## Statistics 512: Homework#3 Solutions

1. Consider the following data set that describes the relationship between the rate of an enzymatic reaction (V) and the substrate concentration (C). A common model used to describe the relationship between rate and concentration is the Michaelis-Menten model  $V = \frac{\theta_1 C}{\theta_2 + C}$ , where  $\theta_1$  is the maximum rate of the reaction and  $\theta_2$  describes how quickly the reaction will reach its maximum rate. With this mode,  $\frac{1}{V}$  can be written as a linear model with explanatory variable  $\frac{1}{C}$ :

$$\frac{1}{V} = \frac{1}{\theta_1} + \frac{\theta_2}{\theta_1} \frac{1}{C}$$

(a) Generate a scatterplot of V vs C. Comment on the shape.

**Solution:** The relationship between V and C does not appear to be linear. V increases rapidly for small values of C but tends to level off as C approaches 1. See Figure 1.



Figure 1: Scatterplot for Problem 1a

(b) Define new variables for  $\frac{1}{V}$  and  $\frac{1}{C}$  in SAS, and generate a scatterplot of the new variables. Does the fit appear linear? Do any assumptions appear to be violated?

**Solution:** The relationship between 1/V and 1/C appears to be linear. However, the assumption of constant variance seems to be violated, since the observations with 1/C = 50 have much larger residuals than do the other observations. See Figure 2.

(c) How is the distribution of  $\frac{1}{C}$  different from the distribution of C? Are there any points that may be more influential in determining the fit?



Figure 2: Scatterplot for Problem 1b

**Solution:** The variable C has mean 0.345, median 0.0165, and standard deviation 0.4, with values ranging from 0.02 and 1.1. The variable 1/C has mean 13.83, median 6.82, and standard deviation 17.8, with values ranging from 0.91 to 50. Neither variable seems to come from a normal distribution.

	Th	e UNIVARIA	TE Proced	ure	
		Varia	ble: c		
	Ba	sic Statis	tical Mea	sures	
Loca	tion		Vari	ability	
Mean	0.345000	Std D	eviation	0.39862	
Median	0.165000	Varia	nce		0.15890
Mode	0.020000	Range			1.08000
		Inter	quartile	Range	0.50000
		Extreme O	bservatio	ns	
Lowes		st	H	ighest	
	Value	Obs	Value	Obs	
	0.02	2	0.22	8	
	0.02	1	0.56	9	
	0.06	4	0.56	10	

## The UNIVARIATE Procedure Variable: cinv Basic Statistical Measures

1.10

1.10

11

12

3

6

0.06

0.11

	Location	Variabilit	2y
Mean	13.83297	Std Deviation	17.77123



		I a a a g	
	Extreme	Observations	
Lowest		Highest	
Value	Obs	Value	Obs
0.909091	12	9.09091	6
0.909091	11	16.66667	3
1.785714	10	16.66667	4
1.785714	9	50.00000	1

(d) Determine the least squares regression line for  $\frac{1}{V}$  vs  $\frac{1}{C}$ . Save the residuals and predicted values. Does the residual plot suggest any problems?

8

4.545455

**Solution:** The least squares line is 1/V = 0.00511 + 0.000247(1/C). The residual plot (Figure 3) suggests that the variance is not constant (heterscedasticity).

50.00000

2

## Parameter Estimates

		Parameter	Standard		
Variable	DF	Estimate	Error	t Value	Pr >  t
Intercept	1	0.00511	0.00070400	7.25	<.0001
cinv	1	0.00024722	0.00003210	7.70	<.0001

For the next 3 questions, use the grade point average data described in the text with Problem 1.19 (CH01PR19.DAT).

2. Describe the distribution of the explanatory variable. Show the plots and output that were helpful in learning about this variable.

Solution: Using proc univariate, we see there are 120 observations ranging between 14 and 35 with a mean of 24.725 and median of 24; their standard deviation



Figure 4: Residual Plot for Problem 1d

is 4.472 . There are no extreme observations (i.e., ones far away from the others) as shown in the histogram (Figure 5). The distribution appears to be reasonably symmetric but not completely so; it has a slight skew to the right.

		The UNIVARI	ATE Procedu	ure	
		Variable:	testscore	е	
		Мо	ments		
N		120	Sum Weig	ghts	120
Mean		24.725	Sum Obse	ervations	2967
Std Deviatio	n	4.47206549	Variance	е	19.9993697
Skewness		-0.1363553	Kurtosi	3	-0.5596968
Uncorrected	SS	75739	Correcte	ed SS	2379.925
Coeff Variat	ion	18.0872214	Std Erro	or Mean	0.40824186
	E	Basic Statis	tical Meas	ires	
Loca	tion		Varia	ability	
Mean	24.725	500 Std	Deviation		4.47207
Median	25.000	000 Vari	ance		19.99937
Mode	24.000	000 Rang	е		21.00000
		Inte	rquartile l	Range	7.00000
		Extreme O	bservation	5	
	Lo	west	H:	ighest	
	Value	Obs	Value	Obs	
	14	2	32	84	
	15	48	32	104	
	16	119	33	15	
	16	52	34	80	

	16	32	35	106
Stem	Leaf		#	Boxplot
35	0		1	Ì
34	0		1	I
33	0		1	I
32	0000		4	Ι
31	0000		4	I
30	0000000		7	I
29	0000000		7	I
28	000000000		10	++
27	000000000		10	
26	000000000		10	
25	000000000		10	**
24	00000000000		12	+
23	00000		5	
22	0000		4	
21	00000000		9	++
20	000000000		10	I
19	000		3	I
18	0000000		7	I
17				I
16	000		3	I
15	0		1	I
14	0		1	I
	++	++		
Dataset of GPA	and Test Scores			Dataset of GPA and Test Scores
			35 -	



Figure 5: Graphs for Problem 2

3. Run the linear regression to predict GPA from the entrance test score, and obtain the residuals (do not include a list of the residuals in your solution).

