

Erik Thompson  
DSCI 449  
Spring 2024 -HW 1  
Data Examination

### **Research Problem**

This study will focus on the ratings of technical support offered by HBAT industries as reported by respondents or purchasing managers based on the length of relationship. A random and independent sample size of 100 with a confidence level of 95% will be used in the examination. Three values will be used to categorize the customer types, 1= less than 1 year; 2=between 1 and 2 years; 3=longer than 5 years.

My variable, Tech\_support is approximately normally distributed. These findings are supported by mean of 5.37 < 5.4 median, indicating nearly symmetrical distribution. The skew is slightly negative, but is nearly symmetrical with a skewness value of -0.20, and a kurtosis value of -0.55 indicating a distribution that is slightly flat, but still normal. Shapiro-Wilk indicates normality with  $W=0.9863$  and p-value of  $0.39 > \alpha$  of .05. There is a positive linear relationship in the normal probability plot, the whiskers of the box plot are similar in length, with the mean and median close to the center, and the histogram appears to be multi-modal validating a normal distribution. There are no outliers in the box plot, no data points above or below the upper and lower fence. My variable does not require transformation because normality can be assumed without any adjustments to the model.

Regression output shows no outliers as defined as standardized residuals with absolute values > 3. The largest is observation 31 with a student residual of 2.16. The absence of bivariate outliers is supported by the lack of values outside the ellipse in the scatterplot. Univariate analysis of residuals is shown to be normal with a mean close to zero, nearly symmetric skewness with a value of -0.18, and a kurtosis value of -0.62. Shapiro-Wilk indicates normality with  $W=0.9876$  and a p-value of  $0.48 > \alpha$  of .05. Regression output shows no relationship between the variables with  $F(1,98)=0.92$  and a P-value of 0.3387. Homoscedasticity is valid because there is no apparent pattern in the scatterplot or residual plots. Lots of scatter, but linear relationship is possible. Homogeneity of variances appears to be valid, Levene's test is not significant,  $F(2,97)=0.73$  and P-value 0.3841. Boxplots appear to be similar among the 3 customer type values. There are no visible outliers in the boxplots for the 3 customer types.

### **Research Study Background Information**

The HBAT database has one hundred observations with 23 variables and is assumed to be based on surveys of HBAT customers completed on a secure Web site managed by an established business research company. The research company contacts purchasing managers and encourages them to participate. Managers log onto the Web site and complete the survey. The survey results are supplemented by other information, compiled, and stored in HBAT's data warehouse and accessible through its decision support system.

## *The SAS System*

### *The MEANS Procedure*

HBAT sells paper products to two market segments: the newsprint industry and the magazine industry. Paper products are sold to these market segments either directly to the customer or indirectly through a broker. The survey questionnaire consisted of two sections. The first section addressed perceptions of HBAT's performance on thirteen attributes. These attributes were developed from focus groups, pretested, and have been used in previous studies. They are considered the most influential in the selection of suppliers in the paper industry. Respondents were purchasing managers of firms buying from HBAT, and they rated HBAT on each of the 13 attributes using a 0-10 scale, where 10 = "excellent" and 0 = "poor". The second section relates to purchase outcomes and business relationships (e.g., satisfaction with HBAT and whether the firm would consider a strategic alliance or partnership with HBAT). A third type of information was available from HBAT's data warehouse and includes information such as size of customer and length of purchase relationship.

Of the twenty-three variables 17 were metric (quantitative). Thirteen variables measured perceptions of HBAT, four variables measured purchase outcomes (e.g., purchase relationships with HBAT).

Variable	N	N Miss	Mean	Std Dev	Minimum	Maximum
Prod_Qual	100	0	7.8100000	1.3962793	5.0000000	10.0000000
Ecommerce	100	0	3.6720000	0.7005164	2.2000000	5.7000000
Tech_support	100	0	5.3650000	1.5304568	1.3000000	8.5000000
Complaint	100	0	5.4420000	1.2084032	2.6000000	7.8000000
Adv	100	0	4.0100000	1.1269428	1.9000000	6.5000000
Prod_Line	100	0	5.8050000	1.3152850	2.3000000	8.4000000
Sales_Image	100	0	5.1230000	1.0723198	2.9000000	8.2000000
Pricing	100	0	6.9740000	1.5450553	3.7000000	9.9000000
Warranty	100	0	6.0430000	0.8197382	4.1000000	8.1000000
New_Prod	100	0	5.1500000	1.4930479	1.7000000	9.5000000
Ordering	100	0	4.2780000	0.9288398	2.0000000	6.7000000
Price_Flex	100	0	4.6100000	1.2060035	2.6000000	7.3000000
Del_Speed	100	0	3.8860000	0.7344372	1.6000000	5.5000000
Satis	100	0	6.9180000	1.1918393	4.7000000	9.9000000
Recommend	100	0	7.0200000	1.0433048	4.6000000	9.9000000
Future_Purch	100	0	7.7130000	0.9361057	5.5000000	9.9000000
Usage_level	100	0	58.4000000	8.8608776	37.1000000	77.1000000

Tech\_support has a mean of 5.37, falling in the middle of the 0-10 scale. This suggests an average user experience that is neither poor or excellent. Tech\_support does have the lowest minimum value of 1.3, showing that out of the respondents that rated HBAT, Tech\_support received the worst score. There were several variables that received a rating above 9, so the maximum rating of 8.5 for Tech\_support is not an exceptional score.

*The SAS System**The FREQ Procedure*

Customer_type	Frequency	Percent	Cumulative Frequency	Cumulative Percent
1	32	32.00	32	32.00
2	35	35.00	67	67.00
3	33	33.00	100	100.00

Industry_type	Frequency	Percent	Cumulative Frequency	Cumulative Percent
0	52	52.00	52	52.00
1	48	48.00	100	100.00

Firm_size	Frequency	Percent	Cumulative Frequency	Cumulative Percent
0	49	49.00	49	49.00
1	51	51.00	100	100.00

Region	Frequency	Percent	Cumulative Frequency	Cumulative Percent
0	39	39.00	39	39.00
1	61	61.00	100	100.00

Distribution	Frequency	Percent	Cumulative Frequency	Cumulative Percent
0	57	57.00	57	57.00
1	43	43.00	100	100.00

Partner	Frequency	Percent	Cumulative Frequency	Cumulative Percent
0	55	55.00	55	55.00
1	45	45.00	100	100.00

## The SAS System

### The UNIVARIATE Procedure

Variable:  
Tech\_support

Moments			
<b>N</b>	100	<b>Sum Weights</b>	100
<b>Mean</b>	5.365	<b>Sum Observations</b>	536.5
<b>Std Deviation</b>	1.53045679	<b>Variance</b>	2.34229798
<b>Skewness</b>	-0.2032586	<b>Kurtosis</b>	-0.5482262
<b>Uncorrected SS</b>	3110.21	<b>Corrected SS</b>	231.8875
<b>Coeff Variation</b>	28.5266876	<b>Std Error Mean</b>	0.15304568

Basic Statistical Measures			
Location		Variability	
<b>Mean</b>	5.365000	<b>Std Deviation</b>	1.53046
<b>Median</b>	5.400000	<b>Variance</b>	2.34230
<b>Mode</b>	4.600000	<b>Range</b>	7.20000
		<b>Interquartile Range</b>	2.45000

*Note: The mode displayed is the smallest of 5 modes with a count of 4.*

#### Examination of 'Tech\_support'

The metric variable 'Tech\_support' shows a slightly negative skew, but is nearly symmetrical with a skewness value of -0.20.

A kurtosis value of -0.55 indicates that the distribution is slightly flat, but still a normal distribution.

The mean of 5.37 < 5.4 median also indicating a nearly symmetrical distribution.

The data is multi-modal with a mode of 4.6 for the smallest of 5 modes with a count of 4.

