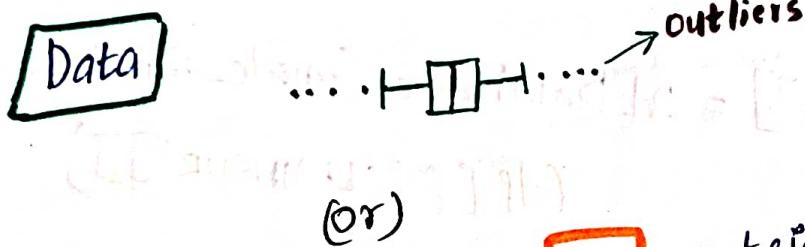


3/4/22 2:08 AM

Outliers → # Applicable Only For
Continuous Data

1 Que :- What is Outlier?

Ans :- An Outlier is a data point in a dataset.
that is distant from all other observation
which is significantly different from the remaining



(or)

- The data point that lies outside the overall distribution of the dataset.

2 Que :- What are Impacts having Outliers in a DataSet?

Ans :- It will cause various problems in statistical analysis. (it may cause a significant impact on mean and standard deviation)

Ex :- $\bar{x} = \text{mean (affect)}$

$$\frac{\sum (x - \bar{x})^2}{n-1} = \text{variance}$$

$$\sqrt{\frac{\sum (x - \bar{x})^2}{n-1}} \quad (\text{standard deviation})$$

all are affected by outliers

* In addition, Some Machine learning models are Sensitive to Outliers, which may Decrease There Performance

3 Que : Reasons For Outliers ?

- Ams :
1. Data Entry Errors (Ex:- Salary = 100,000, But 10,000 Entering)
 2. Measurement Errors (Ex:- Measuring in meters instead of KM)
 3. Instrumental Error

4 Que :- Types of Outliers ?

- Ams :-
- * **Univariate Outlier** ---> Identifying Outlier For Single Variable
 - * **Bivariate Outlier** ---> Identifying as Outlier by Analyzing 2 Variables at a Time

Importing Libraries

- # Import numpy as np
- # Import pandas as pd
- # Import matplotlib as plt.
- # Import Seaborn as sns.

Real Time Case study

boston - House Price Dataset

```
# boston = pd.read_csv("D:\\data science \\Shubham\\Krishna class \\6. Data cleaning \\boston.csv")
```

```
boston.head()
```

out

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO
0	0.00632	18.0	2.31	0	0.538	6.575	65.2	4.090	1	29.6	15.3
1	0.02731	0.0	7.07	0	0.469	6.421	78.9	4.967	2	24.2	17.8
2	0.02729	0.0	7.07	0	0.469	7.185	61.1	4.967	2	24.2	17.8
3	0.03237	0.0	2.18	0	0.458	6.998	45.8	6.062	3	22.2	18.7
4	0.06905	0.0	2.18	0	0.458	7.147	54.2	6.062	3	22.2	18.7

B	LSTAT	PRICE
396.90	4.98	24.0
396.90	9.14	21.6
392.83	4.03	34.7
394.63	2.94	33.4
396.90	5.33	36.2

```
# boston.info()
```

506 rows, 14 columns

out

0	CRIM	float64
1	ZN	float64
2	INDUS	float64
3	CHAS	→ int64
4	NOX	float64
5	RM	float64
6	AGE	float64
7	DIS	
8	RAD	→ int64
9	TAX	→ int64
10	PTRATIO	float64
11	B	float64
12	LSTAT	float64
13	PRICE	float64

(float64(11), int64(3))

Various ways of finding The Outlier

Changing Every Value to Z score.

option ①

Z score

$$|Z\text{score}| > 2$$

$$\left[\begin{array}{l} \therefore z < -2 \\ z > +2 \end{array} \right] \mid \text{modulus} \mid$$

