

# TERRA REAL ESTATE AGENCY

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

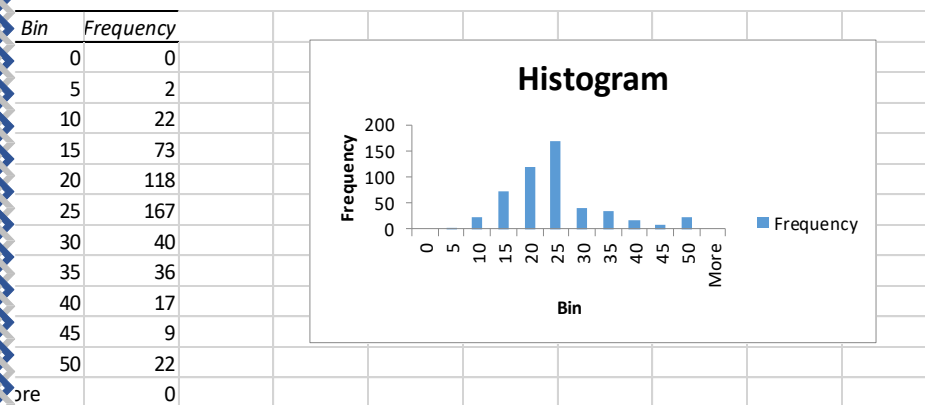
Mean	4.871976	Mean	68.5749	Mean	11.13678	Mean	0.554695	Mean	9.549407	Mean	408.2372	Mean	18.45553	Mean	6.284634	Mean	12.65306	Mean	22.53281
Standard Deviation	0.12986	Standard Deviation	1.25137	Standard Deviation	0.30498	Standard Deviation	0.005151	Standard Deviation	0.387085	Standard Deviation	7.492389	Standard Deviation	0.096244	Standard Deviation	0.031235	Standard Deviation	0.317459	Standard Deviation	0.408861
Median	4.82	Median	77.5	Median	9.69	Median	0.538	Median	5	Median	330	Median	19.05	Median	6.2085	Median	11.36	Median	21.2
Mode	3.43	Mode	100	Mode	18.1	Mode	0.538	Mode	24	Mode	666	Mode	20.2	Mode	5.713	Mode	8.05	Mode	50
Standard Error of the Mean	2.921132	Standard Error of the Mean	28.14886	Standard Error of the Mean	6.860353	Standard Error of the Mean	0.115878	Standard Error of the Mean	8.707259	Standard Error of the Mean	168.5371	Standard Error of the Mean	2.164946	Standard Error of the Mean	0.702617	Standard Error of the Mean	7.141062	Standard Error of the Mean	9.197104
Sample Variance	8.533012	Sample Variance	792.3584	Sample Variance	47.06444	Sample Variance	0.013428	Sample Variance	75.81637	Sample Variance	28404.76	Sample Variance	4.686989	Sample Variance	0.493671	Sample Variance	50.99476	Sample Variance	84.58672
Kurtosis	-1.18912	Kurtosis	-0.96772	Kurtosis	-1.23354	Kurtosis	-0.06467	Kurtosis	-0.86723	Kurtosis	-1.14241	Kurtosis	-0.28509	Kurtosis	1.8915	Kurtosis	0.49324	Kurtosis	1.495197
Skewness	0.021728	Skewness	-0.59896	Skewness	0.295022	Skewness	0.729308	Skewness	1.004815	Skewness	0.669956	Skewness	-0.80232	Skewness	0.403612	Skewness	0.90646	Skewness	1.108098
Range	9.95	Range	97.1	Range	27.28	Range	0.486	Range	23	Range	524	Range	9.4	Range	5.219	Range	36.24	Range	45
Minimum	0.04	Minimum	2.9	Minimum	0.46	Minimum	0.385	Minimum	1	Minimum	187	Minimum	12.6	Minimum	3.561	Minimum	1.73	Minimum	5
Maximum	9.99	Maximum	100	Maximum	27.74	Maximum	0.871	Maximum	24	Maximum	711	Maximum	22	Maximum	8.78	Maximum	37.97	Maximum	50
Sum	2465.22	Sum	34698.9	Sum	5635.21	Sum	280.6757	Sum	4832	Sum	206568	Sum	9338.5	Sum	3180.025	Sum	6402.45	Sum	11401.6
Count	506	Count	506	Count	506	Count	506	Count	506	Count	506	Count	506	Count	506	Count	506	Count	506

