

```

title 'Cyberloafing, Mike Sage'; run;   ODS GRAPHICS ON;
proc univariate plot data=Sage; var Cyberloafing Conscientiousness Age;
Histogram / normal;   PPPLot / normal; run;

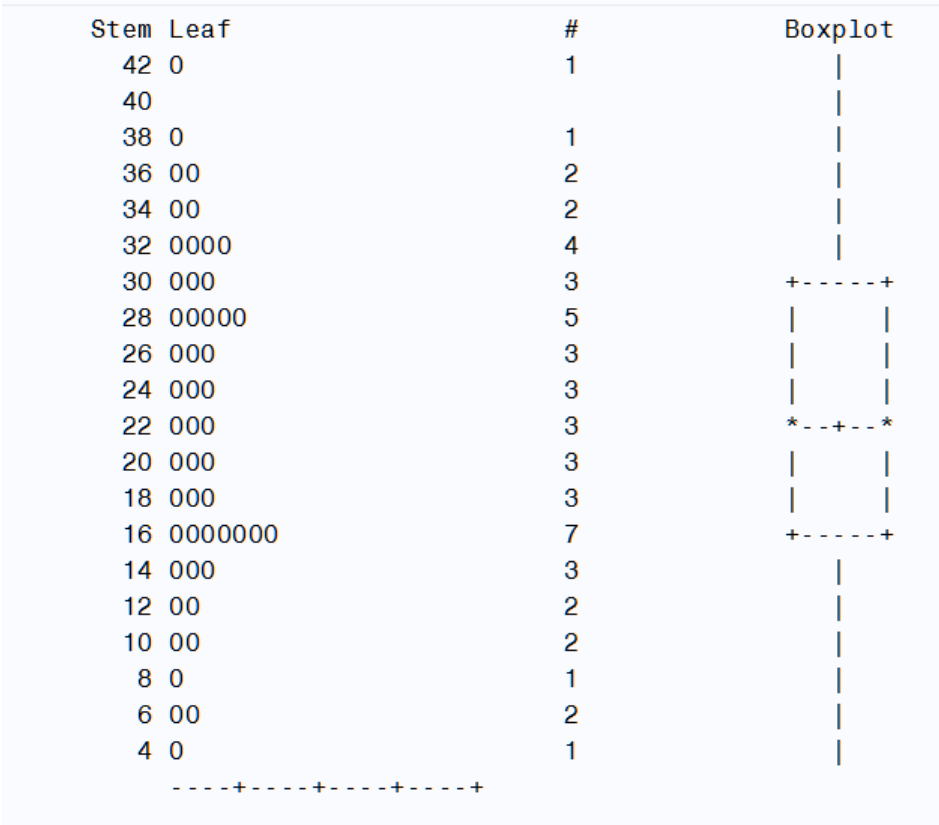
```

The UNIVARIATE Procedure
Variable: **Cyberloafing** (Cyberloafing)

Moments	
Skewness	0.008039
Kurtosis	-0.6908018

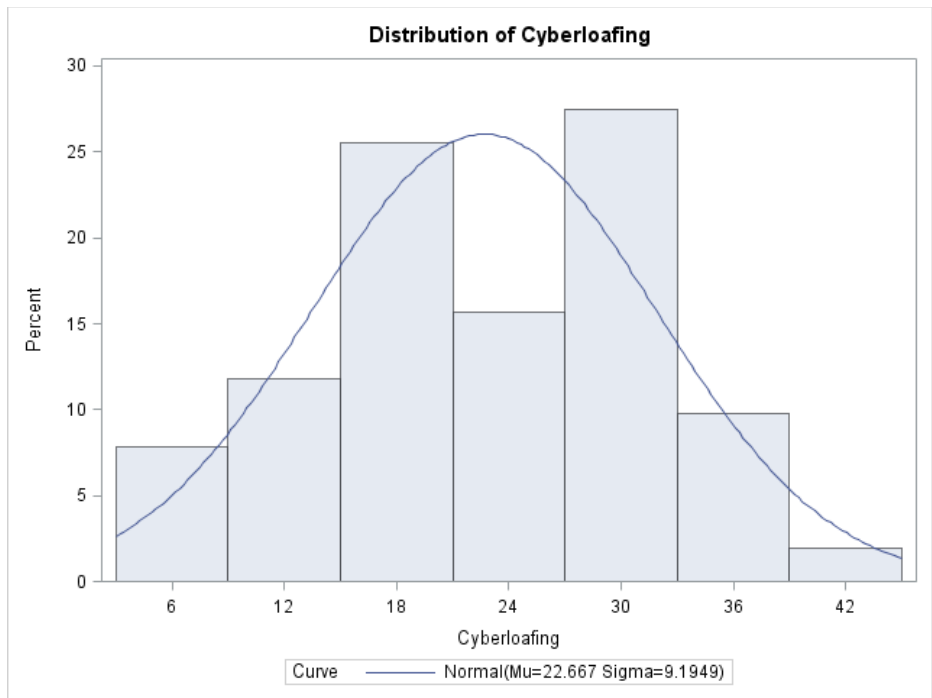
Basic Statistical Measures			
Location		Variability	
Mean	22.66667	Std Deviation	9.19493
Median	23.00000	Variance	84.54667

Extreme Observations			
Lowest		Highest	
Value	Obs	Value	Obs
4	37	35	25
6	38	37	5
7	29	37	14
8	50	38	19
11	49	43	3



I have not provided all of the output here but rather the most important parts. The plot option caused SAS to produce the above stem and leaf plot and boxplot, using old-fashioned methods that were designed for output to be used on a line printer. SAS also produced a normal probability plot in the old fashioned style, but I have not included it in this document. The stem and leaf plot and the boxplot indicate that the data are approximately normally distributed.

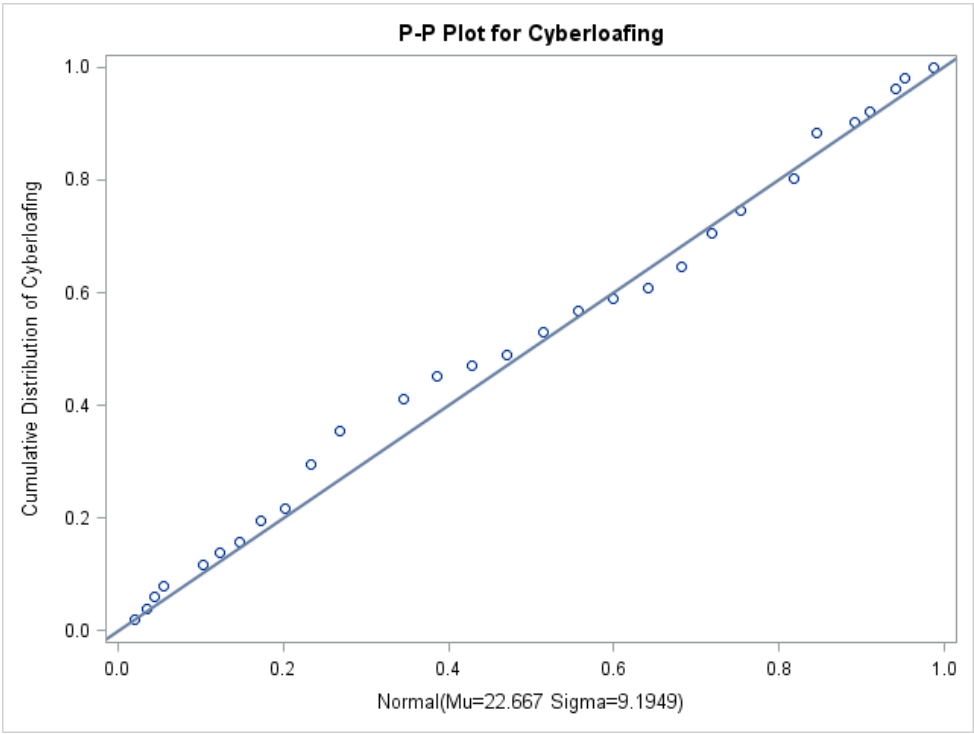
“Histogram / normal” produced a graphic showing a histogram of the observed scores with a overlaid curve of a normal distribution with the same mean and standard deviation as the observed scores. The scores in our sample appear to be close to being normally distributed.



