

# HOUSE PRICE PREDICTION

Date : / /  
Page : /

- Change categorical values into numerical values.

- mainroad (yes/No) : (1/0)
- guestroom (yes/No) : (1/0)
- basement (yes/no) : (1/0)
- hotwaterheating (yes/no) : (1/0)
- A.C. (yes/no) : (1/0)
- prefarea (yes/no) : (1/0)
- furnishingstatus (furnished/semifurnished/unfurnished) : (1/0.5/0)

\* this will help in correlations analysis as well as at the time of model building.

## Descriptive statistics

- mean (avg. of a column)
- median
- mode
- Standard deviation
- variance
- Range

\* mean =  $\frac{x_1+x_2+\dots+x_n}{n}$ ,  $x_i$  = entries,  $n \rightarrow$  total no. of instances

\* median = arrange data in increasing order

{ If  $n \rightarrow$  odd,  $(\frac{n+1}{2})^{\text{th}}$  entry is median

{ If  $n \rightarrow$  even,  $(\frac{n}{2})^{\text{th}}$  entry is median.

\* mode = most frequent entry in the dataset

### Descriptive stats

	Mean	Median	Mode	Range	S.D.
	4766729.24	434000	3500000	11550000	1870440
	5.156.54	3.04600	6000	14550	2170
	2.94≈3	3.0	3.0	N.N.	N.N.
	1.28≈1	1.0	1.0	—	—
	1.80≈2	2.0	2.0	—	—
	0.858≈0.86	1.0	1.0	—	—
	0.177≈0.18	0	0	—	—
	0.345≈0.35	0	0	—	—
	0.045≈0.05	0	0	—	—
A.	0.315	0	0	—	—
parking	0.69	0	0	—	—
pet. area	0.23	0	0	—	—
swimming pool	0.5	0.5	0.5	—	—
studio	N.N.				

(N.N.): does not make sense for the given parameters

### Interpretation

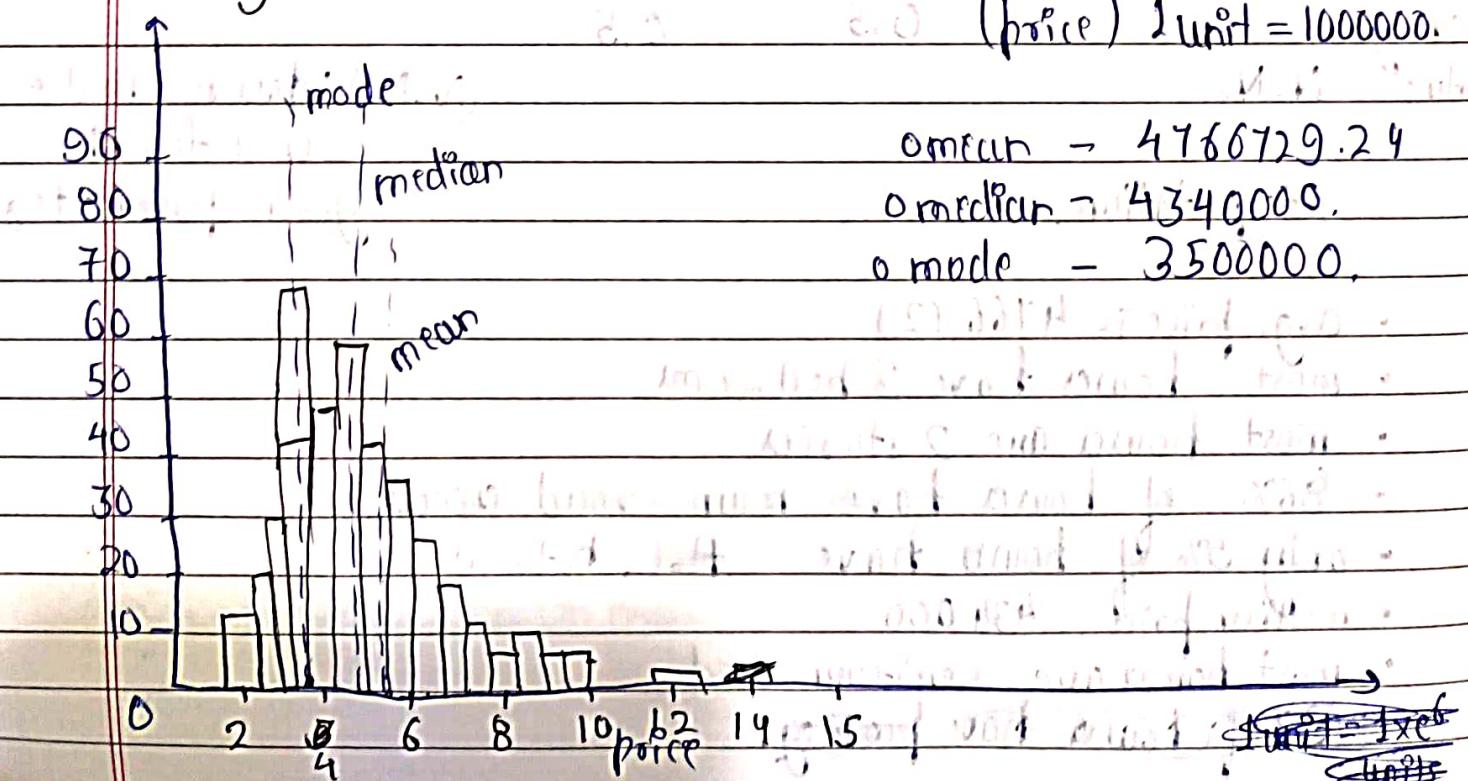
- avg. price is 4766729
- most houses have 3 bedrooms
- most houses are 2 stories
- 86% of houses have main road access
- only 5% of houses have ~~hot~~ hot water
- median price is 434000
- most houses are semifurnished