(a) Verify that the sum of the residuals is zero by running **proc univariate** with the output from the regression.

**Solution:** The given **proc univariate** output shows that the residuals sum to zero.

Moments							
N	120	Sum Weights	120				
Mean	0	Sum Observations	0				
Std Deviation	0.62050134	Variance	0.38502191				
Skewness	-1.0067279	Kurtosis	2.50187662				
Uncorrected SS	45.8176078	Corrected SS	45.8176078				
Coeff Variation		Std Error Mean	0.05664376				

(b) Plot the residuals versus the explanatory variable and briefly describe the plot noting any unusual patterns or points.

**Solution:** There does not appear to be any obvious pattern or outlier in this residual plot (Figure 6). It looks like a random scatter of points, and the variance is reasonably constant.



Figure 6: Scatterplot for Problem 3b

(c) Plot the residuals versus the order in which the data appear in the data file and briefly describe the plot noting any unusual patterns or points.

Solution: There is no obvious pattern over time. See Figure 7.

(d) Examine the distribution of the residuals by getting a histogram and a normal probability plot of the residuals by using the histogram and qqplot statements in proc univariate. What do you conclude?

**Solution:** The residuals appear reasonably normal if somewhat asymmetric, since the histogram appears fairly normal, and the qqplot is fairly linear (Figure



Figure 7: Scatterplot for Problem 3c

8). There is some suggestion of a concave down shape to the qq-plot, but it is not too bad.



Figure 8: Graphs for Problem 3d

4. Change the data set by changing the value of the GPA for the last observation from 2.948 to 29.48 (e.g., a typo). You can do this in a data step. For example,

```
data a2;
set a1;
if _n_ eq 120 then gpa = 29.48;
```

an alternative is simply to edit the data file.

(a) Make a table comparing the results of this analysis with the results of the analysis of the original data. Include in the table the following: fitted equation, *t*-test for the slope, with standard error and *p*-value,  $R^2$ , and the estimate of  $\sigma^2$ . Summarize the differences.

**Solution:** The outlier has a huge impact on these results. The slope becomes nearly double its original value and is no longer significantly different from zero. The changes in  $R^2$  and  $s^2$  are also very extreme. The outlier greatly inflates the estimated variance and makes the  $R^2$  almost zero.

Result	Original	With Outlier					
Equation	Y = 2.114 + 0.0388X	Y = 1.432 + 0.0753X					
<i>t</i> -test	t = 3.04	t = 1.48					
SE	0.0133	0.0509					
<i>p</i> -value	0.0029	0.1414					
Conclusion	Reject H <sub>0</sub>	Do not reject $H_0$					
$R^2$	0.0726	0.0182					
$s^2$	0.388	6.163					

GPA and Test Scores Original:

			Analysis of	Variano	се				
			Sum	of	]	Mean			
Source		DF	Squar	es	Sq	uare	F۷	Value	Pr > F
Model		1	13.507	89	13.5	0789		2.19	0.1414
Error		118	727.289	65	6.1	6347			
Corrected Tot	al	119	740.797	54					
	Root MSE		2.482	263 R-	-Squar	e	0.018	82	
	Dependent	Mean	3.295	515 Ac	dj R-S	q	0.009	99	
	Coeff Var		75.342	206					
		Pa	rameter	Standa	ard				
Variable	DF	E	stimate	Eri	ror	t Va	lue	Pr >	t
Intercep	t 1		1.43243	1.278	850	1	.12	0.20	548
testscor	e 1		0.07534	0.050	089	1	.48	0.14	414

GPA and Test Scores with Outlier:

	A	Inalysis of Var	riance		
		Sum of	Mean		
Source	DF	Squares	Square	F Value	Pr > F
Model	1	13.50789	13.50789	2.19	0.1414
Error	118	727.28965	6.16347		
Corrected Total	119	740.79754			
Root M	SE	2.48263	R-Square	0.0182	
		0 00515			

Dependent Mean	3.29515	Adj R-Sq	0.0099
Coeff Var	75.34206		

		Parameter	: Estimates		
		Parameter	Standard		
Variable	DF	Estimate	Error	t Value	Pr >  t
Intercept	1	1.43243	1.27850	1.12	0.2648
testscore	1	0.07534	0.05089	1.48	0.1414

(b) Repeat parts (b), (c), and (d) from the previous problem and explain how these plots

help you to detect the unusual observation.

**Solution:** The outlier is quite noticeable on the residual plots as a point far away from all the others (Figure 9 left). The sequence plot (Figure 9 right) additionally shows that the outlier is that last observation in the list. The qqplot (Figure 10 right) shows that the data do not fit on the straight line with the estimated mean and standard deviation, and the outlier again appears quite separate from all the others. The histogram (Figure 10 left) shows a distribution that is far from normal, since there is an observation very far out in the tail. All together, these plots clearly show there is one outlying point which should be investigated. (Since GPA cannot be 29.48, this would appear to be a data entry error.)



Figure 9: Graphs for Problem 4b



Figure 10: Graphs for Problem 4b