

Drop the columns "id" and "Unnamed: 0" from axis 1 using the method `drop()`, then use the method `describe()` to obtain a statistical summary of the data. Make sure the `inplace` parameter is set to `True`. Take a screenshot of your code and output. You will need to submit the screenshot for the final project.

```
[16]: #Enter Your Code, Execute and take the Screenshot
df.drop(["id", "Unnamed: 0"], axis=1, inplace=True)

print(df.describe())
```

	price	bedrooms	bathrooms	sqft_living	sqft_lot
count	2.161300e+04	21600.000000	21603.000000	21613.000000	2.161300e+04
mean	5.408881e+05	3.372870	2.115736	2079.899736	1.510697e+04
std	3.671272e+05	0.926657	0.768996	918.448897	4.142851e+04
min	7.500000e+04	1.000000	0.500000	290.000000	5.200000e+02
25%	3.219500e+05	3.000000	1.750000	1427.000000	5.040000e+03
50%	4.500000e+05	3.000000	2.250000	1910.000000	7.618000e+03
75%	6.450000e+05	4.000000	2.500000	2550.000000	1.068800e+04
max	7.700000e+06	33.000000	8.000000	13540.000000	1.651350e+06

	floors	waterfront	view	condition	grade
count	21613.000000	21613.000000	21613.000000	21613.000000	21613.000000
mean	1.404309	0.007542	0.234303	3.409430	7.656873
std	0.539989	0.086517	0.766318	0.650743	1.175459
min	1.000000	0.000000	0.000000	1.000000	1.000000
25%	1.000000	0.000000	0.000000	3.000000	7.000000
50%	1.500000	0.000000	0.000000	3.000000	7.000000
75%	2.000000	0.000000	0.000000	4.000000	8.000000
max	3.500000	1.000000	4.000000	5.000000	13.000000

	sqft_above	sqft_basement	yr_built	yr_renovated	zipcode
count	21613.000000	21613.000000	21613.000000	21613.000000	21613.000000
mean	1788.390601	291.509045	1971.005136	84.402258	98077.939805
std	828.090978	442.575043	20.373411	401.679240	53.505026
min	290.000000	0.000000	1900.000000	0.000000	98001.000000
25%	1190.000000	0.000000	1951.000000	0.000000	98033.000000
50%	1560.000000	0.000000	1975.000000	0.000000	98065.000000
75%	2210.000000	560.000000	1997.000000	0.000000	98118.000000
max	9410.000000	4820.000000	2015.000000	2015.000000	98199.000000

	lat	long	sqft_living15	sqft_lot15
count	21613.000000	21613.000000	21613.000000	21613.000000
mean	47.560053	-122.213096	1986.552492	12768.455652
std	0.138564	0.140828	685.391304	27304.179631
min	47.155900	-122.519000	399.000000	651.000000
25%	47.471000	-122.320000	1490.000000	5100.000000
50%	47.571800	-122.230000	1840.000000	7620.000000
75%	47.670000	-122.125000	2360.000000	10083.000000
max	47.777600	-121.315000	6210.000000	871200.000000

We can see we have missing values for the columns `bedrooms` and `bathrooms`

```
[17]: print("number of NaN values for the column bedrooms :", df['bedrooms'].isnull().sum())
print("number of NaN values for the column bathrooms :", df['bathrooms'].isnull().sum())
```

```
number of NaN values for the column bedrooms : 13
number of NaN values for the column bathrooms : 10
```

We can replace the missing values of the column `'bedrooms'` with the mean of the column `'bedrooms'` using the method `replace()`. Don't forget to set the `inplace` parameter to `True`

```
[18]: mean=df['bedrooms'].mean()
df['bedrooms'].replace(np.nan,mean, inplace=True)
```

We also replace the missing values of the column `'bathrooms'` with the mean of the column `'bathrooms'` using the method `replace()`. Don't forget to set the `inplace` parameter to `True`

```
[19]: mean=df['bathrooms'].mean()
df['bathrooms'].replace(np.nan,mean, inplace=True)
```

```
[20]: print("number of NaN values for the column bedrooms :", df['bedrooms'].isnull().sum())
print("number of NaN values for the column bathrooms :", df['bathrooms'].isnull().sum())
```

```
number of NaN values for the column bedrooms : 0
number of NaN values for the column bathrooms : 0
```

### Question 3

Use the method `value_counts` to count the number of houses with unique floor values, use the method `.to_frame()` to convert it to a data frame. Take a screenshot of your code and output. You will need to submit the screenshot for the final project.

```
[21]: floor_counts = df['floors'].value_counts().to_frame()

floor_counts.columns = ['Number of Houses']

print(floor_counts)
```