Q.

- Max Price of house = ₹ 13300000
- Min Price of house = ₹ 1750000
- Max area = 1650 sq.ft
- Min. area = 16200 sq.ft
- Max bedrooms = 6
- Min bedrooms = 1
- Range for which price varies is ₹ 11550000
- Range for which area varies is sq.ft (14550)

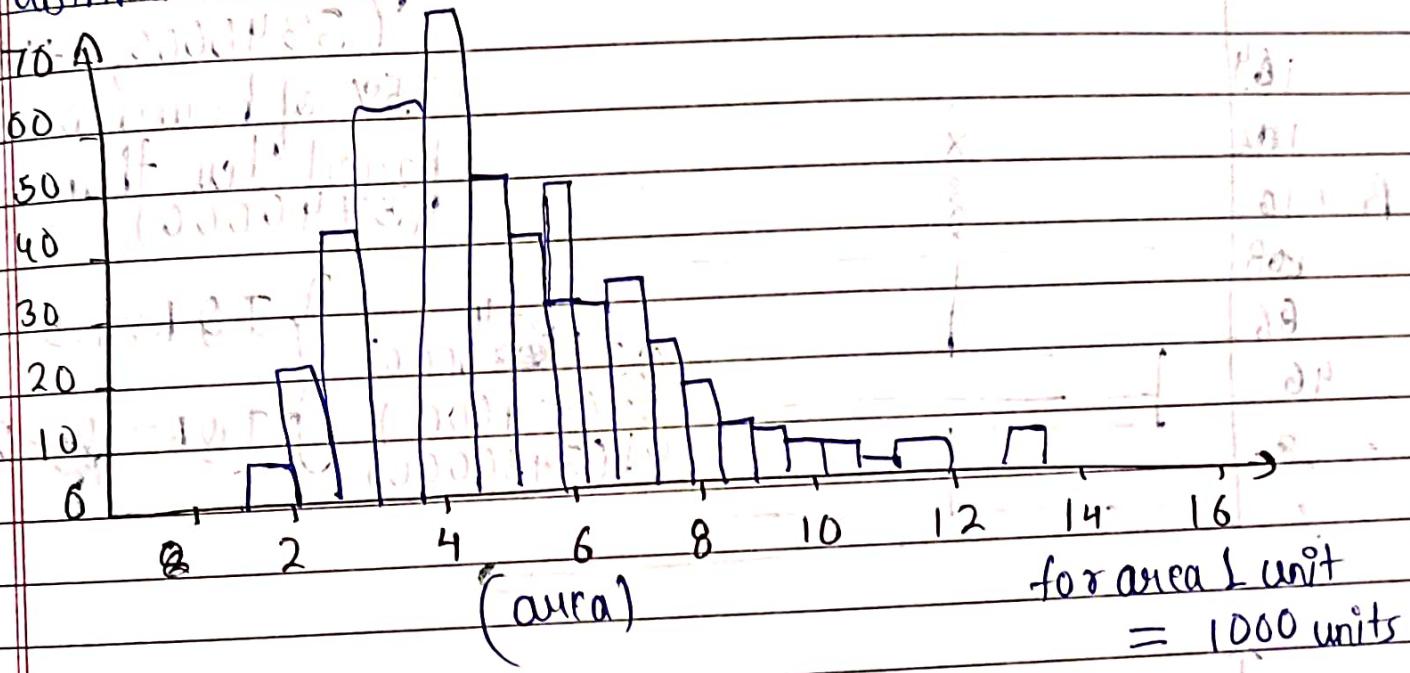
for Price (mean > median > mode) suggest positive skewness