Goodness-of-Fit Tests for Normal Distribution			
Test	Statistic		p Value
Kolmogorov-Smirnov D	0.08408694	Pr > D	>0.150
Cramer-von Mises	W-Sq 0.05210904	Pr > W-Sq	>0.250
Anderson-Darling	A-Sq 0.29443338	Pr > A-Sq	>0.250

“PPPLot / normal;” caused SAS to produce table above and the graphical normal probability plot below. The table includes the results of three tests of the null hypothesis that the data came from a population that is normally distributed. I do not endorse the use of these tests for determining whether or not a normality assumption has been violated. With large sample sizes these tests will have so much power that they will produce “significant” results even when the deviation from normality is too small about which to worry, especially given that with large sample sizes correlation and regression analyses are more robust to violations of the normality assumption. With small samples sizes robustness is less and the goodness-of-fit tests may have too little power to detect a deviation from normality that is large enough to be a problem.

In the plot below, if the data are normally distributed then the line of dots will not deviate much from the solid reference line. Again, the plot indicates that the data are close to normally distributed.



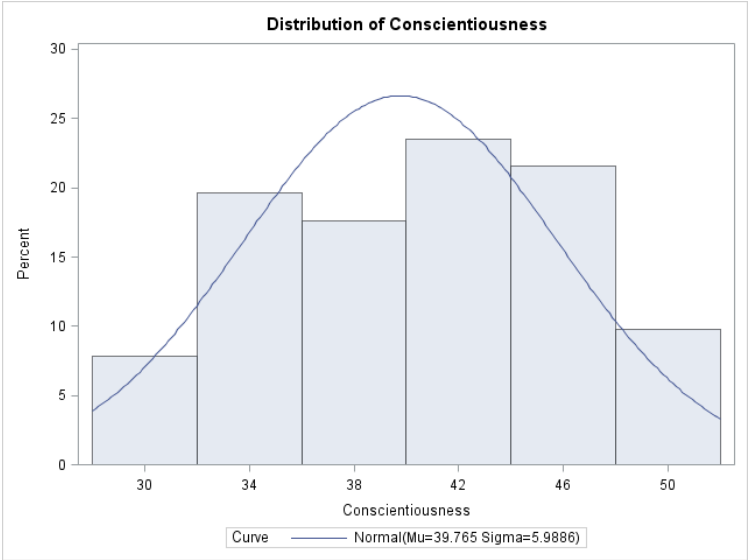
Cyberloafing, Mike Sage

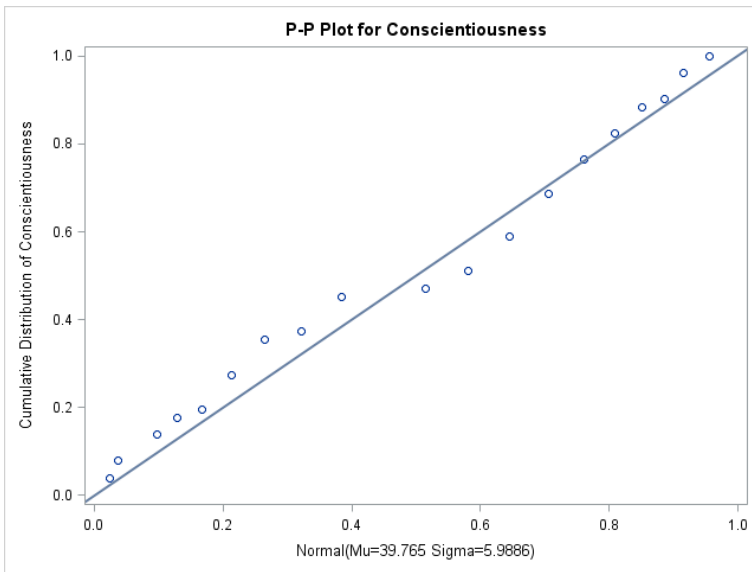
The UNIVARIATE Procedure
 Variable: **Conscientiousness** (Conscientiousness)

Moments
Skewness -0.2692084 **Kurtosis** -0.8823586

Basic Statistical Measures

Location		Variability	
Mean	39.76471	Std Deviation	5.98862
Median	41.00000	Variance	35.86353





The conscientiousness scores are close to normal in their distribution.

Cyberloafing, Mike Sage

The UNIVARIATE Procedure
Variable: Age (Age)

Moments

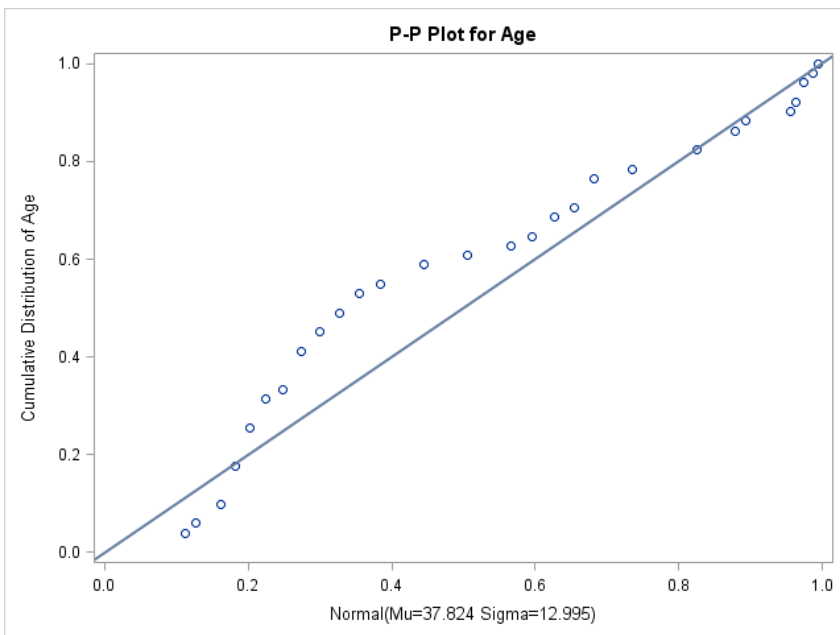
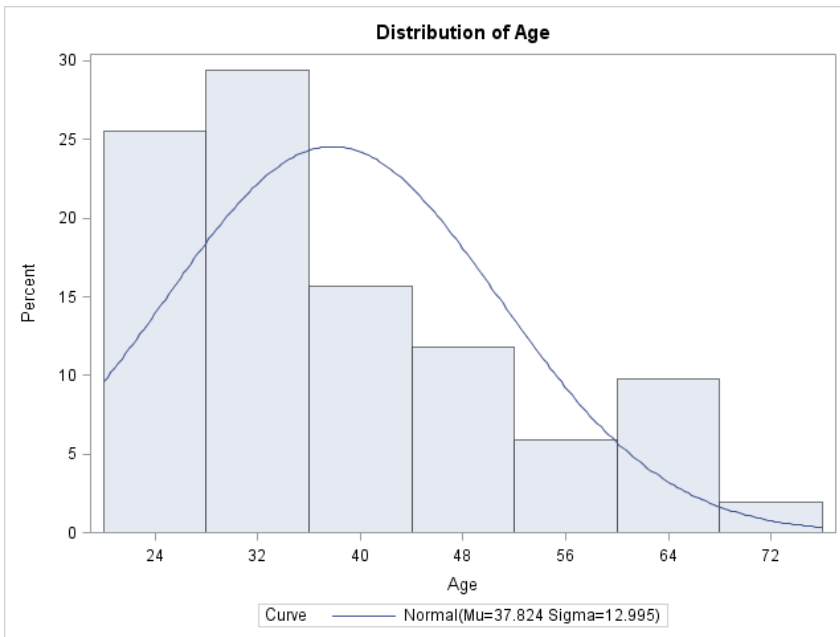
Skewness 0.94090779 **Kurtosis** -0.0914493

Basic Statistical Measures

	Location		Variability
Mean	37.82353	Std Deviation	12.99493
Median	33.00000	Variance	168.86824

Extreme Observations

	Lowest		Highest
	Value	Obs	Value
	22	5	61
	22	3	63
	23	27	63
	25	37	67
	25	4	71



The plots for age reveal a distinct positive skewness. When the absolute value of skewness (g_1) exceeds one, I usually am concerned enough to try transforming the variable to reduce the skewness. Transformations commonly used for that purpose include logarithmic, square roots, negative inverses, and ranks. The table below shows that a log transformation does a good job at reducing the skewness in age.

```
data transform; set Sage;
SR_Age = SQRT(Age); Log_Age = Log10(Age);
Proc Means skewness kurtosis; Var Age SR_Age Log_Age; run;
```

