



Q.3. A regression through the origin is to be fit, relating Y to X for n=5 observations. The X levels are (1,2,3,4,10).

p.3.a. Obtain  $X'X$ ,  $(X'X)^{-1}$ , the P matrix, and the diagonal elements of P ( $v_{ii}$ ).

p.3.b. What do the  $v_{ii}$  elements sum to? Which, if any, elements are potentially influential?

p.3.c. Below is the fit for all cases (top), and with observation 5 dropped (bottom)  $b_1 = \hat{\beta}_1 = \frac{\sum_{i=1}^n X_i Y_i}{\sum_{i=1}^n X_i^2}$  Note:  $\sum_{i=1}^n e_i^2 \neq 0$

X	Y	Y-hat	e		b1
1	2	2.75	-0.75		2.75
2	5	5.51	-0.51		SS(Res)
3	6	8.26	-2.26		
4	7	11.02	-4.02		
10	30	27.54	2.46		
<b>Observation 5 dropped</b>					
X	Y	Y-hat	e		b1
1	2	1.93	0.07		1.93
2	5	3.87	1.13		SS(Res)
3	6	5.80	0.20		
4	7	7.73	-0.73		

All Cases:

SS(Res) =

s =

Observation 5 dropped:

p.3.d. Compute  $\hat{Y}_{5(5)}$ , and using  $s_{(5)}$  as estimate of  $\sigma$ :  $DFFITs_5$ , and  $DFBETAS_{1(5)}$

Q.4. A simple linear regression model is fit, with n = 12 observations (3 each at 4 levels of X). The residual sum of squares from the Regression model is SSResidual = 9042. The 4 fitted values at the distinct X levels are: (40, 80, 120, and 160). The 4 sample means at the distinct X levels are: (60, 70, 80, and 190).

p.4.a. Complete the following ANOVA table (degrees of freedom and sums of squares)

<b>ANOVA</b>		
<b>Source</b>	<b>df</b>	<b>SS</b>
<b>Regression</b>		
<b>Residual</b>		
<b>Lack of Fit</b>		
<b>Pure Error</b>		
<b>Total Corrected</b>		

p.4.b. Conduct the F-test for Lack-of Fit

$$H_0 : E\{Y_{ij}\} = \mu_j = \beta_0 + \beta_1 X_j \quad j = 1, \dots, c; \quad i = 1, \dots, n_j \quad H_A : E\{Y_{ij}\} = \mu_j \neq \beta_0 + \beta_1 X_j$$

Test Statistic: \_\_\_\_\_ Rejection Region: \_\_\_\_\_

Do you conclude that the relationship between  $E\{Y\}$  and X is linear? **Yes** or **No**

Q.5. A simple linear regression model is fit, based on  $n=5$  individuals. The data and the projection matrix are given below:

X		Y	P				
1	2	8	0.344	0.303	-0.129	0.262	0.221
1	4	12	0.303	0.274	-0.035	0.244	0.215
1	25	24	-0.129	-0.035	0.953	0.059	0.153
1	6	14	0.262	0.244	0.059	0.226	0.209
1	8	18	0.221	0.215	0.153	0.209	0.203

p.5.a. Give the leverage values for each observation. Do any exceed twice the average of the leverage values?

Observation 1 \_\_\_\_\_ Obs 2 \_\_\_\_\_ Obs 3 \_\_\_\_\_ Obs 4 \_\_\_\_\_ Obs 5 \_\_\_\_\_

p.5.b. Give  $\beta$ , based on the following results

X'X		X'Y	INV(X'X)	
5	45	76	0.438	-0.026
45	745	892	-0.026	0.003

p.5.c. Compute SSE and  $S_e$  (Note:  $Y'Y = 1304$   $Y'PY = 1282.45$ )

p.5.d. The following table contains the fitted values with and without each observation, residual standard deviation when that observation was not included in the regression, and the regression coefficients when that observation was not included in the regression.