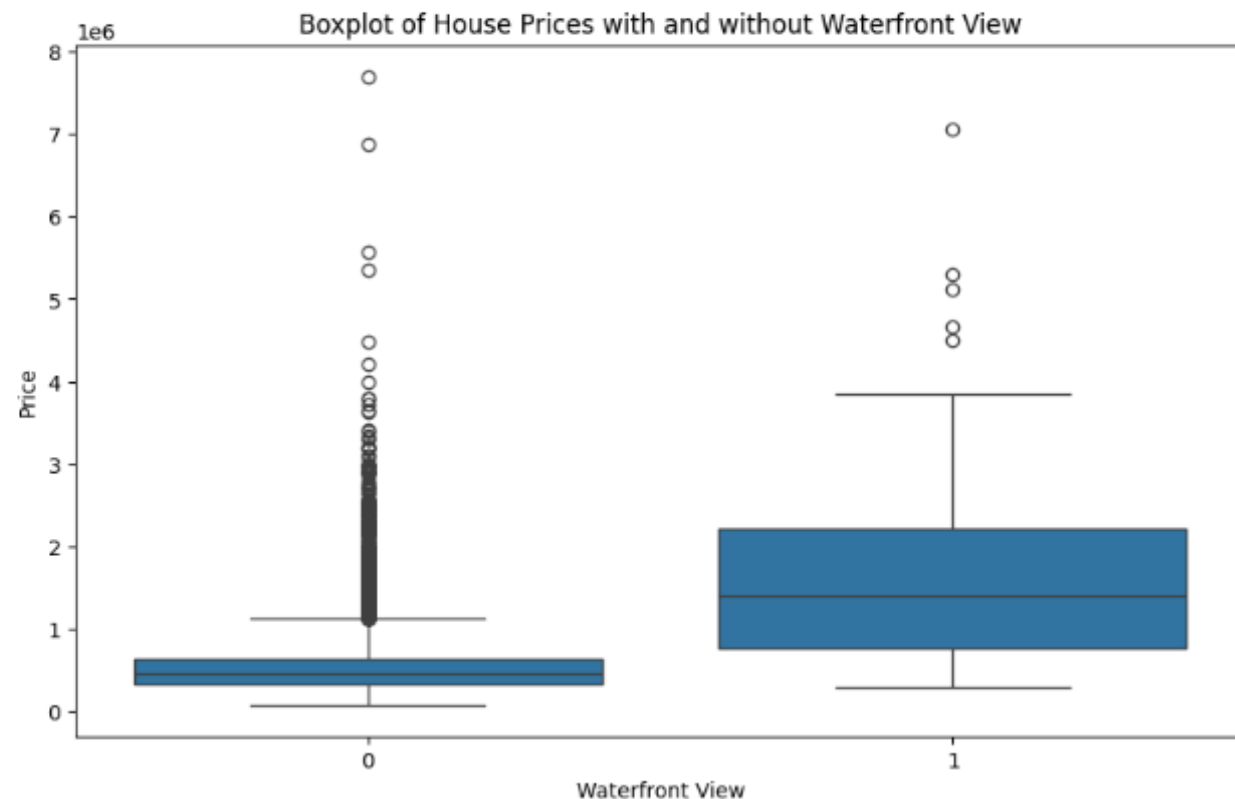
	Number of Houses
1.0	10680
2.0	8241
1.5	1910
3.0	613
2.5	161
3.5	8

## Question 4

Use the function `boxplot` in the seaborn library to determine whether houses with a waterfront view or without a waterfront view have more price outliers. Take a screenshot of your code and boxplot. You will need to submit the screenshot for the final project.

```
[22]: import seaborn as sns

plt.figure(figsize=(10, 6))
sns.boxplot(x='waterfront', y='price', data=df)
plt.title('Boxplot of House Prices with and without Waterfront View')
plt.xlabel('Waterfront View')
plt.ylabel('Price')
plt.show()
```

