

**Testing Assumptions in Multiple Regression:  
Comparison of Procedures Available  
in SAS and SPSSX**

**Paula L. Woehlke**

**Patricia B. Elmore**

**Southern Illinois University Carbondale**

**Debra L. Spearing**

**State University of New York at Albany**

That the use of multiple linear regression requires satisfying several assumptions has seldom been disputed. However, assessing whether one has met important assumptions is not always easy, and given the limited time available to instructors in a typical multiple regression course, the techniques available for checking assumptions are often not taught, or mentioned only briefly. The purpose of this paper is to compare the most easily available techniques for checking assumptions from two of the most popular statistics packages in use today, SAS (SAS Institute, 1985) and SPSS<sup>X</sup> (SPSS, Inc., 1985). It is hoped that the attached examples will make the multiple regression course instructor's job easier by providing concrete examples of computer input and output that illustrate the testing of assumptions.

A condition that should be met for the use of multiple regression, but which is not, strictly speaking, an assumption, is that there be an absence of multicollinearity. Multicollinearity is defined as the existence of substantial correlation among a set of independent variables, and its presence creates three distinct problems:

- the substantial interpretation of partial regression coefficients,
  - the sampling stability of these coefficients
- and
- computational accuracy of the regression analysis.

Thus, although absence of multicollinearity is not a regression assumption, failure to assure that predictor variables are not multicollinear can result in faulty interpretations of analyses, regression equations that cannot be used for prediction, or both.

In terms of actual theoretical assumptions for using multiple regression analyses, errors of the prediction or residuals from estimated values of the regression provide the basis for assessing the adequacy of the model (Cohen & Cohen, 1983). Specifically, it is assumed that errors

- (1) are normally distributed

(2) are independent of one another (that is, errors associated with one observation  $Y_i$  are not correlated with errors associated with any other observation

$Y_j$ )

(3) are identically distributed (that is, are sampled from the same distribution and have constant variances, also known as the assumption of homoscedasticity)

(4) have a mean of zero

and

(5) are uncorrelated with the independent variables ( $X$ 's).

In addition to these assumptions about errors, it is further assumed that

(6) the independent variables, ( $X$ 's) are fixed and measured without error

(7) the regression of  $Y$  on  $X$  is linear

and

(8)  $Y$  is a random variable composed of two components: a fixed component,  $a + bX$ , and a random error  $e_i$ .

Two conditions under which these assumptions about residuals fail to be met occur

when

- the regression of  $Y$  on  $X$  (or  $X$ 's) is curvilinear (so that condition 7 above is not met)

and

- there are one or more extreme residual values, known as "outliers, which not only make relatively large contributions to error or residual variance (thus reducing  $R^2$ ) but also exert a disproportionately strong pull on the regression.

To illustrate the use of SAS and SPSS<sup>X</sup> to test these assumptions, we used the

(in)famous Longley data set. This data set has multicollinearity and some cases of univariate outliers through which to illustrate the diagnostic procedures available in both SAS and

SPSS<sup>X</sup>. The following pages provide annotated output from these two packages, which we will describe in the next section.

### Description of Output

The first assumption about errors is that the residuals are normally distributed. This assumption can be assessed by examining the residual scatterplot in Figure 4.SAS and the normal probability plot and statistical analyses shown in Figure 6.SAS; similar plots and statistics are produced by SPSS<sup>X</sup>, as shown in Figure 4.SPSS<sup>X</sup>, Figure 5.SPSS<sup>X</sup>, and Figure 6.SPSS<sup>X</sup>. If residuals are normally distributed, the plus signs (+'s) and the asterisks (\*'s) will coincide in the SAS normal probability plot (or the asterisks [\*'s] and dots [o's] in the SPSS<sup>X</sup> normal probability plot). Also, a statistical test for normality is provided in SAS in Figure 6.SAS; in this case, W:NORMAL = 0.948682, p=.471. It should be noted that SPSS<sup>X</sup>'s CON- DESCRIPTIVE procedure routinely does not provide a comparable statistical test. All of these plots and tests from both SAS and SPSS<sup>X</sup> indicate that the assumption about normally distributed residuals has been met.

That residuals are independent of one another or errors associated with one observation are not correlated with errors associated with any other observation is the second assumption to be tested. The Durbin-Watson D statistic shown in Figure 3.SAS and Figure 3.SPSS<sup>X</sup> tests for nonindependence of errors when the order of cases is meaningful. For this data set, the Durbin-Watson D statistic is irrelevant. The residual scatterplots in Figure 4.SAS and Figure 5.SPSS<sup>X</sup> show that the residuals are independent.

The third assumption is that residuals are identically distributed. This means that the errors are sampled from the same distribution and have constant variance, also known as homoscedasticity. Examination of the residual scatterplots in Figure 4.SAS and Figure 5.SPSS<sup>X</sup> indicates that the assumption of homoscedasticity has been met.

Assumption 4, that the residuals have a mean of zero, can be determined by examining Figure 6.SAS or Figure 6.SPSS<sup>X</sup>. For this data set, the mean is  $-7.421E-10$  (Figure 6.SAS), which is considered zero for our purposes, or .000 (Figure 6.SPSS<sup>X</sup>).

The correlation matrix showing the correlations between all of the independent variables and the residual should be used to assess assumption 5, that the residuals are uncorrelated with the independent variables. Examination of the correlation matrix for this data set, as found in Figure 7.SAS or Figure 7.SPSS<sup>X</sup> indicates a correlation between each of the six independent variables and the residual equal to zero.

That the regression of Y on X is linear, assumption 7, can be determined by creating bivariate scatterplots for all predictors with the criterion. One example is shown in Figure 5.SAS and another in Figure 1.SPSS<sup>X</sup>; both show the relation between Y and X<sub>1</sub>. All six predictors in this data set are linearly related to the criterion Y.

Figure 1.SAS and Figure 2.SPSS<sup>X</sup> show a check for multicollinearity. Low tolerance value and high condition number with large variance proportion for two or more variables may indicate multicollinearity. Variables X<sub>5</sub> and X<sub>6</sub> in this data set may be multicollinear with previous terms in the model.

Figure 2.SAS has two indices to check for outliers. A studentized residual value in excess of  $\pm 3.00$  may indicate a univariate outlier (Tabachnick & Fidell, 1989, p. 67). Also, a data point with a Cook's Distance value greater than 1.00 is suspected of being an outlier. Cook's Distance is discussed in depth in Tabachnick & Fidell (1989), p. 130, and Kleinbaum, Kupper & Muller (1988), p. 201. Also note that in Figure 3.SPSS<sup>X</sup> a similar casewise plot appears, as well as a listing and a histogram of standardized residuals.