Y-hat	Y-hat(-i)	S <sub>i</sub>	beta0 <sub>i</sub>	beta1 <sub>i</sub>
10.92	12.45	2.07	11.41	0.52
12.14	12.19	3.28	9.76	0.61
24.99	45.00	0.63	5.00	1.60
13.36	9.46	3.24	9.46	0.62
14.59	13.72	1.86	8.72	0.62

p.5.d.i. Compute DFFITS for the fifth observation

p.5.d.ii. Compute DFBETAS0 for the first observation

p.5.d.iii. Compute DFBETAS1 for the third observation

Q.6. A linear regression model was fit, relating energy consumption (Y) to population (X) for n = 183 nations. The following results are for models with and without the United States. The leverage value for the US is  $P_{ii} = .0275$ , and the Population is  $X = 310.0$

<b>INV(X'X)</b>				
<b>0.00590971</b>	<b>-0.00001109</b>			
<b>-0.00001109</b>	<b>0.00000030</b>			
	<b>With US</b>	<b>With US</b>	<b>Without US</b>	<b>Without US</b>
	<b>Coefficients</b>	<b>Standard Error</b>	<b>Coefficients</b>	<b>Standard Error</b>
<b>Intercept</b>	<b>0.6257</b>	<b>0.5912</b>	<b>0.4231</b>	<b>0.3687</b>
<b>popMill</b>	<b>0.0571</b>	<b>0.0042</b>	<b>0.0505</b>	<b>0.0026</b>
<b>MSResidual</b>	<b>59.14</b>		<b>22.98</b>	

p.6.a. Compute the fitted values for the US based on each model.

With US \_\_\_\_\_ Without US \_\_\_\_\_

p.6.b. Compute DFFITS for the US

p.6.c. Compute each of the DFBETAS for the US.

Intercept \_\_\_\_\_ Population \_\_\_\_\_

p.6.d. What is the average leverage value among the 183 nations?

Q.7. An experiment was conducted to measure the subsoil pressure of a steel ground roller. **There were 3 replicates at each of 4 depths (X=5, 10, 15, 20 cm).** The response was measured force (100s of Newtons).

The fitted regression equation is  $\hat{Y} = 49.371 - 2.036X$

We want to test  $H_0: E\{Y_j\} = \beta_0 + \beta_1X_j$   $H_A: E\{Y_j\} = \mu_j \neq \beta_0 + \beta_1X_j$

j	X <sub>j</sub>	Ybar <sub>j</sub>	SD <sub>j</sub>	Yhat <sub>j</sub>	Pure Error	Lack of Fit
1	5	40.38	4.32			
2	10	28.87	6.83			
3	15	16.23	3.76			
4	20	11.15	4.48			

p.7.a. Compute the Pure Error Sum of Squares and Degrees of Freedom. Hint: What is  $SD_j$  equal to?

SSPE = \_\_\_\_\_  $df_{PE}$  = \_\_\_\_\_

p.7.b. Compute the Lack-of-Fit Sum of Squares and Degrees of Freedom.

SSLF = \_\_\_\_\_  $df_{LF}$  = \_\_\_\_\_

p.7.c. Conduct the F-test for Lack-of-Fit

Test Statistic: \_\_\_\_\_ Rejection Region: \_\_\_\_\_ Reject  $H_0$ ? **Yes / No**

Q.8. A simple linear regression model is fit, based on  $n=5$  individuals. The data and the projection matrix are given below:

X		Y	P				
1	0	4	0.402	0.330	-0.175	0.258	0.186
1	2	6	0.330	0.284	-0.041	0.237	0.191
1	16	32	-0.175	-0.041	0.897	0.093	0.227
1	4	7	0.258	0.237	0.093	0.216	0.196
1	6	10	0.186	0.191	0.227	0.196	0.201

p.8.a. Give the leverage values for each observation. Do any exceed twice the average of the leverage values?

Observation 1 \_\_\_\_\_ Obs 2 \_\_\_\_\_ Obs 3 \_\_\_\_\_ Obs 4 \_\_\_\_\_ Obs 5 \_\_\_\_\_

p.8.b. Give  $\beta$ , based on the following results

X'X		X'Y	INV(X'X)	
5	28	59	0.402	-0.036
28	312	612	-0.036	0.006

p.8.c. Compute SSE and  $S_e$ . (Note:  $Y'Y = 1225$   $Y'PY = 1207.144$ )

p.8.d. The following table contains the fitted values with and without each observation, residual standard deviation when that observation was not included in the regression, and the regression coefficients when that observation was not included in the regression.