Variable	Label	Skewness	Kurtosis
Age	Age	0.9409078	-0.0914493
SR_Age		0.7178151	-0.5312523
Log_Age		0.4972690	-0.8392886

PROC CORR data=Sage; var Cyberloafing Conscientiousness Age;
run; quit;

The CORR Procedure

3 Variables: Cyberloafing Conscientiousness Age

Simple Statistics							
Variable	N	Mean	Std Dev	Sum	Minimum	Maximum	Label
Cyberloafing	51	22.66667	9.19493	1156	4.00000	43.00000	Cyberloafing
Conscientiousness	51	39.76471	5.98862	2028	28.00000	50.00000	Conscientiousness
Age	51	37.82353	12.99493	1929	22.00000	71.00000	Age

Pearson Correlation Coefficients, N = 51 Prob > r under H0: Rho=0			
	Cyberloafing	Conscientiousness	Age
Cyberloafing	1.00000	-0.56297	-0.46197
Cyberloafing		<.0001	0.0006
Conscientiousness	-0.56297	1.00000	0.14286
Conscientiousness		<.0001	0.3173
Age	-0.46197	0.14286	1.00000
Age		0.0006	0.3173

Using Cohen's guidelines, the correlations between cyberloafing and both conscientiousness and age are large. Those guidelines are .1 = small, .3 = medium, and .5 = large.

*SAS does not give you the value of t here. The value is $t = \frac{.56297\sqrt{49}}{\sqrt{1-.56297^2}} = 4.77$.

I used the calculator at Vassar to get the confidence interval.

r =	<input type="text" value="-0.56297"/>	<input type="button" value="Reset"/>
n =	<input type="text" value="51"/>	<input type="button" value="Calculate"/>

0.95 and 0.99 Confidence Intervals of rho

	Lower Limit	Upper Limit
0.95	-0.725	-0.341
0.99	-0.765	-0.26

PROC REG data=Sage; **A:** model Cyberloafing = Conscientiousness;
B: model Cyberloafing = Conscientiousness age / **scorr2 pcorr2 stb vif; run; quit;**

Here is a bivariate regression, predicting cyberloafing from Conscientiousness.

The REG Procedure

Model: A

Dependent Variable: Cyberloafing Cyberloafing

Number of Observations Read 51

Number of Observations Used 51

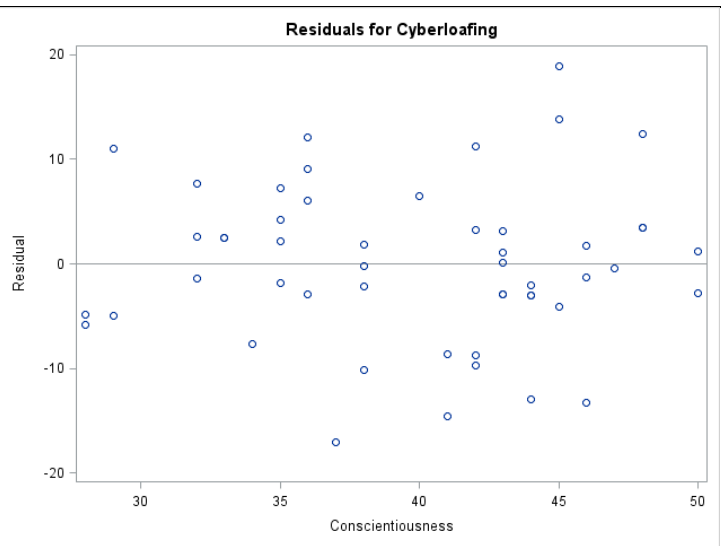
Analysis of Variance					
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	1	1339.80121	1339.80121	22.74	<.0001
Error	49	2887.53213	58.92923		
Corrected Total	50	4227.33333			

Root MSE	7.67654	R-Square	0.3169
Dependent Mean	22.66667	Adj R-Sq	0.3030
Coeff Var	33.86708		

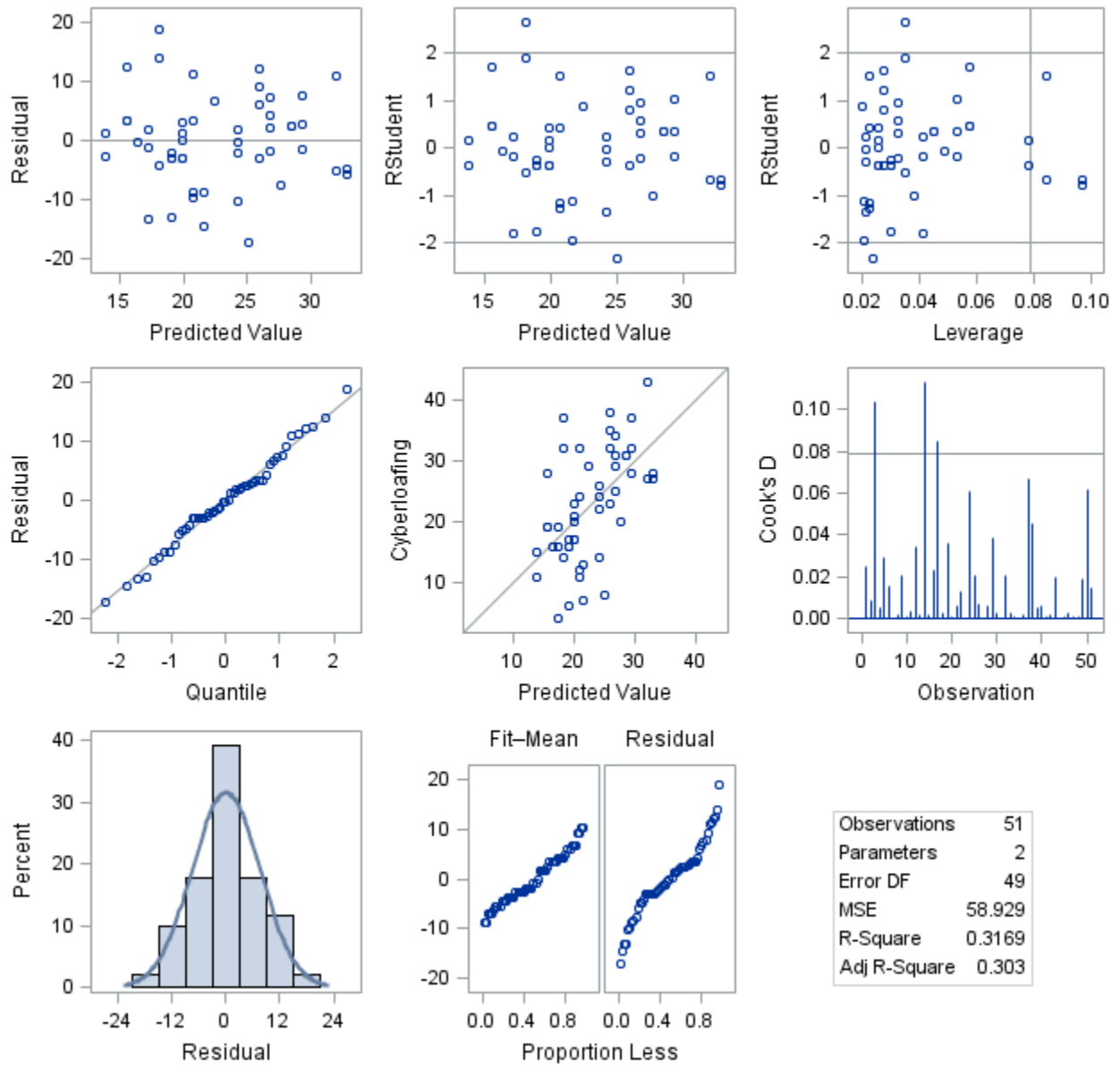
Parameter Estimates						
Variable	Label	DF	Parameter Estimate	Standard Error	t Value	Pr > t
Intercept	Intercept	1	57.03880	7.28832	7.83	<.0001
Conscientiousness	Conscientiousness	1	-0.86439	0.18128	-4.77	<.0001

Notice that the t reported here is that I computed above, by hand, and the F is the square of that t .