## Discussion

Although the output of the regression modules and related descriptive statistics procedures for SAS and SPSS<sup>X</sup> are quite similar, there are a few differences worth noting. First, SPSS<sup>X</sup> includes a histogram of standardized residuals to make the spotting of outliers some-

what easier; the program also has a normal probability plot that is a little easier to read than that provided in SAS. SPSS<sup>X</sup> also provides standard errors for the skewness and kurtosis values for the variables analyzed in the CONDESCRIPTIVE module; these values are not printed in the SAS output. On the other hand, SAS provides a statistical test of normality when requested through PROC UNIVARIATE, as well as stem and leaf diagrams and boxplots of distributions, through the same PROC. It is also easy to obtain Cook's D values through SAS's PROC REG; it is somewhat more difficult to get similar statistics from SPSS<sup>X</sup>, requiring the use of a RESIDUALS subcommand. In most other respects, output is comparable for the data and regression analyses shown here. For more advanced regression applications, it is somewhat easier to obtain leverage (partial regression residual) plots for general linear hypotheses, used in assessing degree of fit, nonfitting points, and multicollinearity (Sall, 1990) from SAS (via an option in PROC REG) than from SPSS<sup>X</sup>, which produces "partial regression plots" through a PARTIALPLOT subcommand. It should be noted, however, that some anomalies recently have been detected in SAS's regression and GLM procedures for models using different types of intercept terms (see Uyar & Erdem, 1990). Finally, although it is somewhat more difficult to obtain several diagnostic statistics from SPSS<sup>X</sup>, the package supplements its regression module with an extensive and flexible MANOVA procedure that allows one to easily build advanced regression models. With these advantages and disadvantages in mind, it should be possible for the reader to choose which computer package is most appropriate for a particular regression analysis.

## References

- Cohen, J. & Cohen, P. (1983). *Applied Multiple Regression/Correlation Analysis for the Behavioral Sciences*, 2nd ed. Hillsdale, NJ: Erlbaum.
- Draper, N. R. & Smith, H. (1981). *Applied Regression Analysis*, 2nd ed. New York: Wiley.
- Kleinbaum, D. G., Kupper, L. L. & Muller, K. E. (1988). *Applied Regression Analysis and Other Multivariable Methods*, 2nd ed. Boston: PWS—Kent.
- Pedhazur, E. J. (1982). *Multiple Regression in Behavioral Research*, 2nd ed. New York: Holt, Rinehart & Winston.
- SAS Institute, Inc. (1985). *SAS User's Guide: Statistics, Version 5 Edition*. Cary, NC: Sas Institute, Inc.
- Sall, J. (1990). Leverage plots for general linear hypotheses. *The American Statistician*, 44(4), 308-315.
- SPSS, Inc. (1986). *SPSS X User's Guide, Edition 2*. New York: McGraw-Hill.
- Tabachnick, B. G. & Fidell, L. S. (1989). *Using Multivariate Statistics*, 2nd ed. New York: Harper & Row.
- Uyar, B. & Erdem, O. (1990). Regression procedures in SAS: Problems? *The American Statistician*, 44(4), 296-301.

SOURCE	DF	SUM OF SQUARES	MEAN SQUARE	F VALUE	PROB>F
MODEL	6	104172402	30695000.32	330.285	0.0001
ERROR	9	834424.06	92736.00616		
C TOTAL	15	185000826			

ROOT MSE	R-SQUARE	ADJ R-SQ	T FOR NO: PARAMETER=0	PROB >  T
DEP MEAN	304.8541			
C.V.	0.65317			
	0.4667301			

PARAMETER ESTIMATES

VARIABLE	DF	PARAMETER ESTIMATE	STANDARD ERROR	T FOR NO: PARAMETER=0	PROB >  T
INTERCEP	1	-3482250.63	89420.30	-3.911	0.0036
X1	1	15.06187236	04.91492570	0.177	0.8631
X2	1	-0.03501910	0.03349101	-1.070	0.3127
X3	1	-2.92022900	0.40039960	-4.136	0.0025
X4	1	-1.03322607	0.21427416	-4.822	0.0009
X5	1	-0.05110411	0.22607320	-0.226	0.8262
X6	1	1829.15146	455.47850	4.016	0.0030

STANDARDIZED ESTIMATE

VARIABLE	DF	STANDARDIZED ESTIMATE	TOLERANCE	VARIANCE INFLATION
INTERCEP	1	0		
X1	1	0.04620202	0.007378307	135.53244
X2	1	-1.01374635	0.00559124	1780.51340
X3	1	-0.53754250	0.02974510	33.6100960
X4	1	-0.29474949	0.27643456	3.50093019
X5	1	-0.10122111	0.002505317	399.15102
X6	1	2.47966430	0.001317557	750.90060

COLLINEARITY DIAGNOSTICS

NUMBER	EIGENVALUE	CONDITION NUMBER	VAR PROP INTERCEP	VAR PROP X1	VAR PROP X2	VAR PROP X3
1	6.061393	1.000000	0.0000	0.0000	0.0000	0.0000
2	0.002103	9.101721	0.0000	0.0000	0.0000	0.0143
3	0.045601	12.255735	0.0000	0.0000	0.0000	0.0003
4	0.010600	25.336607	0.0000	0.0000	0.0003	0.0011
5	0.0001292	230.424	0.0000	0.0000	0.0560	0.0157
6	6.294E-06	1040.000	0.0001	0.5046	0.5204	0.2253
7	3.644E-09	43275.045	0.9999	0.0303	0.6546	0.6893

FIGURE 1.SAS MULTICOLLINEARITY  
 Low tolerance value and high condition number with large variance proportion for two or more variables may indicate multicollinearity. Variables X5 and X6 may be multicollinear with previous terms in the model.

NOTE: (PUB) VERSION = FF SERIAL = 023435 MODEL = 3001 .  
 you have questions, call the Computing Affairs Help Desk at 453-4361

1 DATA AERA; INPUT Y X1 X2 X3 X4 X5 X6;  
 2 OPTIONS NODATE LS=74;  
 3 CARDS;

NOTE: DATA SET WORK.AERA HAS 16 OBSERVATIONS AND 7 VARIABLES.  
 NOTE: THE DATA STATEMENT USED 0.06 SECONDS AND 24K.

20 PROC REG;  
 21 MODEL Y=X1 X2 X3 X4 X5 X6/STB P R CLI BW TOL VIF COLLIN;  
 22 OUTPUT OUT=NEW  
 23 P=PREB  
 24 R=RES;

NOTE: THE DATA SET WORK.NEW HAS 16 OBSERVATIONS AND 9 VARIABLES.  
 NOTE: THE PROCEDURE REG USED 0.18 SECONDS AND 44K  
 AND PRINTED PAGES 1 TO 2.

25 PROC PLOT;  
 26 PLOT PREDRES Y=1 Y=2 Y=3 Y=4 Y=5 Y=6;  
 NOTE: THE PROCEDURE PLOT USED 0.23 SECONDS AND 52K  
 AND PRINTED PAGES 3 TO 9.