### Visualization



- distribution visualization shows positive skewness i.e. (most of houses are priced at lower ends only few houses are very expensive).

- distribution for aura.



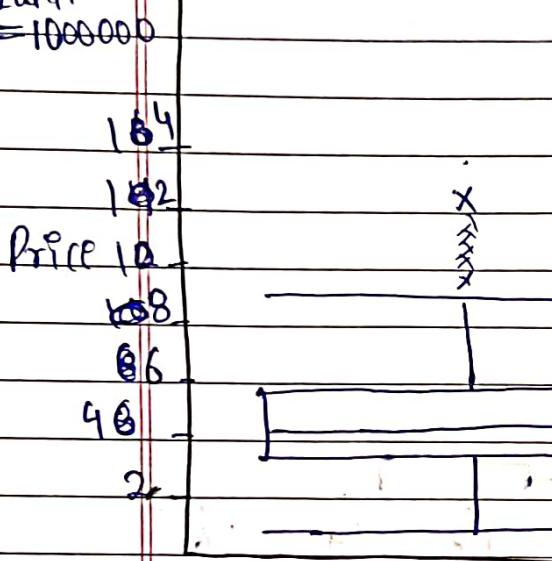
from visualization we can tell it is also +ve skewed.

⇒ most of house area is around the smaller

## Visualization for S.O. and IQR.

$$\text{IQR} = Q_3 - Q_1$$

(Price)  
unit  
= 1000000



- 50% of houses are priced b/w (3340000, 5740000)
- 75% of houses are priced less than (5740000)

$$Q_3(1200) \quad } \text{I.Q.R} \times 1.5 \\ Q_2(1000) \quad } \text{I.Q.R} = Q_3 - Q_1 \\ Q_1(800)$$

Final Box Plot.

Time

Time (x) = ?

Time (y) = ?

Time (z) = ?

Time (w) = ?

Time (v) = ?

Time (u) = ?

Time (t) = ?

Time (s) = ?

Time (r) = ?

Time (q) = ?

Time (p) = ?

Time (n) = ?

Time (m) = ?

Time (l) = ?

Time (k) = ?

Time (j) = ?

Time (i) = ?

Time (h) = ?

Time (g) = ?

Time (f) = ?

Time (e) = ?

Time (d) = ?

Time (c) = ?

Time (b) = ?

Time (a) = ?

Time (0) = ?

Time (-) = ?

