

Introduction to Data Science

Regression Dataset 2

House Prices

By: Pujesh Pradhan

Table of Content

- Objective and Data Description
- Statistical Numerical and Graphical Summaries
- Algorithm Implementation and Statistical Tests
- Performance Improvement
- Conclusion

Objective and Data Description

Overview

There is always a supply and demand in the housing sector and predicting the price of the house at the right amount is very crucial. If the price of the house is stated low than the market value, people will tend to lose the gain they could have made. Similarly, if a house is priced high, there is less likely chance of selling the house or getting any buyer. Thus, pricing the house correctly is very important.

Objective

Predict the price of the house based on the most important feature

About the Dataset

The dataset consists of house prices from King County area in the US State of Washington, this data also covers Seattle.

The dataset was obtained from [Kaggle](#)

The dataset contains 21 fields out of which the first field is an id field and the second field is the price of houses being sold. The description of fields is as follows:

1. **id**: ID of the house being sold.
2. **date**: Date of the house being sold.
3. **price**: Price of the house being sold.
4. **bedrooms**: No. of bedrooms in the house.
5. **bathrooms**: No. of bathrooms in the house.
6. **sqft_living**: The square foot area of the living space of the house.
7. **sqft_lot**: The square foot area of the entire land.
8. **floors**: No. of floors in the house.
9. **waterfront**: If a view of the waterfront is present or not.
10. **view**: Number of views. (1 indicates were out property and 5 great)
11. **condition**: Condition of the house. (Based on King County grading system - 1 poor and 13 excellent)
12. **grade**: Grade of the house.
13. **sqft_above**: The square foot area apart from the basement.
14. **sqft_basement**: The square foot area of the basement.
15. **yr_built**: The year of house built.
16. **yr_renovated**: The year when house is renovated.
17. **zipcode**: Zipcode.
18. **lat**: Latitude co-ordinates of the house.
19. **long**: Longitude co-ordinates of the house.
20. **sqftliving15**: Living room area in 2015.
21. **sqftlot15**: The lot size area in 2015.

```
In [1]: # All the required libraries are imported here.
import numpy as np
import pandas as pd
from pandas.plotting import scatter_matrix

import matplotlib.pyplot as plt
import seaborn as sns

from sklearn.preprocessing import MinMaxScaler
from sklearn.metrics import train_test_split, GridSearchCV, RepeatedStratifiedKFold
from sklearn.linear_model import LinearRegression, RidgeClassifier, LogisticRegression
from sklearn.metrics import classification_report, confusion_matrix, mean_squared_error
from sklearn.feature_selection import SelectKBest, chi2

import math

from scipy import stats
```

```
In [2]: # Describing the size, shape and dimension of the dataset.
house_data = pd.read_csv('house_data.csv')
print('The total size of the dataset is {} bytes.'.format(house_data.size))
print('There are {} number of records and {} number of fields with a {} shape.'.format(house_data.shape[0], house_data.shape[1], house_data.shape))
print('The dataset has a {} dimensional structure.'.format(house_data.ndim))

The total size of the dataset is 4528973 bytes.
There are 21613 number of records and 21 number of fields with a (21613, 21) shape.
The dataset has 4 2 dimensional structure.
```

```
In [3]: print('The top 5 records of the dataset are: \n{}'.format(house_data.head()))

The top 5 records of the dataset are:
   id      date      price  bedrooms  bathrooms  sqft_living  \
0  7129300529  20141813T090000  221900.0      3      1.09      1180
1  6414100102  20141209T080000  538000.0      3      2.25      2578
2  5633500409  20150223T080000  189000.0      2      1.08      778
3  2487200875  20141209T080000  604800.0      4      3.00      1980
4  1954400518  20150218T080000  510000.0      3      2.00      1680

   sqft_lot  floors  waterfront  view  ...  grade  sqft_above  sqft_basement  \
0    5650      1.0      0      0 ...      7      1180      0
1    7242      1.0      0      0 ...      7      1270      498
2    15080      1.0      0      0 ...      6      778      0
3    9090      1.0      0      0 ...      8      1959      918
4    8080      1.0      0      0 ...      8      1680      0

   yr_built  yr_renovated  zipcode      lat  long  sqft_living15  \
0      1991      1991      98115  47.5112 -122.257      1340
1      1951      1991      98115  47.7137 -122.319      1200
2      1953      0      98028  47.7379 -122.233      2720
3      1985      0      98136  47.5209 -122.393      2360
4      1987      0      98074  47.6168 -122.045      1860

   sqft_lot15
0      5650
1      7639
2      8902
3      5900
4      7991
[5 rows x 21 columns]
```