Y-hat	Y-hat(-i)	S <sub>i</sub>	beta0 <sub>i</sub>	beta1 <sub>i</sub>
1.64	12.45	2.07	11.41	0.52
5.27	12.19	3.28	9.76	0.61
30.67	45.00	0.63	5.00	1.60
8.90	9.46	3.24	9.46	0.62
12.53	13.72	1.86	8.72	0.62

p.8.d.i. Compute DFFITS for the fifth observation

p.8.d.ii. Compute DFBETAS0 for the first observation

p.8.d.iii. Compute DFBETAS1 for the third observation

Q.9. For the F-test for Lack-of-Fit, where:

$$H_0 : E\{Y_{ij}\} = \beta_0 + \beta_1 X_j \quad H_A : E\{Y_{ij}\} = \mu_j \neq \beta_0 + \beta_1 X_j \quad j=1, \dots, c; \quad i=1, \dots, n_j$$

with: Residual  $\equiv Y_{ij} - \hat{Y}_j$  Pure Error  $\equiv Y_{ij} - \bar{Y}_j$  Lack of Fit  $\equiv \bar{Y}_j - \hat{Y}_j$   $\hat{Y}_j = \hat{\beta}_0 + \hat{\beta}_1 X_j$   $\bar{Y}_j = \frac{1}{n_j} \sum_{i=1}^{n_j} Y_{ij}$

Show: 
$$\sum_{j=1}^c \sum_{i=1}^{n_j} (Y_{ij} - \hat{Y}_j)^2 = \sum_{j=1}^c \sum_{i=1}^{n_j} (Y_{ij} - \bar{Y}_j)^2 + \sum_{j=1}^c \sum_{i=1}^{n_j} (\bar{Y}_j - \hat{Y}_j)^2$$

Q.10. A simple linear regression model was fit, relating Weight (Y, in pounds) to Height (X, in inches) among a sample of n=9 Women's NBA basketball players. There were 3 distinct height levels, with 3 players per height. Use the table below to conduct the F-test for Lack of Fit.  $H_0 : E\{Y_{ij}\} = \beta_0 + \beta_1 X_j$  vs  $H_A : E\{Y_{ij}\} = \mu_j \neq \beta_0 + \beta_1 X_j$

Height	Weight	Ybar(Grp)	Yhat(Grp)	PE	LF		
68	130						
68	165						
68	155						
72	180						
72	165						
72	150						
76	200						
76	185						
76	230						

Coefficients	
Intercept	-321.667
Height	6.875

Q.11. A random sample of n = 20 cricket players from India was obtained, and a simple linear regression was fit relating the number of runs (Y, in 1000s) and number of completed innings (X, in 100s) in international matches. Player 20 is Sachin Tendulkar who played in 411 completed innings (X = 4.11) and accounted for 18426 runs (Y = 18.426). Regression models were fit with and without Sachin, and results are given below.

X'X		X'Y		X'X_(20)		X'Y_(20)	
20.0000	22.6600	77.1320		19.0000	18.5500	58.7060	
22.6600	44.7782	172.7793		18.5500	27.8861	97.0484	
INV(X'X)		Beta-hat		INV(X'X_(20))		Beta-hat_(20)	
0.1172	-0.0593	-1.2074		0.1501	-0.0999	-0.8785	
-0.0593	0.0523	4.4696		-0.0999	0.1023	4.0646	
Y'Y	Y'PY			Y'Y_(20)	Y'PY_(20)		
693.2798	679.1187			353.7623	342.8852		

p.11.a. Obtain SSE and MSE and the estimated error standard deviation for each model.

SSE<sub>1</sub> = \_\_\_\_\_ MSE<sub>1</sub> = \_\_\_\_\_ s = \_\_\_\_\_

SSE<sub>2</sub> = \_\_\_\_\_ MSE<sub>2</sub> = \_\_\_\_\_ s<sub>(20)</sub> = \_\_\_\_\_

p.11.b. Obtain the fitted value and residual for Satchin for each model.

