SASEG 8A – 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

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**.

Produce the default graphics.

- 1. Create a new process flow and rename it **SASEG8A**.
- 2. Open the **FITNESS** data set.
- 3. Select <u>Analyze</u> \Rightarrow <u>Regression</u> \Rightarrow <u>Linear Regression...</u>.

Anal	Analyze 👻 Export 👻 Send To 💌 📝						
	ANOVA F		🛛 😥 Oxygen_Consumption 😥 Run_F				
	Regression +		🗾 Linear Regression				
	Multivariate	•	Monlinearkegression				
	Survival Analysis	•	Logistic Regression				
	Capability	×	Generalized Linear Models				
	Control Charts	•	48.67				
lín	Pareto Chart		82 45.44				
	Time Series		-08 50.55				
	Time Denes		P1 46.67				
-	Model Scoring		45.31				

4. Drag **Oxygen_Consumption** to the dependent variable task role and **RunTime** to the explanatory variables task role.

5. In order to visually check the assumption of constant variance, we use the **<u>Plots</u>**.

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>Histogram Plot of the residuals</u>, <u>Residuals by predicted values plot</u>, <u>Residual plots</u>, and <u>Scatter</u> <u>Plot with residual line</u>

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:	
	Select all Enables you to choose which plots to include in the output. Scatter plot with regression line creates a plot of the regression line overlayed on a scatter plot of the data.	
Preview code	Run 🔻 Save Cancel Help	,

6. Click Run

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	Squares	Square	F Value	Pr > F	
Model	1	633.01458	633.01458	84.00	<.0001	
Error	29	218.53997	7.53586			
Corrected Total	30	851.55455				

Root MSE	2.74515	R-Square	0.7434
Dependent Mean	47.37581	Adj R-Sq	0.7345
Coeff Var	5.79442		

Parameter Estimates							
Veriable	DE	Parameter	Standard	4. V I	D- 5 141		
variable	DF	Estimate	Error	t value	Pr > t		
Intercept	1	82.42494	3.85582	21.38	<.0001		
RunTime	1	-3.31085	0.36124	-9.17	<.0001		

The plots produced are displayed below.

Note Chapter (14.8) discussion as to how these plots can be used to detect problems (if any).



The histogram of residuals helps you to **find outliers** and is one method to help with assessing the **normality assumption**.



The **Residual by Predicted Value 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 of constant error variance. 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 discussions, but on a smaller scale.

- The plot of the **residuals by the predictor values** (same as the plot on p. 13) in the model shows no patterns or trends. Again, this lends support to the <u>validity of the constant variance assumption</u> for this regression model. 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.
- The plot of the **standardized residuals (RStudent) by the Predicted Value** in the model also shows no patterns or trends; we expect 95% of the residuals to be between -2 and +2. This lends support to the validity of the assumption that the **error terms are distributed normally** for this regression model.

- The plot of the **residuals by the Quartile** in the model shows the normal probability plot; we expect the plotted points to cluster closely around the 45 degree line to indicate a normal distribution. The plot of the residuals against the normal quantiles is the 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. In this plot we see them closely clustered to the line; this lends support to the validity of the assumption that the **error terms are distributed normally** for this regression model.

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.



The plot of the **Residuals by the independent variable (RunTime)** does indicate randomness with no patterns or trends (see discussion on p. 7). 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.



The plot of the **Dependent Variable (Oxygen_Consumption) by independent variable (RunTime)** shows us a plot of the values in addition to the regression line. This plot allows us to not only check the form of the model (linear in this case) but allows us to check for **outliers**. Observations that are far away from the bulk of your data are outliers. There appears to be no outliers by viewing this plot.