27 PROC UNIVARIATE NORMAL PLOT;  
 NOTE: THE PROCEDURE UNIVARIATE USED 0.23 SECONDS AND 64K  
 AND PRINTED PAGES 10 TO 10.

28 PROC CORR;  
 NOTE: THE PROCEDURE CORR USED 0.11 SECONDS AND 76K  
 AND PRINTED PAGES 19 TO 21.  
 NOTE: SAS USED 76K MEMORY.



NUMBER	VAR PROF X4	VAR PROF X5	VAR PROF X6
1	0.0004	0.0000	0.0000
2	0.0919	0.0000	0.0000
3	0.0636	0.0000	0.0000
4	0.4267	0.0000	0.0000
5	0.1154	0.0097	0.0000
6	0.0000	0.0306	0.0002
7	0.3020	0.1597	0.9998

OBS	ACTUAL	PREDICT VALUE	STB ERR PREDICT	LOWESSZ PREDICT	UPPER95Z PREDICT	RESIDUAL
1	60323.0	60055.7	198.6	59252.6	60078.8	267.3
2	61122.0	61216.0	229.1	60353.3	62078.7	-94.0139
3	60171.0	60129.7	183.4	59519.9	60929.6	66.2072
4	61187.0	61597.1	186.0	60789.3	62405.8	-610.1
5	63221.0	62911.3	259.2	62039.7	63787.8	399.7
6	63639.0	63008.3	185.3	63001.2	64095.4	-249.3
7	64909.0	65153.0	213.7	64310.8	65995.3	-144.8
8	63761.0	63774.2	216.6	62928.2	64629.1	-13.1004
9	64019.0	64004.7	286.1	63172.2	64837.2	16.3448
10	67857.0	67401.6	175.3	64406.1	68197.1	655.4
11	68169.0	68106.3	182.9	67382.1	68990.5	-17.2409
12	64813.0	64832.1	211.9	65712.2	67991.9	-39.0550
13	60655.0	60616.5	186.5	60002.1	60619.0	-155.5
14	69544.0	67649.7	145.7	60008.3	70019.0	-85.6713
15	69331.0	60909.1	186.2	60181.0	69797.1	341.9
16	70551.0	70757.8	253.0	69061.6	71653.9	-206.8

OBS	STB ERR RESIDUAL	STUDENT RESIDUAL	COOK'S D
1	231.3	1.1560	0.141
2	201.1	-0.4676	0.041
3	243.5	0.1901	0.003
4	241.5	-1.0779	0.244
5	189.0	1.6309	0.614
6	242.1	-1.0300	0.009
7	217.4	-0.7547	0.079
8	210.6	-0.0614	0.001
9	224.6	0.0637	0.000
10	209.4	1.0258	0.235
11	243.9	-0.0708	0.000
12	219.2	-0.1782	0.004
13	241.1	-0.6451	0.034
14	267.8	-0.3199	0.004
15	241.4	1.0163	0.170
16	170.1	-1.2154	0.467

LM OF RESIDUALS -1.18744E-08  
 LM OF SQUARED RESIDUALS 836424.1  
 PREDICTED RESID SS (PRESS) 2006493  
 DURBIN-WATSON D 2.559  
 FOR NUMBER OF OBS. 16  
 1ST ORDER AUTOCORRELATION -0.348

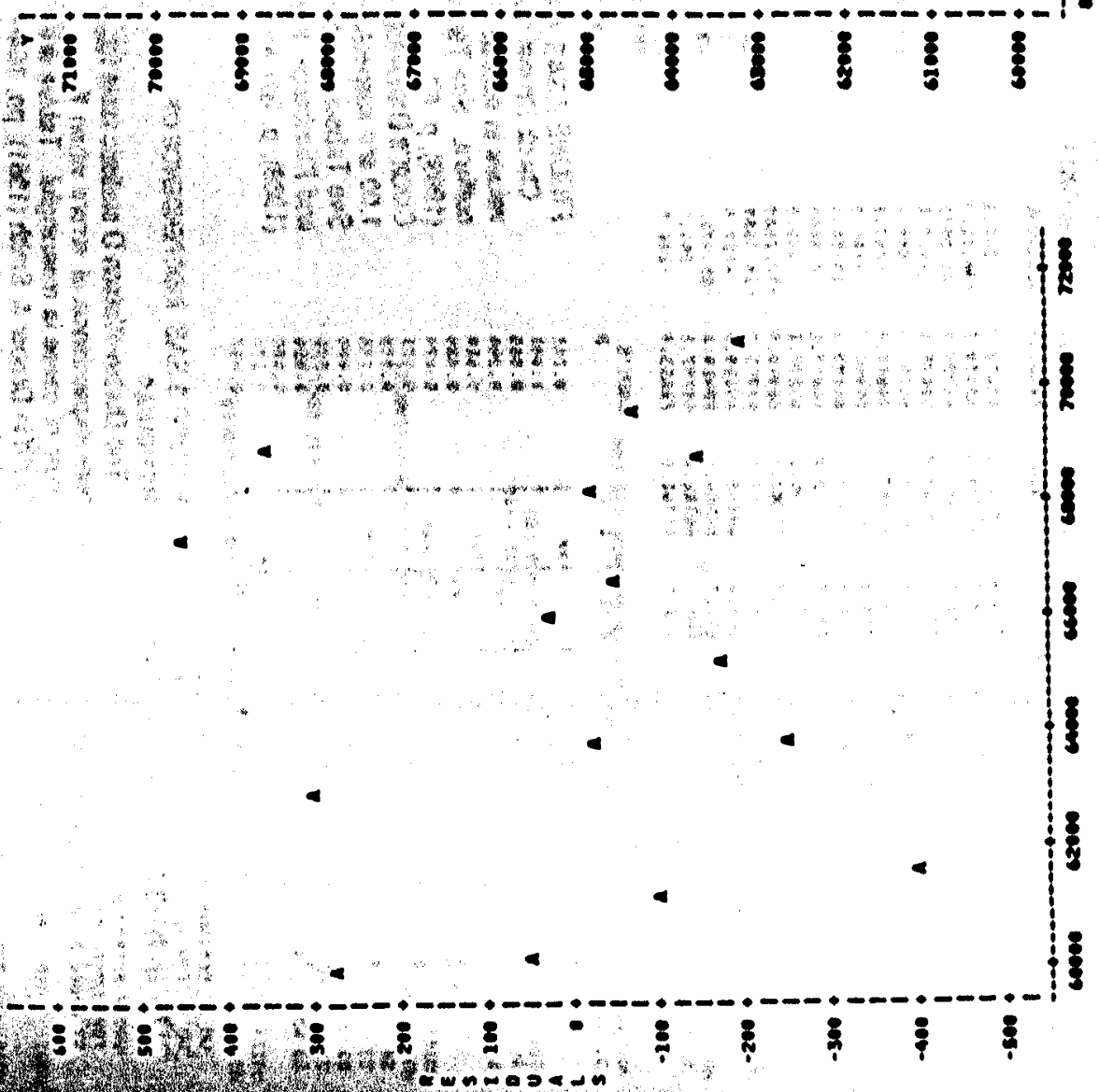
FIGURE 3.SAS INDEPENDENCE OF RESIDUALS

The Durbin-Watson D statistic tests for nonindependence of errors when the order of cases is meaningful. Tables are found in Draper & Smith (1981), pp. 164-166. It is irrelevant for this data set.