Fitted<sub>1</sub> = \_\_\_\_\_ Residual<sub>1</sub> = \_\_\_\_\_ Fitted<sub>2</sub> = \_\_\_\_\_ Residual<sub>2</sub> = \_\_\_\_\_

p.11.c. Obtain the leverage value for Satchin from the full data model, and his studentized residual.

v<sub>20,20</sub> = \_\_\_\_\_ r<sub>20</sub><sup>\*</sup> = \_\_\_\_\_

p.11.d. Compute DFFITS<sub>20</sub> and DFBETAS<sub>1(20)</sub> for Satchin. Note that both DFFITS and DFBETAS make use of s<sub>(20)</sub> to estimate σ .

DFFITS<sub>20</sub> = \_\_\_\_\_ DFBETAS<sub>1(20)</sub> = \_\_\_\_\_

p.11.e. What is the average of the leverage values for the full data model?

Q.12. An experiment was conducted relating springiness in berries (Y, in mm) to sugar equivalent (X, in g/L). There were c = 4 distinct sugar equivalent groups, with n<sub>j</sub> = 5 berries per group. The lack-of-fit test for a linear relation is:

$$H_0 : E\{Y_{ij}\} = \mu_j = \beta_0 + \beta_1 X_j \quad i = 1, \dots, 5; j = 1, \dots, 4 \quad H_A : E\{Y_{ij}\} = \mu_j \neq \beta_0 + \beta_1 X_j$$

ANOVA																										
	df	SS	MS	F	Significance F	j	X <sub>j</sub>	n <sub>j</sub>	Yhat <sub>j</sub>	Ybar <sub>j</sub>	s <sup>2</sup> <sub>j</sub>															
Regression	1	245.74	245.74	57.40	0.0000	1	176.5	5	21.89	21.77	0.38															
Residual	18	77.06	4.28			2	209.3	5	18.04	18.77	2.90															
Total	19	322.80				3	225	5	16.20	15.41	11.30															
						4	259.5	5		12.32	3.18															
<table border="1"> <thead> <tr> <th></th> <th>Coefficients</th> <th>Standard Error</th> <th>t Stat</th> <th>P-value</th> </tr> </thead> <tbody> <tr> <td>Intercept</td> <td>42.6025</td> <td>3.4019</td> <td>12.5233</td> <td>0.0000</td> </tr> <tr> <td>sugCont</td> <td>-0.1174</td> <td>0.0155</td> <td>-7.5763</td> <td>0.0000</td> </tr> </tbody> </table>													Coefficients	Standard Error	t Stat	P-value	Intercept	42.6025	3.4019	12.5233	0.0000	sugCont	-0.1174	0.0155	-7.5763	0.0000
	Coefficients	Standard Error	t Stat	P-value																						
Intercept	42.6025	3.4019	12.5233	0.0000																						
sugCont	-0.1174	0.0155	-7.5763	0.0000																						

Note:  $s_j^2 = \frac{\sum_{i=1}^{n_j} (Y_{ij} - \bar{Y}_j)^2}{n_j - 1}$   $n_j > 1$ , 0 otherwise

p.12.a. Give the fitted value for the linear regression for the 4<sup>th</sup> group (X<sub>4</sub> = 259.5).

p.12.b. Compute the Pure Error Sum of Squares, degrees of freedom and Mean Square.

SS<sub>PE</sub> = \_\_\_\_\_ df<sub>PE</sub> = \_\_\_\_\_ MS<sub>PE</sub> = \_\_\_\_\_

p.12.c. Compute the Lack-of-Fit Sum of Squares, degrees of freedom and Mean Square.

SS<sub>LF</sub> = \_\_\_\_\_ df<sub>LF</sub> = \_\_\_\_\_ MS<sub>LF</sub> = \_\_\_\_\_

p.12.d. Give the Test Statistic, Rejection Region, and P-value relative to .05 for the Lack-of-Fit test.

