

1)Generate the summary statistics for each variable in the table. (Use Data analysis tool pack). Write down your observation.

CRIME_RATE	
Mean	4.871976285
Standard Error	0.129860152
Median	4.82
Mode	3.43
Standard Deviation	2.921131892
Sample Variance	8.533011532
Kurtosis	-1.189122464
Skewness	0.021728079
Range	9.95
Minimum	0.04
Maximum	9.99
Sum	2465.22
Count	506

*The sharp ratio of Crime Rate is 1.67

*Crime Rate is Flat Curve with Positive skewness

AGE	
Mean	68.57490119
Standard Error	1.251369525
Median	77.5
Mode	100
Standard Deviation	28.14886141
Sample Variance	792.3583985
Kurtosis	-0.967715594
Skewness	-0.59896264
Range	97.1
Minimum	2.9
Maximum	100
Sum	34698.9
Count	506

*The sharp ratio of Age is 2.43

*Age is Flat Curve with Negative skewness

INDUS	
Mean	11.13677866
Standard Error	0.304979888
Median	9.69
Mode	18.1
Standard Deviation	6.860352941
Sample Variance	47.06444247
Kurtosis	-1.233539601
Skewness	0.295021568
Range	27.28
Minimum	0.46
Maximum	27.74
Sum	5635.21
Count	506

*The sharp ratio of Indus is 1.62

*Indus is Flat Curve with Positive skewness

NOX	
Mean	0.554695059
Standard Error	0.005151391
Median	0.538
Mode	0.538
Standard Deviation	0.115877676
Sample Variance	0.013427636
Kurtosis	-0.064667133
Skewness	0.729307923
Range	0.486
Minimum	0.385
Maximum	0.871
Sum	280.6757
Count	506

*The sharp ratio of NOX is 4.78

*NOX is Flat Curve with Positive skewness

<i>DISTANCE</i>	
Mean	9.549407115
Standard Error	0.387084894
Median	5
Mode	24
Standard Deviation	8.707259384
Sample Variance	75.81636598
Kurtosis	-0.867231994
Skewness	1.004814648
Range	23
Minimum	1
Maximum	24
Sum	4832
Count	506

*The sharp ratio of Distance is 1.10

*Distance is Flat Curve with Positive skewness

<i>TAX</i>	
Mean	408.2371542
Standard Error	7.492388692
Median	330
Mode	666
Standard Deviation	168.5371161
Sample Variance	28404.75949
Kurtosis	-1.142407992
Skewness	0.669955942
Range	524
Minimum	187
Maximum	711
Sum	206568
Count	506

*The sharp ratio of Tax is 2.42

*Tax is Flat Curve with Positive skewness

<i>PTRATIO</i>	
Mean	18.4555336
Standard Error	0.096243568
Median	19.05
Mode	20.2
Standard Deviation	2.164945524
Sample Variance	4.686989121
Kurtosis	-0.285091383
Skewness	-0.802324927
Range	9.4
Minimum	12.6
Maximum	22
Sum	9338.5
Count	506

*The sharp ratio of PTRATIO is 8.52

*PTRATIO is Flat Curve with Negative skewness

<i>AVG_ROOM</i>	
Mean	6.284634387
Standard Error	0.031235142
Median	6.2085
Mode	5.713
Standard Deviation	0.702617143
Sample Variance	0.49367085
Kurtosis	1.891500366
Skewness	0.403612133
Range	5.219
Minimum	3.561
Maximum	8.78
Sum	3180.025
Count	506

*The sharp ratio of Avg_room is 8.94

*Avg_room is sharp with Positive skewness

LSTAT	
Mean	12.65306324
Standard Error	0.317458906
Median	11.36
Mode	8.05
Standard Deviation	7.141061511
Sample Variance	50.99475951
Kurtosis	0.493239517
Skewness	0.906460094
Range	36.24
Minimum	1.73
Maximum	37.97
Sum	6402.45
Count	506

*The sharp ratio of LSTAT is 1.77

*LSTAT is sharp with Positive skewness

AVG_PRICE	
Mean	22.53280632
Standard Error	0.408861147
Median	21.2
Mode	50
Standard Deviation	9.197104087
Sample Variance	84.58672359
Kurtosis	1.495196944
Skewness	1.108098408
Range	45
Minimum	5
Maximum	50
Sum	11401.6
Count	506

*The sharp ratio of Avg_price is 2.44

*Avg_price is sharp with Positive skewness

2)Plot a histogram of the Avg_Price variable. What do you infer?