#### Location Test

$H_0: \mu=0$

$H_1: \mu \neq 0$

$t(99) = 35.05$

A [p-value below .0001] <  $\alpha$  of .05. Reject the null. The mean of 'Tech\_support' is statistically different from 0. It is much greater than 0.

Tests for Location: Mu0=0				
Test	Statistic		p Value	
<b>Student's t</b>	<b>t</b>	35.05489	<b>Pr &gt;  t </b>	<.0001
<b>Sign</b>	<b>M</b>	50	<b>Pr &gt;=  M </b>	<.0001
<b>Signed Rank</b>	<b>S</b>	2525	<b>Pr &gt;=  S </b>	<.0001

#### Normality Test

$H_0$ : Data is normal.

$H_1$ : Data is not normal.

$W=0.9863$

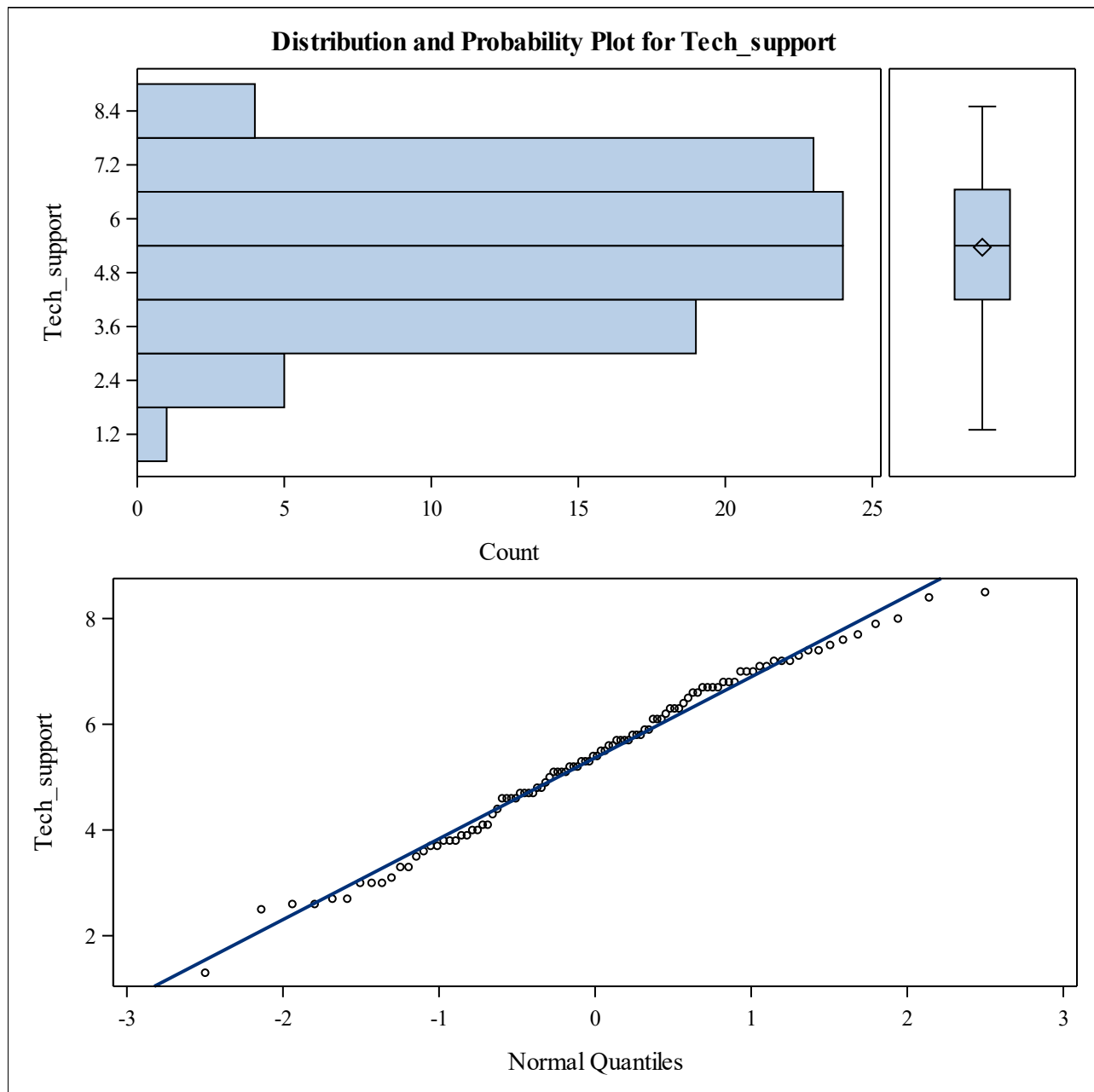
A p-value of 0.39 >  $\alpha$  of .05. Fail to reject the null. The variable, 'Tech\_support' is normally distributed.

Tests for Normality				
Test	Statistic		p Value	
<b>Shapiro-Wilk</b>	<b>W</b>	0.986264	<b>Pr &lt; W</b>	0.3900
<b>Kolmogorov-Smirnov</b>	<b>D</b>	0.060152	<b>Pr &gt; D</b>	>0.1500
<b>Cramer-von Mises</b>	<b>W-Sq</b>	0.051364	<b>Pr &gt; W-Sq</b>	>0.2500
<b>Anderson-Darling</b>	<b>A-Sq</b>	0.370975	<b>Pr &gt; A-Sq</b>	>0.2500

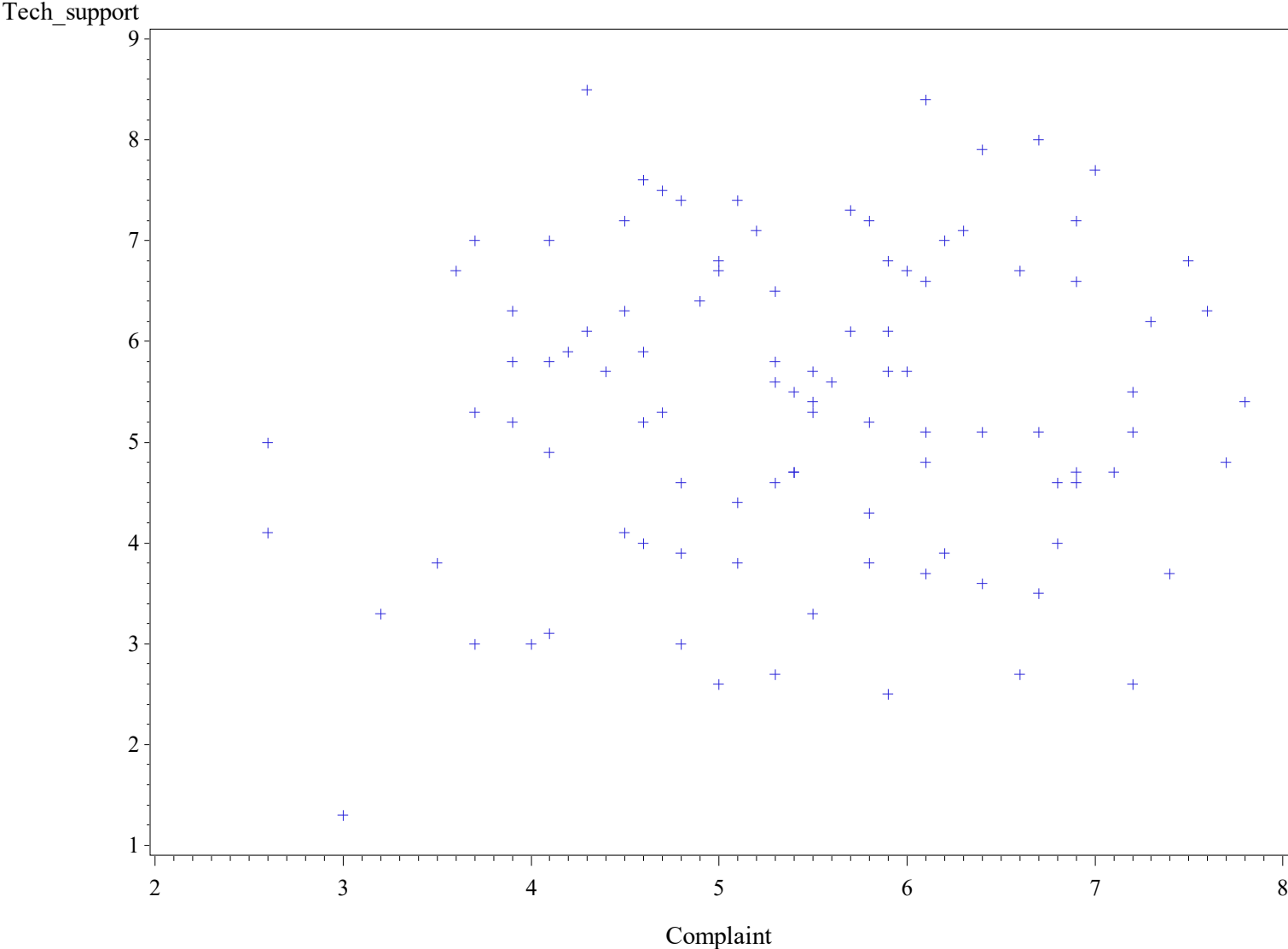
*The SAS System**The UNIVARIATE Procedure**Variable:**Tech\_support*

