

General Linear Models

Summary

The **General Linear Models** procedure is designed to construct a statistical model describing the impact of one or more factors X on one or more dependent variables Y . The factors may be:

1. quantitative or categorical
2. crossed or nested
3. fixed or random

Errors are assumed to follow a normal distribution. Weights may be supplied if a weighted least squares solution is desired. The output includes a wide variety of tables and graphs, including response surface plots, residual plots, and a MANOVA if more than one dependent variable is entered.

Many different types of experimental studies may be analyzed using this procedure. It includes as special cases models that can be estimated by the *Multiple Regression*, *Oneway ANOVA*, *Multifactor ANOVA*, and *Variance Components* procedures. In addition, it can analyze mixed models that cannot be handled by any of the above procedures.

Sample StatFolio: *glm.sgp*

Sample Data:

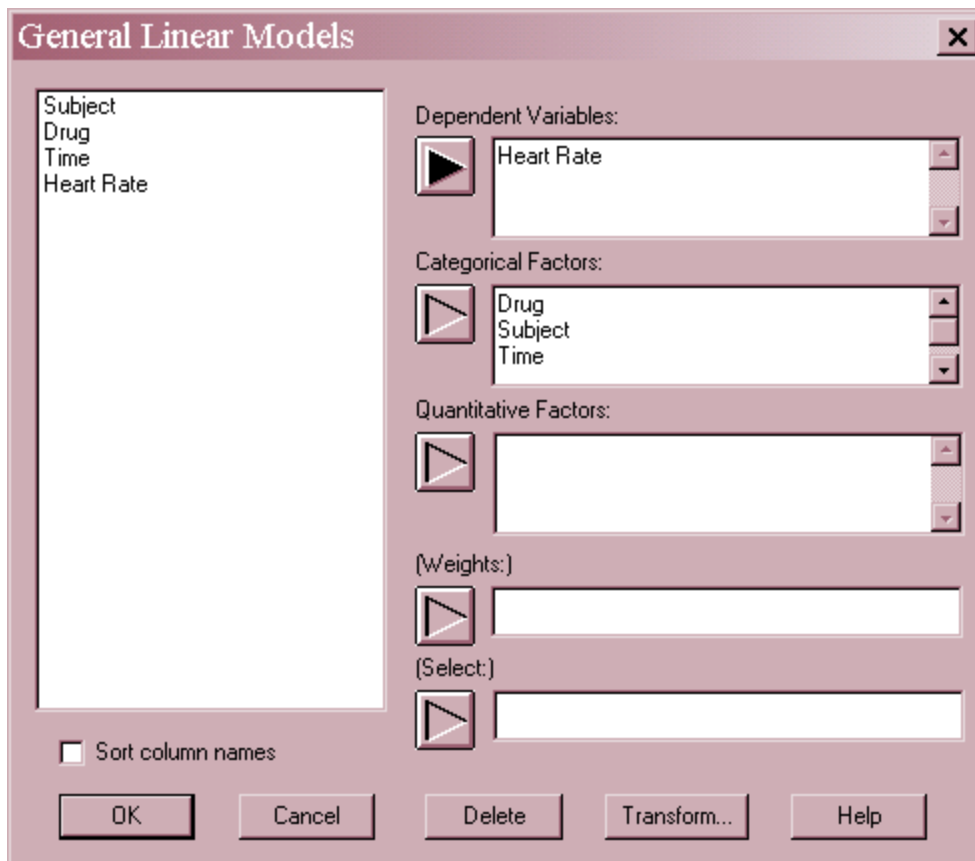
The sample data that will be analyzed is a repeated measures study from Milliken and Johnson (1996). In this study, 2 experimental drugs and a control were each administered to 8 subjects (for a total of 24 subjects). The heart rates of the subjects were measured at 4 different times after the drugs were administered. The data are contained in the file *heartrate.sgd*, a portion of which is shown below:

<i>Subject</i>	<i>Drug</i>	<i>Time</i>	<i>Heart Rate</i>
1	AX23	T1	72
1	AX23	T2	86
1	AX23	T3	81
1	AX23	T4	77
2	BWW9	T1	85
2	BWW9	T2	86
2	BWW9	T3	83
2	BWW9	T4	80
3	CONTROL	T1	69
3	CONTROL	T2	73
3	CONTROL	T3	72
3	CONTROL	T4	74
4	AX23	T1	78
4	AX23	T2	83
4	AX23	T3	88
4	AX23	T4	81
...

Since each of the subjects was given a different drug, *Subject* is said to be “nested” within *Drug*. It is a “repeated measures” experiment since measurements were taken for each subject-drug combination at multiple times.

Data Input

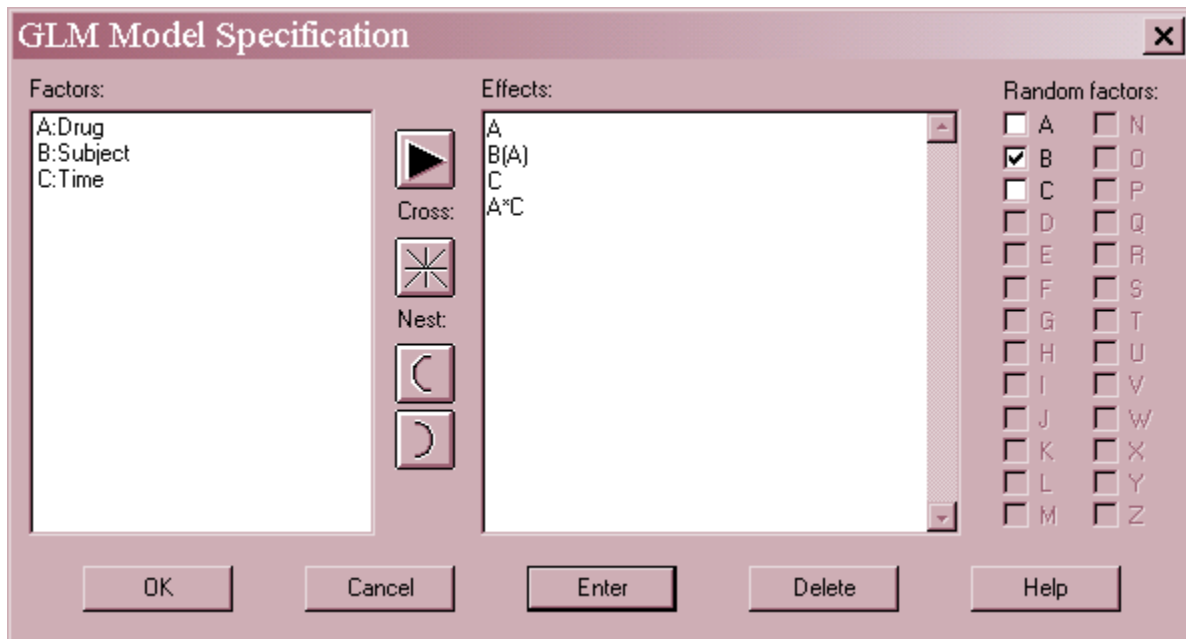
The first of two data input dialog boxes requests the names of the columns containing the dependent variables *Y* and the independent variables *X*:



- **Y:** one or more numeric columns containing the n observations for the dependent variables *Y*. If more than one column is entered, separate models will be fit for each one. In addition, a MANOVA may be requested.
- **Categorical factors:** numeric or non-numeric columns containing the n levels of any non-quantitative factors *X*.
- **Quantitative factors:** numeric columns containing the n values of any quantitative factors *X*.
- **Weight:** an optional numeric column containing n weights w_i to be applied to the squared residuals when performing a weighted least squares fit. In cases where the variance of *Y* is known to vary, the weights should be inversely proportional to those variances. If nothing is specified in this field, all $w_i = 1$.
- **Select:** subset selection.

In the sample study, there is one response and three categorical factors.

The second dialog box is used to specify the model to be fit to the data:



- **Factors:** Each of the categorical and quantitative factors is assigned a letter between A and Z.
- **Effects:** The effects to be included in the model are specified using the letters assigned to the factors. Effects are entered as follows:
 1. *Main effects for crossed factors* - Enter a single letter such as A.
 2. *Interactions between crossed factors* - Enter a term such as A*C to include the interaction between factors A and C or A*B*C to specify a 3-factor interaction.
 3. *Effects of nested factors* - Enter a term such as B(A) if factor B is nested within factor A or C(B A) if factor C is nested within combinations of factors A and B.
 4. *First order effects of quantitative factors* - Enter a single letter such as A.
 5. *Second order effects of quantitative factors* - Enter a term such as A*A for the quadratic effect of A or A*B for a cross-product.
- **Random Factors:** Categorical factors may be either *Fixed* or *Random*. A factor is *Random* if its levels consist of a random sample of levels from a population of possible levels. A factor is *Fixed* if its levels are selected by a nonrandom process or if its levels consist of the entire population of possible levels.

The effects specified on the dialog box above are:

A: the main effects of *Drug*. *Drug* is a fixed factor, since the effects of the specific drugs tested are to be estimated.

B(A): the effects of *Subject*, nested within *Drug*. *Subject* is nested within *Drug*, since different subjects were given each drug. *Subject* is also a random factor, since the 24 subjects selected are a random sample from the population of interest, which consists of everyone who might take those drugs in the future.

C: the main effects of *Time*. *Time* is a fixed factor, since the effects at specific times are to be estimated.

A*C: the interactions between *Drug* and *Time*. This term will allow the *Time* effect to be different for the 3 levels of *Drug*.

Analysis Summary

The *Analysis Summary* shows information about the fitted model. The top section of the output is shown below:

<u>General Linear Models</u>					
Number of dependent variables: 1					
Number of categorical factors: 3					
Number of quantitative factors: 0					
Analysis of Variance for Heart Rate					
<i>Source</i>	<i>Sum of Squares</i>	<i>Df</i>	<i>Mean Square</i>	<i>F-Ratio</i>	<i>P-Value</i>
Model	4487.94	32	140.248	18.83	0.0000
Residual	469.219	63	7.44792		
Total (Corr.)	4957.16	95			
Type III Sums of Squares					
<i>Source</i>	<i>Sum of Squares</i>	<i>Df</i>	<i>Mean Square</i>	<i>F-Ratio</i>	<i>P-Value</i>
Drug	1333.0	2	666.5	5.99	0.0088
Subject(Drug)	2337.91	21	111.329	14.95	0.0000
Time	289.615	3	96.5382	12.96	0.0000
Drug*Time	527.417	6	87.9028	11.80	0.0000
Residual	469.219	63	7.44792		
Total (corrected)	4957.16	95			

Included in the output are:

- **Analysis of Variance:** a decomposition of the sum of squares for *Y* into components for the model and for the residuals. The F-test tests the statistical significance of the model as a whole. A small P-value (less than 0.05 if operating at the 5% significance level) indicates that at least one factor in the model is significantly related to the dependent variable. In the current example, the model is highly significant.
- **Type III Sums of Squares:** decomposition of the model sum of squares into components for each factor. Based on the settings specified on the *Analysis Options* dialog box, either *Type III* or *Type I* sums of squares are displayed. Type III sums of squares test the marginal significance of each factor, assuming it was the last to be entered into the model. Type I sums of squares test the significance of the effects in the order they were added to the model. Small P-values indicate significant effects. In this example, all four effects are highly significant.