Avg_price is sharp with Positive skewness.

3) Compute the covariance matrix. Share your observations

	CRIME_RATE	AGE	INDUS	NOX	DISTANCE	TAX	PTRATIO	AVG_ROOM
CRIME_RATE	8.51614787							
AGE	0.56291522	790.79247						
INDUS	-0.1102152	124.26783	46.97143					
NOX	0.00062531	2.3812119	0.605874	0.0134011				
DISTANCE	-0.2298605	111.54996	35.47971	0.6157102	75.6665313			
TAX	-8.2293224	2397.9417	831.7133	13.020502	1333.11674	28348.6236		
PTRATIO	0.06816891	15.905425	5.680855	0.0473037	8.74340249	167.820822	4.6777263	
AVG_ROOM	0.05611778	-4.742538	-1.884225	-0.024555	-1.2812774	-34.515101	-0.5396945	0.4926952
LSTAT	-0.8826804	120.83844	29.52181	0.4879799	30.3253921	653.420617	5.7713002	-3.073655
AVG_PRICE	1.16201224	-97.39615	-30.4605	-0.454512	-30.50083	-724.82043	-10.090676	4.4845656

* **Positive value** denotes, both the x and y values are above or below their averages.

* **Negative value** denotes, both the x and y values are mostly on opposite sides of their averages.

4) Create a correlation matrix of all the variables (Use Data analysis tool pack)

	CRIME_RATE	AGE	INDUS	NOX	DISTANCE	TAX	PTRATIO
CRIME_RATE	1						
AGE	0.006859463	1					
INDUS	-0.005510651	0.644778511	1				
NOX	0.001850982	0.731470104	0.763651447	1			
DISTANCE	-0.009055049	0.456022452	0.595129275	0.611440563	1		
TAX	-0.016748522	0.506455594	0.72076018	0.6680232	0.910228189	1	
PTRATIO	0.010800586	0.261515012	0.383247556	0.188932677	0.464741179	0.460853035	1
AVG_ROOM	0.02739616	-0.240264931	-0.391675853	-0.302188188	-0.209846668	-0.292047833	-0.35550
LSTAT	-0.042398321	0.602338529	0.603799716	0.590878921	0.488676335	0.543993412	0.37404
AVG_PRICE	0.043337871	-0.376954565	-0.48372516	-0.427320772	-0.381626231	-0.463535934	-0.50778

A) Which are the top 3 positively correlated pairs B) Which are the top 3 negatively correlated pairs.

*Distance and tax

*Index and Nox

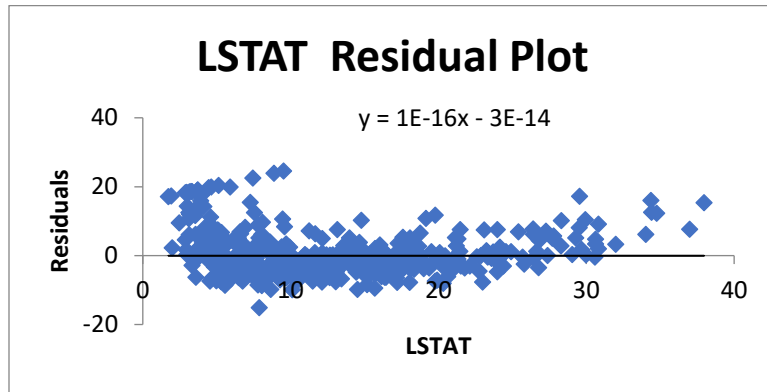
*Age and nox

*LSTAT and Avg price

*Avg room and LSTAT

*PTRATIO and Avg price

5) Build an initial regression model with AVG_PRICE as 'y' (Dependent variable) and LSTAT variable as Independent Variable. Generate the residual plot.



A) What do you infer from the Regression Summary output in terms of variance explained, coefficient value, Intercept, and the Residual plot?

R Square 0.544146298

Coefficients of LSTAT -0.950049354

Intercept 34.55384088

R Square:

There are 54% changes for LSTAT and Avg price

Coefficients of LSTAT:

Coefficient of LSTAT is -0.95005. It is inferred that if the Average price is increase, there will be a 0.95% decrease in population.

Intercept:

It is inferred that the Intercept value is 34.5538.

Residual plot:

* It is inferred that all the values are equally distributed

* Linear equation is Avg price = -0.95+34.554

B) Is LSTAT variable significant for the analysis based on your model?