#### EXPLANATION:

We can observe that the mean of age is high, in price the mode is \$50. from the given data we can find the kurtosis values:

- ❖ CRIME\_RATE is -1.1891, so the curve is not sharp, it is FLAT CURVE & POSITIVE SKEWNESS.
- ❖ AGE is -0.9677, so the curve is not sharp, it is FLAT CURVE & NEGATIVE SKEWNESS.
- ❖ INDUS is -1.23353, so the curve is not sharp, it is FLAT CURVE & POSITIVE SKEWNESS.
- ❖ NOX is -0.0646, so the curve is not sharp, it is FLAT CURVE & POSITIVE SKEWNESS.
- ❖ DISTANCE is -0.86723, so the curve is not so sharp, it is FLAT CURVE & POSITIVE SKEWNESS.
- ❖ TAX is -1.1424, so the curve is not so sharp, it is FLAT CURVE & POSITIVE SKEWNESS.
- ❖ PTRATIO is -0.2850, so the curve is not sharp, it is FLAT CURVE & NEGATIVE SKEWNESS.
- ❖ AVG\_ROOM is 1.8915, so the curve is not sharp, it is FLAT CURVE & POSITIVE SKEWNESS.
- ❖ LSTAT is 0.4932, so the curve is not sharp, it is FLAT CURVE & POSITIVE SKEWNESS.
- ❖ AVG\_PRICE is 1.49519, so the curve is not sharp, it is FLAT CURVE & POSITIVE SKEWNESS.

Plot a histogram of the Avg\_Price variable. What do you infer?



**EXPLANATION:**

From the Histogram, it is inferred that Average Price has a **positive skewers**.

Compute the covariance matrix. Share your observations.

	CRIME_RATE	AGE	INDUS	NOX	DISTANCE	TAX	PTRATIO	VG_ROOM	LSTAT	AVG_PRICE
CRIME_RATE	8.516148									
AGE	0.562915	790.7925								
INDUS	-0.11022	124.2678	46.97143							
NOX	0.000625	2.381212	0.605874	0.013401						
DISTANCE	-0.22986	111.55	35.47971	0.61571	75.66653					
TAX	-8.22932	2397.942	831.7133	13.0205	1333.117	28348.62				
PTRATIO	0.068169	15.90543	5.680855	0.047304	8.743402	167.8208	4.677726			
VG_ROOM	0.056118	-4.74254	-1.88423	-0.02455	-1.28128	-34.5151	-0.53969	0.492695		
LSTAT	-0.88268	120.8384	29.52181	0.48798	30.32539	653.4206	5.7713	-3.07365	50.89398	
AVG_PRICE	1.162012	-97.3962	-30.4605	-0.45451	-30.5008	-724.82	-10.0907	4.484566	-48.3518	84.41956

**EXPLANATION:**

**negative value** Denotes, both the X and Y values are mostly on opposite sides of their average.

create a correlation matrix of all the variables (Use Data analysis tool pack).

	CRIME_RATE	AGE	INDUS	NOX	DISTANCE	TAX	PTRATIO	AVG_ROOM	LSTAT	AVG_PRICE
CRIME_RATE	1									
AGE	0.006859	1								
INDUS	-0.00551	0.644779	1							
NOX	0.001851	0.73147	0.763651	1						
DISTANCE	-0.00906	0.456022	0.595129	0.611441	1					
TAX	-0.01675	0.506456	0.72076	0.668023	0.910228	1				
PTRATIO	0.010801	0.261515	0.383248	0.188933	0.454741	0.460853	1			
AVG_ROOM	0.027396	-0.24026	-0.39168	-0.30219	-0.20985	-0.29205	-0.3555	1		
LSTAT	-0.0424	0.602339	0.6038	0.590879	0.488676	0.543993	0.374044	-0.61381	1	
AVG_PRICE	0.043338	-0.37695	-0.48373	-0.42732	-0.38163	-0.46854	-0.50779	0.59536	-0.73766	1

Which are the top 3 positively correlated pairs

3 positively correlated pairs

< vs DISTANCE **0.910228**

< vs INDUS **0.763651**

X vs AGE **0.73147**

Which are the top 3 negatively correlated pairs.