FIGURE 2.SAS OUTLIERS  
 Check STUDENTIZED RESIDUAL for values in excess of ±3.00 for univariate outliers. See Tabachnick & Fidell (1989), p. 67.  
 COOK'S DISTANCE values greater than 1.00 are suspected of being outliers. See Tabachnick & Fidell (1989), p. 130 and Kleinbaum, Kupper, & Muller (1988), p. 201.

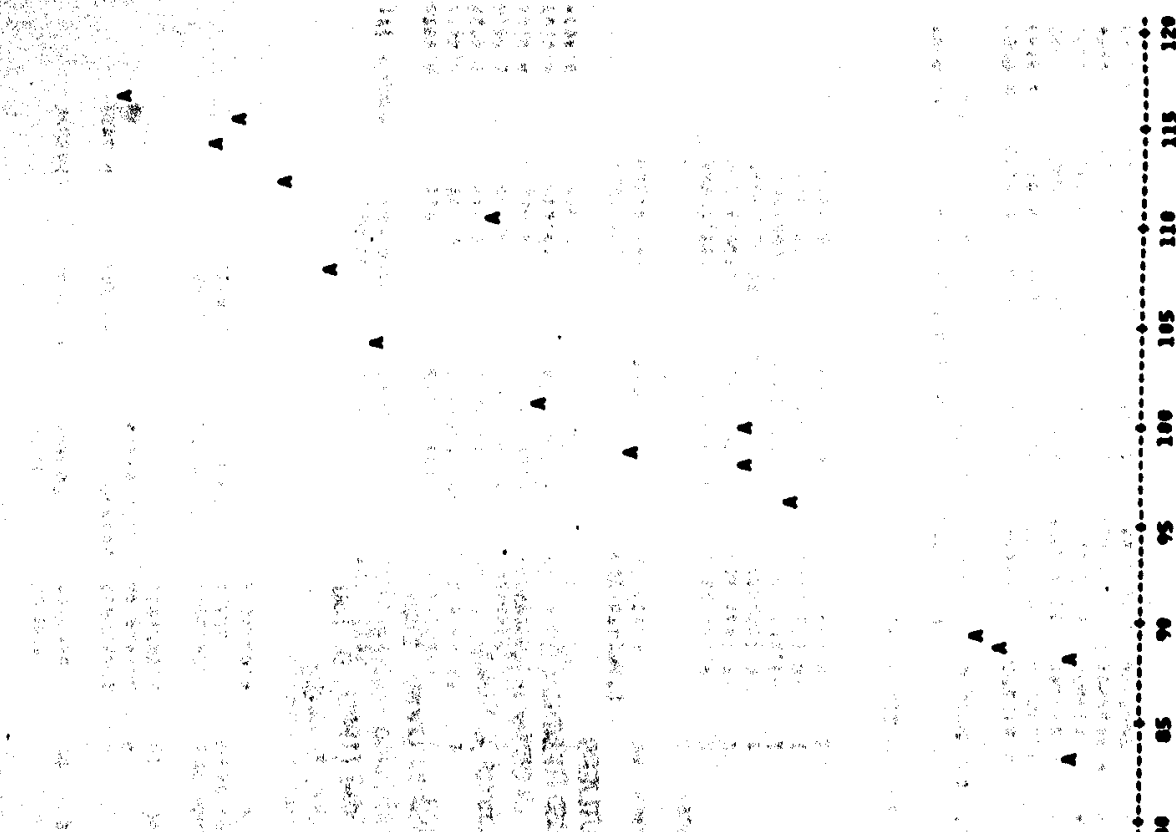
SAS

PLOT OF RESIDUALS  
LEGEND: A = 1 OBS, B = 2 OBS, ETC.



**FIGURE 4. SAS NORMALITY, LINEARITY, HOMOSEDASTICITY, AND INDEPENDENCE OF RESIDUALS**  
Examine the residual scatterplot to assess all four assumptions. All four assumptions are met in this data set.

PLOT OF Y=XI  
LEGEND: A = 1 OBS, B = 2 OBS, ETC.



**FIGURE 5. SAS LINEARITY**  
Use bivariate scatterplots to assess linearity of predictor - criterion association.

UNIVARIATE

VARIABLE=RES

RESIDUALS

MOMENTS

N	16	SUM WGT3	16
MEAN	-7.421E-10	SUM	-1.187E-08
STD DEV	236.139	VARIANCE	55761.6
SKEWNESS	0.464739	KURTOSIS	-0.298894
USS	836424	CS3	836424
CV	-.99999	STD MEAN	59.8347
T:MEAN=0	-1.257E-11	PROB> T	1
SEM BAK	-6	PROB> S	0.776105
NUM = 0	16		
W:NORMAL	0.948682	PROB<W	0.471

QUANTILES (DEF=4)

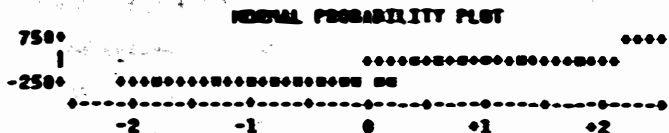
100% MAX	455.394
75% Q3	212.077
50% MED	-20.162
25% Q1	-161.924
0% MIN	-410.115
RANGE	865.509
Q3-Q1	374.001
MODE	-410.115

EXTREMES

99%	455.394	LOWEST	HIGHEST
95%	455.394	-410.115	46.2872
90%	375.97	-249.311	267.34
10%	-297.552	-206.758	309.715
5%	-410.115	-164.049	341.932
1%	-410.115	-155.55	435.394

STEM LEAF	0	0000000000000000
0 5	1	1
0 335	3	0000000
0 422221100000	12	0000000000000000

MULTIPLY STEM LEAF BY 1000000



SR

FIGURE 6.SAS NORMALITY OF RESIDUALS

If residuals are normally distributed, the plus signs (+) and asterisks (\*) should coincide in the normal probability plot. A statistical test for normality is also provided. In this case, WNORMAL = 0.948682, p = .471.