We see that some points are outside the  $((5 \times IQR) \text{ whiskers})$

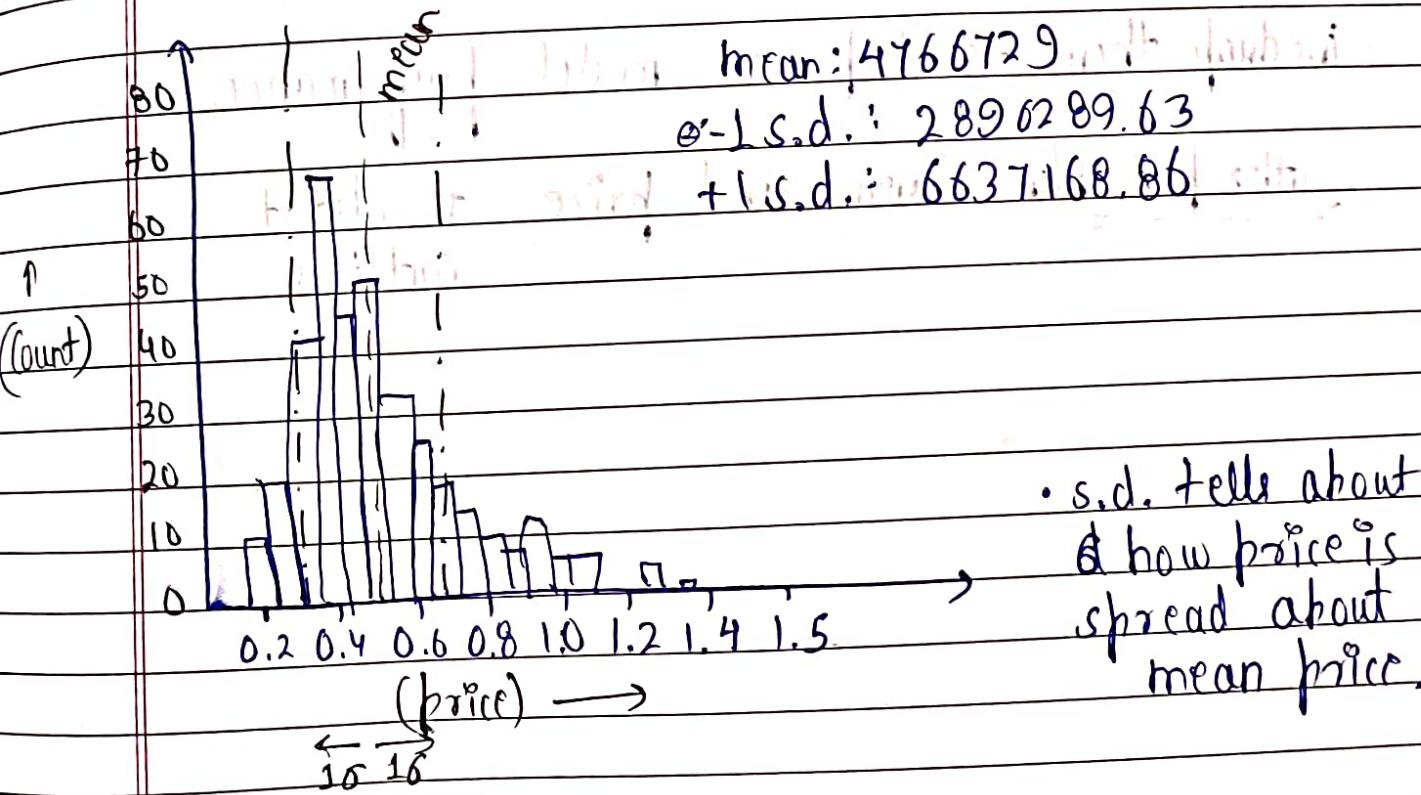
$\Rightarrow$  1-potential outliers.

which can disrupt model performance.

S.D. - tells about the spread of data about the mean.

$$\text{price} = 1870440$$

$$\text{area} = 2170 \text{ bin}(+) \text{ price}$$



## Outlier detection

### Finding Extreme outliers

- we find ( $> 3 \times IQR$ ) and remove them.