Here is a sample record of a house that was sold in different dates and prices but the data records the house twice. We will remove the duplicates and keep the latest transactions. To do that, we need to find and remove the duplicates.

```
In [5]: # We will find the duplicate records based on id, bedrooms, bathrooms, sqft_living, sqft_lot, zipcode, lat and longitude.

print('There are {} number of duplicate values in the dataset. \n'.format(house_data.duplicated(subset=['id', 'bedrooms', 'bathrooms', 'sqft_living', 'sqft_lot', 'zipcode', 'lat', 'long']).sum()))

There are 377 number of duplicate values in the dataset.
```

```
In [6]: # We will find the list of duplicate values as follows: \n
house_data[!house_data.duplicated(subset=['id', 'bedrooms', 'bathrooms', 'sqft_living', 'sqft_lot', 'zipcode', 'lat', 'long'], keep=False), :]

The list of duplicate values are as follows:
```

```
Out[6]:
   id      date      price  bedrooms  bathrooms  sqft_living  sqft_lot  floors  waterfront  view  ...  grade  sqft_above  sqft_basement  yr_built  yr_renovated  zipcode      lat  long  sq
93  6021501535  20141207T000000      43000.0      3      1.50      1580      5000      1.0      0      0 ...      8      1290      290      1939      0      98117  47.6870 -122.386
94  6021501535  20141223T000000      70000.0      3      1.50      1580      5000      1.0      0      0 ...      8      1290      290      1939      0      98117  47.6870 -122.386
313  4139480000  20141037T000000      138400.0      4      3.25      4290      12103      1.0      0      3 ...      11      2690      1600      1997      0      98006  47.5503 -122.102
314  4139480000  20141037T000000      140000.0      4      3.25      4290      12103      1.0      0      3 ...      11      2690      1600      1997      0      98006  47.5503 -122.102
324  7520000020  20141009T000000      232000.0      2      1.00      1240      12082      1.0      0      0 ...      6      960      280      1922      1984      98146  47.4957 -122.352
...
29770  8564060270  20150307T000000      502000.0      4      2.50      2680      5539      2.0      0      0 ...      8      2680      0      2013      0      98048  47.4789 -121.734
29779  8564060270  20150307T000000      240000.0      4      1.00      1200      2171      1.0      0      0 ...      7      1200      0      1993      0      98019  47.7079 -122.342
29988  8564060270  20150307T000000      30000.0      4      1.00      1200      2171      1.0      0      0 ...      7      1200      0      1993      0      98019  47.7079 -122.342
29989  7653403110  20141037T000000      59406.0      3      3.00      2780      6000      2.0      0      0 ...      9      2780      0      2013      0      98006  47.5184 -121.896
21583  7653403110  20150307T000000      62500.0      3      3.00      2780      6000      2.0      0      0 ...      9      2780      0      2013      0      98006  47.5184 -121.896
353 rows x 21 columns
```

```
In [7]: # Now deleting the first occurrence of the duplicate values.
house_data = house_data.drop_duplicates(subset=['id', 'bedrooms', 'bathrooms', 'sqft_living', 'sqft_lot', 'zipcode', 'lat', 'long'], keep='last')
```

```
In [8]: # Printing the data types
print('The datatype of the dataset are: \n{}'.format(house_data.dtypes))

The datatype of the dataset are:
id              int64
date            object
price           float64
bedrooms        int64
bathrooms       int64
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 can see that the data values is of type object. Since, we know date will not play much part in predicting the price of the house, we can either drop the field or leave it as it is.