SKEWNESS can be detected by observing the stem and leaf and the boxplot as well as the skewness index of .464739 indicating the residuals are positively skewed. The kurtosis value of -.298894 indicates the residuals are platykurtic. Statistical tests of skewness and kurtosis are discussed in Tabachnick & Fidell (1989), pp. 72-73. For these data,

$$Z_{\text{SKEWNESS}} = \frac{S - 0}{\sqrt{\frac{6}{N}}} = \frac{.464739}{\sqrt{\frac{6}{16}}} = .76$$

$Z_{\text{SKEWNESS}} < 1.96 \therefore$  not skewed.

$$Z_{\text{KURTOSIS}} = \frac{K - 0}{\sqrt{\frac{24}{N}}} = \frac{-0.298894}{\sqrt{\frac{24}{16}}} = -.24$$

$Z_{\text{KURTOSIS}} > -1.96 \therefore$  no kurtosis

SAS

VARIABLE	N	MEAN	STD DEV	SUM	MINIMUM	MAXIMUM
Y	16	65317.0	3511.97	1045072	60171.0	70551.0
X1	16	101.7	10.79	1627	83.0	116.9
X2	16	387690.4	99394.94	6203175	234209.0	564894.0
X3	16	3193.3	934.46	51093	1870.0	4806.0
X4	16	2606.7	695.92	41707	1456.0	3594.0
X5	16	117424.0	6956.10	1878704	107600.0	130081.0
X6	16	1954.5	4.76	51272	1947.0	1962.0
PRED	16	65317.0	3504.02	1045072	40055.7	70757.8
RES	16	-7.421E-10	236.14	-1.107E-08	-410.1	455.4

SAS

PEARSON CORRELATION COEFFICIENTS / PROB &gt; |r| UNDER H0:RHO=0 / N = 16

	Y	X1	X2	X3	X4	X5
Y	1.00000	0.97090	0.98355	0.50250	0.45731	0.96039
	0.0000	0.0001	0.0001	0.0473	0.0749	0.0001
X1	0.97090	1.00000	0.99159	0.62063	0.46476	0.97916
	0.0001	0.0000	0.0001	0.0103	0.0697	0.0001
X2	0.98355	0.99159	1.00000	0.60426	0.44644	0.99109
	0.0001	0.0001	0.0000	0.0132	0.0030	0.0001
X3	0.50250	0.62063	0.60426	1.00000	-0.17742	0.68655
	0.0473	0.0103	0.0132	0.0000	0.5109	0.0033
X4	0.45731	0.46476	0.44644	-0.17742	1.00000	0.36442
	0.0749	0.0697	0.0030	0.5109	0.0000	0.1652
X5	0.96039	0.97916	0.99109	0.68655	0.36442	1.00000
	0.0001	0.0001	0.0001	0.0033	0.1652	0.0000
X6	0.97133	0.99115	0.99527	0.64826	0.41725	0.99395
	0.0001	0.0001	0.0001	0.0047	0.1078	0.0001
PRED PREDICTED VALUE	0.99774	0.97310	0.98570	0.50364	0.45834	0.96257
	0.0001	0.0001	0.0001	0.0467	0.0742	0.0001
RES RESIDUALS	0.06724	-0.00000	-0.00000	-0.00000	-0.00000	-0.00000
	0.0046	1.0000	1.0000	1.0000	1.0000	1.0000
		X6	PRED	RES		
Y	0.97133	0.99774	0.06724			
	0.0001	0.0001	0.0046			
X1	0.99115	0.97310	-0.00000			
	0.0001	0.0001	1.0000			
X2	0.99527	0.98570	-0.00000			
	0.0001	0.0001	1.0000			
X3	0.64826	0.50364	-0.00000			
	0.0047	0.0467	1.0000			
X4	0.41725	0.45834	-0.00000			
	0.1078	0.0742	1.0000			
X5	0.99395	0.96257	-0.00000			
	0.0001	0.0001	1.0000			
X6	1.00000	0.97353	-0.00000			
	0.0000	0.0001	1.0000			
PRED PREDICTED VALUE	0.97353	1.00000	-0.00000			
	0.0001	0.0000	1.0000			
RES RESIDUALS	-0.00000	-0.00000	1.00000			
	1.0000	1.0000	0.0000			

### FIGURE 7.SAS CORRELATION OF ERRORS AND INDEPENDENT VARIABLES

Use the correlation matrix to determine association between each independent variable and the residual from the multiple regression equation.

Try the new SPSS-X Release 3.1 features:

- Interactive SPSS-X command execution
- Online Help
- Nonlinear Regression
- Time Series and Forecasting (TREMS)
- Macro Facility
- The new EMX procedure
- Improvements in:
  - REWRT and TABLES
  - SIMPLIFIED SYNTAX
  - MATRIX I/O

See SPSS-X User's Guide, Third Edition, for more information on these features.

```

1 0 DATA LIST FREE / Y XI X2 X3 X4 X5 X6
2 0 SET WIDTH =00
3 BEGIN DATA
19 END DATA
    
```

Proceeding task required .02 seconds (CPU time); .03 seconds elapsed.

```
20 PLOT Y WITH XI/
```

There are 1,749,224 bytes of memory available.  
 The largest contiguous area has 1,749,224 bytes.

PLOT requires 15000 bytes of workspace for execution.

\*\*\*\*\* P L O T \*\*\*\*\*

Data Information

16 unweighted cases accepted.