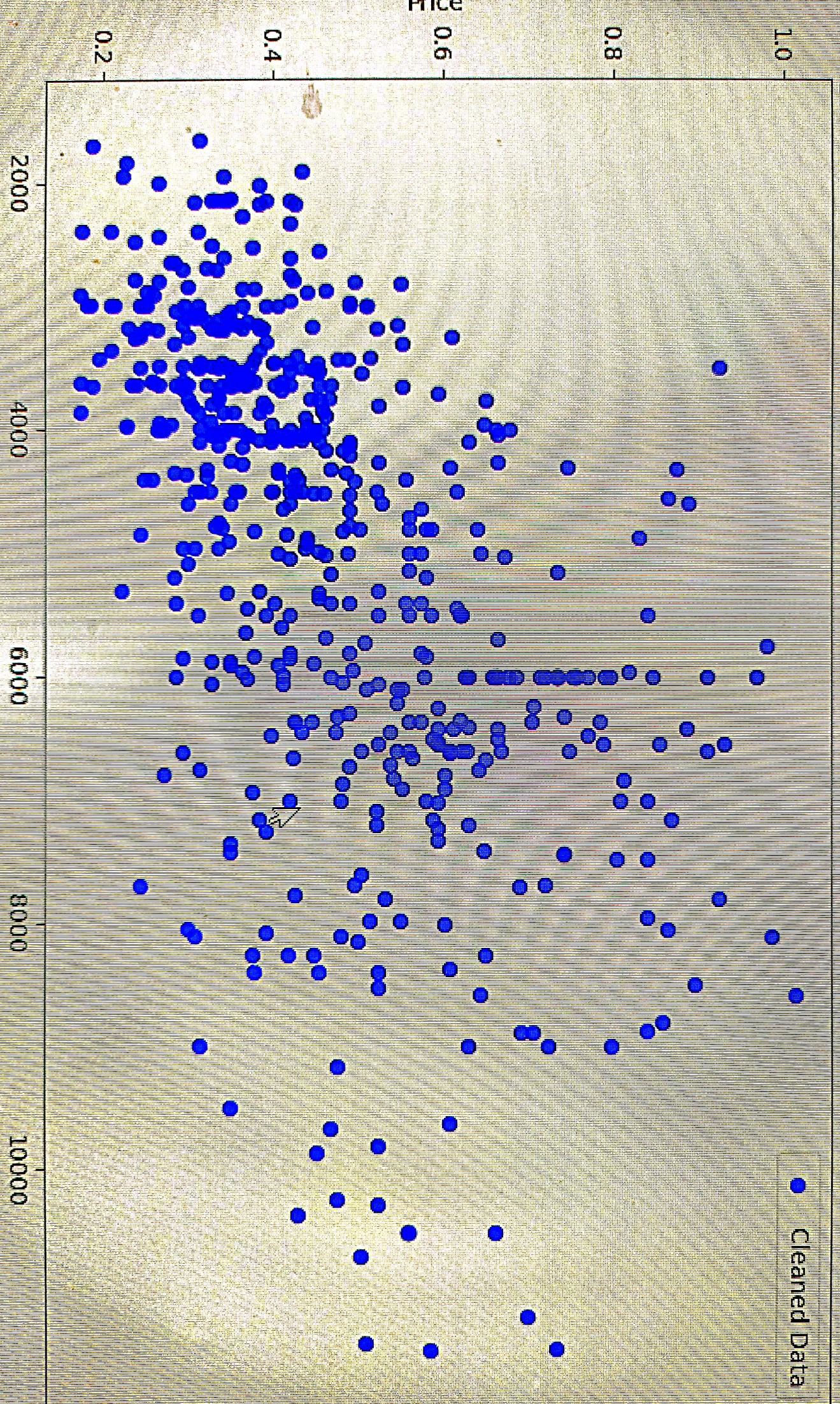
We see that there are few outliers in area.

for Price: the house with price 213300000 is an outlier

for Area in: (16200 sq. ft) and (15608 sq. ft)

We drop them to improve model performance.

# Scatter Plot of Area vs Price (Cleaned Data)



Using Z-score

we calculate  $Z$ -score for potential outliers in price and area.

and for  $Z > 3$  we can say that the points are statistically far enough to consider them outliers.

$$Z = \frac{(X - \text{mean})}{\text{standard deviation}}$$

$$\text{e.g. } Z = \frac{13300000 - 4766729.24}{1870940} = 4.56$$

$$= 4.56$$

$Z > 3$ , so, it can be considered an outlier.

similarly drop these columns: (points, bedrooms)

Now, total rows are 532 now.

→ 13 rows are dropped now.

(12, 2, 1, min, max) and max removed from

## Correlation matrix

$$\rho(x_i, y_j) = \frac{\sum (x_i - \bar{x})(y_j - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2} \sqrt{\sum (y_j - \bar{y})^2}}$$

