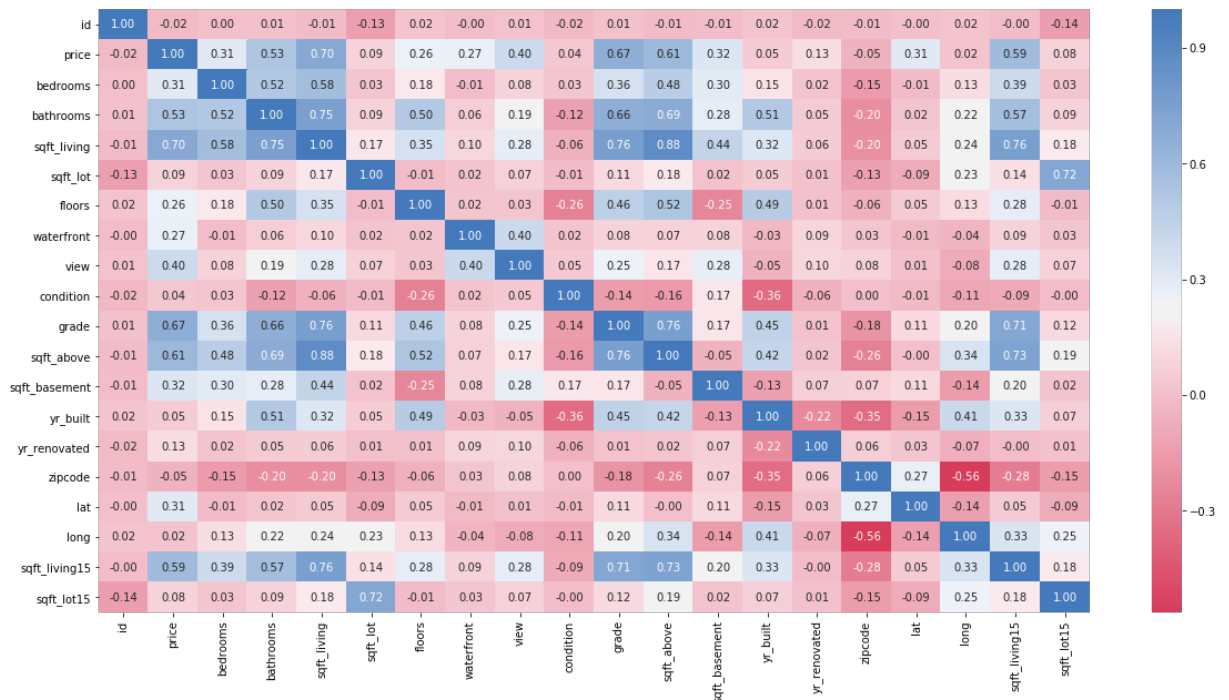


```

In [2]: 1 import numpy as np
        2 import pandas as pd
        3 from sklearn.model_selection import train_test_split
        4 from sklearn import linear_model
        5 from sklearn.neighbors import KNeighborsRegressor
        6 import matplotlib.pyplot as plt
        7 from sklearn import metrics
        8 import seaborn as sns
        9
       10 df = pd.read_csv("kc_house_data.csv")
       11
       12 # Plotting correlation heatmap to get the factors highly correlated
       13 correlation = df.corr()
       14 fig, ax = plt.subplots(figsize=(20, 10))
       15 cm = sns.diverging_palette(5, 250, as_cmap=True)
       16 sns.heatmap(correlation, cmap=cm, annot=True, fmt=".2f")
       17 plt.show()
       18 plt.show()

```



```

In [3]: 1 # Linear regression using sqft_living as X
        2
        3 import numpy as np
        4 import pandas as pd
        5 from sklearn.model_selection import train_test_split
        6 from sklearn import linear_model
        7 import matplotlib.pyplot as plt
        8 from sklearn import metrics
        9
       10 df = pd.read_csv("kc_house_data.csv")
       11
       12 rf = df.drop(['id', 'date', 'price'], 1)