Ex :-

X	Z _{score}
6	4.56
4	-1
9	9 - 5.6
7	7 - 5.6
2	2 - 5.6

$$Z = \frac{X - \mu}{\sigma}$$

$$I = \frac{6 - 5.6}{1}$$

$$\therefore \mu = \frac{6 + 4 + 9 + 7 + 2}{5}$$

$$\mu = 28/5 = 5.6$$

$$EX/\sigma = 1$$

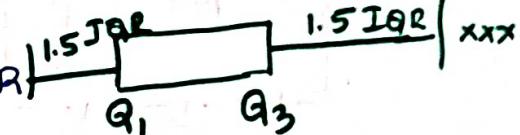
option ②

IQR

(Inter Quartile Range)

any value $< Q_1 - 1.5 \text{ IQR}$

any value $> Q_3 + 1.5 \text{ IQR}$



any value Less than $Q_1 - 1.5 \text{ IQR}$ is Outlier

any value More Than $Q_3 + 1.5 \text{ IQR}$ is Outlier

option ③

Visualization

* Box plot, Histogram For Univariate Outliers

* Scatter plot For bivariate Outliers.

* ...

option③

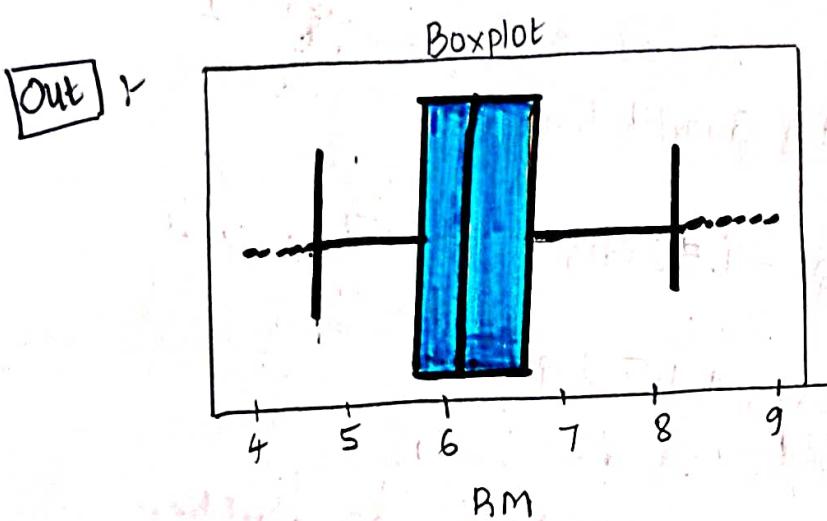
Detection of Outliers of "RM" column (Based on Boxplot)

sns. boxplot (boston.RM)

plt. title ("Boxplot")

plt. show()

* Overall view of Outlier



Outliers in Both Tails of "RM".

option①

Detection of Outliers of "RM" column (Based on Zscore)

* For Exact Outliers in Data

$$Z\text{score} = \frac{x - \mu}{\sigma}$$

boston["RM"].mean(), boston["RM"].std()

[out] (6.284, 0.702)

boston["RM_Zscore"] = ([boston["RM"] - boston["RM"].mean()]) / boston["RM"].std()

creating new column
and storing Zscores

$$\frac{\underline{[RM].mean()}}{\underline{\mu}} \quad \underline{\%} \quad \underline{\sigma}$$

Output :-

clarity :-

$$\frac{x - \mu}{\sigma} = \frac{(\text{boston}['RM']) - \text{boston}['RM'].mean())}{\text{boston}['RM'].std()}$$

$$\frac{x - \mu}{\sigma}$$

$$\frac{x - \mu}{\sigma}$$

E Answers Ami

$$\text{boston}['RM_Zscore'] =$$

↑ Endulo store age + thage.