for each pair of features calculate pearson's coefficient

we see that ~~coefficients with~~ pairs with largest coefficients are

$$*(\text{price}, \text{area}) = \cancel{0.53}$$

$$*(\text{price}, \text{bathroom}) = \cancel{0.46}$$

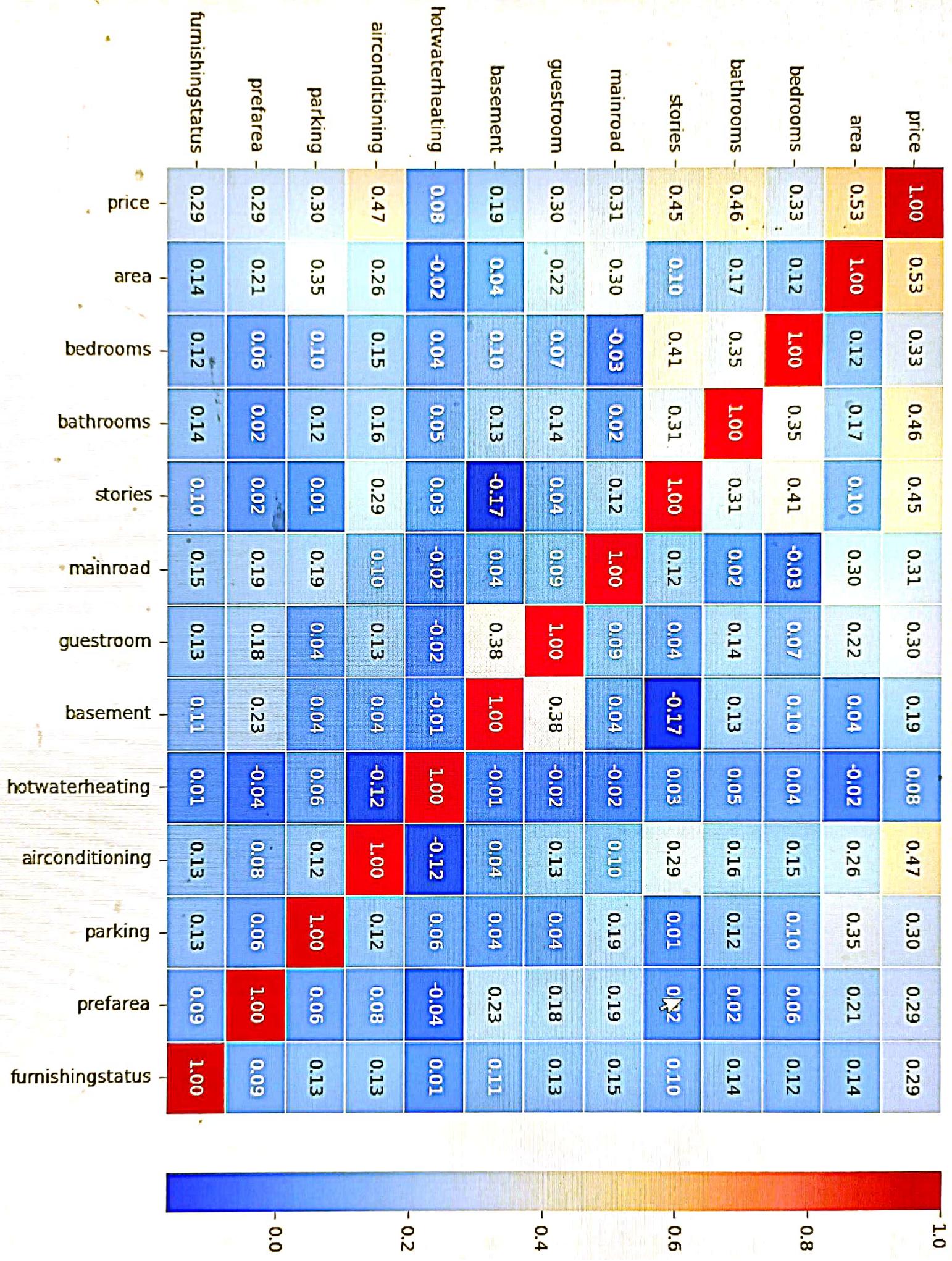
$$*(\text{price}, \text{storing}) = \cancel{0.45}$$

$$*(\text{price}, \text{A.C.}) = \cancel{0.47}$$

$$*(\text{bathrooms}, \text{storing}) = 0.41$$

we see that correlation found is not strong  
most prominent ones are ( $\text{price}$ ,  $\text{area}$ ,  $\text{A.C.}$ ,  $\text{storing}$ )

Correlation Matrix of Cleaned Dataset



## 1.3 Feature Engineering

### feature scaling :

Target : "price" ( $y$ )

$$z = \frac{x - \mu}{\sigma}, \mu \rightarrow \text{mean}, \sigma \rightarrow \text{S.D.}$$

### normalize the features

more about it:

### feature selection :

\* Price ( $y$ ): target variable

\* area ( $X_1$ )

\* bedrooms ( $X_2$ )

\* bathrooms ( $X_3$ )

\* stories ( $X_4$ )