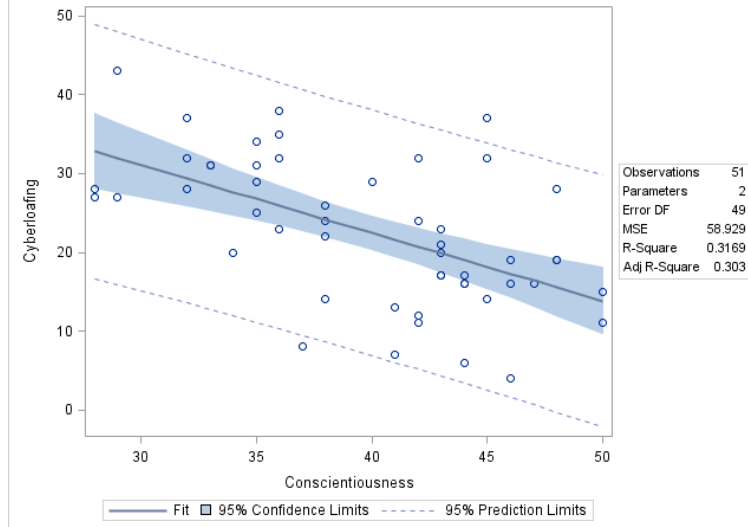
When using t or F with a regression analysis, we assume that the residuals (score minus predicted score) are normally distributed at every value of predicted score and that the variance in the residuals is constant across levels of the predicted score. The plots to the right and below help you evaluate those assumptions as well as identify outliers that may be having great influence on the solution.



Fit Diagnostics for Cyberloafing



Fit Plot for Cyberloafing



The fit plot shows the bivariate scores, the regression line, and the confidence limits. If the relationship were nonlinear, that would be evident here.

B: model Cyberloafing = Conscientiousness age / scorr2 pcorr2 stb vif; run; quit;

Now we add a second predictor, age.

Analysis of Variance					
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	2	1968.02919	984.01459	20.91	<.0001
Error	48	2259.30414	47.06884		
Corrected Total	50	4227.33333			

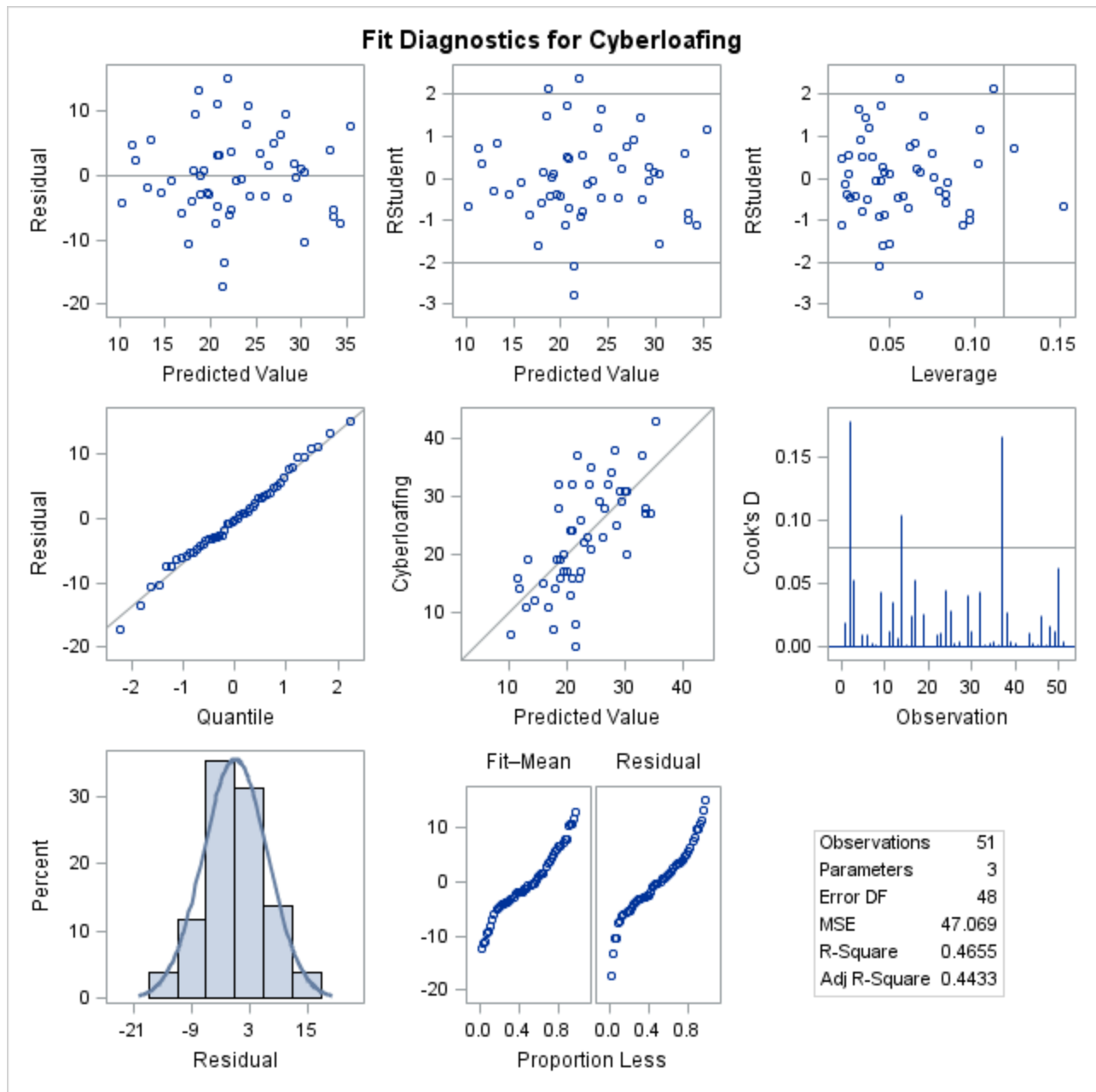
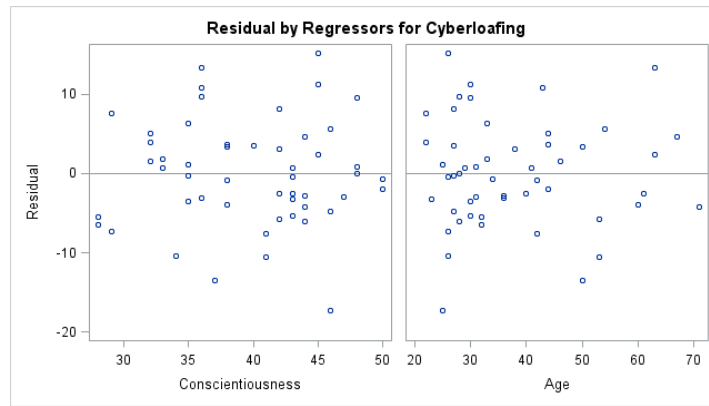
Root MSE	6.86067	R-Square	0.4655
Dependent Mean	22.66667	Adj R-Sq	0.4433
Coeff Var	30.26768		

Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t	Standardized Estimate	Squared Semi-partial Corr Type II	Squared Partial Corr Type II	Variance Inflation
Intercept	1	64.06561	6.79175	9.43	<.0001	0	.	.	0
Conscientiousness	1	-0.77895	0.16369	-4.76	<.0001	-0.50733	0.25213	0.32054	1.02083
Age	1	-0.27560	0.07544	-3.65	0.0006	-0.38950	0.14861	0.21757	1.02083

When we had only one predictor, Conscientiousness, the r^2 was .3169. Adding age to the model increased the R^2 by .4655 - .3169 = .1486; this is the value of the squared semipartial correlation coefficient. That increase is significant, $t(48) = 3.65$, $p < .001$.

Magnitude of the Semipartial Correlations. Some people apply Cohen's guidelines for rho to beta weights, but it is better to use the semipartial correlations. Taking square roots to get the semipartials here, the sr for conscientiousness is .504, a large effect, and the sr for age is .385, a medium to large effect.