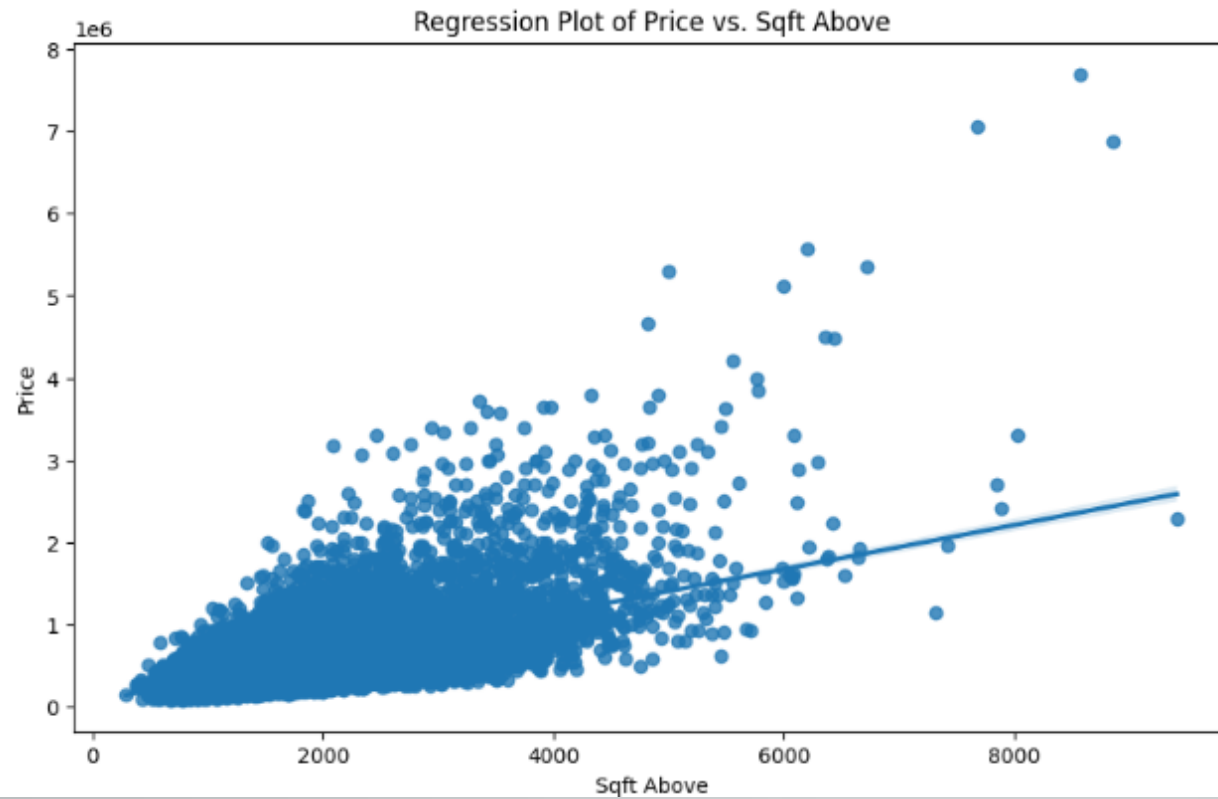


## Question 5

Use the function `regplot` in the seaborn library to determine if the feature `sqft_above` is negatively or positively correlated with price. Take a screenshot of your code and scatterplot. You will need to submit the screenshot for the final project.

```
[23]: import seaborn as sns

plt.figure(figsize=(10, 6))
sns.regplot(x='sqft_above', y='price', data=df)
plt.title('Regression Plot of Price vs. Sqft Above')
plt.xlabel('Sqft Above')
plt.ylabel('Price')
plt.show()
```



## Question 6

Fit a linear regression model to predict the 'price' using the feature 'sqft\_living' then calculate the  $R^2$ . Take a screenshot of your code and the value of the  $R^2$ . You will need to submit it for the final project.

```
[28]: from sklearn.linear_model import LinearRegression
```

```
X = df[['sqft_living']]
```

```
Y = df['price']
```

```
lm = LinearRegression()
```

```
lm.fit(X, Y)
```

```
r2_score = lm.score(X, Y)
```

```
print(f"R^2: {r2_score}")
```

```
R^2: 0.4928532179037931
```

## Question 7

Fit a linear regression model to predict the 'price' using the list of features:

```
[29]: features = ["floors", "waterfront","lat" ,"bedrooms" ,"sqft_basement" ,"view" ,"bathrooms","sqft_living15","sqft_above","grade","sqft_living"]
```

Then calculate the  $R^2$ . Take a screenshot of your code and the value of the  $R^2$ . You will need to submit it for the final project.

```
[30]: X = df[features]
      Y = df['price']

      lm = LinearRegression()

      lm.fit(X, Y)

      r2_score = lm.score(X, Y)
      print(f"R^2: {r2_score}")
```

R^2: 0.6576890354915759