ZSCORE

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B	... RM SCORE
0	0.0632	18.0	2.31		6.515								0.413263
1	0.02731	0.0	7.07		6.421								0.194082
2	0.02729	0.0	2.18		7.185								1.281446
..													
505	0.04741	0.0	11.93		6.030								-0.362408

506 x 15 columns

Now, Identify Outliers in this column ↑ OR

$$\# \text{Outlier} = (\text{boston}['RM_Zscore'] > +2) \quad (\text{boston}['RM_Zscore'] < -2)$$

2 conditions

Greater Than +2
Less Than -2

[Out] :- Outlier

	CRIM	RM SCORE
97		2.539
98		2.185
162		3.47
163		-2.121
66		-3.055

3/4/22 4:30 AM

Time
1:00pm

Outlier. Shape

[Out] :- (31, 15)

option 2
Detection of Outlier of RM (based on IQR)

- Calculate First (Q_1) and third quartile (Q_3)
 - Find interquartile range ($Q_3 - Q_1$)
 - Find Lower bound $Q_1 - 1.5$ & Find upper bound $Q_3 + 1.5$
- for Exact no. of range + - of outliers

$Q_3 = \text{boston}["RM"].quantile(0.75)$

$Q_1 = \text{boston}["RM"].quantile(0.25)$

$$Q_1 = 0.25(w)$$

$$Q_3 = 0.75(w)$$

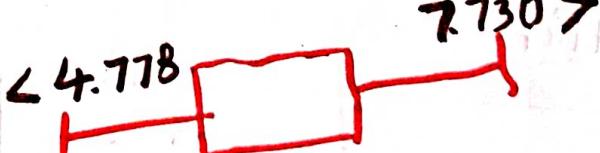
$$IQR = Q_3 - Q_1$$

$$\text{Lower-limit} = Q_1 - (IQR * 1.5) \quad IQR = Q_3 - Q_1$$

$$\text{Upper-limit} = Q_3 + (IQR * 1.5) \quad \Rightarrow Q_3 + IQR * 1.5$$
$$\Rightarrow Q_1 - IQR * 1.5$$

Lower-limit, upper-Limit

[Out] :- 4.778, 7.730



- * Which is having less than 4.778 is outlier
- * Which is having greater than 7.730 is outlier

Methods of deal The Outliers

- 3R → * Remove → Deleting The Outliers
- Technique * Replace → changing values / Data Manipulation
- * Retain → Treat Them Separately.
Work without outlier / with outlier
"Work" Separately.

Outliers should be detected and "REMOVED" Only From Training data set, NOT from The Test Set. So we should first divide our data set into train and test. and remove Outliers in the train set, but keep those in test set, and Measure how well our model is doing.

option ① Removing * (Let's trimm The dataset) 74.77
boston_trimmed = boston[(boston["RM"] > Lower-limit) & (boston["RM"] < Upper-limit)] 7.73

boston_trimmed
out :- which ever Data is Greater than 4.77 and Less Than 7.73 Consider That data Only, Remove Remaining data.

	RM	RM - Zscore
0	6.515	0.413263
1	6.421	0.194082
2	7.185	1.2814416
3	6.794	0.724955
4	6.030	-0.362408

∴ From 506 records, Removed 30 records

that Means 476 having (or) there.

⇒ it's about 6% of the data is Removed

⇒ it is Only in One Variable (1 column) = 6%.

⇒ If Total Variable (all columns) $\frac{1,2,3,4,\dots}{(14) \rightarrow [10] \text{ continuous data}}$ Remove May Delete 100 rows
20% of Data

So, In real Life, we don't use outliers "Remove"

Even, Though, If we Want Remove, still it Contains Outliers.