Multicollinearity exists when a predictor can be nearly perfectly predict from a weighted linear combination of the other predictors – that is, when $R^2_{i \cdot 12 \dots (i) \dots p}$ is very large. In that case, the partial statistics would be unstable, that is, they would tend to vary wildly among samples drawn from the same population. The usual solution here is to drop variables from the model to eliminate the problem with multicollinearity. The **tolerance** statistic is computed as $1 - R^2_{i \cdot 12 \dots (i) \dots p}$. So multicollinearity is present when a tolerance is very low. The **variance inflation factor (VIF)** is computed as $1/\text{tolerance}$, so high values of VIF indicate a problem. So, how high must VIF get before we get worried? Some say 10, some say 5, and a few say 2.5. Clearly we have no problem with these data.



Notice that the residuals plot here is residual (score value minus predicted score value) versus predicted score. This plot is used to evaluate the homoscedasticity assumption and the assumption that the residuals are normally distributed at every value of the predicted scores. The residuals histogram compares the observed residuals with those that would be expected if they are normally distributed. There are no problems apparent here. The leverage and Cook's D plots are used to detect outliers. We shall discuss these statistics later, when I cover regression diagnostics.

```
PROC GLM data=Sage; model Cyberloafing = Conscientiousness / EFFECTSIZE ALPHA=.1; run;
quit;
```

The GLM Procedure

PROC GLM was used here to get a confidence interval for ρ^2 predicting cyberloafing from conscientiousness.

Proportion of Variation Accounted for	
Eta-Square	0.32
90% Confidence Limits	(0.14,0.46)

And here for the confidence interval for ρ^2 predicting cyberloafing from conscientiousness and age and confidence interval for the unique effects of each predictor. GLM calls the squared semipartial correlation coefficient the semipartial eta-squared. I'll explain later why this statistic is preferable to the squared partial correlation coefficient.

```
PROC GLM data=Sage; model Cyberloafing = Conscientiousness Age / EFFECTSIZE ALPHA=.1;
*If the leading * were removed, the below statement would compute confidence intervals for predicted
Cyberloafing for each subject;
*print clm cli;
run; quit;
```

The GLM Procedure

Proportion of Variation Accounted for	
Eta-Square	0.47
90% Confidence Limits	(0.27,0.58)

Type III (unique) Effect Sizes

Source	DF	Partial Variation Accounted For					
		Semipartial Eta-Square	Conservative 90% Confidence Limits		Partial Eta-Square	90% Confidence Limits	
Conscientiousness	1	0.2521	0.0916	0.3993	0.3205	0.1423	0.4580
Age	1	0.1486	0.0272	0.2955	0.2176	0.0649	0.3630

I also used my SAS program to get the confidence interval for R^2 -- [Conf-Interval-R2-Regr.sas](#)

Compute 90% Confidence Interval for R-squared, eta-squared, fixed effect ANOVA/regression

Obs	eta_squared	eta2_lower	eta2_upper
1	0.46560	0.27247	0.57713

Presenting the Results, APA-Style

Multiple correlation/regression analysis was used to predict participants' cyberloading scores from their scores on conscientiousness and age. The data were screened for possible problems with the assumptions of homoscedasticity and normality. No problems were found. The optimally weighted linear combination of conscientiousness and age was significantly associated with cyberloafing, $R^2 = .466$, $F(2, 48) = 20.91$, $p < .001$, 90% CI [.272, .577]. The partial effects of both predictors were significant. For conscientiousness, $t(48) = 4.76$, $p < .001$, $\beta = -.507$, $s^2 = .252$, 90% CI [.092, .399], and for age, $t(48) = 3.65$, $p < .001$, $\beta = -.390$, $s^2 = .149$, 90% CI [.065, .363].

```
*****;
title 'Importance of Plotting Your Data'; run;
*From PSYPARC gopher, Phil Wood, modified by K. Wuensch;
*Originally from Anscombe (1973), American Statistician, pp 17-21;

data PW; input x1 y1 x2 y2 x3 y3 x4 y4; cards;
10 8.04      10 9.14      10 7.46      8 6.58
 8 6.95      8 8.14      8 6.77      8 5.76
13 7.58      13 8.74      13 12.74     8 7.7
 9 8.81      9 8.77      9 7.11      8 8.84
11 8.33      11 9.26      11 7.81      8 8.47
14 9.96      14 8.10      14 8.84      8 7.04
 6 7.24      6 6.13      6 6.08      8 5.25
 4 4.26      4 3.10      4 5.39      19 12.50
12 10.84     12 9.13      12 8.15      8 5.56
 7 4.82      7 7.26      7 6.42      8 7.91
 5 5.68      5 4.74      5 5.73      8 6.89
;
proc reg simple; A: model y1 = x1;
  B: model y2 = x2;;
  C: model y3 = x3;
  D: model y4 = x4; run; quit;
```

Importance of Plotting Your Data

Descriptive Statistics						
Variable	Sum	Mean	Uncorrected SS	Variance	Standard Deviation	
Intercept	11.00000	1.00000	11.00000	0	0	
x1	99.00000	9.00000	1001.00000	11.00000	3.31662	
y1	82.51000	7.50091	660.17270	4.12727	2.03157	
x2	99.00000	9.00000	1001.00000	11.00000	3.31662	
y2	82.51000	7.50091	660.17630	4.12763	2.03166	
x3	99.00000	9.00000	1001.00000	11.00000	3.31662	
y3	82.50000	7.50000	659.97620	4.12262	2.03042	
x4	99.00000	9.00000	1001.00000	11.00000	3.31662	
y4	82.50000	7.50000	659.97840	4.12284	2.03048	

Notice that for each of the four X,Y sets, the means are 99 and 72.5 and the standard deviations are 3.3 and 2.0.

The REG Procedure
Model: A
Dependent Variable: y1

Analysis of Variance					
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	1	27.51000	27.51000	17.99	0.0022
Error	9	13.76269	1.52919		
Corrected Total	10	41.27269			

Root MSE 1.23660 R-Square 0.6665
Dependent Mean 7.50091 Adj R-Sq 0.6295

Parameter Estimates					
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t
Intercept	1	3.00009	1.12475	2.67	0.0257
x1	1	0.50009	0.11791	4.24	0.0022

The REG Procedure
Model: B
Dependent Variable: y2

Analysis of Variance					
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	1	27.50000	27.50000	17.97	0.0022
Error	9	13.77629	1.53070		
Corrected Total	10	41.27629			

Root MSE 1.23721 R-Square 0.6662

Parameter Estimates					
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t
Intercept	1	3.00091	1.12530	2.67	0.0258
x2	1	0.50000	0.11796	4.24	0.0022

The REG Procedure
 Model: C
 Dependent Variable: y3

Analysis of Variance					
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	1	27.47001	27.47001	17.97	0.0022
Error	9	13.75619	1.52847		
Corrected Total	10	41.22620			

Root MSE	1.23631	R-Square	0.6663
Dependent Mean	7.50000	Adj R-Sq	0.6292

Parameter Estimates					
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t
Intercept	1	3.00245	1.12448	2.67	0.0256
x3	1	0.49973	0.11788	4.24	0.0022

The REG Procedure
 Model: D
 Dependent Variable: y4

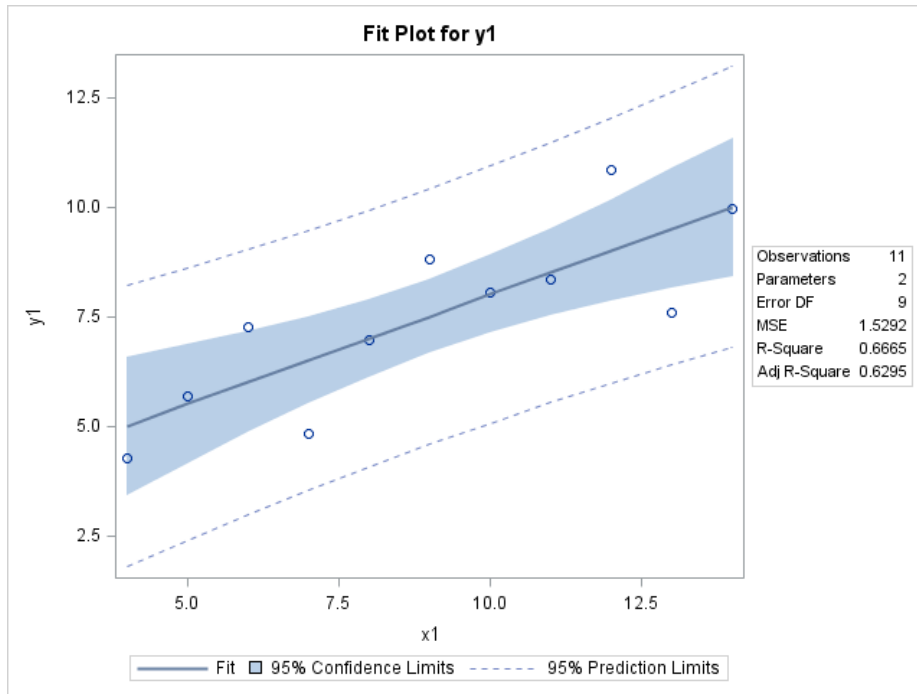
Analysis of Variance					
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	1	27.50000	27.50000	18.03	0.0022
Error	9	13.72840	1.52538		
Corrected Total	10	41.22840			