The second section of the analysis is important if the experiment contains any random factors.

Expected Mean Squares	
Source	EMS
Drug	(5)+4.0(2)+Q1
Subject(Drug)	(5)+4.0(2)
Time	(5)+Q2
Drug*Time	(5)+Q3
Residual	(5)

F-Test Denominators			
Source	Df	Mean Square	Denominator
Drug	21.00	111.329	(2)
Subject(Drug)	63.00	7.44792	(5)
Time	63.00	7.44792	(5)
Drug*Time	63.00	7.44792	(5)

Variance Components	
Source	Estimate
Subject(Drug)	25.9702
Residual	7.44792

It includes:

- Expected Mean Squares:** The expected mean square for each factor is determined using Hartley’s (1967) *synthesis* method. The mean squares in the earlier *Sums of Squares* table are labeled from top to bottom as (1) for *Drug*, (2) for *Subject* within *Drug*, and so on through (5) for the *Residuals*. A term such as *Q1* indicates a quantity unique to the factor in which it appears. The expected mean squares are important in constructing proper F-tests for models involving random factors.
- F-Test Denominators:** the mean square used as the denominator of the F-test for each factor, together with its degrees of freedom and how it was determined. For example, the F-test for *Drug* uses mean square (2) in its denominator, which equates to using *Subject(Drug)* as the error term.
- Variance Components:** for models with random factors, estimates the variance component σ_j of each random effect. The components are derived by equating the mean squares with their expected values, which is referred to as the method of moments. Variance components measure the variability in the response induced by variation in the random factors. For example, the variance in heart rate among persons given the same drug at the same time is estimated to be approximately 26.0.

The final section of the table shows statistics calculated from the fitted model:

R-Squared = 90.5345 percent
 R-Squared (adjusted for d.f.) = 85.7267 percent
 Standard Error of Est. = 2.72909
 Mean absolute error = 1.78841
 Durbin-Watson statistic = 2.23373 (P=0.1049)

Residual Analysis

	Estimation	Validation
N	96	
MSE	7.44792	
MAE	1.78841	
MAPE	2.38762	
ME	3.70074E-16	
MPE	-0.0906573	

The output displays:

- **Statistics:** summary statistics for the fitted model, including:

R-squared - represents the percentage of the variability in Y which has been explained by the fitted regression model, ranging from 0% to 100%. It is calculated by:

$$R^2 = 100 \left(1 - \frac{SS_{error}}{SS_{total}} \right) \% \tag{1}$$

For the sample data, the regression has accounted for about 90.5% of the variability in the heart rates. The remaining 9.5% is attributable to deviations from the model, which may be due to other factors, to measurement error, or to a failure of the current model to fit the data adequately.

Adjusted R-Squared – the R-squared statistic, adjusted for the number of coefficients in the model:

$$R^2_{adj} = 100 \left[1 - \left(\frac{n-1}{n-p} \right) \frac{SS_{error}}{SS_{total}} \right] \% \tag{2}$$

where *p* is the number of estimated model coefficients. This value is often used to compare models with different numbers of coefficients.

Standard Error of Est. – the estimated standard deviation of the residuals (the deviations around the model):

$$\hat{\sigma} = \sqrt{MSE} \tag{3}$$

This value is used to create prediction limits for new observations.

Mean Absolute Error – the average absolute value of the residuals:

$$MAE = \frac{\sum_{i=1}^n |e_i|}{n} \tag{4}$$

This value indicates the average error in predicting the response using the fitted model.

Durbin-Watson Statistic – a measure of serial correlation in the residuals:

$$DW = \frac{\sum_{i=1}^{n-1} (e_{i+1} - e_i)^2}{\sum_{i=1}^n e_i^2} \tag{5}$$

If the residuals vary randomly, this value should be close to 2. A small P-value indicates a non-random pattern in the residuals. For data recorded over time, a small P-value could indicate that some trend over time has not been accounted for. In the current example, the P-value is greater than 0.05, so there is not a significant correlation at the 5% significance level.

- **Residual Analysis:** if a subset of the rows in the datasheet have been excluded from the analysis using the *Select* field on the data input dialog box, the fitted model is used to make predictions of the *Y* values for those rows. This table shows statistics on the prediction errors, defined by

$$e_i = y_i - \hat{y}_i \tag{6}$$

Included are the mean squared error:

$$MSE = \frac{\sum_{i=1}^n e_i^2}{n - 1} \tag{7}$$

the mean absolute error:

$$MAE = \frac{\sum_{i=1}^n |e_i|}{n} \tag{8}$$

the mean absolute percentage error:

$$MAPE = \frac{100 \sum_{i=1}^n |e_i| / y_i}{n} \% \tag{9}$$

the mean error:

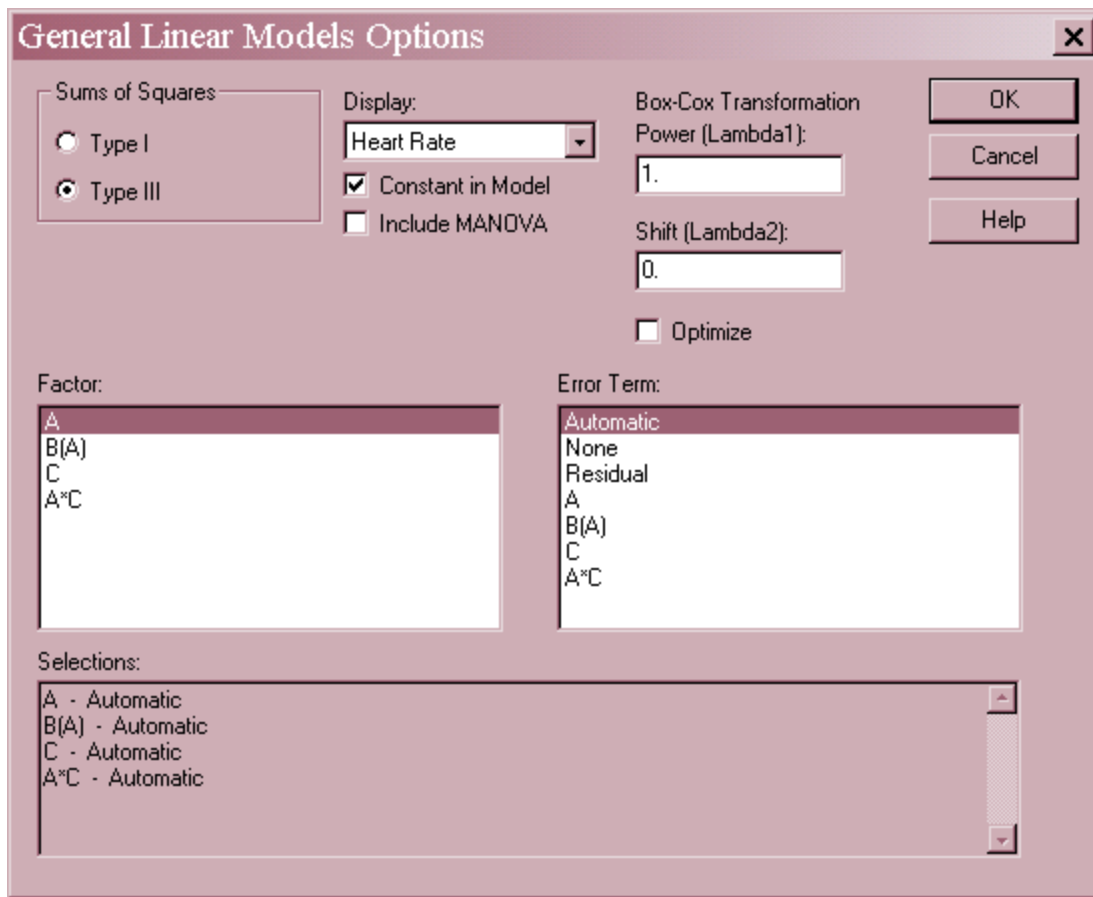
$$ME = \frac{\sum_{i=1}^n e_i}{n} \tag{10}$$

and the mean percentage error:

$$MPE = \frac{100 \sum_{i=1}^n e_i / y_i}{n} \% \tag{11}$$

The validation statistics can be compared to the statistics for the fitted model to determine how well that model predicts observations outside of the data used to fit it.

Analysis Options



- **Sums of Squares:** the sums of squares to display. Type I sums of squares measure the contribution of each variable to the model when added in the order indicated. Type III sums of squares measure the marginal contribution of each effect, assuming it was added last.
- **Display:** if more than one dependent variable has been specified, the variable to use when creating plots and table that display only a single variable.

- **Constant in model:** If this option is not checked, the constant term β_0 will be omitted from the model. Removing the constant term allows for regression through the origin.
- **Include MANOVA:** If more than one dependent variable has been specified, checking this box will cause a multivariate analysis of variance to be included in the *Analysis Summary*. For more information, see the example late in this document.
- **Box-Cox Transformation:** If selected, a Box-Cox transformation will be applied to the dependent variable(s). Box-Cox transformations are a way of dealing with situations in which the deviations from the regression model do not have a constant variance. You may specify the Box-Cox parameters or request that the program automatically find the optimal power. For details, see the *Box-Cox Transformations* documentation.
- **Factor and Error Term:** the denominator to be used for each factor when creating an F-test. The *Automatic* option causes the program to select the proper denominator automatically. You can override the program's selections by clicking on a factor and then clicking on the desired error term. The current error terms are displayed in the *Selections* field.