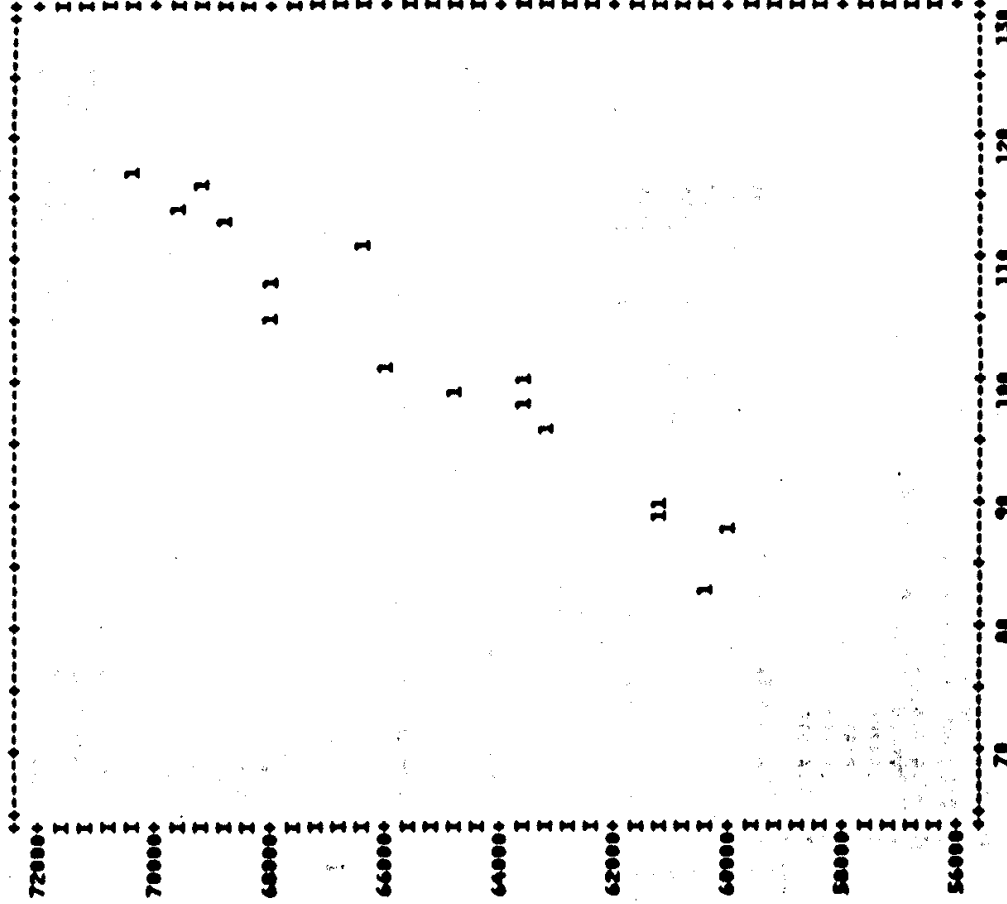
Size of the plots

Horizontal size is 65  
 Vertical size is 40

Frequencies and symbols used (not applicable for control or overlay plots)

1	1	11	B	21	L	31	V
2	2	12	C	22	M	32	W
3	3	13	D	23	N	33	X
4	4	14	E	24	O	34	Y
5	5	15	F	25	P	35	Z
6	6	16	G	26	Q	36	+
7	7	17	H	27	R		
8	8	18	I	28	S		
9	9	19	J	29	T		
10	10	20	K	30	U		

PLOT OF Y WITH XI



16 cases plotted.

FIGURE 1. SPSSX. LINEARITY  
 Use bivariate scatterplots to assess  
 linearity of predictor-criterion  
 association.

Proceeding task required .85 seconds CPU time; .44 seconds elapsed.

```

21 REGRESSION VARS=Y X1 X2 X3 X4 X5 X6/STATISTICS=ALL F/DEP=Y
22 /METHOD=ENTER
23 /RESID=DEFAULTS
24 /CASEWISE=ALL DEFAULTS
25 /SCATTERPLOT=(R2RESID, ZPRED)
26 /SAVE PRED (WHAT) RESID (ERR)
    
```

There are 1,749,184 bytes of memory available.  
The largest contiguous area has 1,749,832 bytes.

\*\*\*\*\* MULTIPLE REGRESSION \*\*\*\*\*

Listwise Deletion of Missing Data

Equation Number 1 Dependent Variable... Y

Beginning Block Number 1. Method: Enter

Variable(s) Entered on Step Number

- 1... X6
- 2... X4
- 3... X3
- 4... X1
- 5... X5
- 6... X2

```

Multiple R      .97774
R Square        .99548
Adjusted R Square .97247
Standard Error  304.85487
R Square Change .99548
F Change       330.28534
Signif F Change .0000
    
```

```

Analysis of Variance
DF      Sum of Squares      Mean Square
Regression      6      104172481.94449      34695488.32448
Residual        9      836429.86351       92936.66617
F = 330.28534      Signif F = .0000
    
```

```

AIC      187.82804
PC        .01155
CP        7.06000
SBC      193.23696
    
```

Var-Cov Matrix of Regression Coefficients (B)  
Below Diagonal: Covariances Above: Correlations

	X6	X4	X3	X1	X5	X2
X6	207468.643	-.54957	-.82918	-.18428	.30816	-.88168
X4	-.531674	.04591	.61857	-.34081	-.18091	.46860
X3	-.18332591	.06473	.23853	-.53300	-.73826	.94561
X1	7204.91263	-6.34671	-25.01719	7218.54462	.65918	-.64942
X5	39.96948	-.08372	.01722	12.65424	.05111	-.83321
X2	-.1222010	.00000	.00000	.00000	.00000	.00000

XTX Matrix

	X6	X4	X3	X1	X5	X2
X6	758.98068	-28.67217	-131.6399	59.74668	213.64593	-934.8346
X4	-28.67217	3.50893	6.79454	-7.69386	-7.15817	37.54356
X3	-131.6399	6.79454	33.61889	-37.44332	-87.83677	231.87218
X1	59.74668	-7.69386	-37.44332	153.53244	153.51803	-319.7367
X5	213.64593	-7.15817	-87.83677	153.51803	399.15102	-703.9988
X2	-934.8346	37.54356	231.87218	-319.7367	-703.9988	1788.5135
Y	2.47966	-.20474	-.53754	.04628	-.10122	-1.01375

	Y
X6	-2.47966
X4	-.20474
X3	-.53754
X1	-.04628
X5	-.10122
X2	1.01375
Y	.00452

Equation Number 1 Dependent Variable... Y

Variables in the Equation

Variable	B	SE B	95% Confidence Interval B	Beta
X6	1829.151665	455.478499	798.788692 2869.514237	2.479664
X4	-1.833227	.214274	-1.517948 -.548506	-.294741
X3	-2.832230	.488400	-3.125065 -.537943	-.537943
X1	15.861872	84.914926	-177.828816 207.152540	.046282
X5	-.831104	.226873	-.542517 -.460308	-.181221
X2	-.838819	.833491	-.111881 -.899943	-.1013746
(Constant)	-3482258.635	876429.3836	-5496527.183 -1467799.084	

Variables in the Equation

Variable	SE Beta	Correl Part Cor	Partial Tolerance	VIF	F
X6	.617943	.971329	.090087	.001318	798.901
X4	.042448	.457307	-.188874	.276435	3.589
X3	.129953	.582498	-.092789	.827945	33.619
X1	.264926	.978099	.063973	.007378	138.532
X5	.447788	.948391	-.085046	.075137	399.151
X2	.947855	.983862	-.025971	.335804	5.571E-04
(Constant)					1788.513
					1.144
					15.294

in

Variable Sig F

X6	.0030
X4	.0009
X3	.0025
X1	.0431
X5	.0262
X2	.3127
(Constant)	.0036

Number	Eigenval	Cond Index	Variances Constant	Proportions X1	X2	X3	X4
1	6.86139	1.000	.00000	.00000	.00000	.00004	.00035
2	.08210	9.142	.00000	.00000	.00001	.01428	.09191
3	.04568	12.256	.00000	.00000	.00026	.00004	.06357
4	.01869	25.337	.00000	.00034	.00107	.06464	.42672
5	.00013	238.424	.00000	.45677	.01566	.00559	.11848
6	.00001	1048.000	-.00015	.50456	.52839	.22334	.00000
7	.00000	43275.85	.97985	.03033	.65463	.68726	.30285

	X6
1	.00000
2	.00000
3	.00000
4	.00000
5	.00000
6	.00016
7	.97984

End Block Number 1 All requested variables entered.