```
In [9]: # Statistical description of the data
print(house_data.describe())

id              count      mean      std      min      max
price           21436.0  5.45090e+05  5.67134e+05  1.000102e+06  7.50000e+06
bedrooms        21436.0  3.37357e+00  0.250347e+01  0.000000e+00  6.00000e+00
bathrooms       21436.0  2.17349e+00  0.699128e+01  0.000000e+00  6.00000e+00
sqft_living     21436.0  2.08270e+03  5.15140e+03  2.00000e+02  5.20000e+02
sqft_lot        21436.0  1.51356e+04  4.15382e+04  0.00000e+00  5.20000e+02
floors          21436.0  1.48510e+00  5.40338e+01  0.00000e+00  6.00000e+00
waterfront      21436.0  7.08493e+03  0.68709e+02  0.00000e+00  6.00000e+00
view            21436.0  1.25151e+01  1.67929e+01  0.00000e+00  6.00000e+00
condition       21436.0  3.41039e+00  0.56235e+01  0.00000e+00  6.00000e+00
grade           21436.0  7.41217e+00  0.348937e+01  0.00000e+00  6.00000e+00
sqft_above      21436.0  1.78959e+03  0.29825e+02  0.00000e+00  6.00000e+00
sqft_basement   21436.0  2.61476e+02  4.47820e+02  0.00000e+00  6.00000e+00
yr_built        21436.0  1.97199e+03  0.23852e+04  1.00000e+03  6.00000e+03
yr_renovated    21436.0  0.47298e+03  0.82431e+02  0.00000e+00  6.00000e+00
zipcode         21436.0  0.88778e+04  0.348937e+01  0.00000e+00  6.00000e+00
lat             21436.0  4.75092e+01  1.38601e+01  0.00000e+00  6.00000e+00
long            21436.0  1.27359e+02  4.48986e+01  0.00000e+00  6.00000e+00
sqft_living15   21436.0  1.98831e+03  0.85692e+02  0.00000e+00  6.00000e+00
sqft_lot15      21436.0  1.27859e+04  7.73757e+04  6.51000e+02  6.51000e+02
```

```
In [10]: # We can see that the data values is of type object. Since, we know date will not play much part in predicting the price of the house, we can either drop the field or leave it as it is.

print('The top 5 records of the dataset are: \n{}'.format(house_data.head()))

Now our dataframe with the updated index looks like:
   id      date      price  bedrooms  bathrooms  sqft_living  \
0  7129300529  20141813T090000  221900.0      3      1.09      1180
1  6414100102  20141209T080000  538000.0      3      2.25      2578
2  5633500409  20150223T080000  189000.0      2      1.08      778
3  2487200875  20141209T080000  604800.0      4      3.00      1980
4  1954400518  20150218T080000  510000.0      3      2.00      1680

   sqft_lot  floors  waterfront  view  condition  grade  sqft_above  \
0    5650      1.0      0      0      7      1180
1    7242      1.0      0      0      7      1270
2    15080      1.0      0      0      6      778
3    9090      1.0      0      0      8      1959
4    8080      1.0      0      0      8      1680

   yr_built  yr_renovated  zipcode      lat  long  sqft_living15  \
0      1991      1991      98115  47.5112 -122.257
1      1951      1991      98115  47.7137 -122.319
2      1953      0      98028  47.7379 -122.233
3      1985      0      98136  47.5209 -122.393
4      1987      0      98074  47.6168 -122.045

   sqft_lot15
0      5650
1      7639
2      8902
3      5900
4      7991
[21613 rows x 21 columns]
```

Now, since the id field is unique we can create that as an index column as the field does not make any contribution on linear regression or classification.

```
In [15]: # Now, since the id field is unique we can create that as an index column as the field does not make any contribution on linear regression or classification.

house_data = house_data.set_index('id')

print('Now our dataframe with the updated index looks like: \n{}'.format(house_data.head()))

Now our dataframe with the updated index looks like:
   date      price  bedrooms  bathrooms  sqft_living  \
0  20141813T090000  221900.0      3      1.09      1180
1  20141209T080000  538000.0      3      2.25      2578
2  20150223T080000  189000.0      2      1.08      778
3  20141209T080000  604800.0      4      3.00      1980
4  20150218T080000  510000.0      3      2.00      1680

   sqft_lot  floors  waterfront  view  condition  grade  sqft_above  \
0    5650      1.0      0      0      7      1180
1    7242      1.0      0      0      7      1270
2    15080      1.0      0      0      6      778
3    9090      1.0      0      0      8      1959
4    8080      1.0      0      0      8      1680

   yr_built  yr_renovated  zipcode      lat  long  sqft_living15  \
0      1991      1991      98115  47.5112 -122.257
1      1951      1991      98115  47.7137 -122.319
2      1953      0      98028  47.7379 -122.233
3      1985      0      98136  47.5209 -122.393
4      1987      0      98074  47.6168 -122.045