<b>Quantiles (Definition 5)</b>	
<b>Level</b>	<b>Quantile</b>
<b>100% Max</b>	8.50
<b>99%</b>	8.45
<b>95%</b>	7.65
<b>90%</b>	7.25
<b>75% Q3</b>	6.65
<b>50% Median</b>	5.40
<b>25% Q1</b>	4.20
<b>10%</b>	3.20
<b>5%</b>	2.70
<b>1%</b>	1.90
<b>0% Min</b>	1.30

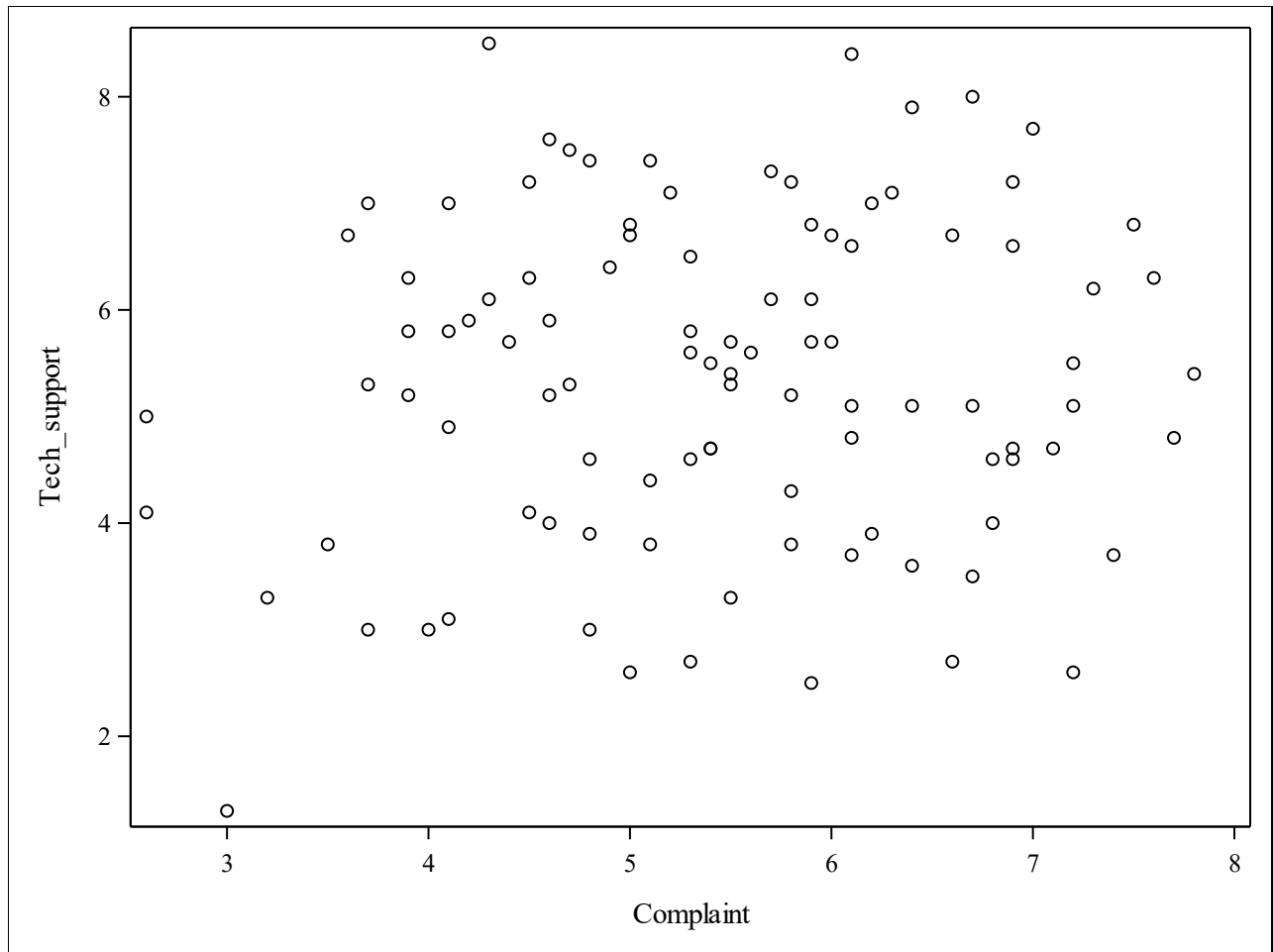
<b>Extreme Observations</b>			
<b>Lowest</b>		<b>Highest</b>	
<b>Value</b>	<b>Obs</b>	<b>Value</b>	<b>Obs</b>
1.3	87	7.7	90
2.5	1	7.9	67
2.6	97	8.0	88
2.6	24	8.4	61
2.7	73	8.5	31

*The SAS System**The UNIVARIATE Procedure*

There is a positive linear relationship in the normal probability plot. The skew is nearly symmetrical and the data appears normally distributed. The whiskers of the box plot are similar in length, with the mean and median close to the center. The histogram appears to be multi-modal and validates a normal distribution. There are no outliers in the box plot.



The above and below scatter plots appears to be randomized with no noticeable pattern, suggesting homoscedasticity.