Model Coefficients

Underlying the analysis is a linear statistical model of the form

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_{p-1} X_{p-1} + \varepsilon_i \quad (12)$$

where Y is the dependent variable, the X's carry information about each of the effects in the model, and the ε 's are assumed to be normally and independently distributed with a mean of 0. The *Model Coefficients* pane displays the estimated coefficients, standard errors, lower and upper confidence limits, and variance inflation factors:

95.0% confidence intervals for coefficient estimates (Heart Rate)					
		Standard			
Parameter	Estimate	Error	Lower Limit	Upper Limit	V.I.F.
CONSTANT	76.4063	0.278536	75.8496	76.9629	
Drug	-0.125	0.39391	-0.912167	0.662167	1.33333
Drug	4.625	0.39391	3.83783	5.41217	1.33333
Subject(Drug)	2.71875	1.27641	0.168036	5.26946	1.75
Subject(Drug)	2.46875	1.27641	-0.0819642	5.01946	1.75
Subject(Drug)	0.09375	1.27641	-2.45696	2.64446	1.75
Subject(Drug)	6.21875	1.27641	3.66804	8.76946	1.75
Subject(Drug)	1.96875	1.27641	-0.581964	4.51946	1.75
Subject(Drug)	-4.90625	1.27641	-7.45696	-2.35554	1.75
Subject(Drug)	0.96875	1.27641	-1.58196	3.51946	1.75
Subject(Drug)	-7.53125	1.27641	-10.082	-4.98054	1.75
Subject(Drug)	15.3438	1.27641	12.793	17.8945	1.75
Subject(Drug)	0.46875	1.27641	-2.08196	3.01946	1.75
Subject(Drug)	1.71875	1.27641	-0.831964	4.26946	1.75
Subject(Drug)	5.59375	1.27641	3.04304	8.14446	1.75
Subject(Drug)	-4.28125	1.27641	-6.83196	-1.73054	1.75
Subject(Drug)	-0.28125	1.27641	-2.83196	2.26946	1.75
Subject(Drug)	-1.15625	1.27641	-3.70696	1.39446	1.75
Subject(Drug)	3.21875	1.27641	0.668036	5.76946	1.75
Subject(Drug)	1.46875	1.27641	-1.08196	4.01946	1.75
Subject(Drug)	-8.65625	1.27641	-11.207	-6.10554	1.75
Subject(Drug)	-5.53125	1.27641	-8.08196	-2.98054	1.75
Subject(Drug)	-0.28125	1.27641	-2.83196	2.26946	1.75
Subject(Drug)	-1.65625	1.27641	-4.20696	0.894464	1.75
Time	-1.40625	0.482439	-2.37033	-0.442171	1.5
Time	2.55208	0.482439	1.588	3.51616	1.5
Time	0.635417	0.482439	-0.328663	1.5995	1.5
Drug*Time	-4.375	0.682272	-5.73841	-3.01159	2.0
Drug*Time	1.66667	0.682272	0.303253	3.03008	2.0
Drug*Time	4.08333	0.682272	2.71992	5.44675	2.0
Drug*Time	2.125	0.682272	0.761586	3.48841	2.0
Drug*Time	0.416667	0.682272	-0.946747	1.78008	2.0
Drug*Time	-3.04167	0.682272	-4.40508	-1.67825	2.0

The model can get quite complicated, particularly when categorical factors are involved. It includes a term for each degree of freedom associated with the effects. Except in simple cases, it is not expected that the user will calculate values using the model, since the *Reports* pane will create predictions for any combination of the factors.

- **Parameter:** the estimated model coefficients. The columns of X are defined as follows:
 1. *Constant:* X contains a columns of 1's.
 2. *Main effect of a quantitative factor:* X contains the value of the independent variable.
 3. *Main effects of a categorical factor:* For a factor with k levels, X contains k-1 indicator variables. The first variable equals 1 when the factor is at its first level, -1 when the factor is at its last level, and 0 otherwise. The second variable equals 1 when the factor is at its second level, -1 when the factor is at its last level, and 0 otherwise. Etc.
 4. *Interactions between factors:* X contains the product of the columns created for those factors.

For example, the equation for the first subject that was given the first drug at the first time in the above table is:

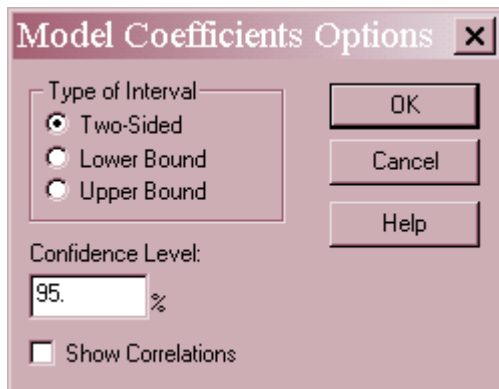
$$Time = 76.4063 - 0.125(1) + 2.71875(1) - 1.40625(1) - 4.375(1) = 73.2188$$

The equation for the first subject that was given the last drug at the first time is:

$$Time = 76.4063 - 0.125(-1) + 4.625(-1) + 0.09375(1) - 1.40625(1) - 4.375(-1) + 2.125(-1) = 72.8438$$

- **Standard errors:** estimated standard errors for each of the model coefficients.
- **Confidence Limits:** two-sided confidence limits or one-sided confidence bounds for the model coefficients.
- **V.I.F.:** variance inflation factors. The variance inflation factors measure how large the variance of the coefficients is compared to what it would be if the independent variables were uncorrelated. Values greater than 10.0 usually indicate serious multicollinearity amongst the predictor variables, which leads to imprecise estimates of the model coefficients.

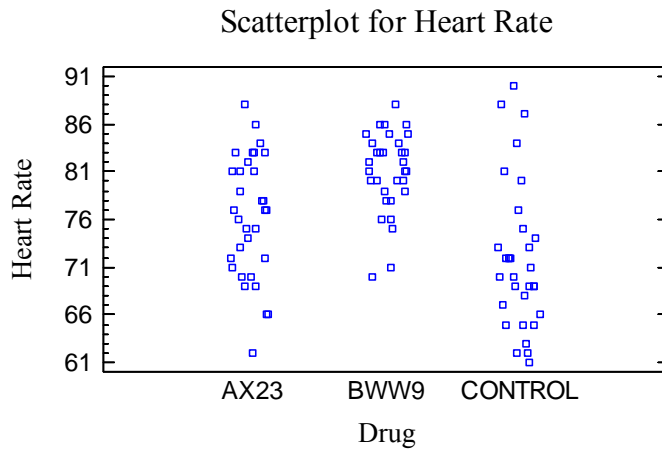
Pane Options



- **Type of Interval:** Select either two-sided confidence limits or one-sided confidence bounds.
- **Confidence Level:** percentage used for the limits or bounds.
- **Show Correlations:** If selected, a table of estimated correlations between the model coefficients will be displayed. This table can be helpful in determining how well the effects of different independent variables have been separated from each other.

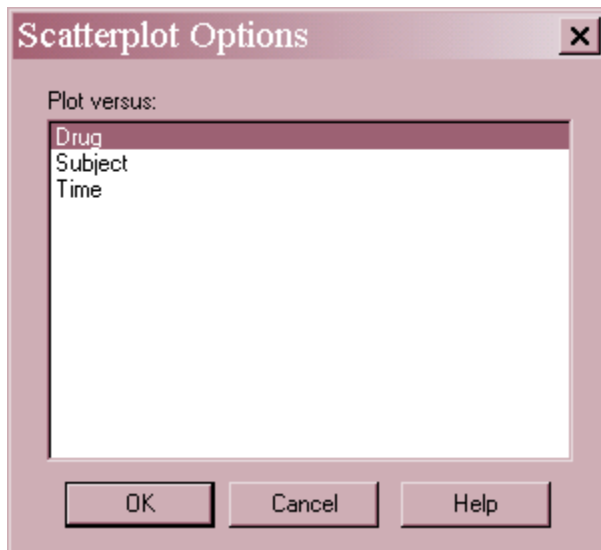
Scatterplot

The *Scatterplot* plots the observations versus any one of the selected factors.



It is often helpful to jitter the points in the horizontal direction by pressing the *Jitter* button on the analysis toolbar, as in the above plot. Jittering offsets each point by a random amount to prevent the points from falling exactly on top of each other.

Pane Options



- **Plot versus:** the factor to plot on the horizontal axis.

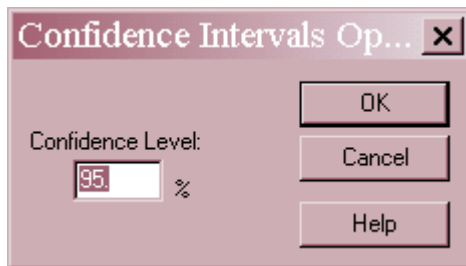
Table of Means

This table displays the least squares means for each level of the factors and for pairs of levels for any included two-factor interactions. Least squares means represent the predicted mean value of *Y* at a specified level of a categorical factor *X* when all quantitative variables are set equal to their observed means and all indicator variables for other categorical factors are set equal to 0. Each mean is shown together with its estimated standard error and a confidence interval:

Table of Least Squares Means for Heart Rate with 95.0 Percent Confidence Intervals					
			<i>Std.</i>	<i>Lower</i>	<i>Upper</i>
<i>Level</i>	<i>Count</i>	<i>Mean</i>	<i>Error</i>	<i>Limit</i>	<i>Limit</i>
GRAND MEAN	96	76.4063	0.278536	75.8496	76.9629
Drug					
AX23	32	76.2813	1.86522	72.4023	80.1602
BWW9	32	81.0313	1.86522	77.1523	84.9102
CONTROL	32	71.9063	1.86522	68.0273	75.7852
Subject within Drug					
1 AX23	4	79.0	1.36454	76.2732	81.7268
2 BWW9	4	83.5	1.36454	80.7732	86.2268
3 CONTROL	4	72.0	1.36454	69.2732	74.7268
4 AX23	4	82.5	1.36454	79.7732	85.2268
5 BWW9	4	83.0	1.36454	80.2732	85.7268
6 CONTROL	4	67.0	1.36454	64.2732	69.7268
7 AX23	4	77.25	1.36454	74.5232	79.9768
8 BWW9	4	73.5	1.36454	70.7732	76.2268
9 CONTROL	4	87.25	1.36454	84.5232	89.9768
10 AX23	4	76.75	1.36454	74.0232	79.4768
11 BWW9	4	82.75	1.36454	80.0232	85.4768
12 CONTROL	4	77.5	1.36454	74.7732	80.2268
13 AX23	4	72.0	1.36454	69.2732	74.7268
14 BWW9	4	80.75	1.36454	78.0232	83.4768
15 CONTROL	4	70.75	1.36454	68.0232	73.4768
16 AX23	4	79.5	1.36454	76.7732	82.2268
17 BWW9	4	82.5	1.36454	79.7732	85.2268
18 CONTROL	4	63.25	1.36454	60.5232	65.9768
19 AX23	4	70.75	1.36454	68.0232	73.4768
20 BWW9	4	80.75	1.36454	78.0232	83.4768
21 CONTROL	4	70.25	1.36454	67.5232	72.9768
22 AX23	4	72.5	1.36454	69.7732	75.2268
23 BWW9	4	81.5	1.36454	78.7732	84.2268
24 CONTROL	4	67.25	1.36454	64.5232	69.9768
Time					
T1	24	75.0	0.557073	73.8868	76.1132
T2	24	78.9583	0.557073	77.8451	80.0716
T3	24	77.0417	0.557073	75.9284	78.1549
T4	24	74.625	0.557073	73.5118	75.7382
Drug by Time					
AX23 T1	8	70.5	0.964878	68.5718	72.4282
AX23 T2	8	80.5	0.964878	78.5718	82.4282
AX23 T3	8	81.0	0.964878	79.0718	82.9282
AX23 T4	8	73.125	0.964878	71.1968	75.0532
BWW9 T1	8	81.75	0.964878	79.8218	83.6782
BWW9 T2	8	84.0	0.964878	82.0718	85.9282
BWW9 T3	8	78.625	0.964878	76.6968	80.5532
BWW9 T4	8	79.75	0.964878	77.8218	81.6782
CONTROL T1	8	72.75	0.964878	70.8218	74.6782
CONTROL T2	8	72.375	0.964878	70.4468	74.3032
CONTROL T3	8	71.5	0.964878	69.5718	73.4282
CONTROL T4	8	71.0	0.964878	69.0718	72.9282

For example, the mean heart rate of subjects given drug AX23 at time T1 is estimated to be between 68.6 and 72.4, with 95% confidence.

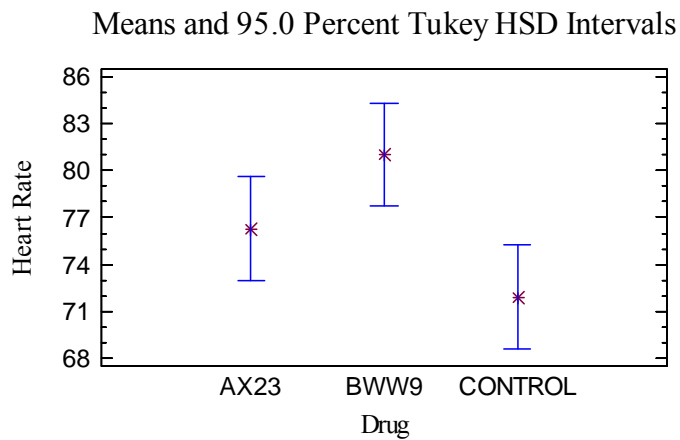
Pane Options



- **Confidence Level:** the level of confidence associated with each interval.

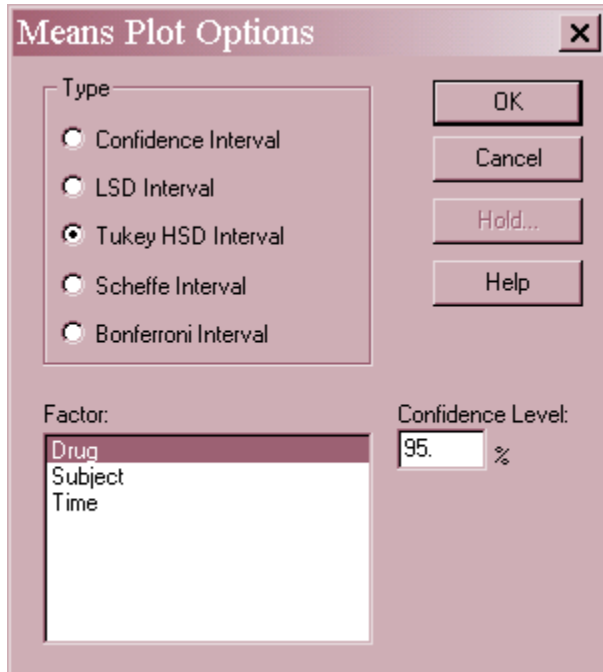
Means Plot

The level means for a selected factor may be plotted using the *Means Plot*.



If the factor plotted on the horizontal axis is categorical, then the plot shows the least squares means with uncertainty intervals. The type of interval displayed depends on the settings in *Pane Options*. If the factor on the horizontal axis is quantitative, the plot displays the fitted model with all other quantitative factors set equal to their observed means and all categorical indicator variables set equal to 0.

Provided all of the sample sizes are the same (or close), the analyst can determine which level means of a categorical factor are significantly different from which others using the LSD, Tukey, Scheffe, or Bonferroni procedure simply by looking at whether or not a pair of intervals overlap in the vertical direction. A pair of intervals that do not overlap indicates a statistically significant difference between the means at the selected confidence level. In this case, note that the interval for drug BWW9 does not overlap the interval for CONTROL, indicating a statistically significant difference between the means at those two levels. The intervals for AX23 and CONTROL overlap, however, so they cannot be declared to be significantly different.

Pane Options

- **Intervals:** the method used to construct the intervals.
- **Factor:** the factor to be plotted.
- **Confidence Level:** the level of confidence associated with each interval.

The type of intervals that may be selected are:

Confidence intervals - displays confidence intervals for the level means using the estimated standard errors.

LSD intervals - designed to compare any pair of means with the stated confidence level.

Tukey HSD Intervals - designed for comparing all pairs of means. The stated confidence level applies to the entire family of pairwise comparisons.

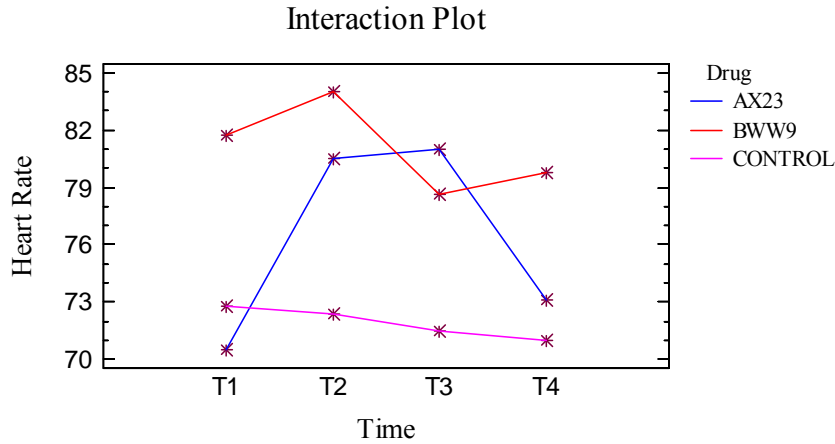
Scheffe Intervals - designed for comparing all contrasts. Not usually relevant here.

Bonferroni Intervals - designed for comparing a selected number of contrasts. Tukey's intervals are usually tighter.

Each of the intervals is formed by adding a multiple of the standard error of the least squares mean to the estimated mean. The multiple depends upon the method used, as described in the *Oneway ANOVA* documentation. The degrees of freedom are those associated with the estimate of the standard error and depend on the structure of the experiment.

Interaction Plot

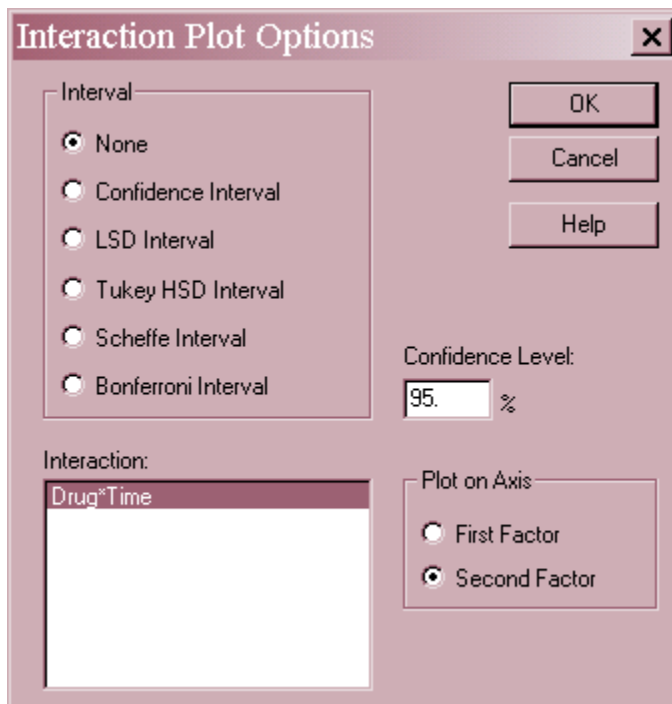
When one or more significant interactions exist amongst the categorical factors, the factors involved should be examined together using the *Interaction Plot*.



The interaction plot displays the least squares means at all combinations of two factors. If the factors do not interact, the lines on the plot should be approximately parallel. If they are not, then the effect of one factor depends upon the level of the other, which is the definition of an interaction.

Notice that the heart rate for the CONTROL group changes very little over time, while it shows significant changes for the other two drugs. In addition, drug BWW9 appears to have a quicker and more sustained effect than drug AX23.

Pane Options



- **Interval:** type of interval to be drawn around each mean. The interaction is treated as a factor with number of levels equal to the total number of plotted points.
- **Interaction:** interaction to plot.
- **Confidence Level:** percentage used to define the intervals.
- **Plot on Axis:** the factor used to define points along the horizontal axis. Lines will be drawn at each level of the other factor.

Multiple Range Tests

For factors that shows significant P-Values in the ANOVA table and which do not interact with other factors, a further analysis can be performed by selecting the *Multiple Range Tests*.

Multiple Comparisons for Heart Rate by Drug				
Method: 95.0 percent LSD				
Drug	Count	LS Mean	LS Sigma	Homogeneous Groups
CONTROL	32	71.9063	1.86522	X
AX23	32	76.2813	1.86522	XX
BWW9	32	81.0313	1.86522	X
Contrast	Sig.	Difference	+/- Limits	
AX23 - BWW9		-4.75	5.48564	
AX23 - CONTROL		4.375	5.48564	
BWW9 - CONTROL	*	9.125	5.48564	

* denotes a statistically significant difference.