Test Statistic: \_\_\_\_\_ Rejection Region: \_\_\_\_\_ P-value > or < .05

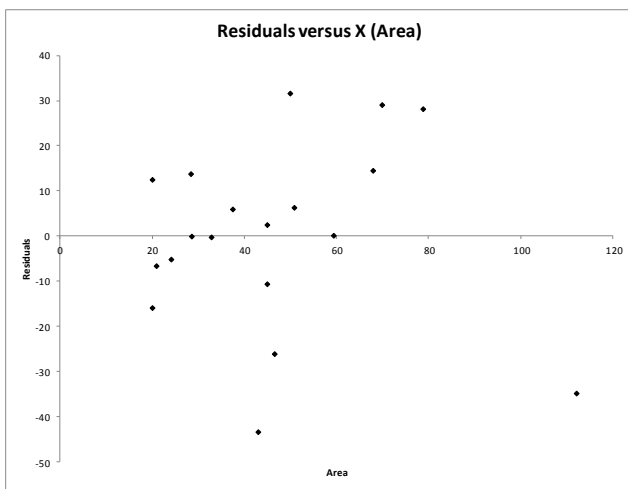


Q.13. A regression model is fit, relating energy consumption (Y) to total area (X) for a sample of  $n = 19$  luxury hotels in Hainan Province, China. The Analysis of Variance for the simple linear regression model is given below.

ANOVA					
	<i>df</i>	<i>SS</i>	<i>MS</i>	<i>F</i>	<i>gnificance</i>
Regressio	1	25521.76	25521.76	57.75	0.0000
Residual	17	7512.94	441.94		
Total	18	33034.70			

p.13.a. A plot of the residuals versus area is given below. It demonstrates which possible violations of assumptions (circle all that apply).

Non-normal Errors      Unequal Variance      Serial Correlation of Errors      Non-linear Relation between Y and X



p.13.b. A second regression model is fit, relating the squared residuals (Y) to area (X). Conduct the Breusch-Pagan test to test whether the equal variance assumption is reasonable. The sums of squares are given below.

ANOVA		
	<i>df</i>	<i>SS</i>
Regressio	1	1239379
Residual	17	3871645
Total	18	5111024

Test Statistic: \_\_\_\_\_ Rejection Region: \_\_\_\_\_ P-value > or < .05

Q.14. A random sample of  $n = 15$  Bollywood movies was obtained, and a simple linear regression was fit relating the log Revenues (Y) the log Budget (X) in international matches. Movie 15 was **Sultan** with log Revenues  $Y_{15} = 5.705$  and log Budget  $X_{15} = 4.500$ . Regression models were fit with and without **Sultan**, and results are given below.

X'X			X'Y		X'X_(15)			X'Y_(15)	
15.000	55.474		53.796		14.000	50.974		48.091	
55.474	214.736		213.087		50.974	194.488		187.417	
INV(X'X)			Beta-hat		INV((X'X)_(15))			Beta-hat_(15)	
1.495	-0.386		-1.871		1.563	-0.410		-1.609	
-0.386	0.104		1.476		-0.410	0.112		1.385	
Y'Y	Y'PY				Y'Y_15	Y'PY_15			
218.04	213.79				185.50	182.26			

p.14.a. Obtain SSE and MSE and the estimated error standard deviation for each model.

$SSE_1 =$  \_\_\_\_\_  $MSE_1 =$  \_\_\_\_\_  $s =$  \_\_\_\_\_

$SSE_2 =$  \_\_\_\_\_  $MSE_2 =$  \_\_\_\_\_  $s_{(15)} =$  \_\_\_\_\_

p.14.b. Obtain the fitted value for **Sultan** for each model and its leverage value based on the full data set.

Fitted<sub>1</sub> = \_\_\_\_\_ Fitted<sub>2</sub> = \_\_\_\_\_ Leverage = \_\_\_\_\_

p.8.c. . Compute  $DFFITS_{15}$  for **Sultan**. Note that DFFITS makes use of  $s_{(15)}$  to estimate  $\sigma$  .

$DFFITS_{15} =$  \_\_\_\_\_