⇒ Outliers in The Trimmed dataset

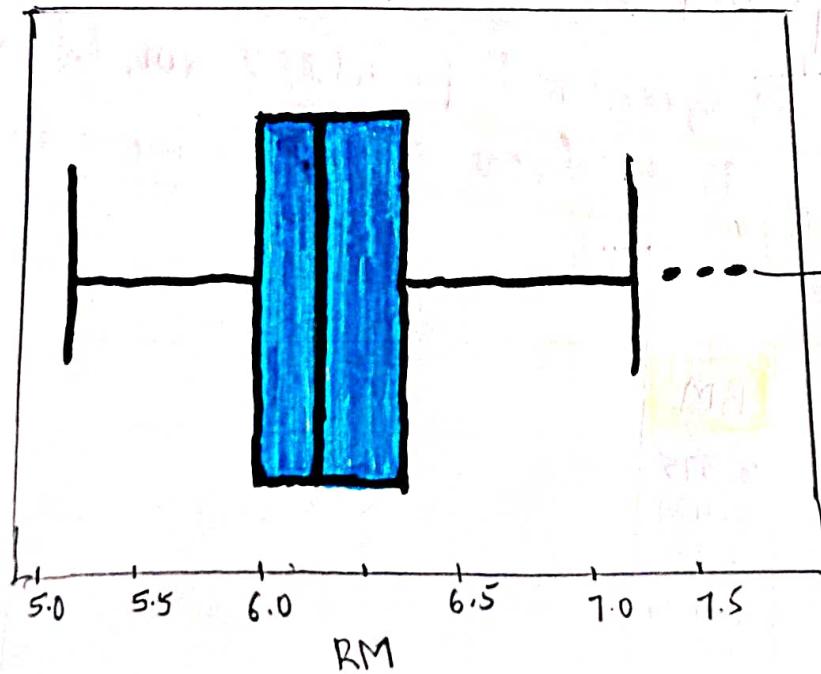
sns.boxplot(boston_trimmed.RM)

plt.title("Boxplot")

plt.show()

Boxplot

Out



still, we have Outliers, after Removing The Outlier records

Q: Why? it showing Outliers after Removing records of outliers

A:- When we remove datapoints from our dataset, all the parameters of the distribution are re-calculated.

For 476 records,
Again, it's calculating (Q_1 , Q_3 , IQR)
(again)

for
506

$$G_1 = 4.77, \quad G_1 = \text{New Value (again)} \\ G_3 = 7.730, \quad G_3 = \text{New value (again)}$$

- * Therefore, In The new-trimmed variables Values that before were not considered Outliers.

That's why, it showing new Outliers after removing some records also.

option ②. Replace

Replace with $\frac{\text{Upper Limit}}{\text{Based IQR}} * \frac{\text{Lower Limit}}{\text{Calculated}}$

Replacing with
 4.7778 (min value)

min. Value
10

Every Outlier

130

$$Q_1 - [IQR * 1.5]$$

$$77.730 = Q_3 + [IGR^* 1.5]$$

↓
Every Outlier Replacing with

7.730 (max.value)

⇒ More Detailness

Replacing with Q_1, Q_3 Values

CRIM	...	RM	...
97	More Than 7.730, so, they are Outliers. Reduce with $\bar{x} = 7.730$	8.069 7.820 7.802 8.375	7.730 7.730 7.730 7.730
98			
162			
163			

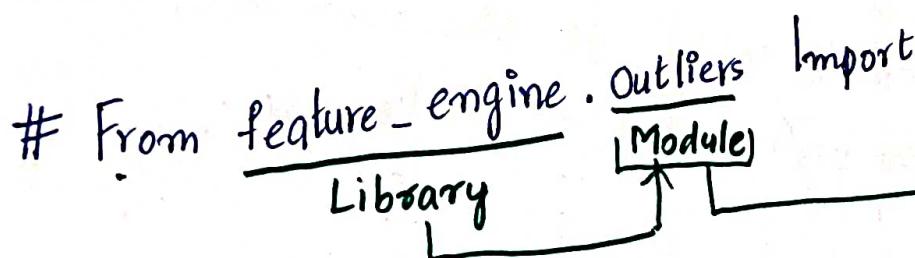
$$Q_3 - IQR * 1.5 < [4.778] \xrightarrow{\text{Lower-value}} \text{less than } 4.778 \xrightarrow{\text{= outlier}}$$

$$Q_1 + IQR * 1.5 > [7.730] \xrightarrow{\text{Upper-value}} \text{greater than } 7.730 \xrightarrow{\text{= outlier}}$$

* Winsorizer

PIP install

feature-engine



Winsorizer

Winsorizer
Function

In Winsorizer (capping method)
Function argument

- * Capping method = 1. IQR
- 2. "Gaussian"

In feature-engine. Outliers Import ArbitraryOutlierCapper

```

graph LR
    FE[feature-engine] --> Outliers[Outliers]
    Outliers --> Import[Import]
    subgraph Library [Library]
        Outliers
    end
    
```

We can give our own Lower, upper values.