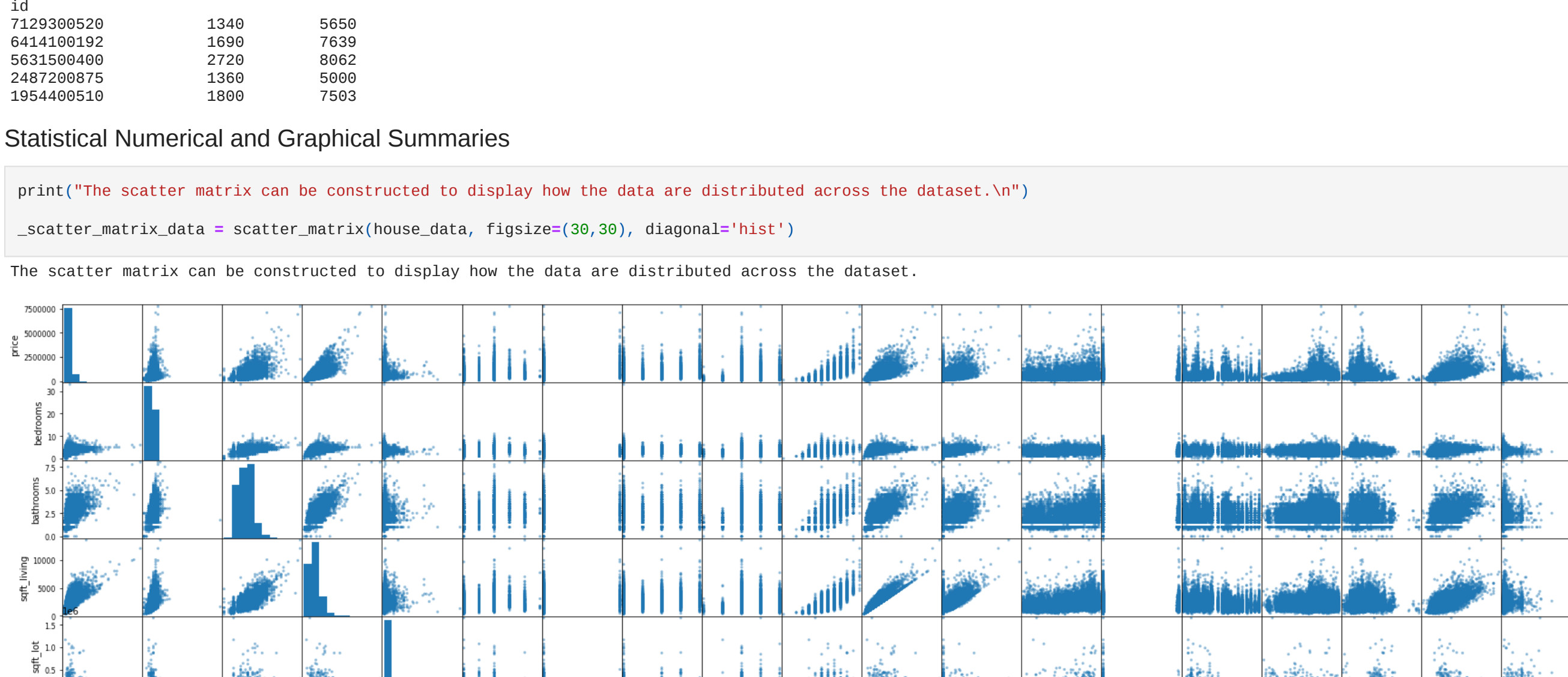
   sqft_lot15
0      5650
1      7639
2      8902
3      5900
4      7991
[21613 rows x 21 columns]
```

Statistical Numerical and Graphical Summaries

```
In [11]: print('The scatter matrix can be constructed to display how the data are distributed across the dataset.\n')

...scatter_matrix(data = scatter_matrix(house_data, figsize=(30,30), diagonal='hist'))