Summary table

Step	Mult R	Rsq	F (Em)	Sig	Variable	Delta In
1					In: X6	.9713
2					In: X4	.0630
3					In: X3	-.3910
4					In: X1	-.0224
5					In: X5	-.5902
6	.9977	.9955	330.285	.000	In: X2	-1.0137

FIGURE 2.SPSSX. MULTICOLLINEARITY  
 Low tolerance value and high condition number with large variance proportion for two or more variables may indicate multicollinearity. Variables X5 and X6 may be multicollinear with previous terms in the model.

\*\*\*\*\* MULTIPLE REGRESSION \*\*\*\*\*

Equation Number 1 Dependent Variable.. Y

Casewise Plot of Standardized Residual

N: Selected N: Missing

Case #	-3.0	0.0	3.0	Y	#PRED	#RESID
1	.	.	.	60025.00	60025.6400	267.3400
2	.	.	.	61122.00	61216.8139	-94.8139
3	.	.	.	60171.00	60124.7120	46.2872
4	.	.	.	61187.00	61577.1146	-410.1146
5	.	.	.	63221.00	62911.2094	309.7146
6	.	.	.	63439.00	63000.3112	-249.3112
7	.	.	.	64909.00	65153.0490	-144.0490
8	.	.	.	63761.00	63774.1004	-13.1004
9	.	.	.	64019.00	64004.6932	14.3048
10	.	.	.	67057.00	67401.6059	455.3941
11	.	.	.	60169.00	60106.2409	-17.2409
12	.	.	.	64513.00	64552.0550	-39.0550
13	.	.	.	60453.00	60010.5500	-155.5500
14	.	.	.	69544.00	69449.6713	-85.6713
15	.	.	.	69331.00	68909.0405	341.9315
16	.	.	.	70531.00	70757.7570	-206.7570

Outliers - Standardized Residual

Case #	#ZRESID
10	1.49381
4	-1.34520
15	1.12162
5	1.01594
1	.87694
6	-.81781
16	-.67022
7	-.53012
13	-.51024
2	-.30839

Histogram - Standardized Residual

Height	N	(N = 1 Cases, . : = Normal Curve)
0	.01	Out
0	.02	3.00
0	.06	2.67
0	.14	2.33
0	.29	2.00
0	.53	1.67
1	.88	1.33
3	1.29	1.00
0	1.70	.67
0	2.00	.33
5	2.12	.00
2	2.00	-.33
4	1.70	-.67
0	1.29	-1.00
1	.88	-1.33
0	.53	-1.67
0	.29	-2.00
0	.14	-2.33
0	.06	-2.67
0	.02	-3.00
0	.01	Out

\*\*\*\*\* MULTIPLE REGRESSION \*\*\*\*\*

Equation Number 1 Dependent Variable.. Y

Residuals Statistics:

	Min	Max	Mean	Std Dev	N
#PRED	60025.6543	70757.7500	65317.0000	3500.0206	16
#RESID	-410.1146	455.3940	.0000	236.1309	16
#ZPRED	-1.5015	1.5527	.0000	1.0000	16
#ZRESID	-1.3453	1.4938	.0000	.7746	16

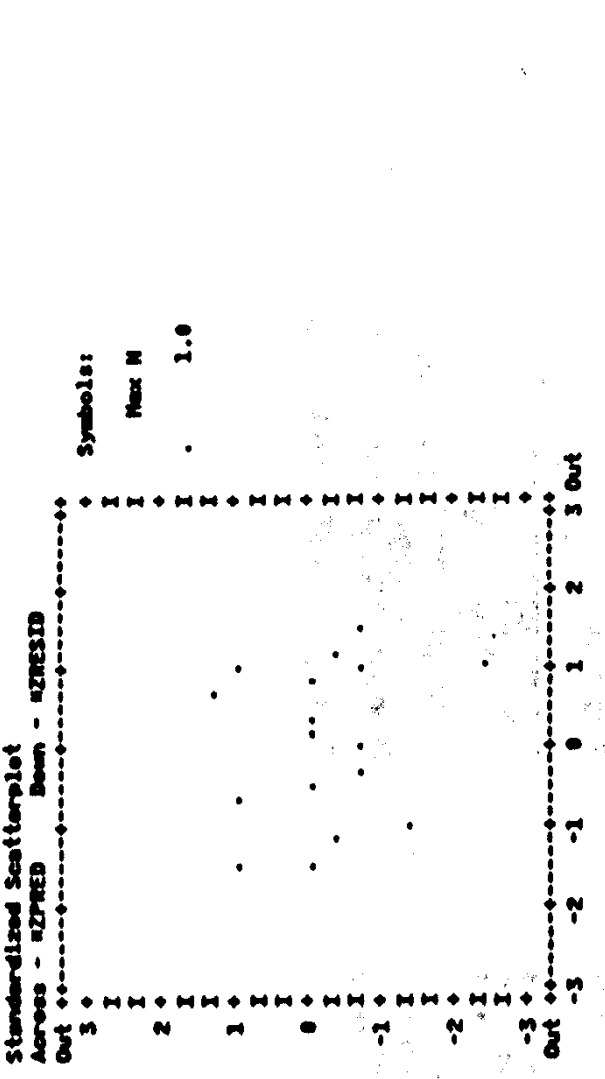
Total Cases = 16

Durbin-Watson Test = 2.59999

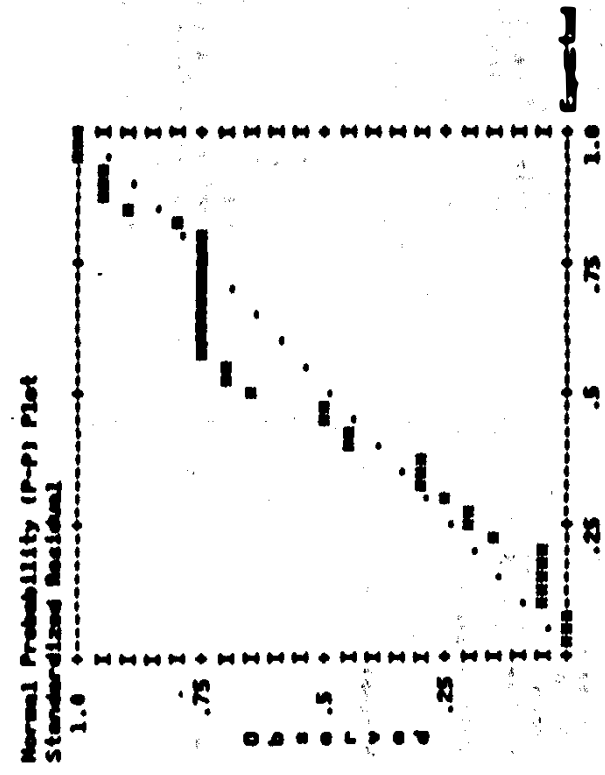
FIGURE 3.SPSSX. OUTLIERS

Check STUDENTIZED RESIDUAL for values in excess of  $\pm 3.00$  for univariate outliers. See Tabachnick & Fidell (1989), p. 67. Also check plot of STANDARDIZED RESIDUAL. Also note that the Durbin-Watson D statistic tests for nonindependence of errors when the order of cases is meaningful. Tables are found in Draper & Smith (1981), pp. 164-166. It is irrelevant for this data set.