1. By IQR method

In Depth of **Boxplot**

Firstly,

we calculate Q_1 Value

$Q_1 = \text{boston}["RM"].quantile(0.25)$

Q_1

[Out] :- 5.885

506

Records

$$\therefore Q_1 = 5.885$$

we calculate Q_3 Value.

$Q_3 = \text{boston}["RM"].quantile(0.75)$

Q_3

[Out] :- 6.623

$$Q_3 = 6.623$$

$$\begin{aligned} IQR &= Q_3 - Q_1 \\ &= 0.737 \end{aligned}$$

Distance between Q_3 and Q_1

IQR Value

$IQR = Q_3 - Q_1$

IQR

0.737

CODE :-

From feature_engine.Outliers import Wimsozier

$wim = \text{Wimsozier}(\text{capping_method} = "iqr", \text{tail} = "both", \text{fold} = 1.5, \text{Variables} = ["RM"])$

$boston_t = wim.fit_transform(\text{boston}[[\text{"RM"}]])$

Print(wim.left_tail_caps_, wim.right_tail_caps_)

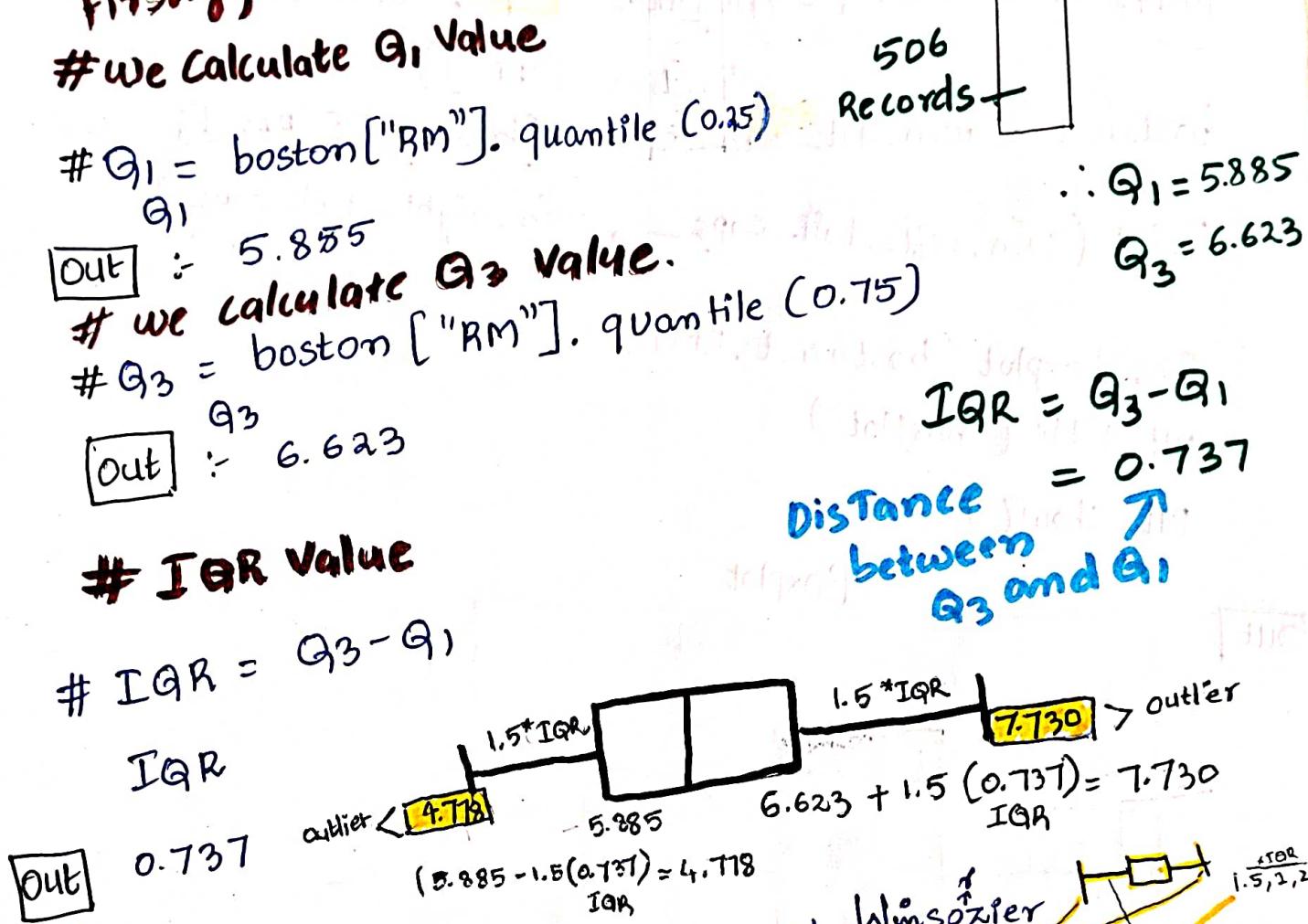
4.778 7.730

sns.boxplot(boston_t.RM)

plt.title("boxplot")

plt.show()