The scatter matrix can be constructed to display how the data are distributed across the dataset.
```



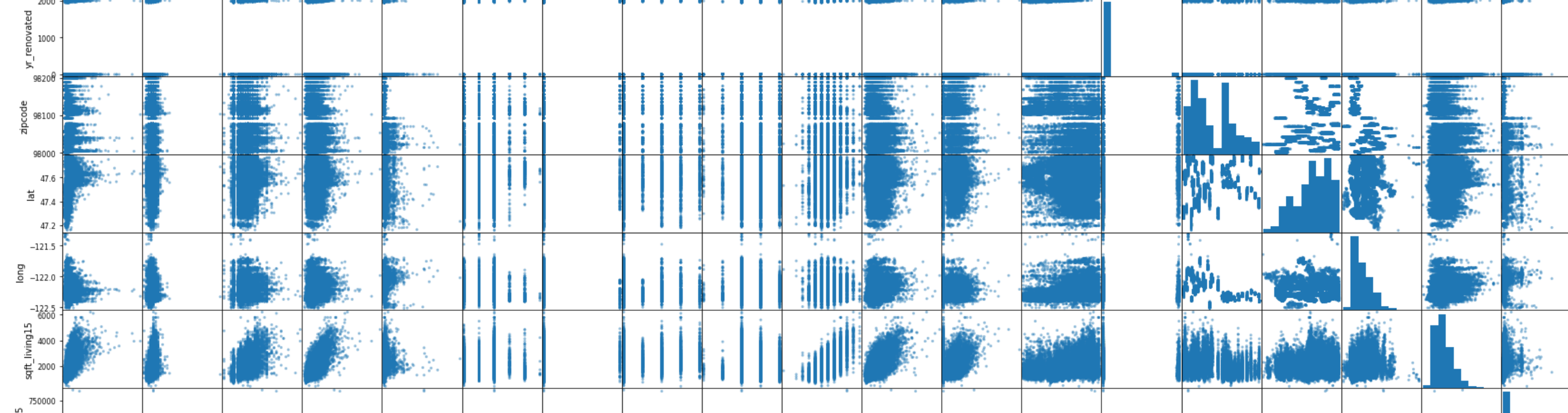
```
In [12]: # From the above graph while considering the histograms diagonally we can see that most of the fields are right skewed. \
# Whereas from the scatter plot we can see there are few outliers in the dataset. \
# Below is the list of all the keys of the dataset and determines whether they are normally distributed or skewed. \n

skewness = {}
for data_keys in house_data.keys():
    if skewness > 0:
        print(data_keys, "is Right Skewed")
    elif skewness < 0:
        print(data_keys, "is Left Skewed")
    else:
        print(data_keys, "is Normally Distributed")

From the above graph while considering the histograms diagonally we can see that most of the fields are right skewed. \
whereas from the scatter plot we can see there are few outliers in the dataset. \
Below is the list of all the keys of the dataset and determines whether they are normally distributed or skewed.
```

price is Right Skewed
bedrooms is Right Skewed
bathrooms is Right Skewed
sqft_living is Right Skewed
sqft_lot is Right Skewed
floors is Right Skewed
waterfront is Right Skewed
view is Right Skewed
condition is Right Skewed
grade is Right Skewed
sqft_above is Right Skewed
sqft_basement is Right Skewed
yr_built is Left Skewed
yr_renovated is Right Skewed
zipcode is Right Skewed
lat is Left Skewed
long is Right Skewed
sqft_living15 is Right Skewed
sqft_lot15 is Right Skewed

```
In [13]: correlations = house_data.corr()
plt.figure(figsize=(20,20))
heatmap = sns.heatmap(correlations, annot=True)
```



The heatmap shows the correlation between each variables. Since, we are trying to find a variable that is strongly related with the price of the house, we can see that sqft_living field is the most strongly correlated with the price of the house. After sqft_living, the grade and the sqft_above are more strongly related to the price.

We will choose sqft_living to predict the price of the house. With this finding we can also say that the larger the sqft of the house the higher the price of the house is to be estimated. Which is true, but as all correlation does not imply causation, we cannot neglect other factors such as sqft_above, the grade of the house, number of bathrooms and bedrooms and other fields, as they could play vital role in determining the prices.

We will only find the correlation between the price of the house and the sqft_living of the house in this dataset.