**FIGURE 5. SPSSX NORMALITY, LINEARITY, HOMOSEDASTICITY, AND INDEPENDENCE OF RESIDUALS**  
 Examine the standardized scatterplot of predicted and residual values to assess all four assumptions. All are met in this data set.



**FIGURE 4. SPSSX. NORMALITY OF RESIDUALS**  
 If residuals are normally distributed, the dots (•) and asterisks (\*) should coincide in the NORMAL PROBABILITY PLOT.

Preceding task required .19 seconds CPU time; 1.83 seconds elapsed.

27 COMBINE STATISTICS YMAT ERR/  
 28 STATISTICS ALL

>Warning 8 11003  
 >The new default column-style printing cannot be used for this DESCRIPTIVES, as  
 >there are too many statistics to print on one line per variable. Old style  
 >printing will be used instead.

There are 1,746,296 bytes of memory available.  
 The largest contiguous area has 1,746,720 bytes.

148 bytes of memory required for the DESCRIPTIVES procedure.  
 4 bytes have already been acquired.  
 144 bytes remain to be acquired.

Number of valid observations (listwise) = 16.00

Variable YMAT Predicted Values

Mean	65317.000	S.E. Mean	876.005
Std Dev	3594.821	Variances	12278160.136
Kurtosis	-1.299	S.E. Kurt	1.891
Skewness	-.101	S.E. Skew	.564
Range	18782.898	Minimum	60055.65977
Maximum	70757.75783	Sum	1065872.000

Valid observations - 16 Missing observations - 0

Variable ERR Residual

Mean	.000	S.E. Mean	59.835
Std Dev	236.139	Variances	55761.606
Kurtosis	-.299	S.E. Kurt	1.891
Skewness	.465	S.E. Skew	.564
Range	845.589	Minimum	-418.11462
Maximum	455.39489	Sum	-7.6161041598E-09

Valid observations - 16 Missing observations - 0

FIGURE 6.SPSSX. NORMALITY OF  
 RESIDUALS

SKWENESS index of .465 indicates the  
 residuals are positively skewed. The  
 KURTOSIS value of -.299 indicates the  
 residuals are platykurtic. Statistical tests  
 of skewness and kurtosis are discussed  
 in Tabachnick & Fidell (1989), pp. 72-73.  
 For these data,

$$Z_{SKWENESS} = \frac{S - 0}{\sqrt{\frac{6}{N} \sqrt{\frac{6}{16}}}} = \frac{.465}{.564} = .82$$

or, using reported SE<sub>SKWEN</sub>.

$$Z_{SKWENESS} = \frac{.465}{.564} = .82$$

Z<sub>SKWENESS</sub> < 1.96 ∴ not skewed.

Z<sub>KURTOSIS</sub> =

$$\frac{K - 0}{\sqrt{\frac{24}{N} \sqrt{\frac{24}{16}}}} = \frac{-0.299}{.1091} = -.27$$

or, using reported SE<sub>KURTOSIS</sub>.

$$Z_{KURTOSIS} = \frac{-.299}{.1091} = -.27$$

Z<sub>KURTOSIS</sub> > -1.96 ∴ no kurtosis

29 PEARSON CORR Y X1 X2 X3 X4 X5 X6 YHAT ERR/  
30 OPTIONS 6

PEARSON CORR problem requires 1,872 bytes of workspace.

--- PEARSON CORRELATION COEFFICIENTS ---

VARIABLE PAIR		VARIABLE PAIR		VARIABLE PAIR		VARIABLE PAIR	
Y WITH X1	.9709 N( 16) SIG .000	Y WITH X2	.9836 N( 16) SIG .000	Y WITH X3	.5025 N( 16) SIG .029	Y WITH X4	.4573 N( 16) SIG .037
Y WITH X5	.9600 N( 16) SIG .000	Y WITH X6	.9713 N( 16) SIG .000	Y WITH YHAT	.9977 N( 16) SIG .000	Y WITH ERR	.0672 N( 16) SIG .402
X1 WITH X2	.9916 N( 16) SIG .000	X1 WITH X3	.6206 N( 16) SIG .005	X1 WITH X4	.4647 N( 16) SIG .035	X1 WITH X5	.9792 N( 16) SIG .000
X1 WITH X6	.9911 N( 16) SIG .000	X1 WITH YHAT	.9731 N( 16) SIG .000	X1 WITH ERR	.0000 N( 16) SIG .500	X2 WITH X3	.6043 N( 16) SIG .007
X2 WITH X4	.4464 N( 16) SIG .042	X2 WITH X5	.9911 N( 16) SIG .000	X2 WITH X6	.9953 N( 16) SIG .000	X2 WITH YHAT	.9888 N( 16) SIG .000
X2 WITH ERR	.0000 N( 16) SIG .500	X3 WITH X4	-.1774 N( 16) SIG .255	X3 WITH X5	.6066 N( 16) SIG .042	X3 WITH X6	.6483 N( 16) SIG .002
X3 WITH YHAT	.5036 N( 16) SIG .023	X3 WITH ERR	.0000 N( 16) SIG .500	X4 WITH X5	.3644 N( 16) SIG .003	X4 WITH X6	.4172 N( 16) SIG .054
X4 WITH YHAT	.4583 N( 16) SIG .037	X4 WITH ERR	.0000 N( 16) SIG .500	X5 WITH X6	.9940 N( 16) SIG .000	X5 WITH YHAT	.9626 N( 16) SIG .000
X5 WITH ERR	.0000 N( 16) SIG .500	X6 WITH YHAT	.9735 N( 16) SIG .000	X6 WITH ERR	.0000 N( 16) SIG .500	YHAT WITH ERR	.0000 N( 16) SIG .500

SIG IS 1-TAILED, "." IS PRINTED IF A COEFFICIENT CANNOT BE COMPUTED.

FIGURE 7.SPSSX. CORRELATION OF ERRORS AND INDEPENDENT VARIABLES  
Use the correlation matrix to determine association between each independent variable and the residual from the multiple regression equation.