\* parking ( $X_5$ )

to make the model simple we can combine some features

we combine bedrooms and bathrooms so total rooms.

$$\text{total rooms} = \text{bedrooms} + \text{bathrooms}$$

## final selected features

- $x_1$  (area)
- $x_2$  (total rooms)
- ~~• age (years)~~

## Linear Regression model -

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon$$

$\beta_0$ : intercept  
 $\beta_1$ : coefficient of area  
 $\beta_2$ : coefficient of total\_rooms  
 $y$ : target variable.  
 $x_1$ : area  
 $x_2$ : total\_rooms  
 $\epsilon$ : error term

## Analytical method:

$$\beta = (X^T X)^{-1} X^T y$$

## Gradient descent method:

### Cost function:

$$J(\beta) = \frac{1}{2m} \sum_{i=1}^m (y_i - \hat{y}_i)^2$$

$J(\beta)$ : cost function, but unrolled over labels

$y_i$ : actual input

$\hat{y}_i$ : predicted output

## Model Training

~~Step 1~~

Initialize parameters

$$\beta_0 = 0, \beta_1 = 0, \beta_2 = 0$$

Compute gradients

$$\frac{\partial L}{\partial \beta_0} = -2 \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_{1,i} + \beta_2 x_{2,i}))$$

$$\frac{\partial L}{\partial \beta_1} = -2 \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_{1,i} + \beta_2 x_{2,i})) x_{1,i}$$

$$\frac{\partial L}{\partial \beta_2} = -2 \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_{1,i} + \beta_2 x_{2,i})) x_{2,i}$$

Iteration one

$$\frac{\partial L}{\partial \beta_0} = 0$$

$$\frac{\partial L}{\partial \beta_1} = -1.0671$$

$$\frac{\partial L}{\partial \beta_2} = -0.9173$$

Update Coefficients

$$\alpha = 0.01$$

$$\beta_0^{(\text{new})} = \beta_0 - (0.01) \frac{\partial L}{\partial \beta_0}$$

$$= 0 - (0.01) \times 0$$

$$\beta_1^{(\text{new})} = \beta_1^{(\text{old})} - (0.01) \times (-0.9173) (-1, 0.67)$$

$$= 0.0107 + 0.01 \times (-0.9173) = 0.0107 - 0.009173 = 0.0092$$

$$\beta_2^{(\text{new})} = \beta_2^{(\text{old})} - (0.01) \times (-0.9173)$$

Iteration 2

$$\frac{\partial L}{\partial \beta_0} = 0$$

$$\frac{\partial L}{\partial \beta_1} = -1.0426$$

$$\frac{\partial L}{\partial \beta_2} = -0.8954$$

$$\frac{\partial L}{\partial \beta_2} = -0.8954$$

Update coefficients:

$$\beta_0^{(\text{new})} = \beta_0^{(\text{old})} - (0.01) \frac{\partial L}{\partial \beta_0}$$

$$= 0 - (0.01) \times 0 \\ = 0$$

$$\beta_1^{(\text{new})} = \beta_1^{(\text{old})} - (0.01)(-1.0426) \\ = 0.0211$$

$$\beta_2^{(\text{new})} = \beta_2^{(\text{old})} - (0.01)(-0.8954) \\ = 0.0181$$

Model evaluation

(After convergence using numpy)

mean absolute error:  $\frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$

RMSE:  $\sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$

model after 2 iterations

$$y = 0 + 0.0211(x_1) + 0.0181(x_2)$$

we will calculate for 53 samples.

$x_1$  (area)

$x_2$  (total room)

(Area) (room)

7420, 6

8960, 8

9360, 5

$$\hat{y}_1 = 0 + 0.0211(7420) + 0.0181(6)$$

$$= 156.6706$$

$$\hat{y}_2 = 0 + 0.0211(8960) + 0.0181(8)$$

$$= 189.2008$$

$$\hat{y}_3 = 202.11$$

now,

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

$$= \frac{1}{3} \left[ (3300000 - 156.6706) + \left( \frac{8960}{12250000} - 189.2008 \right) + (12250000 - 202.11) \right]$$

= ~~if~~ very large value

. it is done for very few iterations the value is extremely large after convergence it will decrease.

Final results (using numpy)

$$\beta_1 = 0.448$$

$$\beta_2 = 0.364$$

$$\beta_0 = -0.024756 \approx 0$$

$$MAE = 0.604$$