```
In [14]: # Linear Regression between sqft_living and price.

price_data = house_data[['sqft_living']]
price_data = house_data[['price']]

X_train, X_test, y_train, y_test = train_test_split(price_data, price_data, test_size=0.25)

# scaler = MinMaxScaler()
# scaler.fit(X_train)

# X_train = scaler.transform(X_train)
# X_test = scaler.transform(X_test)

regression = LinearRegression()
regression.fit(X_train, y_train)

y_pred = regression.predict(X_test)
```

```
In [15]: # Plotting the points and the line of best fits. Here the scatter plot is plotted between sqft_living and price of the house and the best line of fit is created by plot() function. The predicted values are displayed on the x-axis and the explanatory variable values are displayed on the y-axis.

plt.scatter(X_test, y_test, c='grey', alpha=0.5)
plt.plot(X_test, y_pred, color='red', linewidth=2)
plt.xlabel('Price (in millions)')
plt.ylabel('sqft_living')
plt.show()
```



```
In [16]: print('Thus, the line of best fit for the above can be found using : \n y = ( ) + ( ) x. \n'.format(regression.intercept_, regression.coef_[0]))

mse = mean_squared_error(y_test, y_pred)
print('Where the Slope is ( ) and the y intercept is ( ) \n')
print('With mean squared error of ( ) and Root mean squared error of ( )')
print('The Price of the house can be put up as : $({})'.format(regression.intercept_ + (regression.coef_[0] * x)))

Thus, the line of best fit for the above can be found using :
y = [-44825.93903958] + [281.47743554] x.
```

Where the Slope is [281.47743554] and the y intercept is [-44825.93903958]
With mean squared error of 68886.82955.7337 and root mean squared error of 264349.8666179805

```
In [17]: # Pearson coefficient of the variables is r = ( ) \n
pearson_coef, p_value = stats.pearsonr(price_data['sqft_living'], price_data['price'])
print('The Pearson coefficient of the variables is r = ( ) \n'.format(pearson_coef))

The Pearson coefficient of the variables is r = 0.761432395241632
```

Thus, from the above calculation we can see that with increase in 1 square foot, there is an increase of 280.316 in price.

The correlation between the sqft_living and price creates a strong positive pattern. The coefficient correlation between the variables is 0.701.

Using the above equation we can predict price of houses by providing the sqft_living value.

```
In [18]: print('Similarly the coefficient of variation is : (%).format(regression.score(X_test, y_test)*100))

Similarly the coefficient of variation is : 0.6157865273373897%
```

```
In [19]: # Title and Labels
sns.residplot(sqft_data['sqft_living'], price_data['price'], color='magenta')

# Title and Labels
plt.title('Residual plot', size=24)
plt.xlabel('sqft_living', size=18)
plt.ylabel('Price', size=18);
```

Residual plot



The residual plot created above looks like it learns more toward Homoscedastic as the residual plot shows homogeneous spread out of the plot and the price of the house is also axis.

```
In [20]: print('To predict a price of a house with 2570 of sqft_living, we can use the above equation. \n y = ( ) + ( ) x. \n'.format(regression.intercept_, regression.coef_[0]))

mse = mean_squared_error(y_test, y_pred)
print('Where the Slope is ( ) and the y intercept is ( ) \n')
print('With mean squared error of ( ) and Root mean squared error of ( )')
print('The Price of the house can be put up as : $({})'.format(regression.intercept_ + (regression.coef_[0] * x)))

To predict a price of a house with 2570 of sqft_living, we can use the above equation.
```

The equation is given by : y = [-44825.93903958] + [281.47743554] * 2570
y = [679371.07831926]

Thus the price of the house can be put up as : \$[679371.07831926]

Performance Improvement

From the above we can see that we are able to obtain coefficient of variation of 46%. Let us see if we could find better score using the multi-linear regression.