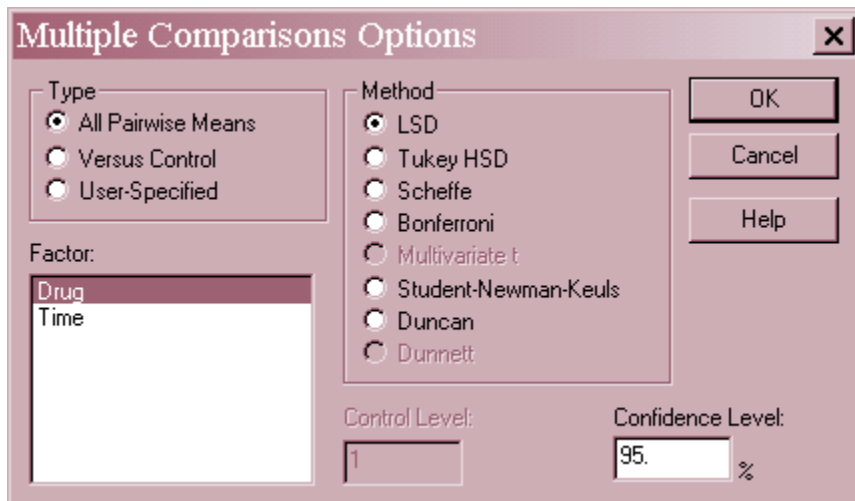
The top half of the table displays each of the estimated least squares means in increasing order of magnitude. It shows:

- **Count** - the number of observations at the specified level of the factor.
- **LS Mean** - the estimated least squares mean. In the case of a balanced design, the least squares mean is equivalent to the average of all observations at the indicated factor level. In unbalanced designs, the least squares mean is the predicted value of the dependent variable when the specified factor is set to a particular level while all other factors are set equal to their mean levels. The least squares means adjust for any imbalance in the data by making predictions at a common level of all the factors.
- **LS Sigma** – the estimated standard error of the least squares mean.
- **Homogeneous groups** - a graphical illustration of which means are significantly different from which others, based on the contrasts displayed in the second half of the table. Each column of X's indicates a group of means within which there are no statistically significant differences. In the example, there are 2 columns, each containing a pair of X's. It indicates that drug AX23 is not significantly different from either the CONTROL or from drug BWW9. However, since CONTROL and BWW9 are not within the same group anywhere, their means are significantly different.

The second half of the table displays a comparison between each pair of level means.

- **Difference** - the difference between the two least squares means.
- **Limits** - an interval estimate of that difference, using the currently selected multiple comparisons procedure.
- **Sig.** - An asterisk is placed next to any difference that is statistically significantly different from 0 at the currently selected significance level, i.e., any interval that does not contain 0.

Pane Options



- **Type:** type of contrasts to be created.
- **Factor:** factor to be analyzed.
- **Method:** the method used to make the multiple comparisons.
- **Control Level:** if *Type* is set to *Versus Control*, the number of the level against which all other levels will be compared.
- **Confidence Level:** the level of confidence used by the selected multiple comparison procedure.

The available methods are:

LSD - forms a confidence interval for each pair of means at the selected confidence level using Student's t distribution. This procedure is due to Fisher and is called the *Least Significant Difference* procedure, since the magnitude of the limits indicates the smallest difference between any two means that can be declared to represent a statistically significant difference. It should only be used when the F-test in the ANOVA table indicates significant differences amongst the level means. The probability of making a Type I error α applies to each pair of means separately. If making more than one comparison, the overall probability

of calling at least one pair of means significantly different when they are not may be considerably larger than α .

Tukey HSD - widens the intervals to allow for multiple comparisons amongst all pairs of means using Tukey's T. Tukey called his procedure the *Honestly Significant Difference* procedure since it controls the experiment-wide error rate at α . If all of the means are equal, the probability of declaring *any* of the pairs to be significantly different in the entire experiment equals α . Tukey's procedure is more conservative than Fisher's LSD procedure, since it makes it harder to declare any particular pair of means to be significantly different.

Scheffe - designed to permit the estimation of all possible contrasts amongst the sample means (not just pairwise comparisons).

Bonferroni - designed to permit the estimation of any preselected number of contrasts. These limits are usually wider than Tukey's limits when all pairwise comparisons are being made.

Multivariate t – designed for sets of linearly independent combinations of the means.

Student-Newman-Keuls - Unlike the previous methods, this method does not create intervals for the pairwise differences. Instead, it sorts the means in increasing order and then begins to separate them into groups according to values of the Studentized range distribution. Eventually, the means are separated into homogeneous groups within which there are no significant differences.

Duncan - similar to the Student-Newman-Keuls procedure, except that it uses a different critical value of the Studentized range distribution when defining the homogeneous groups. A detailed discussion of the Duncan and Student-Newman-Keuls procedures is given by Milliken and Johnson (1992).

Dunnett – designed for pairwise comparisons when one level is a control.

Example – User-Specified Contrasts

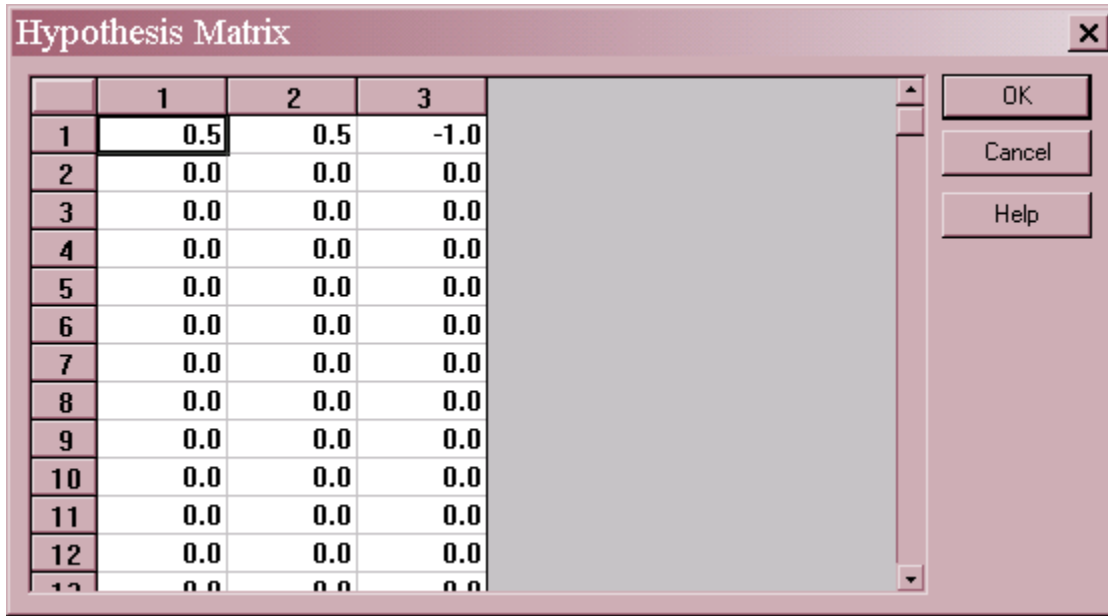
User-specified contrasts may be tested by setting *Type* to *User-Specified*. When *OK* is pressed, a small datasheet will be displayed on which to define the contrasts. Each row of the datasheet specifies the coefficients in the contrast

$$c_1\mu_1 + c_2\mu_2 + \dots + c_k\mu_k \quad (13)$$

where the coefficients c_j must sum to 1. For example, the datasheet below defines a contrast of the form

$$0.5\mu_1 + 0.5\mu_2 - \mu_3 \quad (14)$$

which contrasts the average response of the two experimental drugs to the control.



The resulting output displays each least squares mean and an interval estimate for the contrasts:

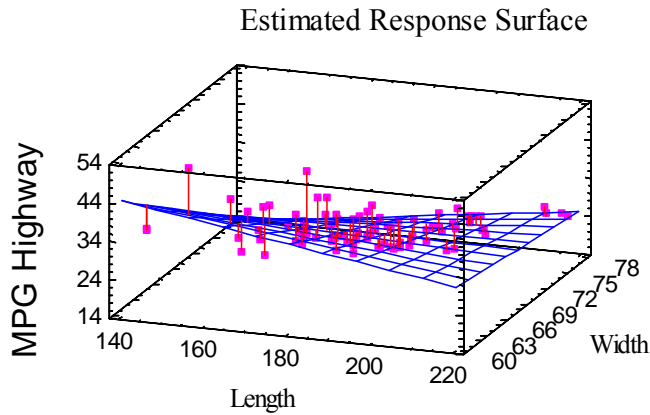
Multiple Comparisons for Heart Rate by Drug			
Method: 95.0 percent LSD			
Drug	Count	LS Mean	
AX23	32	76.2813	
BWW9	32	81.0313	
CONTROL	32	71.9063	
Contrast	Sig.	Estimate	+/- Limits
0.5 0.5 -1.0	*	6.75	4.75071

* denotes a statistically significant estimate.

If *LSD* is selected, the +/- *Limits* correspond to 95% confidence intervals for the desired contrasts.

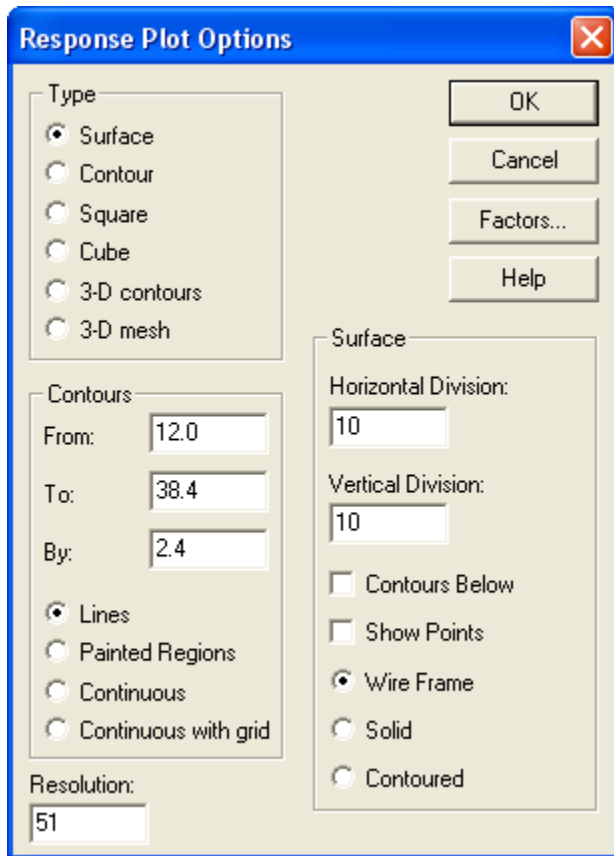
Surface and Contour Plots

If the model involves at least two quantitative factors, surface and contour plots can be created. For example, using the *93cars.sfb* dataset, the following plots displays a model for *MPG Highway* as a function of the *Length* and *Width* of the automobiles in that file.