```

```

13
14 # From the heat map we identified that sqft_living and grade have
15 # So we did two models
16 # -> Using all attributes against price
17 # -> Using sqft_living and grade attributes against price
18
19
20 column_selected = ['sqft_living']
21 predicted_array = []
22
23 train, test = train_test_split(df, train_size =0.90, random_state =
24
25 lm = linear_model.LinearRegression()
26
27 X_train = np.array(train[column_selected]).reshape(-1,1)
28 X_test = np.array(test[column_selected]).reshape(-1, 1)
29
30 Y_train = np.array(train['price']).reshape(-1, 1)
31 Y_test = np.array(test['price']).reshape(-1, 1)
32
33 lm.fit(train[column_selected], train['price'])
34
35 prediction = lm.predict(test[column_selected])
36
37 print("Model using {} as X".format(column_selected))
38 mse = metrics.mean_squared_error(Y_test, prediction)
39 error = np.sqrt(mse)
40 intercept = lm.intercept_
41 accuracy = lm.score(test[column_selected], test['price'])
42
43 print("\nThe root mean squared error is ", np.round(error, 2))
44 print("\nThe coefficient array is ", np.round(lm.coef_,2))
45 print("\nThe intercept is ", np.round(intercept,2))
46 print("\nThe accuracy is given by is ", round(accuracy, 2))
47
48 fig, ax = plt.subplots(figsize= (15, 10))
49 plt.scatter(X_test, Y_test, color= 'blue', label = 'Scattered Data')
50 plt.plot(X_test, lm.predict(X_test), color='black', label= 'Predict')
51 plt.xlabel('Square ft Living')
52 plt.ylabel('Price of the house')
53 plt.legend()
54
55
56
57
58
59

```

Model using ['sqft_living'] as X

The root mean squared error is 249846.53

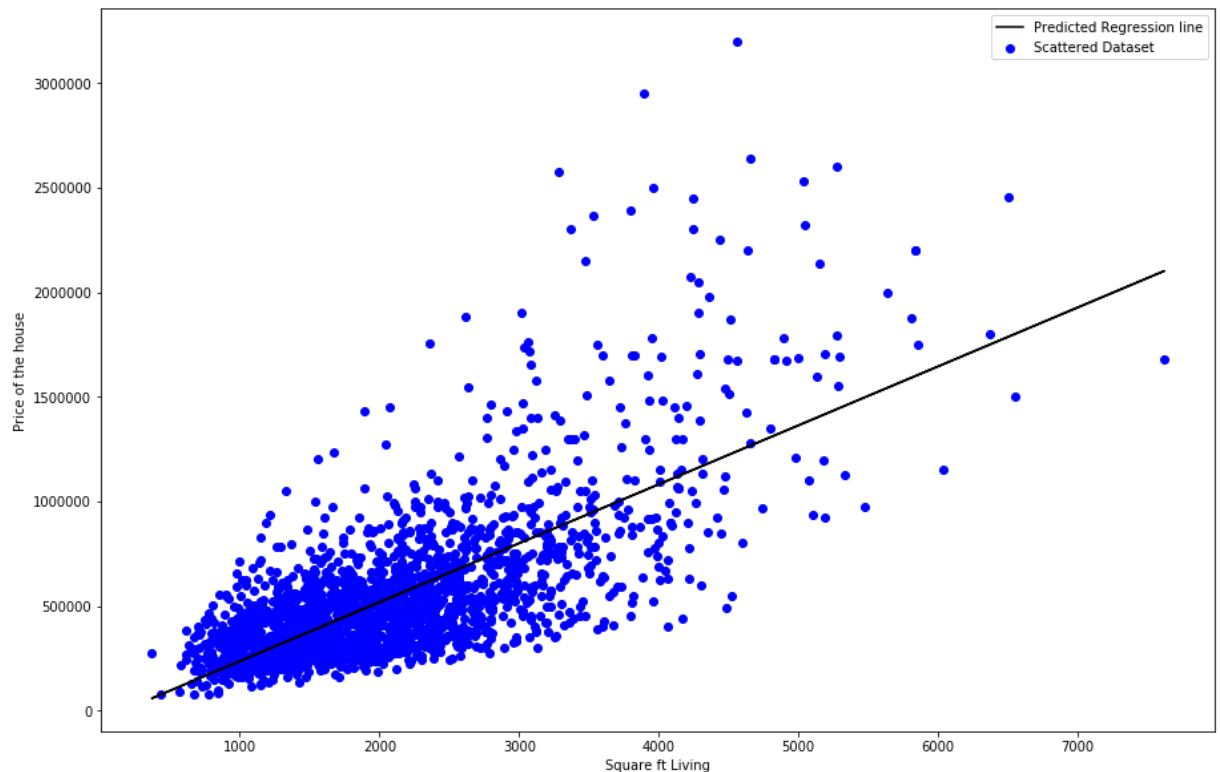
The coefficient array is [282.03]