Root MSE	1.23506	R-Square	0.6670
Dependent Mean	7.50000	Adj R-Sq	0.6300

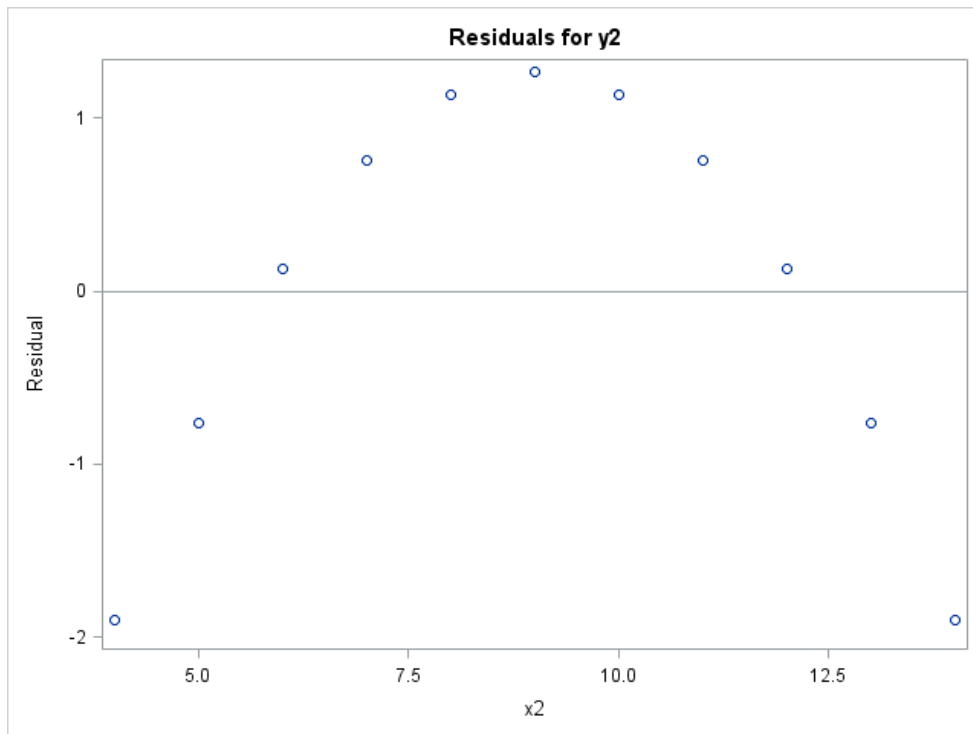
Parameter Estimates					
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t
Intercept	1	3.00000	1.12334	2.67	0.0256
x4	1	0.50000	0.11776	4.25	0.0022

Notice that for each of the X,Y pairs, the r^2 is .67, the intercept is 3, and the slope is .5.

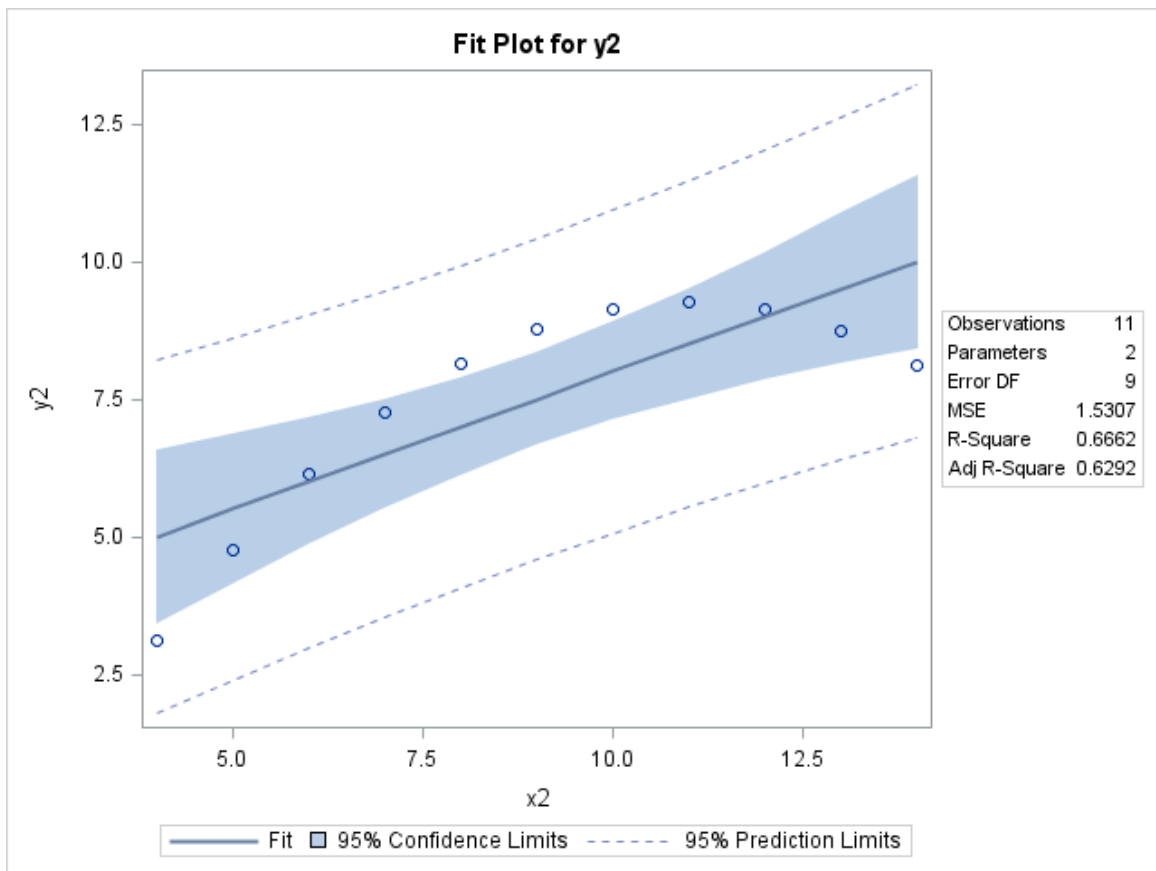
Below is the sort of plot that most researchers would anticipate for an r^2 of .67.



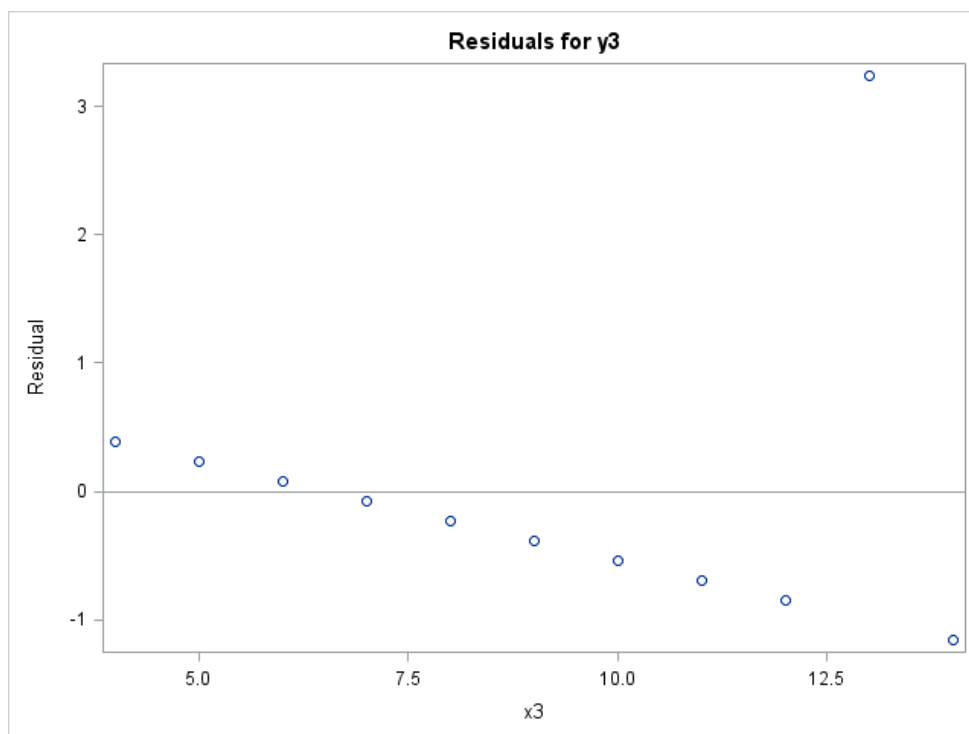
The residuals plot for X2, Y2 is quite revealing. It shows that there is a quadratic effect not included in the model.



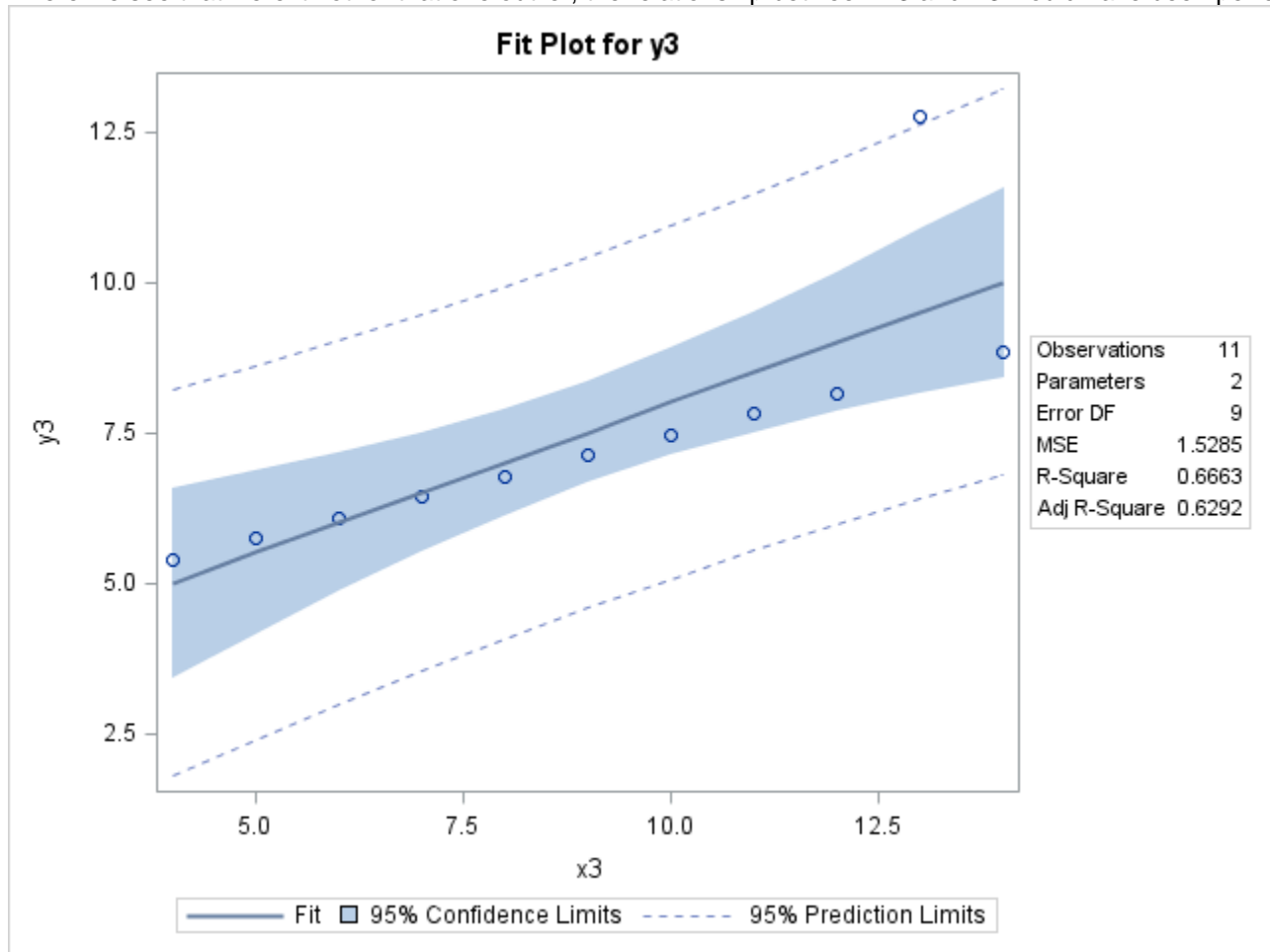
Here you see that the relationship between X2 and Y2 is actually perfect, but it is not linear.



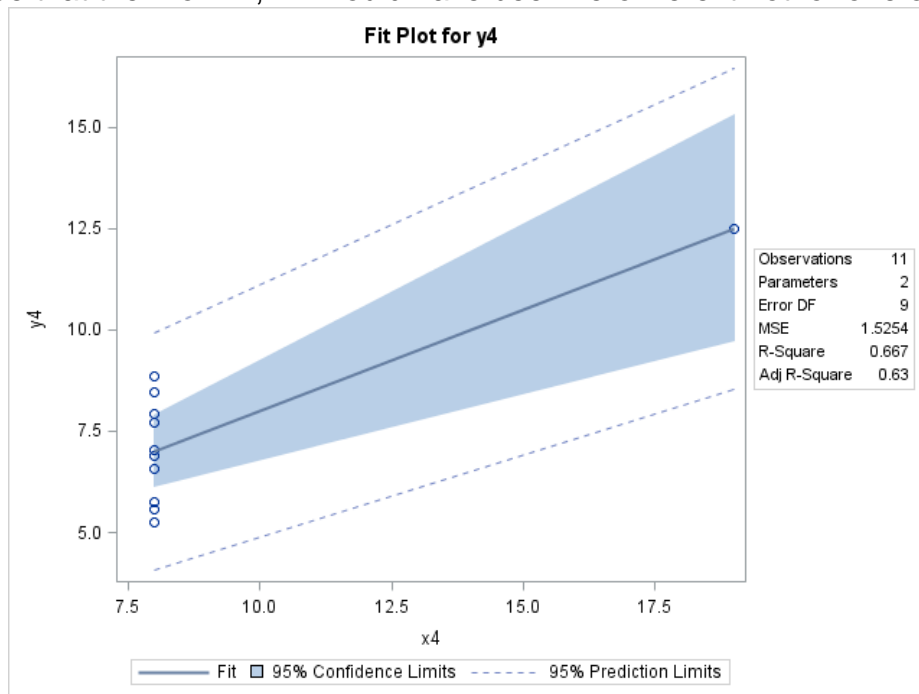
The residuals plot for X3, Y3 is also quite revealing. One observation has a very high residual. Such a case is going to have a large effect on the solution, pulling the regression line toward itself.



Here we see that were it not for that one outlier, the relationship between X3 and Y3 would have been perfect.



Here we see that the r for X4, Y4 would have been zero were it not for one outlier



Moderation Analysis – is the relationship between Y and X₁ modified by the value of X₂?

*****;

```

title 'Ar-Misanth Relationship for Nonidealists versus Idealists.'; run;
proc format; value l 0='NonIdealist' 1='Idealist';
data kevin2; infile 'C:\Users\Vati\Documents\StatData\potthoff.dat'; input ar misanth idealism;
format idealism l.;
proc sort; by idealism;
proc reg simple corr; model ar=misanth;
by idealism; run; quit;
    
```

The REG Procedure
idealism=NonIdealist

Descriptive Statistics					
Variable	Sum	Mean	Uncorrected SS	Variance	Standard Deviation
Intercept	91.00000	1.00000	91.00000	0	0
misanth	216.20000	2.37582	554.44000	0.45319	0.67319
ar	212.82100	2.33869	525.45037	0.30808	0.55505

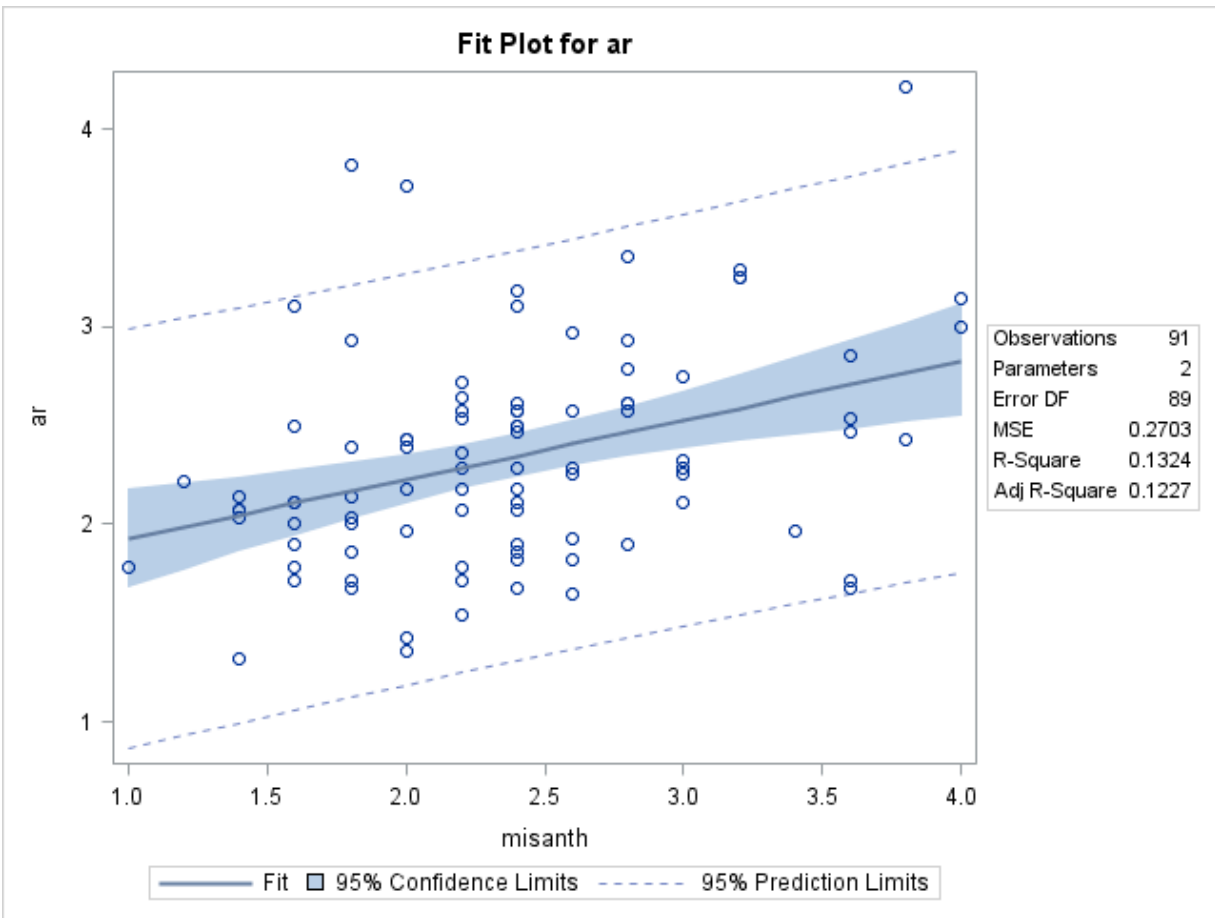
Correlation		
Variable	misanth	ar
misanth	1.0000	0.3639
ar	0.3639	1.0000

Analysis of Variance					
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	1	3.67218	3.67218	13.59	0.0004
Error	89	24.05535	0.27028		
Corrected Total	90	27.72753			

Root MSE	0.51989	R-Square	0.1324
Dependent Mean	2.33869	Adj R-Sq	0.1227
Coeff Var	22.22991		

Parameter Estimates					
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t
Intercept	1	1.62581	0.20094	8.09	<.0001
misanth	1	0.30006	0.08140	3.69	0.0004

Among non-idealists, misanthropy is significantly associated with support for animal rights, $r(n = 91) = .36, p < .001, 95\% \text{ CI } [.17, .53]$.



Now, for the idealists.

The REG Procedure
 Model: MODEL1
 Dependent Variable: ar
 idealism=Idealist

Number of Observations Read 63

Number of Observations Used 63

Descriptive Statistics					
Variable	Sum	Mean	Uncorrected SS	Variance	Standard Deviation
Intercept	63.00000	1.00000	63.00000	0	0
misanth	141.20000	2.24127	344.40000	0.45053	0.67121
ar	153.64900	2.43887	390.42107	0.25308	0.50307

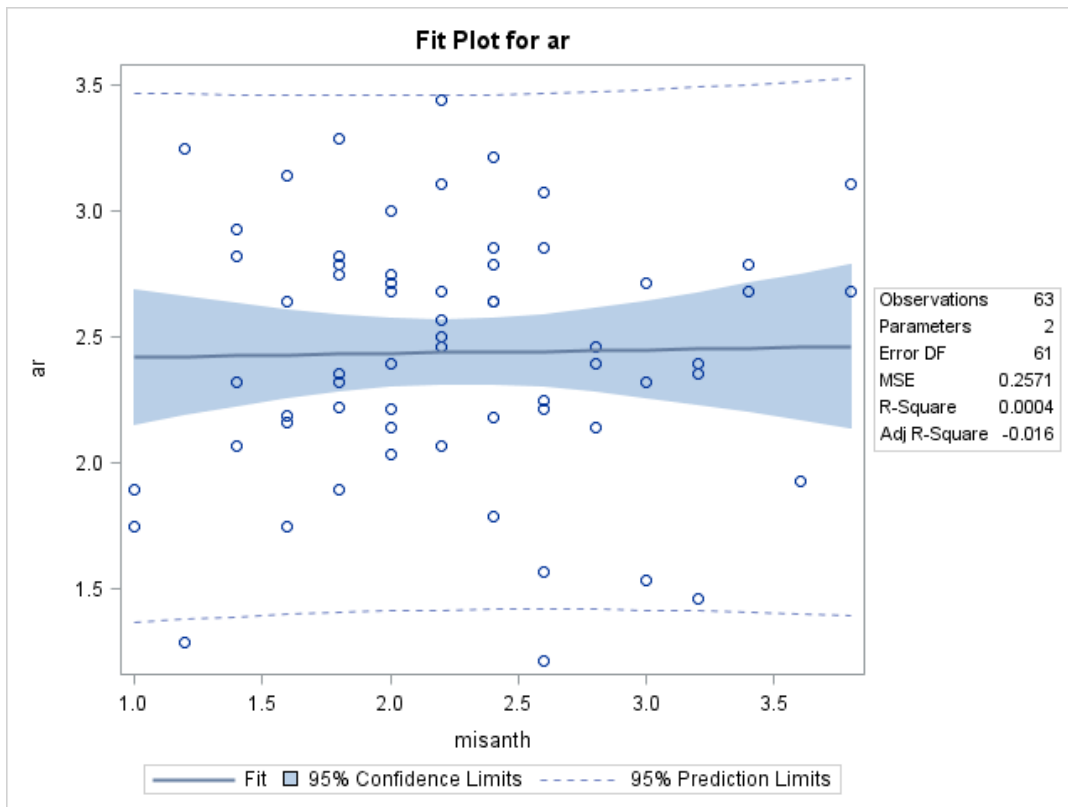
Correlation		
Variable	misanth	ar
misanth	1.0000	0.0205
ar	0.0205	1.0000

Analysis of Variance					
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	1	0.00657	0.00657	0.03	0.8735
Error	61	15.68410	0.25712		
Corrected Total	62	15.69067			

Root MSE	0.50707	R-Square	0.0004
Dependent Mean	2.43887	Adj R-Sq	-0.0160
Coeff Var	20.79102		

Parameter Estimates					
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t
Intercept	1	2.40450	0.22432	10.72	<.0001
misanth	1	0.01533	0.09594	0.16	0.8735

Among the idealists, misanthropy was not significantly associated with support for animal rights, $r(n = 63) = .02$, $p = .87$, 95% CI [-.23, .27].



Later I shall show you more sophisticated methods of moderation analysis.

[Return to my SAS Lessons page.](#)