The fitted model includes the main effects of both factors together with their interaction. Lines have been dropped from each point perpendicularly to the estimated model.

Pane Options



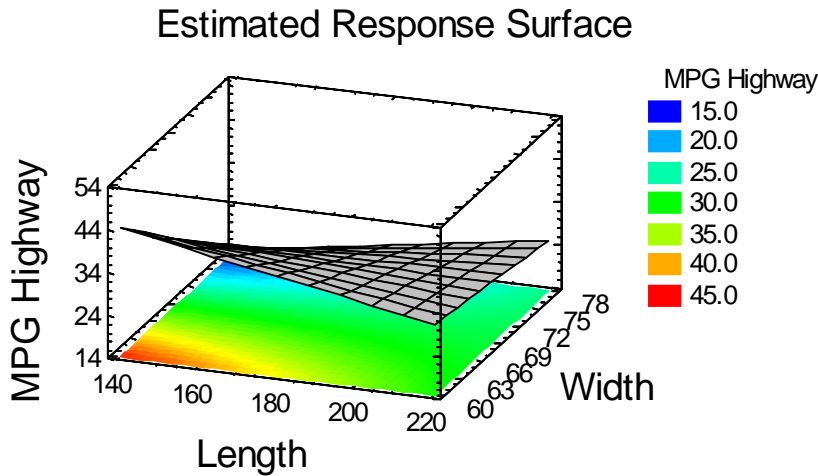
- **Type:** type of plot to display. The fitted model may be plotted as a 3-D *Surface* plot, a 2-D *Contour* plot, at each corner of a square, or at each corner of a cube (given at least 3 quantitative factors).
- **Contours From, To, and By:** defines the contour regions when contours are added to the plot. The contours may be drawn as solid *Lines*, *Painted Regions* of solid color, using a *Continuous* range of colors, or using as *Continuous with grid*.

- **Resolution:** the number of X and Y locations at which the function is evaluated when creating the plot. A larger resolution results in a smoother plot. You can set the default resolution using the *Preferences* selection on the *Edit* menu.
- **Surface Horizontal and Vertical Divisions:** the number of intervals between the grid lines along the X and Y axes.
- **Contours Below:** draws contours in the base of the cube when creating a surface plot.
- **Draw Points:** plots each observation and drops a vertical line to the surface.
- **Type:** the type of surface to be drawn:
 - **Wire frame:** a surface defined by grid lines only.
 - **Solid:** a surface defined by grid lines with a solid color between the lines.
 - **Contoured:** a surface with colored regions showing the value of the function.
 - **3-D contours:** a cube defined by 3 factors with contours on 3 faces.
 - **3-D mesh plot:** a cube in which the response is evaluated on a mesh throughout the cube.
- **Factors:** Press this button to set the limits for the factors on the plot and the values at which to fix other factors. The following dialog box will be displayed:

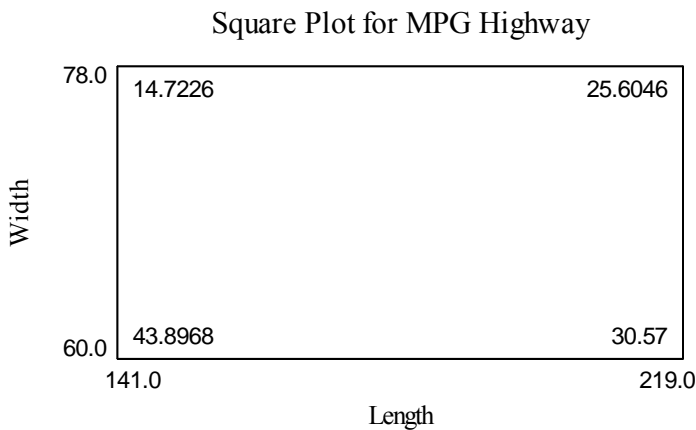
	Low	High	Hold
<input checked="" type="checkbox"/> Length	140.	220.	180.
<input checked="" type="checkbox"/> Width	60.	78.	69.
<input type="checkbox"/>			
<input type="checkbox"/>			
<input type="checkbox"/>			
<input type="checkbox"/>			
<input type="checkbox"/>			
<input type="checkbox"/>			
<input type="checkbox"/>			
<input type="checkbox"/>			
<input type="checkbox"/>			
<input type="checkbox"/>			
<input type="checkbox"/>			
<input type="checkbox"/>			
<input type="checkbox"/>			
<input type="checkbox"/>			
<input type="checkbox"/>			
<input type="checkbox"/>			
<input type="checkbox"/>			

- **Low and High:** plotting limits for the selected factors.
- **Hold:** values to fix other factors at when evaluating the fitted model.

Example: Surface Plot with Continuous Contours Below



Example: Square Plot



The values displayed at each corner of the square are the predicted values \hat{Y} .

Reports

The *Reports* pane displays predictions from the fitted least squares model. By default, the table includes a line for each row in the datasheet that has complete information on the X variables and a missing value for the Y variable. This allows you to add rows to the bottom of the datasheet corresponding to levels at which you want predictions without affecting the fitted model.

For example, suppose you wished to display the estimated values for each of the two experimental drugs at the four time periods. Additional rows would be added to the bottom of the datasheet as follows:

<i>Row</i>	<i>Subject</i>	<i>Drug</i>	<i>Time</i>	<i>Heart Rate</i>
97	0	AX23	T1	
98	0	AX23	T2	
99	0	AX23	T3	
100	0	AX23	T4	
101	0	BWW9	T1	
102	0	BWW9	T2	
103	0	BWW9	T3	
104	0	BWW9	T4	

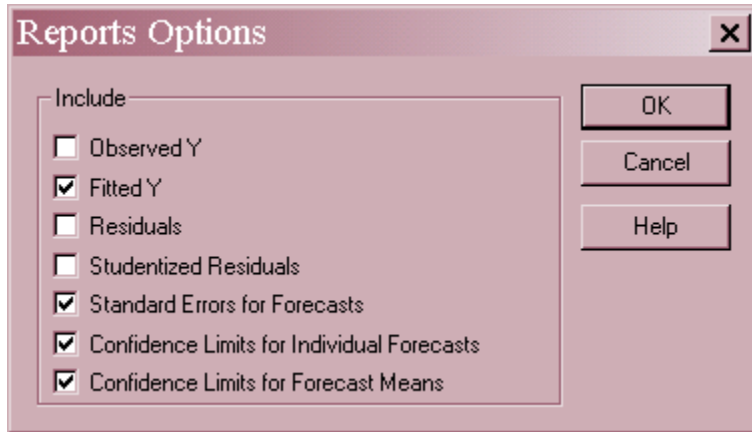
Subject is set to 0 so that all of the indicator variables for that factor will be set to 0, effectively averaging across all subjects. The resulting table is shown below:

Regression Results for Heart Rate						
	<i>Fitted</i>	<i>Std. Error</i>	<i>Lower 95.0% CL</i>	<i>Upper 95.0% CL</i>	<i>Lower 95.0% CL</i>	<i>Upper 95.0% CL</i>
<i>Row</i>	<i>Value</i>	<i>for Forecast</i>	<i>for Forecast</i>	<i>for Forecast</i>	<i>for Mean</i>	<i>for Mean</i>
97	70.5	2.89463	64.7155	76.2845	68.5718	72.4282
98	80.5	2.89463	74.7155	86.2845	78.5718	82.4282
99	81.0	2.89463	75.2155	86.7845	79.0718	82.9282
100	73.125	2.89463	67.3405	78.9095	71.1968	75.0532
101	81.75	2.89463	75.9655	87.5345	79.8218	83.6782
102	84.0	2.89463	78.2155	89.7845	82.0718	85.9282
103	78.625	2.89463	72.8405	84.4095	76.6968	80.5532
104	79.75	2.89463	73.9655	85.5345	77.8218	81.6782

The table displays:

- **Row** - the row number in the datasheet.
- **Fitted Value** - the predicted value of the dependent variable \hat{Y} using the fitted model.
- **Standard Error for Forecast** - the estimated standard error for predicting a single new observation.
- **Confidence Limits for Forecast** - prediction limits for new observations at the selected level of confidence.
- **Confidence Limits for Mean** - confidence limits for the mean value of Y at the selected level of confidence.

For example, an additional subject given drug BWW9 is likely to have a heart rate at time T1 between 76.0 and 87.5 (row #101). The 95% confidence interval for the mean heart rate of many subjects given that drug at that time runs from 79.8 to 83.7.

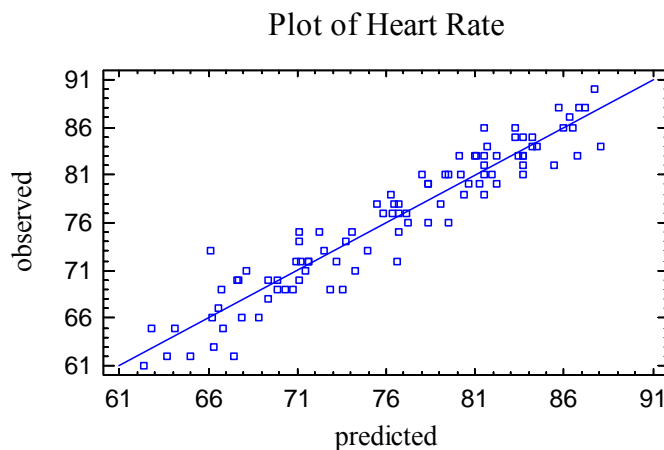
Pane Options

You may include:

- *Observed Y* – the observed values of the dependent variable.
- *Fitted Y* – the predicted values from the fitted model.
- *Residuals* – the ordinary residuals (observed minus predicted).
- *Studentized Residuals* – the Studentized deleted residuals as described below.
- *Standard Errors for Forecasts* – the standard errors for new observations at values of the independent variables corresponding to each row of the datasheet.
- *Confidence Limits for Individual Forecasts* – confidence intervals for new observations.
- *Confidence Limits for Forecast Means* – confidence intervals for the mean value of Y at values of the independent variables corresponding to each row of the datasheet.