*The SAS System**The CORR Procedure*

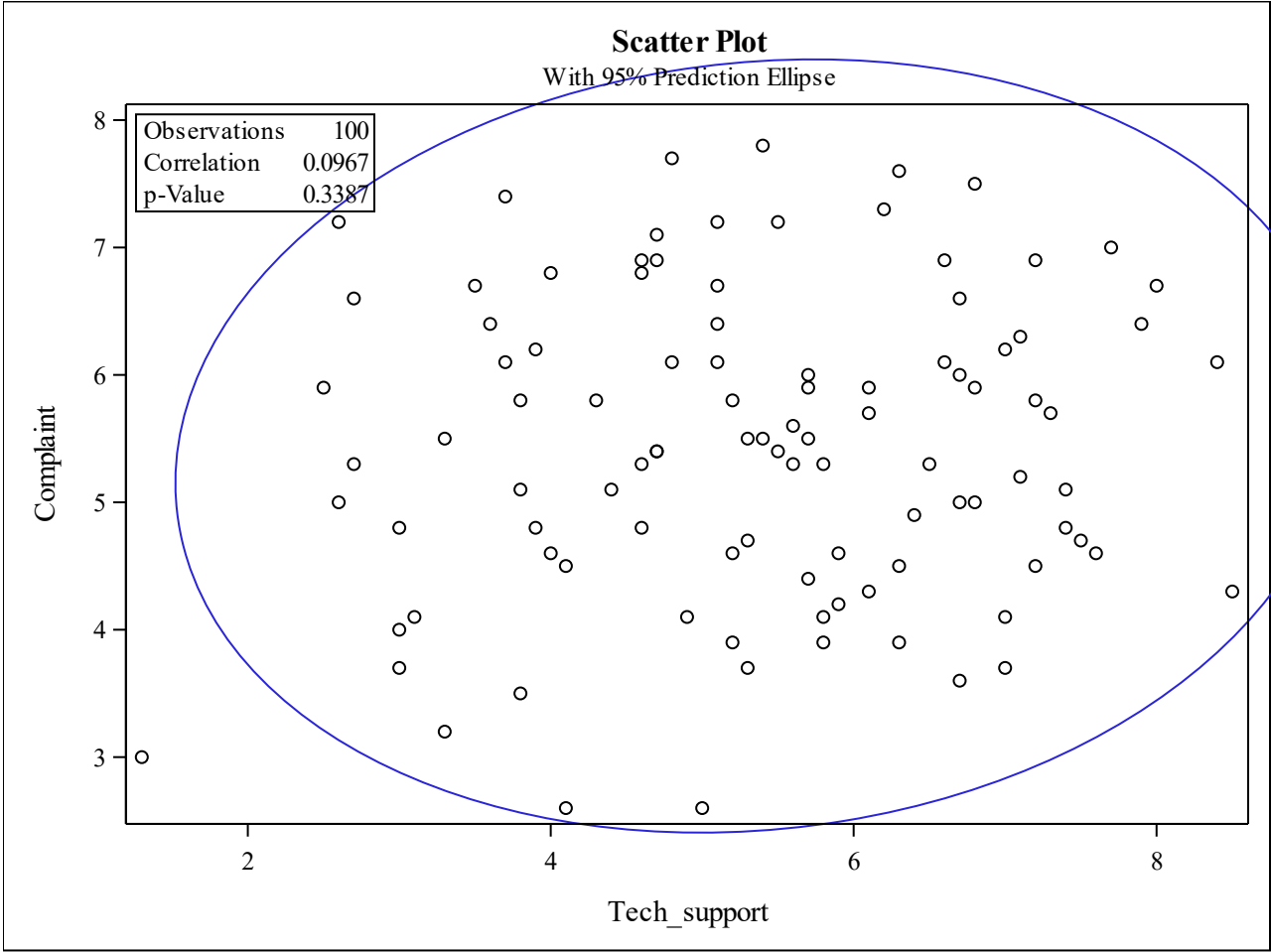
<b>2</b> <b>Variables:</b>	Tech_support Complaint
-------------------------------	------------------------

Simple Statistics						
Variable	N	Mean	Std Dev	Sum	Minimum	Maximum
<b>Tech_support</b>	100	5.36500	1.53046	536.50000	1.30000	8.50000
<b>Complaint</b>	100	5.44200	1.20840	544.20000	2.60000	7.80000

Pearson Correlation Coefficients, N = 100 Prob >  r  under H0: Rho=0		
	Tech_support	Complaint
Tech_support	1.00000	0.09666 0.3387
Complaint	0.09666 0.3387	1.00000

*The SAS System*

*The CORR Procedure*



There are no values outside the ellipse in the above scatterplot. No bivariate outliers are observed.

*The SAS System**The ANOVA Procedure*

Class Level Information		
Class	Levels	Values
Customer_type	3	1 2 3

Number of Observations Read	100
Number of Observations Used	100

*The SAS System**The ANOVA Procedure**Dependent Variable: Tech\_support*

Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
<b>Model</b>	2	4.3031870	2.1515935	0.92	0.4031
<b>Error</b>	97	227.5843130	2.3462300		
<b>Corrected Total</b>	99	231.8875000			

R-Square	Coeff Var	Root MSE	Tech_support Mean
0.018557	28.55062	1.531741	5.365000

Source	DF	Anova SS	Mean Square	F Value	Pr > F
<b>Customer_type</b>	2	4.30318696	2.15159348	0.92	0.4031

**Testing:** Means of the 3 customer types. Where 1= less than one year, 2 = between 1&5 years, 3= longer than 5 years.  $F(2,97)=0.92$

**H<sub>0</sub>:**  $\mu_1 = \mu_2 = \mu_3$

**H<sub>1</sub>:** Not all  $\mu_j$  are equal.

The P-Value of 0.4031 > than  $\alpha = 0.05$  Fail reject the null hypothesis. There is insufficient evidence that not all means are equal.

*The SAS System**The ANOVA Procedure*

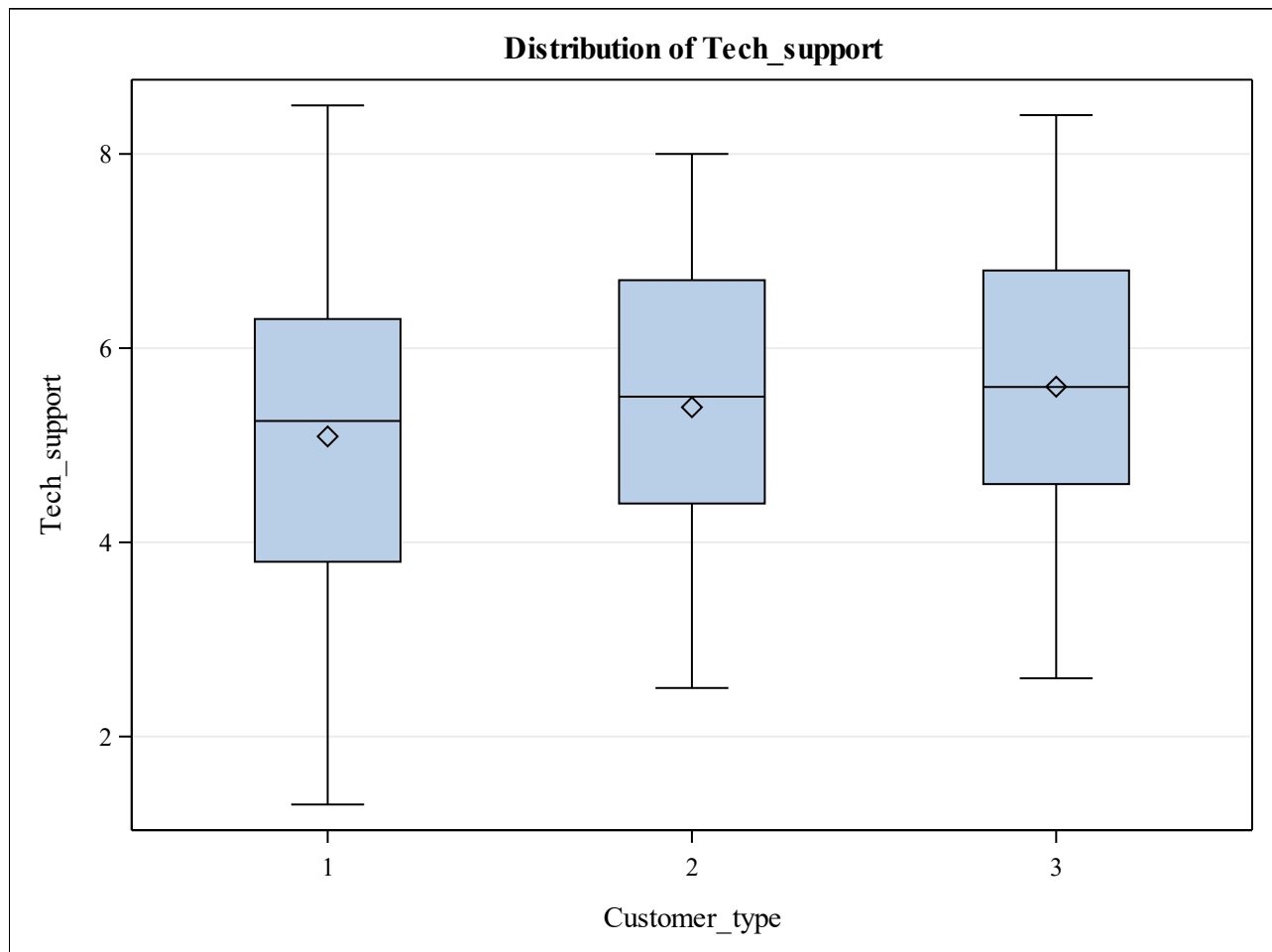
Levene's Test for Homogeneity of Tech_support Variance ANOVA of Squared Deviations from Group Means					
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Customer_type	2	10.4526	5.2263	0.73	0.4831
Error	97	691.5	7.1292		