Edited (replaced IQR). RM



again

From feature_engine.Outliers import Winsorizer

win = Winsorizer(capping_method="lqr", fold=1.5, Tail="both")

Variables = ["RM"])

boston_t = win.fit_transform(boston[["RM"]])

Print(win.left_tail_caps_, win.right_tail_caps_)

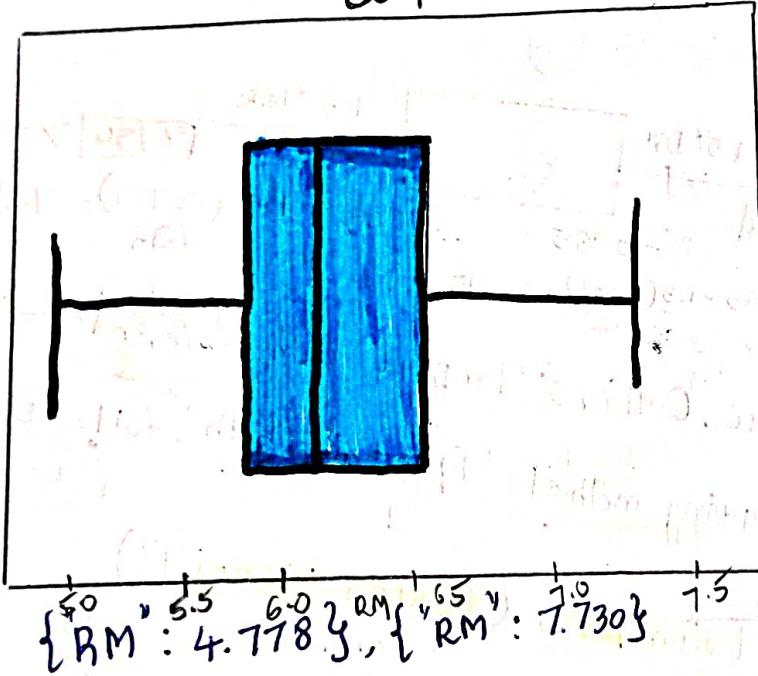
Sns.boxplot(boston_t.RM)

plt.title("boxplot")

plt.show()

Boxplot

Out



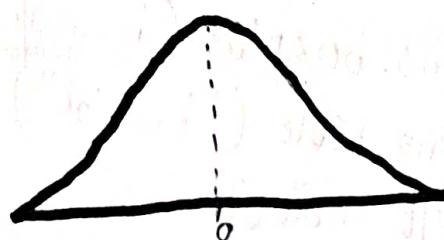
2, Replacing outliers using Winsorizer (min, max, automatically taken by "Gaussian method")

Gaussian method

= "Normal"

Same

Distribution method



```
# from feature_engine.outliers import Winsorizer
```

```
# Wim = Winsorizer(capping_method = "gaussian", tail = "both")
```

```
fold = 1.5, variables = ['RM'])
```

```
boston_t = Wim.fit_transform(boston[['RM']])
```

```
print(Wim.left_tail_caps_, Wim.right_tail_caps_)
```