Observed versus Predicted

The *Observed versus Predicted* plot shows the observed values of Y on the vertical axis and the predicted values \hat{Y} on the horizontal axis.



If the model fits well, the points should be randomly scattered around the diagonal line. Any change in variability from low values of Y to high values of Y might indicate the need to transform the dependent variable before fitting a model to the data.

Residual Plots

As with all statistical models, it is good practice to examine the residuals. In a regression, the residuals are defined by

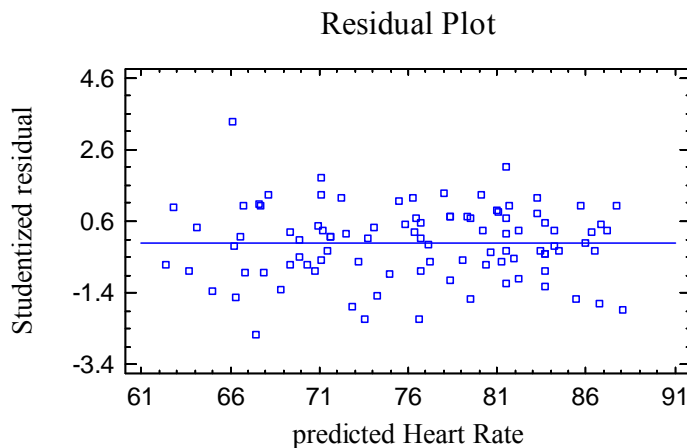
$$e_i = y_i - \hat{y}_i \quad (15)$$

i.e., the residuals are the differences between the observed data values and the fitted model.

The *General Linear Models* procedure creates various types of residual plots, depending on the settings in *Pane Options*.

Scatterplot versus Predicted Values

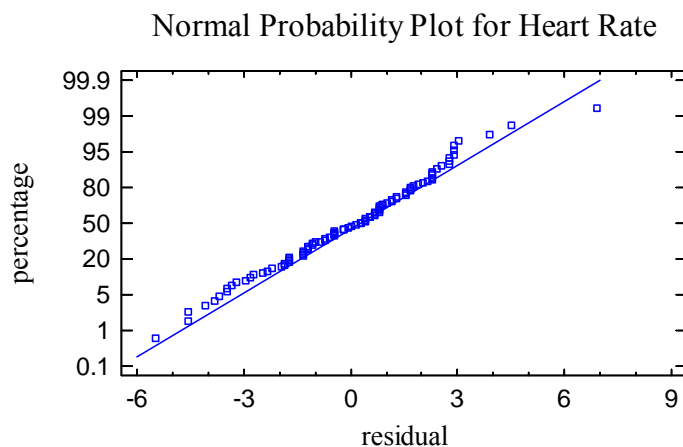
This plot is helpful in visualizing any possible dependence of the residual variance on the mean, which might necessitate a weighted least squares fit.



The above plot shows a fairly constant variance, although one possible outlier is evident.

Normal Probability Plot

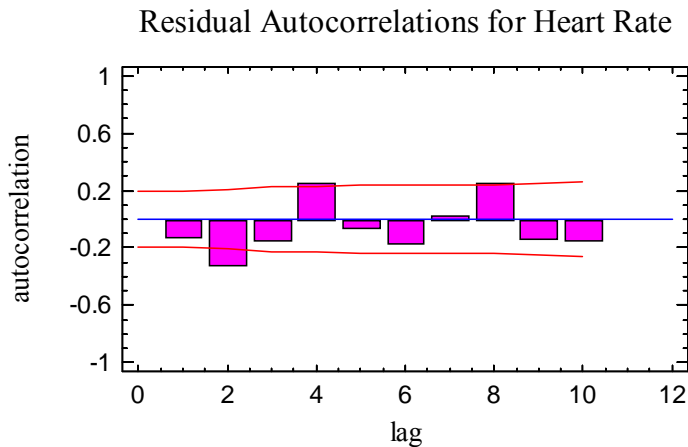
This plot can be used to determine whether or not the deviations around the line follow a normal distribution, which is the assumption used to form the prediction intervals.



If the deviations follow a normal distribution, they should fall approximately along a straight line. In the above plot, the points fall fairly close to the line.

Residual Autocorrelations

This plot calculates the autocorrelation between residuals as a function of the number of rows between them in the datasheet.



It is only relevant if the data have been collected sequentially. Any bars extending beyond the probability limits would indicate significant dependence between residuals separated by the indicated “lag”, which would violate the assumption of independence made when fitting the regression model.

Pane Options

Residual Plots Options

Plot

Residuals

Studentized Residuals

Direction

Horizontal

Vertical

OK

Cancel

Help

Type

Scatterplot

Normal Probability Plot

Autocorrelation Function

Fitted Line

None

Using Quartiles

Using Least Squares

Plot versus:

Predicted values

Row number

Drug

Subject

Time

Number of Lags:

Confidence Level:

- **Plot:** the type of residuals to plot:

1. *Residuals* – the residuals from the least squares fit.
2. *Studentized residuals* – the difference between the observed values y_i and the predicted values \hat{y}_i when the model is fit using all observations except the i -th, divided by their estimated standard error. These residuals are sometimes called

externally deleted residuals, since they measure how far each value is from the fitted model when that model is fit using all of the data except the point being considered. This is important, since a large outlier might otherwise affect the model so much that it would not appear to be unusually far away from the line.

- **Type:** the type of plot to be created. A *Scatterplot* is used to test for curvature. A *Normal Probability Plot* is used to determine whether the model residuals come from a normal distribution. An *Autocorrelation Function* is used to test for dependence between consecutive residuals.
- **Plot Versus:** for a *Scatterplot*, the quantity to plot on the horizontal axis.
- **Number of Lags:** for an *Autocorrelation Function*, the maximum number of lags. For small data sets, the number of lags plotted may be less than this value.
- **Confidence Level:** for an *Autocorrelation Function*, the level used to create the probability limits.

Unusual Residuals

Once the model has been fit, it is useful to study the residuals to determine whether any outliers exist that should be removed from the data. The *Unusual Residuals* pane lists all observations that have Studentized residuals of 2.0 or greater in absolute value.

		<i>Predicted</i>		<i>Studentized</i>
<i>Row</i>	<i>Y</i>	<i>Y</i>	<i>Residual</i>	<i>Residual</i>
22	62.0	67.4687	-5.46875	-2.58
24	73.0	66.0938	6.90625	3.37
40	69.0	73.5938	-4.59375	-2.14
48	72.0	76.5938	-4.59375	-2.14
53	86.0	81.4688	4.53125	2.10

Studentized residuals greater than 3 in absolute value correspond to points more than 3 standard deviations from the fitted model, which is a rare event for a normal distribution. Row #24 is more than 3.3 standard deviations from the fitted model, which is a very rare event if the deviations follow a normal distribution.

Note: Points can be removed from the fit while examining the *Scatterplot* by clicking on a point and then pressing the *Exclude/Include* button on the analysis toolbar. Excluded points are marked with an X.

Influential Points

In fitting a regression model, all observations do not have an equal influence on the parameter estimates in the fitted model. Points located at extreme values of X have greater influence than those located nearer to the center of the experimental region. The *Influential Points* pane displays any observations that have high influence on the fitted model:

Influential Points for Heart Rate				
		Mahalanobis		Cook's
Row	Leverage	Distance	DFITS	Distance
9	0.34375	48.2486	-1.27936	0.0479806
14	0.34375	48.2486	-1.23576	0.0449106
22	0.34375	48.2486	-1.86911	0.0971251
24	0.34375	48.2486	2.43976	0.154896
33	0.34375	48.2486	-1.3672	0.0544249
40	0.34375	48.2486	-1.54576	0.0685315
48	0.34375	48.2486	-1.54576	0.0685315
53	0.34375	48.2486	1.52322	0.0666794
81	0.34375	48.2486	1.30124	0.0495536

Average leverage of single data point = 0.34375

Points are placed on this list for one of the following reasons:

- **Leverage** – measures how distant an observation is from the mean of all n observations in the space of the *independent* variables. The higher the leverage, the greater the impact of the point on the fitted values \hat{y} . Points are placed on the list if their leverage is more than 3 times that of an average data point.
- **Mahalanobis Distance** – measures the distance of a point from the center of the collection of points in the multivariate space of the independent variables. Since this distance is related to *leverage*, it is not used to select points for the table.
- **DFITS** – measures the difference between the predicted values \hat{y}_i when the model is fit with and without the i -th data point. Points are placed on the list if the absolute value of DFITS exceeds $2p/\sqrt{n}$, where p is the number of coefficients in the fitted model.
- **Cook's Distance** – an overall measure of the influence of the i -th observation on the estimated coefficients. Points are placed on the list if the value is beyond the 50th percentile of an F distribution with p and $n - p$ degrees of freedom.