**Levene's Test:** Equal Variance of the 3 customer types. Where 1= less than one year, 2 = between 1&5 years, 3= longer than 5 years.  $F(2,97)=0.73$

**H<sub>0</sub>:**  $\sigma_1^2 = \sigma_2^2 = \sigma_3^2$

**H<sub>1</sub>:** Not all  $\sigma_j^2$  are equal.

The P-Value of 0.4831 > than  $\alpha = 0.05$  Fail reject the null hypothesis. There is insufficient evidence that not all variances are equal. The assumption of equal variances appears to be valid.

*The SAS System**The ANOVA Procedure*

Level of Customer_type	N	Tech_support	
		Mean	Std Dev
1	32	5.09062500	1.67473818
2	35	5.39142857	1.50555275
3	33	5.60303030	1.40945132

A boxplot of the data is used to approximate normality for each of the 3 Customer\_types. The median is approximately centered, the whiskers are approximately equal, and the variability is similar on both sides of the IQR. There is a negative skew in number 1 and 2, with number 3 being nearly symmetric. Approximate normality is validated.

*The SAS System**The REG Procedure**Model: MODEL1**Dependent Variable: Tech\_support*

<b>Number of Observations Read</b>	100
<b>Number of Observations Used</b>	100

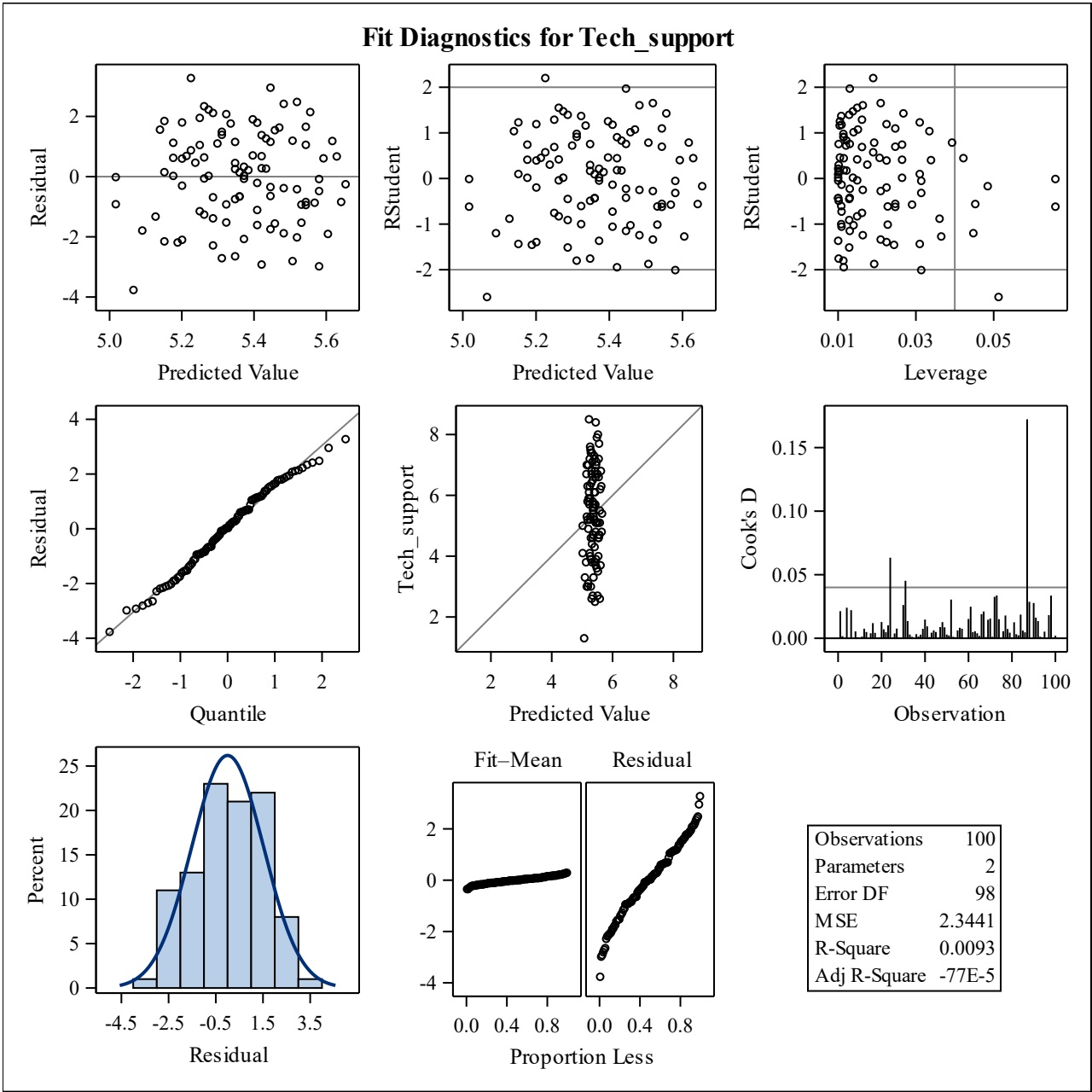
<b>Analysis of Variance</b>					
<b>Source</b>	<b>DF</b>	<b>Sum of Squares</b>	<b>Mean Square</b>	<b>F Value</b>	<b>Pr &gt; F</b>
<b>Model</b>	1	2.16641	2.16641	0.92	0.3387
<b>Error</b>	98	229.72109	2.34409		
<b>Corrected Total</b>	99	231.88750			

<b>Root MSE</b>	1.53104	<b>R-Square</b>	0.0093
<b>Dependent Mean</b>	5.36500	<b>Adj R-Sq</b>	-0.0008
<b>Coeff Var</b>	28.53761		

<b>Parameter Estimates</b>					
<b>Variable</b>	<b>DF</b>	<b>Parameter Estimate</b>	<b>Standard Error</b>	<b>t Value</b>	<b>Pr &gt;  t </b>
<b>Intercept</b>	1	4.69881	0.70969	6.62	<.0001
<b>Complaint</b>	1	0.12242	0.12734	0.96	0.3387