3 negatively correlated pairs

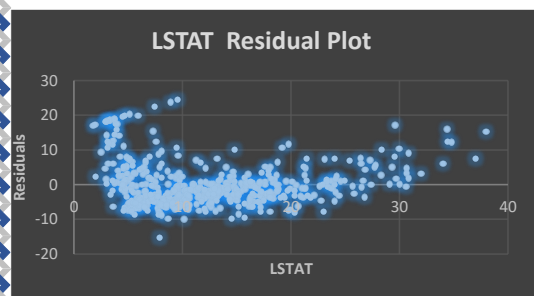
G\_PRICE vs LSTAT **-0.73766**

AT vs AVG\_ROOM **-0.61381**

G\_PRICE vs PTRATIO **-0.50779**

Build an initial regression model with AVG\_PRICE as 'y' (Dependent variable) and LSTAT variable as Independent variable. Generate the residual plot.

SUMMARY OUTPUT				
Regression Statistics				
Multiple R	0.737663			
R Square	0.544146			
Adjusted R Square	0.543242			
Standard Error	6.21576			
Observations	506			
ANOVA				
	df	SS	MS	F
Regression	1	23243.91	23243.91	601.6179
Residual	504	19472.38	38.63568	
Total	505	42716.3		
	Coefficient	Standard Error	t Stat	P-value
Intercept	34.55384	0.562627	61.41515	3.7E-236
LSTAT	-0.95005	0.038733	-24.5279	5.08E-88



R Square 0.544146

Coefficient of LSTAT -0.95005

Intercept 34.55384

#### R square

R Square is just above 0.5 so this value is not significant. Square has to be near to 1. **Coefficient of LSTAT**

Coefficient of LSTAT is -0.95005. It is inferred that for each \$1000 increase in Average price, there will be a 0.95% decrease in population.

#### Intercept

It is inferred that the Intercept value is 34.5538.

#### Residual plot

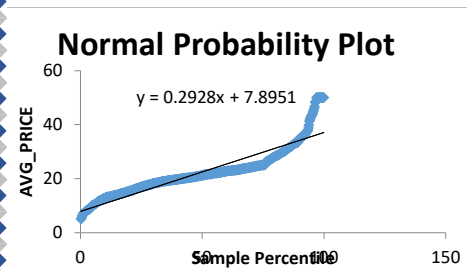
It is inferred that all the values are equally distributed

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

The p-value for LSTAT variable is 5.08110339438E-88. It is less than 0.05. So it is inferred that LSTAT variable is **significant** for the analysis.

Build a new Regression model including LSTAT and AVG\_ROOM together as Independent variables and AVG\_PRICE as dependent variable.

SUMMARY OUTPUT				
<b>Regression Statistics</b>				
Multiple R	0.7991			
Adjusted R Square	0.637124			
Standard Error	5.540257			
Observations	506			
<b>ANOVA</b>				
	<i>df</i>	<i>SS</i>	<i>MS</i>	<i>F</i>
Regression	2	27276.99	13638.49	444.3309
Residual	503	15439.31	30.69445	
Total	505	42716.3		
<b>Coefficients</b>				
	<i>Standard Error</i>	<i>t Stat</i>	<i>P-value</i>	
Intercept	-1.35827	3.172828	-0.4281	0.668765
AVG_ROOM	5.094788	0.444466	11.46273	3.47E-27
LSTAT	-0.64236	0.043731	-14.6887	6.67E-41



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?

b)  $AVG\_PRICE = \text{Intercept} + (\text{Coefficient of } AVG\_ROOM * \text{value of } AVG\_ROOM) +$

c)  $(\text{Coefficient of } LSTAT * \text{value of } LSTAT)$

d)  $AVG\_PRICE = -1.35827281187456 + (5.09478798433655 * 7) + (-0.642358334244129 * 20)$

e)  $AVG\_PRICE = 21.4581$

f) It is inferred that the Average price is **\$21.4581**. But the company quoting a value of 30000 USD for this locality. By the result, it is concluded that the company is **overcharging**.

g) R Square = **0.637124475470123** (Qn. 6)

h) R Square = **0.543241825954707** (Qn. 5)

c) 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.