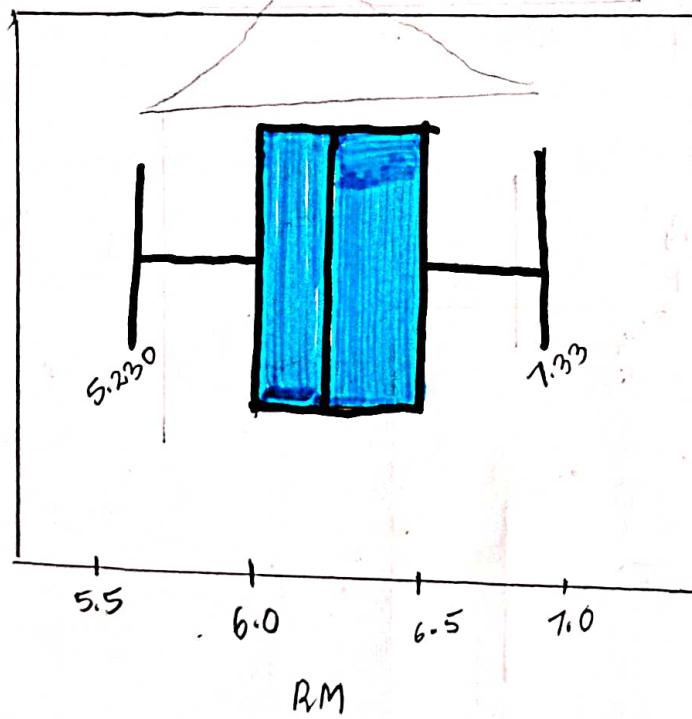
```
Sns.boxplot(boston_t.RM)
```

```
plt.title("Boxplot")
```

```
plt.show()
```

```
{'RM': 5.230} {'RM': 7.33}
```

BOXPLOT



Out

3. Replace arbitrary Outlier Capper ("The min, max, by user")
~~~~~  
# we get the **min** and **max** values based on "Domain Expert"  
(or) by **Our Own Research**

# from feature\_engine.Outliers import ArbitraryOutlierCapper

Capper = ArbitraryOutlierCapper (**max\_capping\_dict** = { "RM": 7.5 },  
**min\_capping\_dict** = { "RM": 4.8 })

boston\_c = Capper.fit\_transform(boston[["RM"]])

Print (Capper.right\_tail\_caps\_, Capper.left\_tail\_caps\_)

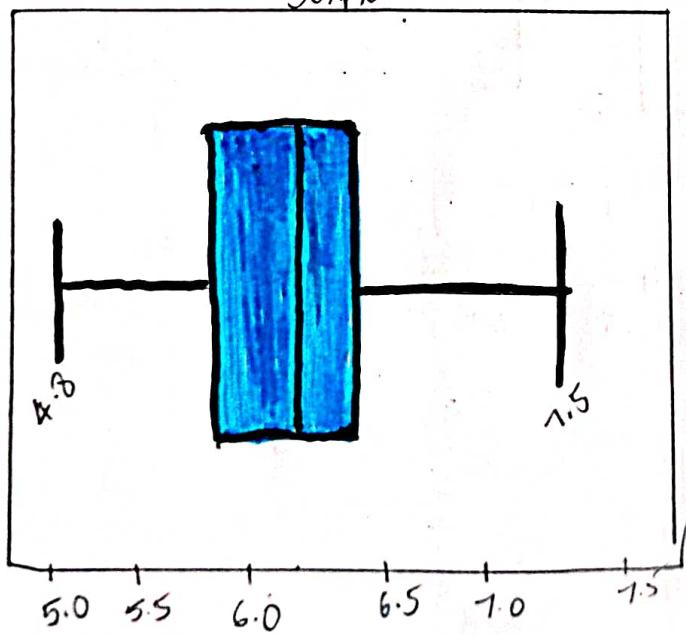
Sns.boxplot(boston\_c.RM)

plt.title("Boxplot")

plt.show()

{ "RM": 4.8 }, { "RM": 7.5 }  
Boxplot

Out



3/4/2022 5:20 PM