The SAS System

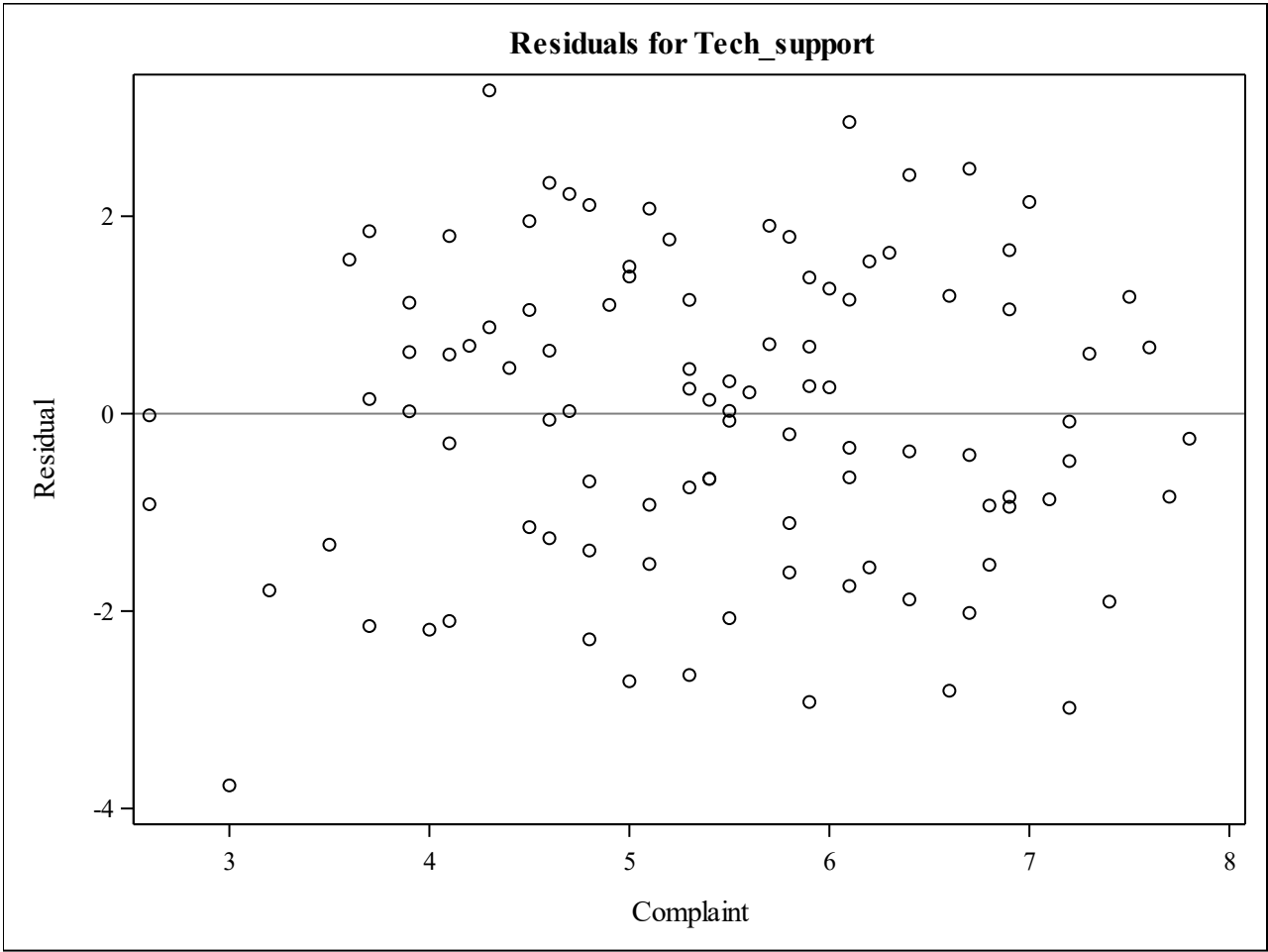
The REG Procedure  
Model: MODEL1  
Dependent Variable: Tech\_support



From this histogram, it appears the residuals are normal, validating the assumption of normality. Residual versus predicted shows two potential outliers and appears to be randomized with no noticeable pattern, suggesting homoscedasticity. Residual versus Quantile shows a positive linear relationship with a slope close to 1, also validating an assumption of normal distribution.

*The SAS System*

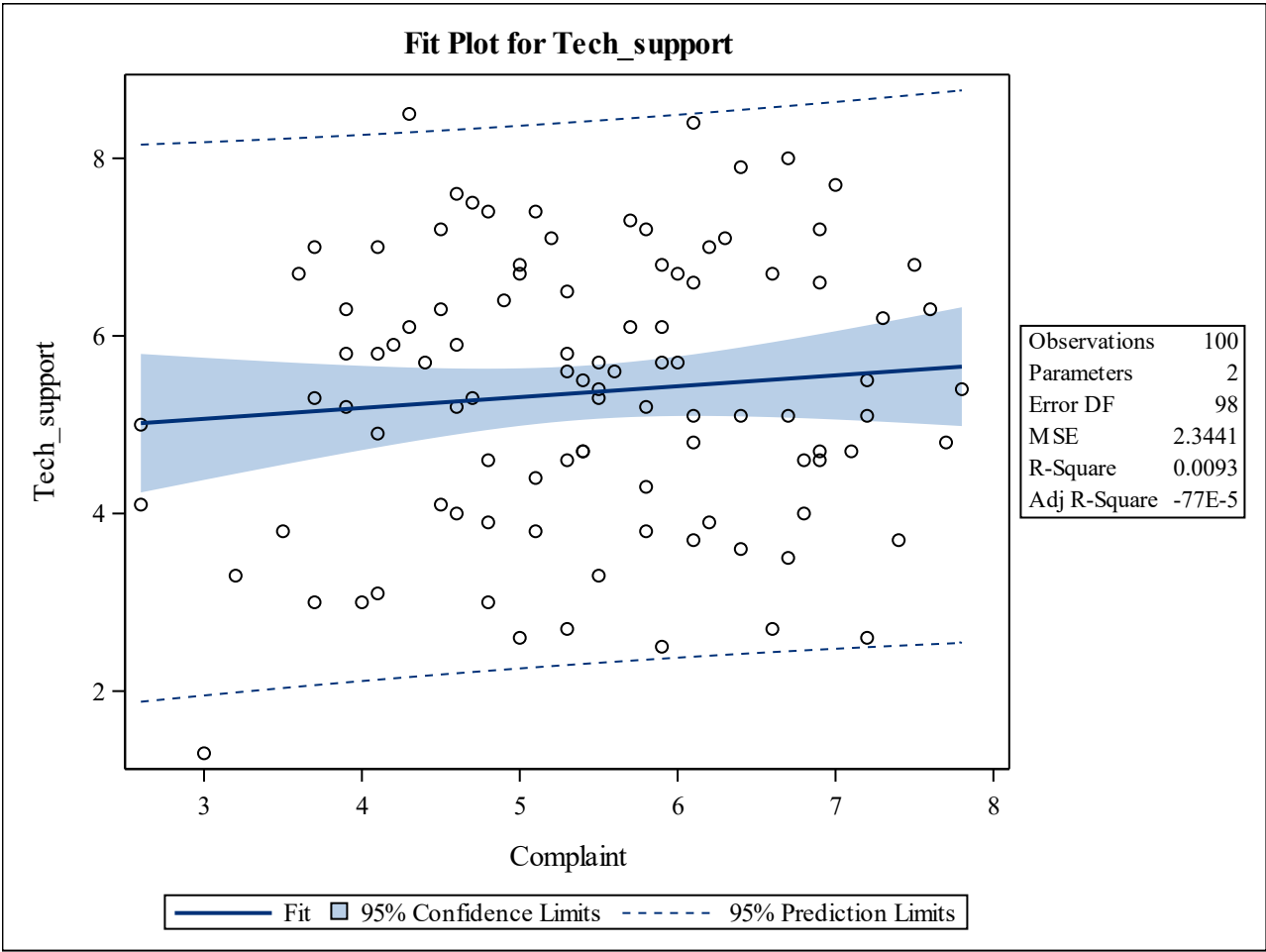
*The REG Procedure  
Model: MODEL1  
Dependent Variable: Tech\_support*



Residuals are normally distributed with about 68% of the residuals should be within 1 standard deviation of the mean of 0; 95% within 2 standard deviation, and about 99% within 3 standard deviations.