It is inferred that the value of R Square is close to **1**, if the count of independent variable increases.

Based on the analysis, the **performance** of this model is **better** than the previous model.(Qn. 5)

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.

SUMMARY OUTPUT				
<b>Regression Statistics</b>				
Multiple R	0.83298			
Adjusted R Square				
Standard Error	5.13476			
Observations	506			
<b>ANOVA</b>				
	<i>df</i>	<i>SS</i>	<i>MS</i>	<i>F</i>
Regression	9	29638.9	3293.21	124.905
Residual	496	13077.4	26.3658	
Total	505	42716.3		
<b>Coefficients</b>				
	<i>Coefficient</i>	<i>Standard Error</i>	<i>t Stat</i>	<i>P-value</i>
Intercept	29.2413	4.81713	6.07028	2.5E-09
CRIME_RATE	0.04873	0.07842	0.62135	0.53466
AGE	0.03277	0.0131	2.502	0.01267
INDUS	0.13055	0.06312	2.06839	0.03912
POX	-10.3212	3.89404	-2.65051	0.00829
DISTANCE	0.26109	0.06795	3.8426	0.00014
MAX	-0.0144	0.00391	-3.68774	0.00025
TRATIO	-1.07431	0.1336	-8.0411	6.6E-15
AVG_ROOMS	4.12541	0.44276	9.3175	3.9E-19
STAT	-0.60349	0.05308	-11.3691	8.9E-27



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:

Interpret the output of this model.

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? 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? d) Write the regression equation from this model.

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?

Write the regression equation from this model.

#### Regression Statistics

Multiple R	0.8328
Adjusted R Square	0.6936
Standard Error	0.6887
Observations	506

#### ANOVA

	df	SS	MS	F
Regression	8	29629	3703.6	140.64
Residual	497	13088	26.333	
Total	505	42716		

	Coefficient	Standard Error	t Stat	P-value
Intercept	29.428	4.8047	6.1249	2E-09
AGE	0.0329	0.0131	2.5166	0.0122
INDUS	0.1307	0.0631	2.0722	0.0388
NOX	-10.27	3.8908	-2.64	0.0085
DISTANCE	0.2615	0.0679	3.8512	0.0001
TAX	-0.014	0.0039	-3.704	0.0002
PTRATIO	-1.072	0.1335	-8.031	7E-15
AVG_ROOM	4.1255	0.4425	9.3234	4E-19
LSTAT	-0.605	0.053	-11.42	5E-27

Adjusted R Square = 0.68868

Adjusted R Square = 0.6886836818 (Qn.8)

Adjusted R Square = 0.6882986468 (Qn.7)

the result, Adjusted R square for this model is greater comparing to the previous model.  
it is concluded that this model performs better than previous model.

#### Coefficients

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

is inferred that if the value of NOX is more in a locality in this town, the value of the average price will be reduced.

$G\_PRICE = \text{Intercept} + (\text{coefficient of Age} * \text{value of Age}) + (\text{coefficient of Indus} * \text{value of Indus}) + (\text{coefficient of NOX} * \text{value of NOX}) + (\text{coefficient of Distance} * \text{value of Distance}) + (\text{coefficient of Tax} * \text{value of Tax}) + (\text{coefficient of PTRATIO} * \text{value of PTRATIO}) + (\text{coefficient of Avg\_room} * \text{value of Avg\_room}) + (\text{coefficient of LSTAT} * \text{value of LSTAT})$

Commented [E1]:

