SASEG 9B – Regression Assumptions

(Fall 2015)

Sources (adapted with permission)-

T. P. Cronan, Jeff Mullins, Ron Freeze, and David E. Douglas Course and Classroom Notes Enterprise Systems, Sam M. Walton College of Business, University of Arkansas, Fayetteville Microsoft Enterprise Consortium IBM Academic Initiative SAS[®] Multivariate Statistics Course Notes & Workshop, 2010 SAS[®] Advanced Business Analytics Course Notes & Workshop, 2010 Microsoft[®] Notes Teradata[®] University Network

Copyright © 2013 ISYS 5503 Decision Support and Analytics, Information Systems; Timothy Paul Cronan. For educational uses only - adapted from sources with permission. No part of this publication may be reproduced, stored in a retrieval system, or transmitted, in any form or by any means, electronic, mechanical, photocopying, or otherwise, without the prior written permission from the author/presenter.

Examining Residuals



Recall that the model for the linear regression has the form $Y=\beta_0+\beta_1X+\epsilon$. When you perform a regression analysis, several assumptions about the error terms must be met to provide valid tests of hypothesis and confidence intervals. The assumptions are that the error terms

- have a mean of 0 at each value of the predictor variable
- are normally distributed at each value of the predictor variable
- have the same variance at each value of the predictor variable
- are independent.





In the first plot, a regression line adequately describes the data.



In the second plot, a simple linear regression model is not appropriate because you are fitting a straight line through a curvilinear relationship.



In the third plot, there seems to be an outlying data value that is affecting the regression line. This outlier is an influential data value in that it is substantially changing the fit of the regression line.



In the fourth plot, the outlying data point dramatically changes the fit of the regression line. In fact the slope would be undefined without the outlier.

The four plots illustrate that relying on the regression output to describe the relationship between your variables can be misleading. The regression equations and the R² statistics are the same even though the relationships between the two variables are different. Always produce a scatter plot before you conduct a regression analysis.



To verify the assumptions for regression, you can use the residual values from the regression analysis. Residuals are defined as

 $r_{i} = Y_{i} - \hat{Y}_{i}$

where \hat{Y}_{i} is the predicted value for the *i*th value of the dependent variable.

You can examine two types of plots when verifying assumptions:

- the residuals versus the predicted values
- the residuals versus the values of the independent variables



The graphs above are plots of residual values versus predicted values or predictor variable values for four models fit to different sets of data. If model assumptions are valid, then the residual values should be randomly scattered about a reference line at 0. Any patterns or trends in the residuals might indicate problems in the model.

- 1. The model form appears to be adequate because the residuals are randomly scattered about a reference line at 0 and no patterns appear in the residual values.
- 2. The model form is incorrect. The plot indicates that the model should take into account curvature in the data. One possible solution is to add a quadratic term as one of the predictor variables.
- 3. The variance is not constant. As you move from left to right, the variance increases. One possible solution is to transform your dependent variable.
- 4. The observations are not independent. For this graph, the residuals tend to be followed by residuals with the same sign, which is called *autocorrelation*. This problem can occur when you have observations that have been collected over time. A possible solution is to use the Regression Analysis with Autoregressive Errors task.



Besides verifying assumptions, it is also important to check for outliers. Observations that are far away from the bulk of your data are outliers. These observations are often data errors or reflect unusual circumstances. In either case, it is good statistical practice to detect these outliers and find out why they have occurred.

Residual Plots

Using the **FITNESS** data set, invoke the Linear Regression task to test the regression model of **Oxygen_Consumption** against the predictor variables of **RunTime**, **Age**, **Run_Pulse** and **Maximum_Pulse** (the model that was best based on Mallows' Cp prediction criterion). Produce the default graphics.

- 1. Create a new project and name it **SASEG 9B Demos**.
- 2. Open the **FITNESS** data set.

🜀 SAS Enterprise Guide - EGBS.egp			
File Edit View Tasks Program To	ols	Help 🛛 🗎 🕶 🚰 🕶	🚳 📇
Project Tree 👻 👻	FIT	NESS 👻	
🗇 🖧 Data Creation	7	Filter and Sort 🕮 🤇	Query Builde
eabs00d01		🔌 Name	🔌 Ge
⊕ Seg Chapter 1 Demos	1	Donna	F
🕀 🗞 😌 Chapter 1 Exercises	2	Gracie	F
🕀 🗞 🔁 Chapter 2 Demos	3	Luanne	F
🕂 🚱 Chapter 2 Exercises	4	Mimi	F
🕀 စိုဗ္ဗရူ Chapter 3 Demos	5	Chris	М
🕀 စိုဗ္ဗရူ Chapter 3 Exercises	6	Allen	М
E to the second	7	Nancy	F
En terre de Exercises	8	Patty	F
Chapter 5 Demos	9	Suzanne	F
FILNESS	10	Teresa	F

3. Select <u>Analyze</u> \Rightarrow <u>Regression</u> \Rightarrow <u>Linear Regression...</u>.

Analyze 👻 Export 👻 Send To 👻 📝						
	ANOVA	۲	🔞 Охуд	jen_Consumption 🔞	Run_F	
	Regression	×	📈 Linear	Regression		
	Multivariate	۲	🦗 Nonlin	earRegression		
	Survival Analysis	۲	🔢 Logist	Logistic Regression		
	Capability 🕨 🕨		🕂 Gener	alized Linear Models		
	Control Charts	۲	2	43.67		
m	Pareto Chart		2	45.44		
	Time Series		8	50.55		
	Time Denes		1	46.67		
B	Model Scoring		7	45.31		

4. Drag Oxygen_Consumption to the dependent variable task role and RunTime, Age, Run_Pulse, and Maximum_Pulse to the explanatory variables task role.

🖉 Linear Regressio	on9 for Local:SASUSER.FITNESS		
Data Model Statistics Plots Predictions Titles Properties	Data Data source: Local:SASUSER.FITNESS Task filter: None		
, isponto	Variables to assign: Name Sender RunTime Age Weight Oxygen_Consumption Run_Pulse Rest_Pulse Maximum_Pulse Performance	Task roles: Dependent variable (Limit: 1) 0 Oxygen_Consumption Explanatory variables 8 Run_Pulse 8 Group analysis by Frequency count (Limit: 1) Relative weight (Limit: 1)	수 수

5. Click

Run

Linear Regression Results										
	Dep	Mod ende	T lel: nt V	he Line aria	REG Pro ear_Regr ble: Oxy	es ge	edure ssion_Moo en_Consu	iel mpt	tion	
	1	Numl	ber (of O) bservati	or	ns Read	31	Ī	
	1	Num	ber (of O)bservati	or	ns Used	31]	
			Δ	nalv	vsis of V	ar	iance			
					Sum of		Mean			
Sourc	е		DF	Squares		Square		Fν	alue	Pr > F
Model			4	711.45087		177.86272		3	3.01	<.0001
Error			26	140.10368		5.38860				
Correc	cted To	tal	30	85	1.55455					
	Root I	NSE			2.321	34	R-Square	e 0	.835	5
	Depen	dent	Me	an	47.3758	B1	Adj R-Sq	0	.810	2
	Coeff	Var			4.8998	84				_
			Р	araı	neter Es	tir	nates			
			Ť	P	aramete	r	Standard			
Variat	Variable		DF	:	Estimate	9	Error	t V	alue	Pr > t
Interc	Intercept		1	1	97.16952	2	11.65703		8.34	<.0001
RunTime		1	1	-2.77576	5	0.34159	-	8.13	<.0001	
Age			1	1	-0.18903	3	0.09439	-	2.00	0.0557
Run_F	Pulse		1	1	-0.34568	3	0.11820	-	2.92	0.0071

0.27188

0.13438

2.02 0.0534

The histogram of residuals helps you to find outliers and assess the normality assumption.

1

Maximum_Pulse

Note – review SASEG 8A (pp. 12 - 16) regarding interpretation of the plots – much of the information regarding interpretation for a one variable model will be the same for the multiple variable model.

The plot of the residuals versus the values of the independent variables, **Runtime**, **Age**, **Run_Pulse**, and **Maximum_Pulse** are produced by SASEG. They show no obvious trends or patterns in the residuals. Recall that independence of residual errors (no trends) is an assumption for linear regression, as is constant variance across all levels of all predictor variables (and across all levels of the predicted values, which is seen below).

The diagnostic plots shown above will be described later in greater detail.

The plot of the residuals against the normal quantiles is shown above left (quantile-quantile plot, *also known as the Q-Q Plot*). If the residuals are normally distributed, the plot should appear to follow closely a straight, diagonal line. If the plot deviates substantially from the reference line, then there is evidence against normality.

The plot shows little deviation from the expected pattern. Thus, you can conclude that the residuals do not significantly violate the normality assumption. If the residuals did violate the normality assumption, then a transformation of the response variable or a different model might be warranted.

More diagnostic plots and plots are included by default, as well as a box and whisker plot for residuals.

6. In order to visually check the assumption of constant variance, you can reopen the last task by rightclicking it and modifying it.

🗄 🖧 Chapter 5 Demos	
🖻 🊟 FITNESS	
Linear Re 🧭	Open •
►	Run Linear Regression9
	Modify Linear Regression9
	Select Input Data
- and - and - and - and -	Publish
	Add as Code Template
1	Create Task Template
and the second s	Create Stored Process

 With <u>Plots</u> selected at the left, click the radio button next to <u>Custom list of plots</u>. The box next to <u>Diagnostic plots</u> should already be checked. In addition, check the boxes next to <u>Residuals by predicted values plot</u> and <u>Residual plots</u>.

🔟 Linear Regressio	Linear Regression9 for Local:SASUSER.FITNESS						
Linear Regression	Plots Plots Show plots for regression analysis All appropriate plots for the current data selection Custom list of plots Custom plots: Histogram plot of the residuals Residuals by predicted values plot Studentized residuals by predicted values plot Dbserved by Predicted values plot Plot Cook's D statistic Studentized residuals by leverage plot Normal quantile plot of the residuals Residual-Fit plot						
	 Residual-Fit plot Box plot of the residuals Diagnostic plots DFFITS plots DEFETAS plots 						
	Residual plots Conter plot with regression line Select all						

Run

and do not replace the results from the previous run.

The plots produced are displayed below:

The Residual by Predicted plot shows no pattern of residuals around the residual mean of 0. One of the assumptions of linear regression is constant variance across all levels of all predictors. This plot, along with the plots of residuals against predictors, helps you to assess that assumption. In this case, there is no clear pattern, indicating no strong evidence against the assumption of constant variance.

The Fit Diagnostics panel plot displays many of the plots seen in the previous part of the demonstration, but on a smaller scale.

The plots of the residuals by each of the predictor variables in the model show no patterns or trends. Again, this lends support to the validity of the constant variance assumption for this regression model.

Influential Observations (Any Outliers?) – Going Beyond

Recall in the previous section that you saw examples of data sets where the simple linear regression model fits were essentially the same. However, plotting the data revealed that the model fits were different.

One of the examples showed a highly influential observation like the example above.

Identifying influential observations in multiple linear regression is more complex because you have more predictors to consider.

For our purposes, to detect outliers we will use the **Studentized Residuals**, **Cook's D statistic**, and the **RSTUDENT residuals**. Note that there are others...

Studentized Residuals - One way to check for outliers is to use the **studentized residuals**. These are calculated by dividing the residual values by their standard errors. For a model that fits the data well and has no outliers, most of the studentized residuals should be close to 0. In general, studentized residuals that have an absolute value less than 2.0 could have easily occurred by chance. Studentized residuals that are between an absolute value of 2.0 to 3.0 occur infrequently and could be outliers. Studentized residuals that are larger than an absolute value of 3.0 occur rarely by chance alone and should be investigated.

Cook's D statistic - To detect influential observations, you can also use **Cook's D statistic**. This statistic measures the change in the parameter estimates that results from deleting each observation.

Identify observations above the cutoff and investigate the reasons they occurred.

Cook's D Statistic
Cook's D statistic is a measure of the simultaneous change in the parameter estimates when an observation is deleted from the analysis.
A suggested cutoff is $D_i > \frac{4}{n}$, where <i>n</i> is the sample size.
If the above condition is true, then the observation might have an adverse effect on the analysis.
22

RSTUDENT Residuals - Recall that studentized residuals are the ordinary residuals divided by their standard errors. **The RSTUDENT residuals** are similar to the studentized residuals except that they are calculated after deleting the i^{th} observation. In other words, the RSTUDENT residual is the difference between the observed Y and the predicted value of Y excluding this observation from the regression.

If the RSTUDENT residual is different from the studentized residual for a specific observation, that observation is likely to be influential. A suggested cutoff for |RSTUDENT| residuals is greater than 3.

An Exercise - Looking for Influential Observations

model.

Generate the **RStudent** and **Cook's D** influence statistics and plots for the prediction

Save the statistics to an output data set and create a data set with only observations that exceed the suggested cutoffs of the influence statistics.

Refer to the last task (linear model where you used the **FITNESS** data set, the regression model of **Oxygen_Consumption** against the predictor variables of **RunTime**, **Age**, **Run_Pulse** and **Maximum_Pulse**).

- 1. Modify the last task by right-clicking the Project and selecting Modify....
- 2. With <u>Plots</u> selected at the left, check the boxes shown checked below in the Custom plots area.

RSTUDENT residuals are referred to as Studentized residuals in the task windows.

Data Model	Plots
Statistics Plots Predictions Titles Properties	 Show plots for regression analysis All appropriate plots for the current data selection Custom list of plots
	Custom plots: Histogram plot of the residuals Residuals by predicted values plot Studentized residuals by predicted values plot Observed by Predicted values plot Plot Cook's D statistic Studentized residuals by leverage plot Normal quantile plot of the residuals Box plot of the residuals Diagnostic plots DFFITS plots DFBETAS plots Residual plots

- 3. With <u>Predictions</u> selected at the left:
 - a. Check the box for **Original sample** under Data to predict.
 - b. Check Predictions and Diagnostic statistics under Save output data.
 - c. Check the box for <u>Residuals</u> under Additional statistics.

You can change the name and library of the data set where the diagnostic statistic variables will be stored by clicking Browse... in the Save output data area.

L	Linear Regression911 for Local:SASUSER.FITNESS							
	Data Model Statistics	Predictions						
	Plots	Data to predict	Save output data					
	Predictions Titles	Original sample	✓ Predictions					
	Properties	Additional data	Diagnostic statistics					
		Browse	LocatSASUSER.PREDLINR Browse					
		Additional statistics Residuals Prediction limits	Display output and plots					

4. Click Run and do not replace the results from the previous run.

Linea	ar Regression911 👻								
	📰 Input Data 🗒 Code 📋 Log 📓 Output Data 😰 Results								
\$5	属 Modify Task 🏼 🐺 Filter and	Sort 🛄 Query Build	er Data 🕶 Des	cribe 👻 Graph 👻 Analyze	e → Export → Se	end To 👻 📝			
	Oxygen_Consumption	🕽 Run_Pulse 🔞	Rest_Pulse	😟 Maximum_Pulse 😡	Performance 🤅	predicted_0xygen_Consumption	stdp_0xygen_Consumption		
1	59.57	166	40	172	90	55.9332897	0.91043968		
2	60.06	170	48	186	94	57.8362043	1.6123022		
3	54.3	156	45	168	83	56.7811803	1.07752127		

Along with the other output from the task, a tab for the Output Data table appears. Select that tab to see the data set created with all variables from the **Fitness** data set, along with several new variables containing values for the diagnostic statistics and residuals, along with relevant standard errors.

Return to the Results tab.

Linear Regression Results

The REG Procedure Model: Linear_Regression_Model Dependent Variable: Oxygen_Consumption

Number of Observations Read31Number of Observations Used31

Analysis of Variance							
Sum of Mean							
Source	DF	Squares	Square	F Value	Pr > F		
Model	4	711.45087	177.86272	33.01	<.0001		
Error	26	140.10368	5.38860				
Corrected Total	30	851.55455					

Root MSE	2.32134 R-Square	0.8355
Dependent Mean	47.37581 Adj R-Sq	0.8102
Coeff Var	4.89984	

Parameter Estimates								
		Parameter	Standard					
Variable	DF	Estimate	Error	t Value	Pr > t			
Intercept	1	97.16952	11.65703	8.34	<.0001			
RunTime	1	-2.77576	0.34159	-8.13	<.0001			
Age	1	-0.18903	0.09439	-2.00	0.0557			
Run_Pulse	1	-0.34568	0.11820	-2.92	0.0071			
Maximum_Pulse	1	0.27188	0.13438	2.02	0.0534			

The RStudent by Predicted plot shows only two values outside the range of [-2,2] and no values outside the range of [-3,3]. These values are not different from what one would normally expect by chance from a normally distributed population.

A horizontal reference line is drawn at the critical value of Cook's D. Only one observation's Cook's D value exceeded that cutpoint and merits further investigation.

5. Right-click the previous task and select Add as a Code Template.

6. Double-click the node for the code in order to edit it and find the PROC REG section of the code.

```
TITLE;
TITLE1 "Linear Regression Results";
FOOTNOTE ;
FOOTNOTE1 "Generated by the SAS System (& SASSERVERNAME, &SYSSCPL) on
%TRIM(%QSYSFUNC(DATE(), NLDATE20.)) at %TRIM(%SYSFUNC(TIME(),
NLTIMAP20.))";
PROC REG DATA=WORK.SORTTempTableSorted
        PLOTS (ONLY) = RSTUDENTBYPREDICTED
        PLOTS (ONLY) = COOKSD
        PLOTS (ONLY) = DFFITS
        PLOTS (ONLY) = DFBETAS
    Linear Regression Model: MODEL Oxygen Consumption = RunTime Age
Run Pulse Maximum Pulse
        /
                 SELECTION=NONE
    ;
    OUTPUT OUT=SASUSER.PREDLINREGPREDICTIONSFITNES 0001(LABEL="Linear
regression predictions and statistics for SASUSER.FITNESS")
        PREDICTED=predicted Oxygen Consumption
        RESIDUAL=residual Oxygen Consumption
        STUDENT=student Oxygen Consumption
        RSTUDENT=rstudent Oxygen Consumption
        COOKD=cookd Oxygen Consumption
        DFFITS=dffits Oxygen Consumption
        H=h Oxygen Consumption
        STDI=stdi Oxygen Consumption
        STDP=stdp Oxygen Consumption
        STDR=stdr Oxygen Consumption ;
RUN;
QUIT;
```

- 7. Make the following changes:
 - a. Add the option (LABEL) at the end of each PLOTS(ONLY) line.

```
PROC REG DATA=WORK.SORTTempTableSorted
    PLOTS (ONLY) =RSTUDENTBYPREDICTED (LABEL)
    PLOTS (ONLY) =COOKSD (LABEL)
    PLOTS (ONLY) =DFFITS (LABEL)
    pLOTS (ONLY) =DFBETAS (LABEL)
;
```

b. Add the statement ID NAME; immediately above the OUTPUT statement.

```
ID NAME;
OUTPUT
OUT=SASUSER.PREDLINREGPREDICTIONSFITNES_0001(LABEL="Linear
regression predictions and statistics for SASUSER.FITNESS")
```

8. Click \triangleright Run above the code window.

The RStudent plot shows two observations beyond 2 standard errors from the mean of 0. Those are identified as Sammy and Jack. Because you expect 5% of values to be beyond 2 standard errors from the mean (remember that these RStudent residuals are assumed to be normally distributed), the fact that you have 2 that far out gives no cause for concern (5% of 31 is 1.55 expected observations). William and Gracie have the most extreme "leverage" values, which mean that they are most extreme in the predictor variable space.

The Cook's D plot shows Gracie to be an influential point.

If the unusual data are erroneous, correct the errors and reanalyze the data.

Another possibility is that the observation, although valid, could be unusual. If you had a larger sample size, there might be more observations like the unusual ones.

You might have to collect more data to confirm the relationship suggested by the influential observation.

In general, we try not to exclude data. In many circumstances, some of the unusual observations contain important information. However, if you do choose to exclude some observations, include a description of the types of observations you exclude and provide an explanation. Also discuss the limitation of your conclusions, given the exclusions, as part of your report or presentation.

Collinearity

In the **Fitness** data set example, **RunTime** and **Oxygen_Consumption** have a strong linear relationship. **Performance** and **Oxygen_Consumption** also have a strong linear relationship. In addition, **RunTime** and **Performance** are linearly related to a large degree.

The goal of multiple linear regression with two predictor variables is to find a best fit plane through the data to predict **Oxygen_Consumption**. This perspective shows a very strong relationship between the

predictor variables **RunTime** and **Performance**. You can imagine that the prediction plane you are trying to build is like a tabletop, where the observations guide the angle of the tabletop, relative to the floor, like legs for the table. If the legs line up with one another, then the plane built atop will tend to be unstable.

Here is another way of looking at the three dimensions of two predictor variables and a response variable. Where should the prediction plane be placed? The slopes of the prediction plane relative to each X and the Y are the parameter coefficient estimates.

 X_1 and X_2 almost follow a straight line $X_1 = X_2$ in the (X_1, X_2) plane.

Why is this a problem? Two reasons exist.

- 1. Neither might appear to be significant when both are in the model; however, either might be significant when only one is in the model. Thus, collinearity can hide significant effects. (The reverse can be true as well: collinearity can increase the apparent significance of effects.)
- 2. Collinearity also increases the variance of the parameter estimates and consequently increases prediction error.

This is a representation of a best-fit plane through the data.

However, the removal of just one data point (or even just moving the data point) results in a very different prediction plane (as represented by the lighter plane). This illustrates variability of the parameter estimates when there is extreme collinearity.

When collinearity is a problem, the estimates of the coefficients are unstable. This means that they have a large variance. Consequently, the true relationship between Y and the Xs might be quite different from that suggested by the magnitude and sign of the coefficients.

Collinearity is **not** a violation of the assumptions of linear regression.

Example of Collinearity

Generate a regression with Oxygen_Consumption as the dependent variable and Performance, Runtime, Age, Weight, Run_Pulse, Rest_Pulse, and Maximum_Pulse as the independent variables. Compare this model with the Mallows prediction model from the previous section.

- 1. With the Fitness data set active, select <u>Tasks</u> \Rightarrow <u>Regression</u> \Rightarrow <u>Linear Regression...</u>.
- 2. Drag **Oxygen_Consumption** to the dependent variable role and all other numeric variables to the explanatory variables role.

3. With <u>Plots</u> selected at the left, uncheck the box for <u>Show plots for regression analysis</u>.

Linear Regression Results

The REG Procedure
Model: Linear_Regression_Model
Dependent Variable: Oxygen_Consumption

Number of Observations Read 31									
Number of Observations Used 31									
	Analysis of Variance								
			Sum of		Mean				
Source	DF	S	quares		Square	F	Value	Pr > F	
Model	7	722	2.66124	1	03.23732		18.42	<.0001	
Error	23	128	3.89331		5.60406				
Corrected Total	30	851	1.55455						
D (MCC			0.007	00	D.C.		0.040		
Root MSE			2.36/29 R-Squar			e	0.848	0	
Dependent	Me	an	47.37581 Adj R-So			1	0.8020	5	
Coeff Var			4.996	83					
	Р	arar	neter Es	stir	nates				
		P	aramete	r	Standard				
Variable	DF	-	Estimat	e	Error	t	Value	Pr > t	
Intercept	1	1 1	31.7824	9	72.20754		1.83	0.0810	
RunTime	1	1	-3.8601	9	2.93659		-1.31	0.2016	
Age	1	1	-0.4608	2	0.58660		-0.79	0.4401	
Weight	1		-0.0581	2	0.06892	2	-0.84	0.4078	
Run_Pulse	1	1	-0.3620	7	0.12324		-2.94	0.0074	
Rest_Pulse	1		-0.0151	2	0.06817	1	-0.22	0.8264	
Maximum_Pulse	1	1	0.3010	2	0.13981	T	2.15	0.0420	
Performance	1	1	-0.1261	9	0.30097	'	-0.42	0.6789	

For the full model, Model F is highly significant and the \mathbb{R}^2 is large. These statistics suggest that the model fits the data well.

- However, when you examine the *p*-values of the parameters, only **Run_Pulse** and **Maximum Pulse** are statistically significant.
- Recall that the 4-variable prediction model included **Runtime**; however, in the full model, this same variable is not statistically significant (*p*-value=0.2016). The *p*-value for **Age** changed from 0.0557 to 0.4401 between the 4-variable model and the full model.

When you have a highly significant Model *F* but no (or few) highly significant terms, *collinearity is a likely problem*.

VIFprovides a measure of the magnitude of the collinearity (Variance
Inflation Factor).Collinearity Analysisincludes the intercept vector when analyzing the X'X matrix for
collinearity.Collinearity (No Intercept)excludes the intercept vector.

The two Collinearity Analysis options also provide a measure of the magnitude of the problem as well as give information that can be used to identify the sets of Xs that are the source of the problem. They are not described in this course.

You can calculate a VIF for each term in the model.

Marquardt (1990) suggests that a VIF > 10 indicates the presence of strong collinearity in the model.

 $VIF_i = 1/(1 - R_i^2)$, where R_i^2 is the R^2 of X_i , regressed on all the other Xs in the model.

For example, if the model is Y = X1 X2 X3 X4, i = 1 to 4.

To calculate the R² for X3, fit the model X3 = X1 X2 X4. Take the R² from the model with X3 as the dependent variable and replace it in the formula VIF₃ = $1/(1 - R_3^2)$. If VIF₃ is greater than 10, X3 is possibly involved in collinearity.

Collinearity Diagnostics

Invoke the Linear Regression task and use the VIF option to assess the magnitude of the collinearity problem and identify the terms involved in the problem.

- 1. Reopen the previous task by right-clicking it and selecting Modify....
- 2. With <u>Statistics</u> checked at the left, check the box next to <u>Variance inflation values</u> in the Diagnostics area.

🗵 Linear Regressio	on10 for Local:SASUSER.FITNESS	
Linear Regression	Statistics Details on estimates Standardized regression coefficients Sum of squares, Type 1 Sum of squares, Type 2 Correlation matrix of estimates Covariance matrix of estimates Confidence limits for parameter estimates Confidence level: 95%	Diagnostics Collinearity analysis Collinearity analysis without the intercept Tolerance values for estimates Variance inflation values Heteroscedasticity test Asymptotic covariance matrix Durbin-Watson statistic
	Correlations Partial correlations Semi-partial correlations	

3. Click Run and do replace the results from the previous run.

SAS Enter	rprise Guide 🔀
?	Do you want to replace the results from the previous run? Choosing "No" will save the changes to a new task, named "Linear Regression101".
	Yea No Cancel

Partial Output

Parameter Estimates							
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t	Variance Inflation	
Intercept	1	131.78249	72.20754	1.83	0.0810	0	
RunTime	1	-3.86019	2.93659	-1.31	0.2016	88.86251	
Age	1	-0.46082	0.58660	-0.79	0.4401	51.01176	
Weight	1	-0.05812	0.06892	-0.84	0.4078	1.76383	
Run_Pulse	1	-0.36207	0.12324	-2.94	0.0074	8.54498	
Rest_Pulse	1	-0.01512	0.06817	-0.22	0.8264	1.44425	
Maximum_Pulse	1	0.30102	0.13981	2.15	0.0420	8.78755	
Performance	1	-0.12619	0.30097	-0.42	0.6789	162.85399	

The only change in the output from the previous run of the task is the final column of the Parameter Estimates table. There is now a listing of Variance Inflation values for each predictor variable.

Marquardt (1990) suggests that a VIF > 10 indicates the presence of strong collinearity in the model.

Some of the VIFs are much larger than 10. *A severe collinearity problem is present*. At this point there are many ways to proceed. <u>However, it is always a good idea to use some subject-matter expertise</u>. For instance, a quick conversation with the analyst and a view of the data coding scheme turned up this bit of information.

<u>We just happen to know</u> - The variable **Performance** was not a measured variable. The researchers, on the basis of prior literature, created a summary variable, which is a weighted function of the three variables, **RunTime**, **Age**, and **Gender**. This is not at all an uncommon occurrence and illustrates an important point. *If a summary variable is included in a model along with some or all of its composite measures, there is bound to be collinearity. In fact, this can be the source of great problems.*

- If the composite variable has meaning, it can be used as a stand-in measure for all three composite scores and you can remove the variables **RunTime** and **Age** from the analysis.

A decision was made to remove **Performance** from the analysis. Another check of collinearity is warranted.

- 4. Reopen the previous task.
- 5. Remove **Performance** from the list of explanatory variables by highlighting it and clicking

6. Click Run and do not replace the results from the previous run.

	Number of Observations Read 31											
	Number of Observations Used 31											
			٨	nah	eie /	ofVa	rianc	0				_
				inary	515	of	Ianc	loon				
Source			DF	s	auai	res	Sa	uare	F۱	/alue	Pr>	F
Model			6	721	1.676	605 1	20.2	7934	2	22.23	<.000)1
Error			24	129	9.878	351	5.4	1160				_
Correcte	ed To	tal	30	851	1.554	55						
					-							
R	loot N	ISE			2.	32629	R-S	quare	e (0.8475)	
D)epen	den	t Me	an	47.	.37581 Adj R-So		R-Sq		0.8094	ŀ	
C	Coeff \	/ar			4.	91028						
			P	arar	nete	r Esti	mate	s				
			Pa	rame	eter	Stan	dard				V	ariance
Variable		DF	E	stim	ate	E	Fror	t Va	lue	Pr >	t Ir	nflation
Intercept		1	10	1.96	313	12.2	7174	8	.31	<.000)1	0
RunTime		1	-	2.63	994	0.3	8532	-6	.85	<.000)1 1	1.58432
Age		1	-	0.21	848	0.0	9850	-2	22	0.03	63 1	1.48953
Weight		1	-	0.07	503	0.0	5492	-1	.37	0.184	45 1	1.15973
Run_Pulse		1	-	0.36	721	0.1	2050	-3	.05	0.00	55 8	3.46034
Rest_Pulse		1	-	0.01	952	0.0	6619	-0	.29	0.77	06 1	1.41004
Maximum_P	ulse	1		0.30	457	0.1	3714	2	.22	0.03	50 8	3.75535

The greatest VIF values are much smaller now. The variables **Maximum_Pulse** and **Run_Pulse** are also collinear, but for a natural reason. The pulse at the end of a run is highly likely to correlate with the maximum pulse during the run. One might be tempted simply to remove one variable from the model, but the small *p*-values for each indicate that this would adversely affect the model.

- 7. Reopen the previous task.
- 8. Remove Maximum_Pulse from the list of explanatory variables by highlighting it and clicking

SAS Enter	prise Guide
?	Do you want to replace the results from the previous run? Choosing "No" will save the changes to a new task, named "Linear Regression101".
	Yes Nr Cancel

	Number of Observations Read 31										
	Number of Observations Used 31										
		_		^	nah	ala of l	lari				
				μ	mary		an	ance			
C	_			DE		Sum of		Mean	-		
Source	е			DF	3	quares		Square	F	value	PT > F
Model				5	694	1.98323	13	8.99665		22.19	<.0001
Error				25	156	5.57132		6.26285			
Correc	ted 1	Tota	al	30	851	1.55455					
	D	4 8.8	с г			0.500	57		-	0.040	4
	ROO		SE		2.002		:57	57 R-Square		0.816	1
	Dep	end	ent	t Me	an 47.375		81	1 Adj R-Sq		0.7794	4
	Coe	ff Va	ar		5.282		238				
				D	arar	notor E	otim	atoo.			
			-	F	arar	neter E	sun	lates	_		
			Р	aran	nete	r Stand	lard				Variance
Variable		DF		Esti	mate	e E	rroi	t Value	P	r > t	Inflation
Intercept	t	1	1	15.4	611	5 11.40	6893	10.07	' <	.0001	0
RunTime	9	1		-2.71594		4 0.41	1288	-6.58	<	.0001	1.57183
Age		1		-0.27650		0.10)217	-2.71	0	.0121	1.38477
Weight		1		-0.05300		0.05	5811	-0.91	0	.3704	1.12190
Run_Pul	se	1		-0.1	2213	3 0.05	5207	-2.35	i 0	.0272	1.36493
Rest_Pu	se	1		-0.0	2485	5 0.07	7116	-0.35	i 0	.7298	1.40819

With **Maximum_Pulse** removed, all of the VIF values are low, but the R-Square and Adj R-Sq values were reduced and the *p*-value for **Run-Pulse** actually increased!

??Even with collinearity still present in the model, it might be advisable to keep the previous model including **Maximum Pulse**.??

Collinearity can have a substantial effect on the outcome of a stepwise procedure for model selection. Because the significance of important variables can be masked by collinearity, the final model might not include very important variables. This is why it is advisable to deal with collinearity before using any automated model selection tool.

Just FYI - there are other approaches to dealing with collinearity. Two techniques are ridge regression and principle components regression. In addition, re-centering the predictor variables can sometimes eliminate collinearity problems, especially in a polynomial regression and ANCOVA models.

- (1) **Preliminary Analysis** This step includes the use of descriptive statistics, graphs, and correlation analysis.
- (2) Candidate Model Selection This step uses the numerous selection options in the Linear Regression task to identify one or more candidate models.
- (3) Assumption Validation This step includes the plots of residuals and graphs of the residuals versus the predicted values. It also includes a test for equal variances.
- (4) Collinearity and Influential Observation Detection The former includes the use of the VIF statistic, condition indices, and variation proportions; the latter includes the examination of Rstudent residuals, Cook's D statistic, and DFFITS statistics.
- (5) Model Revision If steps (3) and (4) indicate the need for model revision, generate a new model by returning to these two steps.
- (6) Prediction Testing If possible, validate the model with data not used to build the model.

Comprehensive Exercise – but, Optional

1. Assessing Collinearity

Using the **BodyFat2** data set, run a regression of **PctBodyFat2** on all the other numeric variables in the file.

- **a.** Determine whether there is a collinearity problem.
- **b.** If so, decide what you would like to do about that. Will you remove any variables? Why or why not?

Solutions to Exercises

1. Examining Residuals

Assess the model obtained from the final forward stepwise selection of predictors for the **BodyFat2** data set. Run a regression of **PctBodyFat2** on **Abdomen**, **Weight**, **Wrist**, and **Forearm**. Create plots of the residuals by the four regressors and by the predicted values and a normal quantilequantile plot.

Invoke the Linear Regression task to test the regression model of **PctBodyFat2** against the predictor variables of **Abdomen**, **Weight**, **Wrist**, and **Forearm**.

a. Do the residual plots indicate any problems with the constant variance assumption?

• Create a new process flow and rename it Chapter 5 Exercises.

- Open the BodyFat2 data set.
- Select <u>Analyze</u> \Rightarrow <u>Regression</u> \Rightarrow <u>Linear Regression...</u>.

Analy	yze 🕶 🛛 Export 👻 Se	nd To	- 🛃	נ		
ANOVA		•		🕖 Weight	🔞 Height	
	Regression	•	1/1	Linear Regression		
	Multivariate	•	20	NonlinearRegress	sion	
	Survival Analysis			Logistic Regression		
	Capability	•	<u>n</u>	Generalized Linear Models		
	Control Charts	•	24	210.25	69.75	
lín	Pareto Chart		25	176	72.5	
			25	191	74	
	Time Series	•	23	198.25	73.5	
щ.	Model Scoring		26	186.25	74.5	

• Drag PctBodyFat2 to the dependent variable task role and Abdomen, Weight, Wrist, and Forearm to the explanatory variables task role.

Data			
Data source: Local:SASUSER.BODYFAT2 Task filter: None Variables to assign: Name Density Density Density Weight	[Task roles: Dependent variable (Limit: 1) PotBodyFat2 Explanatory variables Abdomen	令
 Height Adioposity FatFreeWt Neck Chest Abdomen Hip Thigh Knee Ankle Biceps Forearm Wrist 	A	 Weight Wrist Forearm Group analysis by Frequency count (Limit: 1) Relative weight (Limit: 1) 	

• With <u>Plots</u> selected at the left, click the radio button next to <u>Custom list of plots</u>. The box next to <u>Diagnostic plots</u> should already be checked. In addition, check the boxes next to <u>Residuals by predicted values plot</u> and <u>Residual plots</u>.

Linear Regression9 for Local:SASUSER.FITNESS						
Data Model	Plots					
Plots Predictions Titles Properties	 Show plots for regression analysis All appropriate plots for the current data selection Custom list of plots 					
	Custom plots: Histogram plot of the residuals Second Studentized residuals by predicted values plot Studentized residuals by predicted values plot Plot Cook's D statistic Studentized residuals by leverage plot Normal quantile plot of the residuals Residual-Fit plot Box plot of the residuals DFFITS plots DFFITS plots DFBETAS plots Residual plots Scatter plot with regression line					
Click Run .						

•

It does not appear that the data violate the assumption of constant variance.

b. Are there any outliers indicated by the evident in any of the residual plots?

There are a few outliers for **Wrist** and **Forearm** and one clear outlier in each of **Abdomen** and **Weight**.

c. Does the quantile-quantile plot indicate any problems with the normality assumption?

The quantile-quantile plot in the center left panel shows that the normality assumption seems to be met.

2. Generating Potential Outliers

Using the BodyFat2 data set, run a regression model of PctBodyFat2 on Abdomen, Weight, Wrist, and Forearm.

a. Use plots to identify potential influential observations based on the suggested cutoff values.

- Reopen the last task by right-clicking in it in the Project Tree and selecting Modify....
- With <u>Plots</u> selected at the left, check the boxes that are checked below in the Custom plots area.

Linear Regression91 for Local:SASUSER.FITNESS					
Data Model	Plots				
Statistics Plots Predictions Titles Properties	 Show plots for regression analysis All appropriate plots for the current data selection Custom list of plots 				
	Custom plots: Histogram plot of the residuals Residuals by predicted values plot Studentized residuals by predicted values plot Observed by Predicted values plot Plot Cook's D statistic Studentized residuals by leverage plot Normal quantile plot of the residuals Residual-Fit plot Box plot of the residuals Diagnostic plots DFFITS plots DFFITS plots Scatter plot with regression line				

- With <u>**Predictions**</u> selected at the left:
 - Check the box for **Original sample** under Data to predict.
 - Check <u>Predictions</u> and <u>Diagnostic statistics</u> under Save output data.
 - Check the box for <u>Residuals</u> under Additional statistics.
- Click Save and do not replace the results from the previous run.

• Right-click the saved task icon in the Project Tree and select Add as Code Template.

Englishing Chapter 5 Exercises	on11	
	2	Open 🕨
	▶	Run Linear Regression111
Task List	R	Modify Linear Regression111
🔁 🕞 🛛 🍪		Select Input Data
Tasks by Category	Ē	Publish
Data		Add as Code Template
Filter and Sort		Create Task vemplate
En Ouer Puilder	Z)	Create Stored Process

• Edit the code template in the PROC REG section by adding the option (LABEL) at the end of each PLOTS(ONLY) line and the statement ID CASE; immediately after the next semi-colon.

• Click \triangleright Run above the code window.

There are only a modest number of observations further than 2 standard error units from the mean of 0.

There are 10 labeled outliers, but observation 39 is clearly the most extreme.

The same observations are shown to be influential by the DFFITS statistic.

DFBETAS are particularly high for observation 39 on the parameters for weight and forearm circumference.

3. Assessing Collinearity

Using the **BodyFat2** data set, run a regression of **PctBodyFat2** on all the other numeric variables in the file.

- **a.** Determine whether there is a collinearity problem.
 - Open the BodyFat2 data set.
 - Select <u>Analyze</u> \Rightarrow <u>Regression</u> \Rightarrow <u>Linear Regression...</u>.

Analyze 👻 Export 👻 Send To 👻 📝						
	ANOVA 🕨) Weight	🔞 Height	
	Regression	•	1	🗵 Linear Regression		
	Multivariate 🕨 🕨			Nonlinear Regression		
	Survival Analysis 🔹 🕨			Logistic Regression		
	Capability		ղի	Generalized Linear Models		
			Z4:	210.23	74.73	
	Control Charts		26	181	69.75	
lín I	Pareto Chart		25	176	72.5	
			25	191	74	
	Time Series	•	23	198.25	73.5	
*	Model Scoring		26	186.25	74.5	

• Drag **PctBodyFat2** to the dependent variable task role and all other continuous variables shown to the explanatory variables task role.

• With <u>Statistics</u> selected at the left, check the box for <u>Variance inflation values</u> in the Diagnostics area.

Data Model	Statistics	
Statistics Plots Predictions Titles Properties	Details on estimates Standardized regression coefficients Sum of squares, Type 1 Sum of squares, Type 2 Correlation matrix of estimates Covariance matrix of estimates Confidence limits for parameter estimates Confidence level: 95% Correlations Partial correlations Semi-partial correlations	Diagnostics Collinearity analysis Collinearity analysis without the intercept Tolerance values for estimates Variance inflation values Heteroscedasticity test Asymptotic covariance matrix Durbin-Watson statistic
Click Run		

Linear Regression Results

The REG Procedure Model: Linear_Regression_Model Dependent Variable: PctBodyFat2

Number of Observations Read252Number of Observations Used252

Analysis of Variance							
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F		
Model	13	13159	1012.22506	54.50	<.0001		
Error	238	4420.06401	18.57170				
Corrected Total	251	17579					

Root MSE	4.30949	R-Square	0.7486
Dependent Mean	19.15079	Adj R-Sq	0.7348
Coeff Var	22.50293		

Parameter Estimates							
		Parameter	Standard			Variance	
Variable	DF	Estimate	Error	t Value	Pr > t	Inflation	
Intercept	1	-21.35323	22.18616	-0.96	0.3368	0	
Age	1	0.06457	0.03219	2.01	0.0460	2.22447	
Weight	1	-0.09638	0.06185	-1.56	0.1205	44.65251	
Height	1	-0.04394	0.17870	-0.25	0.8060	2.93911	
Neck	1	-0.47547	0.23557	-2.02	0.0447	4.43192	
Chest	1	-0.01718	0.10322	-0.17	0.8679	10.23469	
Abdomen	1	0.95500	0.09016	10.59	<.0001	12.77553	
Hip	1	-0.18859	0.14479	-1.30	0.1940	14.54193	
Thigh	1	0.24835	0.14617	1.70	0.0906	7.95866	
Knee	1	0.01395	0.24775	0.06	0.9552	4.82530	
Ankle	1	0.17788	0.22262	0.80	0.4251	1.92410	
Biceps	1	0.18230	0.17250	1.06	0.2917	3.67091	
Forearm	1	0.45574	0.19930	2.29	0.0231	2.19193	
Wrist	1	-1.65450	0.53316	-3.10	0.0021	3.34840	

There seems to be high collinearity with Weight and less so with Hip, Abdomen, Chest, and Thigh.

b. If so, decide what you would like to do about that. Will you remove any variables? Why or why not?

The answer is not so easy. True, **Weight** is collinear with some set of the other variables, but as you have seen before in your model-building process, **Weight** actually ends up as a relatively significant predictor in the "best" models.