The intercept is -46861.8

The accuracy is given by is 0.49

```
/anaconda3/lib/python3.7/site-packages/sklearn/model_selection/_split.py:2026: FutureWarning: From version 0.21, test_size will always complement train_size unless both are specified.
FutureWarning)
```

Out[3]: <matplotlib.legend.Legend at 0x1a2369f7f0>



```
In [5]: 1 # Linear regression using sqft_living as X
2
3 import numpy as np
4 import pandas as pd
5 from sklearn.model_selection import train_test_split
6 from sklearn import linear_model
7 import matplotlib.pyplot as plt
8 from sklearn import metrics
9 import seaborn as sns
10
11
12 df = pd.read_csv("kc_house_data.csv")
13
14 rf = df.drop(['id', 'date', 'price'], 1)
15
16
17 column_selected = ['grade']
18 predicted_array = []
19
20 train, test = train_test_split(df, train_size = 0.90, random_state =
21
```

```

22 lm = linear_model.LinearRegression()
23
24 X_train = np.array(train[column_selected]).reshape(-1,1)
25 X_test = np.array(test[column_selected]).reshape(-1, 1)
26
27 Y_train = np.array(train['price']).reshape(-1, 1)
28 Y_test = np.array(test['price']).reshape(-1, 1)
29
30 lm.fit(train[column_selected], train['price'])
31
32 prediction = lm.predict(test[column_selected])
33
34 print("Model using {} as X".format(column_selected))
35 mse = metrics.mean_squared_error(Y_test, prediction)
36 error = np.sqrt(mse)
37 intercept = lm.intercept_
38 accuracy = lm.score(test[column_selected], test['price'])
39
40 print("\nThe root mean squared error is ", np.round(error, 2))
41 print("\nThe coefficient array is ", np.round(lm.coef_,2))
42 print("\nThe intercept is ", np.round(intercept,2))
43 print("\nThe accuracy is given by is ", round(accuracy, 2))
44
45
46 fig, ax = plt.subplots(figsize= (15, 10))
47 plt.scatter(X_test, Y_test, color= 'blue', label = 'Scattered Datas')
48 plt.plot(X_test, lm.predict(X_test), color='black', label= 'Predicted Regression line')
49 plt.xlabel('Grade')
50 plt.ylabel('Price of the house')
51 plt.legend()
52 plt.show()
53
54
55
56
57
58

```

```

/anaconda3/lib/python3.7/site-packages/sklearn/model_selection/_split.py:2026: FutureWarning: From version 0.21, test_size will always complement train_size unless both are specified.
FutureWarning)