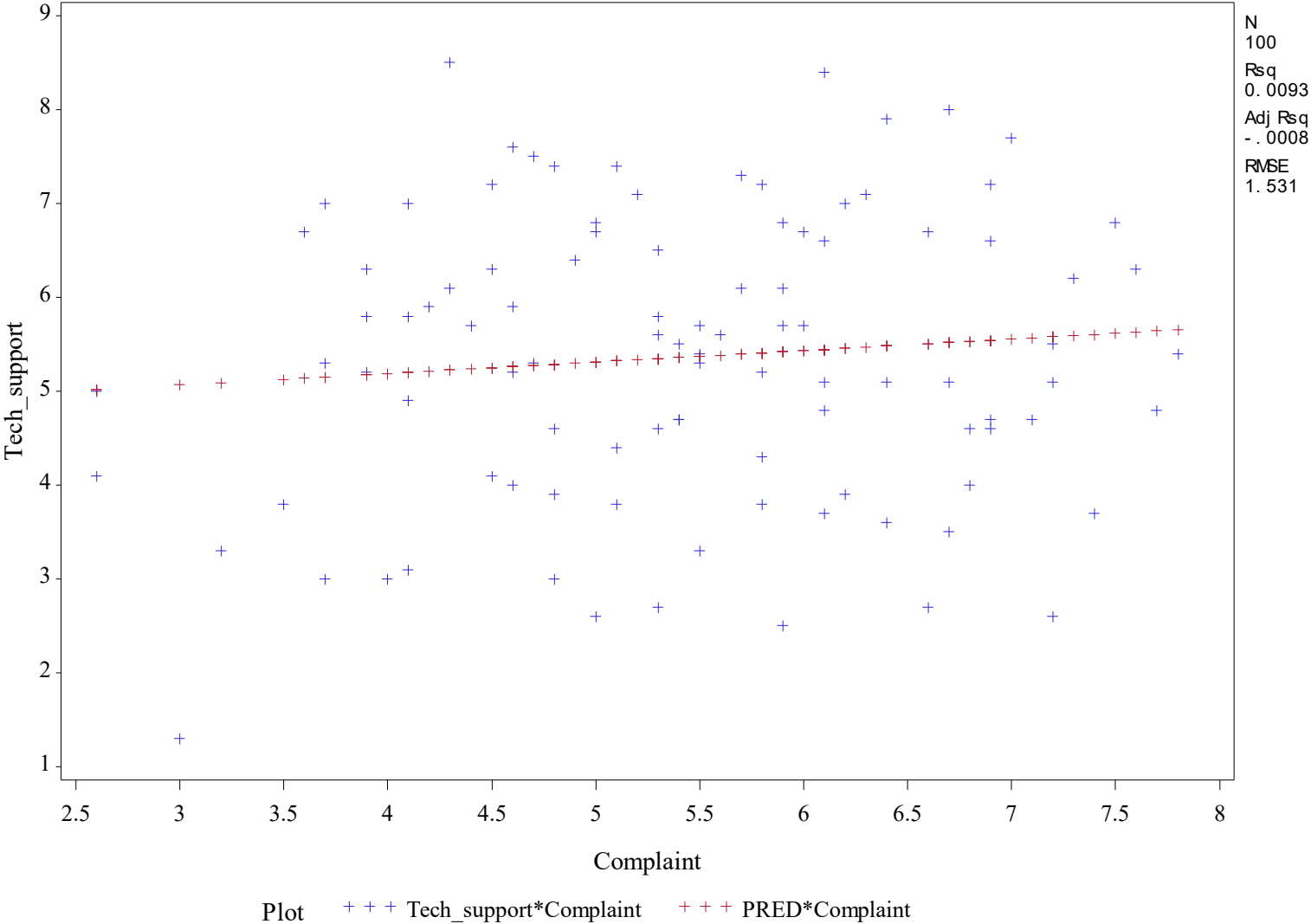
The SAS System

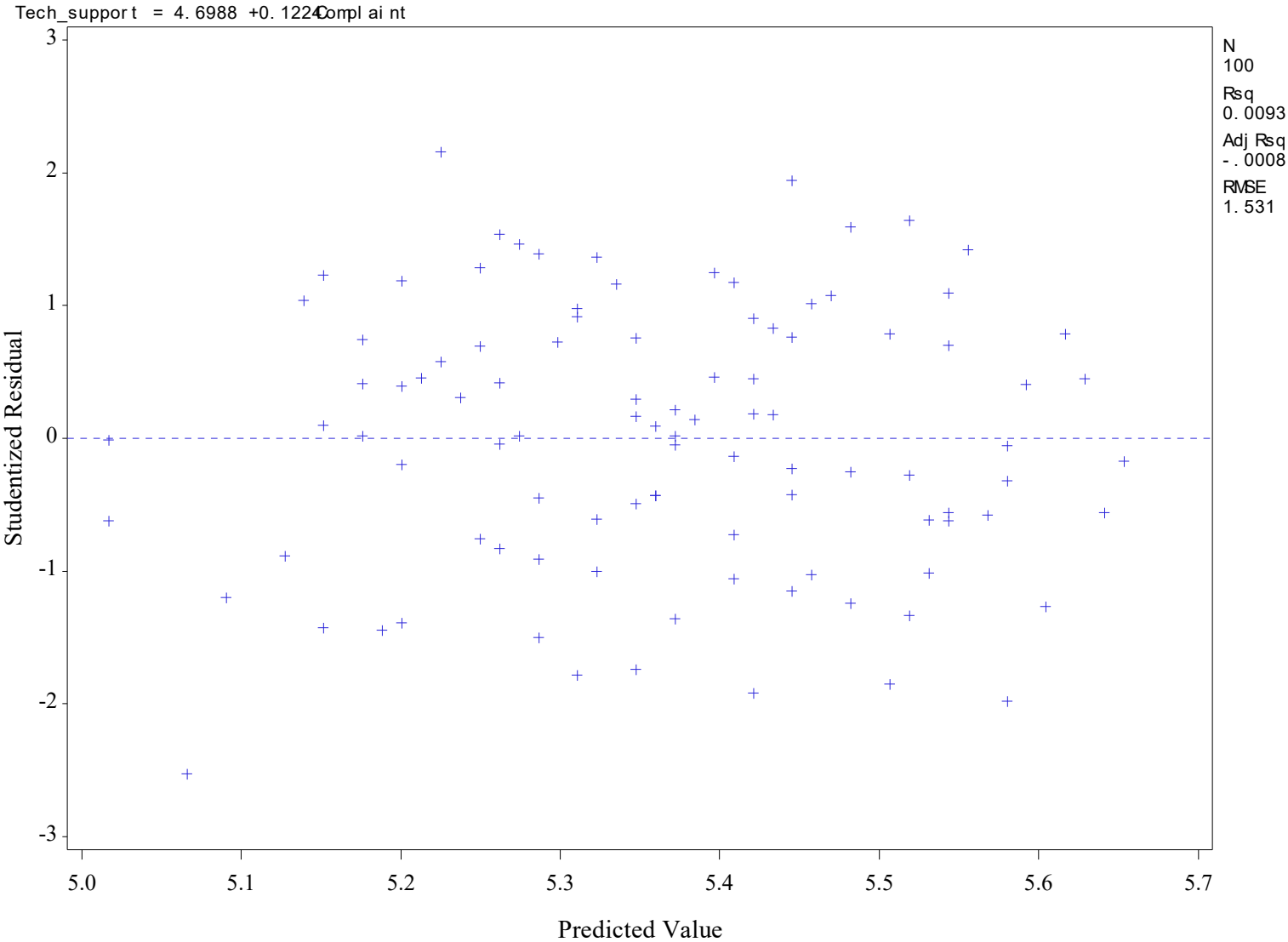
The REG Procedure  
Model: MODEL1  
Dependent Variable: Tech\_support



The fitted plot shows normality with two outliers.

Tech\_support = 4.6988 + 0.1224 \* Complaint





Residuals are homoscedastic, there is no apparent pattern. The distribution of the error terms is assumed to be normal, about 68% of the residuals should be within 1 standard deviation of the mean, about 95% within 2 standard deviations, and 99% within 3.