```
In [21]: # Considering most of the features to find the top features that helps to determine the price.
# Grid search is a method to find the best hyperparameters for a machine learning model. It is used to find the best combination of hyperparameters for a given model.
X = house_data[['bedrooms', 'bathrooms', 'sqft_living', 'sqft_lot', 'floors', 'waterfront', 'view', 'condition', 'grade', 'sqft_above', 'sqft_basement', 'yr_built', 'yr_renovated', 'zipcode']]
y = house_data['price']

select_k_best = SelectKBest(score_func=chi2, k=10)

best_result = select_k_best.fit(X, y)

scores = pd.DataFrame(best_result.score_)
columns = pd.DataFrame(X.columns)

feature_scores = pd.concat([columns, scores], axis=1)
feature_scores.columns = ['Field', 'Score']

print(feature_scores.nlargest(10, 'Score'))
```

```
Field      Score
0  sqft_living  0.930220e+00
1  yr_renovated  0.228790e+00
2  sqft_living15  0.450602e+00
3  sqft_above  0.381121e+00
4  sqft_basement  0.191784e+00
5  waterfront  0.161339e+00
6  bedrooms  0.093892e+00
7  grade  0.236518e+00
8  yr_built  0.007580e+00
9  sqft_lot  0.007580e+00
```

```
In [25]: # Defining the dataset for Logistic Regression.
predictor_data = house_data[['sqft_lot', 'yr_renovated', 'sqft_living']]
target_data = house_data['price']

X_train, X_test, y_train, y_test = train_test_split(predictor_data, target_data)

# Data normalization
scaler = MinMaxScaler()
scaler.fit(X_train)

X_train = scaler.transform(X_train)
X_test = scaler.transform(X_test)

# Defining models and parameters for our grid search
regression = LogisticRegression(n_jobs=-1, C=0.1, solver='newton-cg')
solvers = ['newton-cg', 'lbfgs', 'liblinear']
C = [0.05, 1]
# C = [0.01, 0.1]

param_grid = dict(solvers=solvers, C=C)
GridSearchCV = GridSearchCV(estimator=regression, param_grid=param_grid, n_jobs=-1, cv=cv, scoring='accuracy', error_score=0)
grid_search = GridSearchCV(estimator=regression, param_grid=param_grid, n_jobs=-1, cv=cv, scoring='accuracy', error_score=0)

# Summarizing results
print('Best: %f using %s' % (grid_search.best_score_, grid_search.best_params_))
scores = grid_search.cv_results_['mean_test_score']
std = grid_search.cv_results_['std_test_score']
params = grid_search.cv_results_['params']
for mean_score, params in zip(scores, params):
    print('(%f) using (%s)' % (mean_score, params))

The result is obtained as below. The result is copied after executing the above code in Google Colab:
```

Best: 0.008453 using [C:0.1, solver:'newton-cg']
0.008205 (0.00282) with [C:0.01, solver:'newton-cg']
0.008205 (0.00282) with [C:0.1, solver:'lbfgs']
0.008453 (0.000462) with [C:0.1, solver:'newton-cg']
0.008453 (0.000462) with [C:0.1, solver:'lbfgs']

```
In [26]: # Defining the dataset for Logistic Regression.
predictor_data = house_data[['sqft_lot', 'yr_renovated', 'sqft_living']]
target_data = house_data['price']

X_train, X_test, y_train, y_test = train_test_split(predictor_data, target_data)

# Data normalization
scaler = MinMaxScaler()
scaler.fit(X_train)

X_train = scaler.transform(X_train)
X_test = scaler.transform(X_test)

# Defining models and parameters for our grid search
regression = LogisticRegression(n_jobs=-1, C=0.1, solver='newton-cg')
regression.fit(X_train, y_train)

y_pred = regression.predict(X_test)
```

```
In [28]: print('Thus, the line of best fit for the above can be found using : \n y = ( ) + ( ) x. \n'.format(regression.intercept_, regression.coef_[0]))

mse = mean_squared_error(y_test, y_pred)
print('Where the Slope is ( ) and the y intercept is ( ) \n')
print('With mean squared error of ( ) and Root mean squared error of ( )')
print('The Price of the house can be put up as : $({})'.format(regression.intercept_ + (regression.coef_[0] * x)))

Thus, the line of best fit for the above can be found using :
y = [-6.70372325] + [0.009178372] x.
```

Where the Slope is [0.009178372] and the y intercept is [-6.70372325]
With mean squared error of 146476276489.88504 and root mean squared error of 382722.3523432748

```
In [48]: print('Similarly the coefficient of variation is : (%).format(regression.score(X_test, y_test)*100))

Similarly the coefficient of variation is : 0.6157865273373897%
```

Conclusion

Thus we can conclude that the price of the house increases with the increase in the square