Yes, it is significant for analysis

6) Build a new Regression model including LSTAT and AVG_ROOM together as Independent variables and AVG_PRICE as dependent variable

A) Write the Regression equation. If a new house in this locality has 7 rooms (on an average) and has a value of 20 for L-STAT, then what will be the value of AVG_PRICE? How does it compare to the company quoting a value of 30000 USD for this locality? Is the company Overcharging/ Undercharging?

*Avg price = Coefficient of intercept + (Coefficient of AVG_Room * Avg_Room) +(Coefficient of LSTAT * LSTAT)

*The company is overcharging, and the company can quote the amount as 21458 USD.

B) Is the performance of this model better than the previous model you built in Question 5? Compare in terms of adjusted R-square and explain

Yes, the performance of this model is better than the previous model(Q5). The R square value is improved because we have added AVG_Room for this regression.

7) Build another Regression model with all variables where AVG_PRICE alone be the Dependent Variable and all the other variables are independent. Interpret the output in terms of adjusted R square, coefficient and Intercept values. Explain the significance of each independent variable with respect to AVG_PRICE

	<i>Coefficients</i>
Intercept	29.24131526
CRIME_RATE	0.048725141
AGE	0.032770689
INDUS	0.130551399
NOX	-10.3211828
DISTANCE	0.261093575
TAX	-0.01440119
PTRATIO	-1.074305348
AVG_ROOM	4.125409152
LSTAT	-0.603486589

- *For every \$1000 of avg. price of houses, per capita crime rate by town increases by 0.0487.
- *For every \$1000 of avg. price of houses, proportion of houses built prior to 1940 increases by 0.03%.
- *For every \$1000 of avg. price of houses, proportion of non-retail business acres per town increases by 0.13%.
- *For every \$1000 of avg. price of houses, nitric oxides concentration decreases by 10 million.
- *For every \$1000 of avg. price of houses, distance from highway increases by 0.2610 miles.
- *For every \$1000 of avg. price of houses, full-value property-tax rate decreases by 0.0144.
- *For every \$1000 of avg. price of houses, pupil-teacher ratio by town decreases by 1.0743.
- *For every \$1000 of avg. price of houses, average number of rooms per house increases by 4.12540.
- *For every \$1000 of avg. price of houses, lower status(LSTAT) of the population decreases by 0.603%.

8) Pick out only the significant variables from the previous question. Make another instance of the Regression model using only the significant variables you just picked and answer the questions below:

a) Interpret the output of this model

	<i>Coefficients</i>	<i>P-value</i>
Intercept	29.42847349	1.84597E-09
AGE	0.03293496	0.012162875
INDUS	0.130710007	0.038761669
NOX	-10.27270508	0.008545718
DISTANCE	0.261506423	0.000132887
TAX	-0.014452345	0.000236072
PTRATIO	-1.071702473	7.08251E-15
AVG_ROOM	4.125468959	3.68969E-19
LSTAT	-0.605159282	5.41844E-27

Adjusted R Square = 0.68868

B) Compare the adjusted R-square value of this model with the model in the previous question, which model performs better according to the value of adjusted R-square?

R Square = 0.6886836818 (Qn.8)

R Square = 0.6882986468 (Qn.7)

Adjusted R square value for this model is slightly a good percentage of changes for analysis compared to the previous model.

C) Sort the values of the Coefficients in ascending order. What will happen to the average price if the value of NOX is more in a locality in this town?

<i>Coefficients</i>	
NOX	-10.2727
PTRATIO	-1.0717
LSTAT	-0.60516
TAX	-0.01445
AGE	0.032935
INDUS	0.13071
DISTANCE	0.261506
AVG_ROOM	4.125469
Intercept	29.42847

If the value of NOX is more in a locality in this town, the value of the average price will be **reduced**.

D) Write the regression equation from this model

$$\text{AVG_PRICE} = \text{Intercept} + (\text{coefficient of Age} * \text{value of Age}) + (\text{coefficient of Indus} * \text{value of Indus}) + (\text{coefficient of NOX} * \text{value of NOX}) + (\text{coefficient of Distance} * \text{value of Distance}) + (\text{coefficient of Tax} * \text{value of Tax}) + (\text{coefficient of PTRATIO} * \text{value of PTRATIO}) + (\text{coefficient of Avg_room} * \text{value of Avg_room}) + (\text{coefficient of LSTAT} * \text{value of LSTAT})$$