```

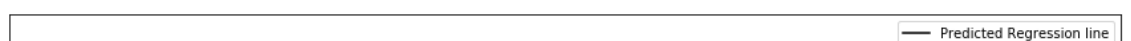
Model using ['grade'] as X

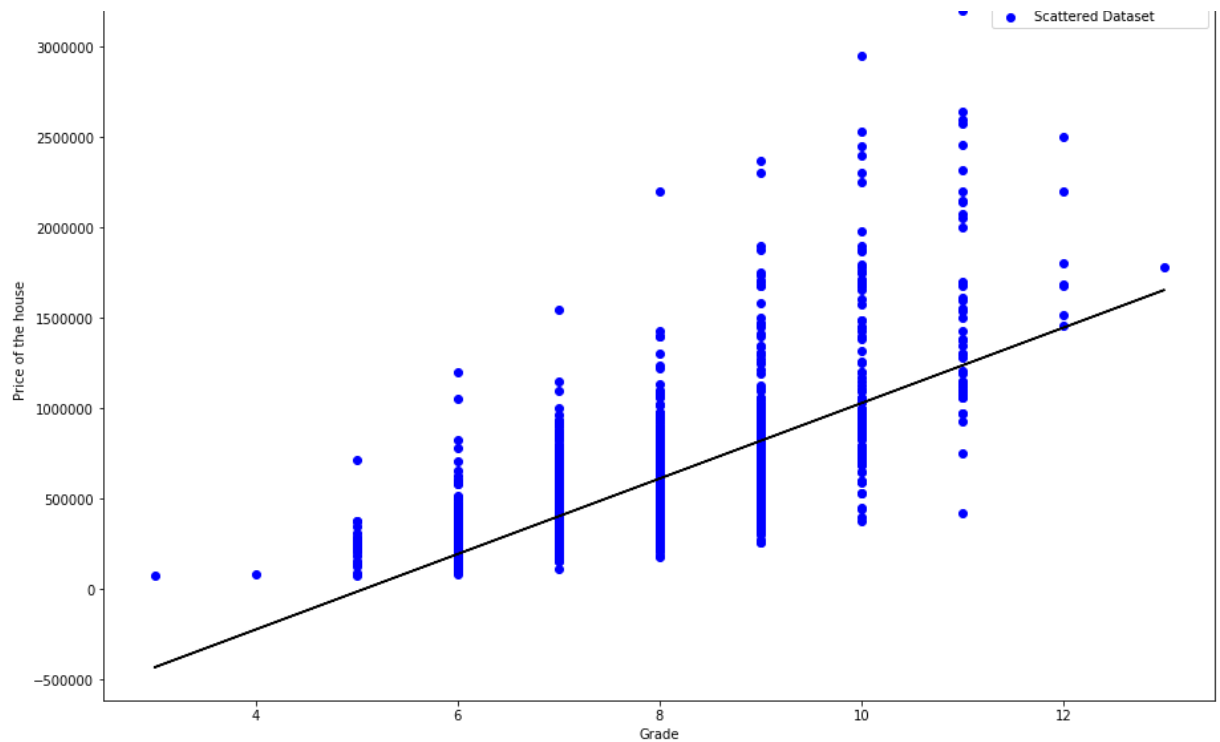
The root mean squared error is 252130.8

The coefficient array is [208591.32]

The intercept is -1057294.06

The accuracy is given by is 0.48





In []: 1

```
In [6]: 1 import numpy as np
2 import pandas as pd
3 from sklearn.model_selection import train_test_split
4 from sklearn import linear_model
5 from sklearn.neighbors import KNeighborsRegressor
6 import matplotlib.pyplot as plt
7
8 df = pd.read_csv("kc_house_data.csv")
9
10 nf = df
11 date_string = nf['date'].tolist()
12 date_int = [int(str[:4]) for str in date_string]
13 df['year'] = date_int
14
15 rf = df.drop(['id', 'date', 'price', 'zipcode'], 1)
16
17 # From the heat map we identified that sqft_living and grade have
18 # So we did two models
19 # -> Using all attributes against price
20 # -> Using sqft_living and grade attributes against price
21
22 column_selected = list(rf.columns.values)
23
24 predicted_array = []
25
26 train, test = train_test_split(df, train_size = 0.90, random_state =
27
28 lm = linear_model.LinearRegression()
29
30 X_train = np.array(train[column_selected]).reshape(-1,1)
```

```

31 X_test = np.array(test[column_selected]).reshape(-1, 1)
32
33 Y_train = np.array(train['price']).reshape(-1, 1)
34 Y_test = np.array(test['price']).reshape(-1, 1)
35
36 lm.fit(train[column_selected], train['price'])
37
38 prediction = lm.predict(test[column_selected])
39
40 print("Model using {} as X".format(column_selected))
41 mse = metrics.mean_squared_error(Y_test, prediction)
42 error = np.sqrt(mse)
43 intercept = lm.intercept_
44 accuracy = lm.score(test[column_selected], test['price'])
45
46 print("\nThe root mean squared error is ", np.round(error, 2))
47 print("\nThe coefficient array is ", np.round(lm.coef_, 2))
48 print("\nThe intercept is ", np.round(intercept, 2))
49 print("\nThe accuracy is given by is ", round(accuracy, 2))
50

```

Model using ['bedrooms', 'bathrooms', 'sqft_living', 'sqft_lot', 'floors', 'waterfront', 'view', 'condition', 'grade', 'sqft_above', 'sqft_basement', 'yr_built', 'yr_renovated', 'lat', 'long', 'sqft_living15', 'sqft_lot15', 'year'] as X

The root mean squared error is 193425.01

The coefficient array is [-3.3978050e+04 4.4056360e+04 1.0924000e+02 1.3000000e-01

9.5241000e+02 6.0859354e+05 4.9747160e+04 3.2698890e+04

9.6312430e+04 7.1300000e+01 3.7940000e+01 -2.4791000e+03

2.3030000e+01 5.6349587e+05 -1.1762568e+05 2.7110000e+01

-3.9000000e-01 2.9429070e+04]

The intercept is -96292145.69

The accuracy is given by is 0.7

/anaconda3/lib/python3.7/site-packages/sklearn/model_selection/_split.py:2026: FutureWarning: From version 0.21, test_size will always complement train_size unless both are specified.
FutureWarning)

In []:

1