## The SAS System

### The UNIVARIATE Procedure

Variable: *stresid* (Studentized Residual)

Moments			
<b>N</b>	100	<b>Sum Weights</b>	100
<b>Mean</b>	-0.0010047	<b>Sum Observations</b>	-0.1004723
<b>Std Deviation</b>	1.00541852	<b>Variance</b>	1.0108664
<b>Skewness</b>	-0.1789027	<b>Kurtosis</b>	-0.6170686
<b>Uncorrected SS</b>	100.075875	<b>Corrected SS</b>	100.075774
<b>Coeff Variation</b>	-100069.27	<b>Std Error Mean</b>	0.10054185

Basic Statistical Measures			
Location		Variability	
<b>Mean</b>	-0.00100	<b>Std Deviation</b>	1.00542
<b>Median</b>	0.01765	<b>Variance</b>	1.01087
<b>Mode</b>	-0.43316	<b>Range</b>	4.68494
		<b>Interquartile Range</b>	1.44908

Tests for Location: $\mu_0=0$				
Test	Statistic		p Value	
<b>Student's t</b>	<b>t</b>	-0.00999	<b>Pr &gt;  t </b>	0.9920
<b>Sign</b>	<b>M</b>	2	<b>Pr &gt;=  M </b>	0.7644
<b>Signed Rank</b>	<b>S</b>	40	<b>Pr &gt;=  S </b>	0.8914

Tests for Normality				
Test	Statistic		p Value	
<b>Shapiro-Wilk</b>	<b>W</b>	0.987642	<b>Pr &lt; W</b>	0.4822
<b>Kolmogorov-Smirnov</b>	<b>D</b>	0.064553	<b>Pr &gt; D</b>	>0.1500
<b>Cramer-von Mises</b>	<b>W-Sq</b>	0.050063	<b>Pr &gt; W-Sq</b>	>0.2500
<b>Anderson-Darling</b>	<b>A-Sq</b>	0.344713	<b>Pr &gt; A-Sq</b>	>0.2500

#### Descriptive Values

Nearly symmetrical with a skewness value of -0.18.

A kurtosis value of -0.62 indicates that the distribution is slightly flat, but still a normal distribution.

The mean of -0.00 < 0.02 median also indicating a nearly symmetrical distribution.

#### Location Test

$H_0: \mu=0$   
 $H_1: \mu \neq 0$

$t(99) = -0.01$

A p-value 0.99 >  $\alpha$  of .05. Fail to reject the null. The mean is statistically zero.

#### Normality Test

$H_0$ : Data is normal.  
 $H_1$ : Data is not normal.

$W=0.9876$

A p-value of 0.48 >  $\alpha$  of .05.  
Fail to reject the null. The variable, Validating normal distribution.

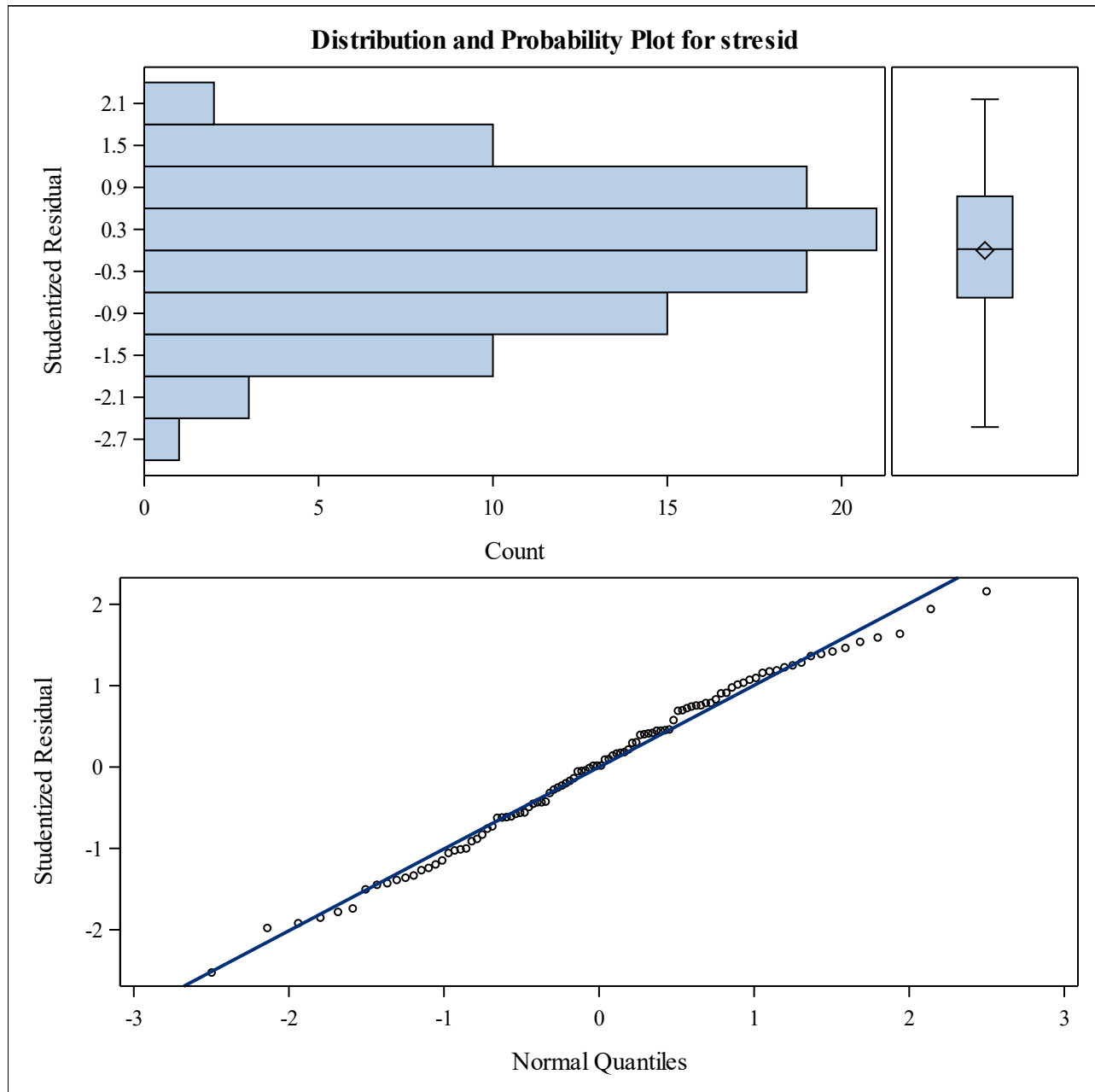
## *The SAS System*

### *The UNIVARIATE Procedure*

*Variable: stresid (Studentized Residual)*

<b>Quantiles (Definition 5)</b>	
<b>Level</b>	<b>Quantile</b>
<b>100% Max</b>	2.1595716
<b>99%</b>	2.0509654
<b>95%</b>	1.5012829
<b>90%</b>	1.2670113
<b>75% Q3</b>	0.7729814
<b>50% Median</b>	0.0176527
<b>25% Q1</b>	-0.6760975
<b>10%</b>	-1.3739895
<b>5%</b>	-1.7594387
<b>1%</b>	-2.2515815
<b>0% Min</b>	-2.5253637

<b>Extreme Observations</b>			
<b>Lowest</b>		<b>Highest</b>	
<b>Value</b>	<b>Obs</b>	<b>Value</b>	<b>Obs</b>
-2.52536	87	1.53862	77
-1.97780	24	1.59220	67
-1.91891	1	1.63771	88
-1.85116	73	1.94236	61
-1.78075	97	2.15957	31

*The SAS System**The ANOVA Procedure*

There is a positive linear relationship in the normal probability plot of studentized residuals. The skew is nearly symmetrical, and the data appears normally distributed. The whiskers of the box plot are similar in length, with the mean and median aligned. The histogram validates a normal distribution. There are no outliers in the box plot.