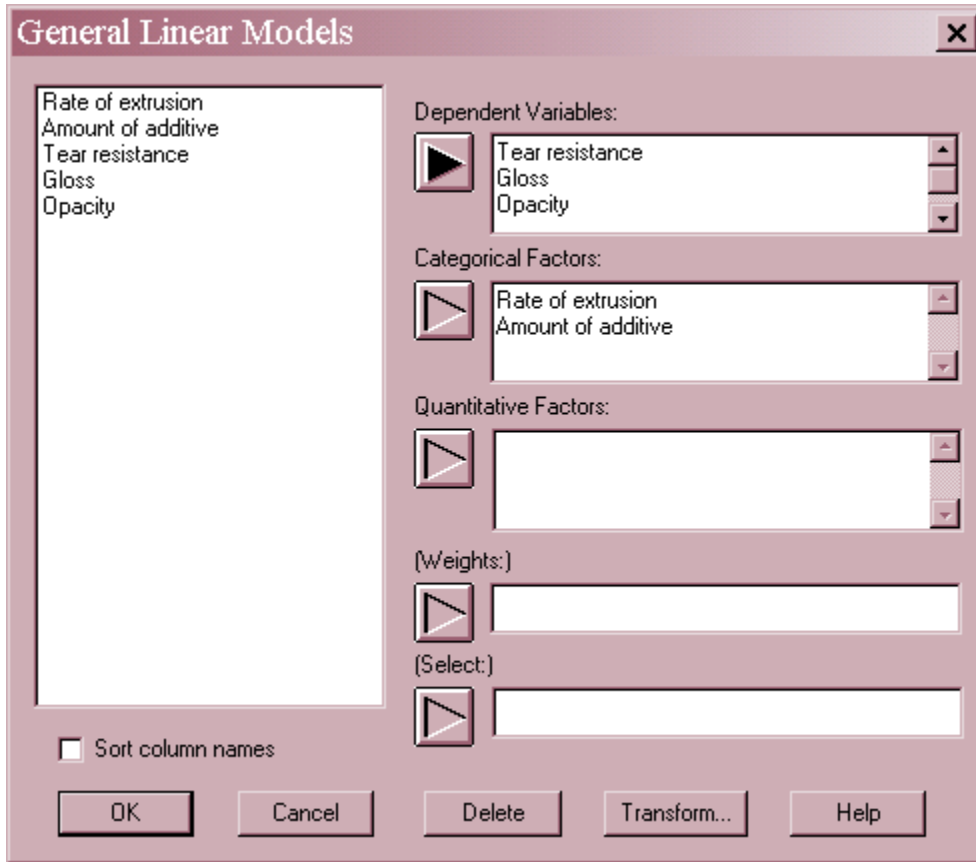
Because of the perfect balance in this design, all leverage values are equal. However, 9 points made the list because of a large value of DFITS, including all of the points previously identified as large residuals.

MANOVA

When more than one dependent variable is specified on the data input dialog box, a multivariate analysis of variance may be included if requested using *Analysis Options*. For example, consider the data from an experiment reported by Johnson and Wichern (2002) performed to determine the optimal conditions for extruding plastic film. Three response variables, *Tear resistance*, *Gloss*, and *Opacity* were measured at different levels of two factors, *Rate of Extrusion* and *Amount of additive*. The data is contained in the file *film.sfb*:

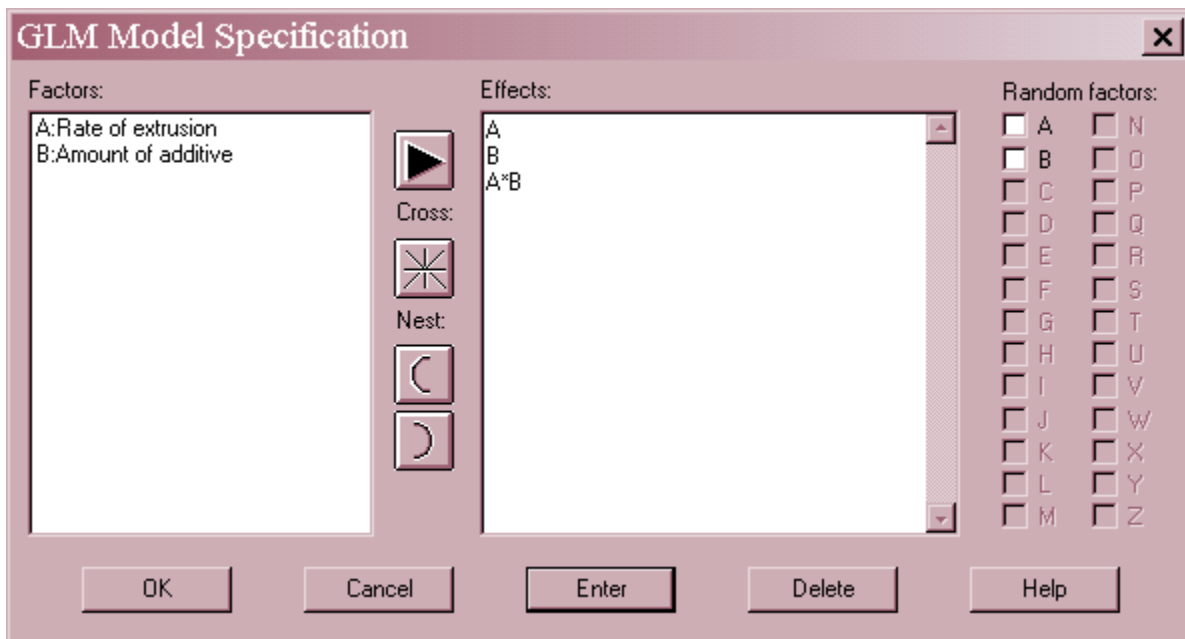
<i>Rate of Extrusion</i>	<i>Amount of Additive</i>	<i>Tear Resistance</i>	<i>Gloss</i>	<i>Opacity</i>
-10	1	6.5	9.5	4.4
-10	1	6.2	9.9	6.4
-10	1	5.8	9.6	3
-10	1	6.5	9.6	4.1
-10	1	6.5	9.2	0.8
-10	1.5	6.9	9.1	5.7
-10	1.5	7.2	10	2
-10	1.5	6.9	9.9	3.9
-10	1.5	6.1	9.5	1.9
-10	1.5	6.3	9.4	5.7
10	1	6.7	9.1	2.8
10	1	6.6	9.3	4.1
10	1	7.2	8.3	3.8
10	1	7.1	8.4	1.6
10	1	6.8	8.5	3.4
10	1.5	7.1	9.2	8.4
10	1.5	7.0	8.8	5.2
10	1.5	7.2	9.7	6.9
10	1.5	7.5	10.1	2.7
10	1.5	7.6	9.2	1.9

The data input dialog box specifies the names of the three response variables and the two factors:



Since the factors are each at only two levels, they can be entered as either categorical factors or quantitative factors.

The specified model includes both main effects and a two-factor interaction:



For multiple dependent variables, the *Analysis Summary* includes separate analyses for each response. If requested on the *Analysis Options* dialog box, a MANOVA will also be performed. The additional output from that analysis is shown below:

MANOVA for A			
Wilks' lambda = 0.381858 F = 7.55427 P-value = 0.00303404			
Pillai trace = 0.618142 F = 7.55427 P-value = 0.00303404			
Hotelling-Lawley trace = 1.61877 F = 7.55427 P-value = 0.00303404			
Roy's greatest root = 1.61877 s = 1 m = 0.5 n = 6.0			
Hypothesis Matrix H			
	<i>Tear resistance</i>	<i>Gloss</i>	<i>Opacity</i>
Tear resistance	1.7405	-1.5045	0.8555
Gloss	-1.5045	1.3005	-0.7395
Opacity	0.8555	-0.7395	0.4205
Error Matrix E			
	<i>Tear resistance</i>	<i>Gloss</i>	<i>Opacity</i>
Tear resistance	1.764	0.02	-3.07
Gloss	0.02	2.628	-0.552
Opacity	-3.07	-0.552	64.924
MANOVA for B			
Wilks' lambda = 0.523035 F = 4.25562 P-value = 0.0247453			
Pillai trace = 0.476965 F = 4.25562 P-value = 0.0247453			
Hotelling-Lawley trace = 0.911918 F = 4.25562 P-value = 0.0247453			
Roy's greatest root = 0.911918 s = 1 m = 0.5 n = 6.0			
Hypothesis Matrix H			
	<i>Tear resistance</i>	<i>Gloss</i>	<i>Opacity</i>
Tear resistance	0.7605	0.6825	1.9305
Gloss	0.6825	0.6125	1.7325
Opacity	1.9305	1.7325	4.9005
Error Matrix E			
	<i>Tear resistance</i>	<i>Gloss</i>	<i>Opacity</i>
Tear resistance	1.764	0.02	-3.07
Gloss	0.02	2.628	-0.552
Opacity	-3.07	-0.552	64.924
MANOVA for A*B			
Wilks' lambda = 0.777106 F = 1.33852 P-value = 0.301782			
Pillai trace = 0.222894 F = 1.33852 P-value = 0.301782			
Hotelling-Lawley trace = 0.286826 F = 1.33852 P-value = 0.301782			
Roy's greatest root = 0.286826 s = 1 m = 0.5 n = 6.0			
Hypothesis Matrix H			
	<i>Tear resistance</i>	<i>Gloss</i>	<i>Opacity</i>
Tear resistance	0.0005	0.0165	0.0445
Gloss	0.0165	0.5445	1.4685
Opacity	0.0445	1.4685	3.9605
Error Matrix E			
	<i>Tear resistance</i>	<i>Gloss</i>	<i>Opacity</i>
Tear resistance	1.764	0.02	-3.07
Gloss	0.02	2.628	-0.552
Opacity	-3.07	-0.552	64.924

For each effect, the table shows four statistics designed to test whether or not there are significant overall effects due to that factor. The statistics are based on the matrices of sums of squares and cross-products attributable to the hypothesized effects (**H**) and to the residuals (**E**). The statistics displayed are:

- *Wilks' lambda*: a statistic based on the ratio of two determinants

$$\Lambda^* = \frac{|E|}{|E + H|} \quad (16)$$

- *Pillai Trace*: a statistic calculated from

$$\text{tr}[H(H + E)^{-1}] \quad (17)$$

- *Hotelling-Lawley Trace*: a statistic calculated from

$$\text{tr}[HE^{-1}] \quad (18)$$

- *Roy's Greatest Root*: a statistic equal to

$$\frac{\eta_1}{1 + \eta_1} \quad (19)$$

where η_1 is the largest eigenvalue of HE^{-1} .

The output line for Roy's statistic also displays the values of s , m , and n , three values used to calculate the F tests for the other statistics. It is worth noting that the tests are exact if $s = 1$ or 2 and approximate otherwise.

The first three statistics are shown together with the results of F-tests. Small P-Values (less than 0.05 if operating at the 5% significance level) indicate significant effects. In the example, the main effects of both factors are statistically significant at the 5% level, but the interaction is not.

Save Results

The following results may be saved to the datasheet:

1. *Predicted Values* – the predicted value of Y corresponding to each of the n observations.
2. *Standard Errors of Predictions* - the standard errors for the n predicted values.
3. *Lower Limits for Predictions* – the lower prediction limits for each predicted value.
4. *Upper Limits for Predictions* – the upper prediction limits for each predicted value.
5. *Standard Errors of Means* - the standard errors for the mean value of Y at each of the n values of X.
6. *Lower Limits for Forecast Means* – the lower confidence limits for the mean value of Y at each of the n values of X.
7. *Upper Limits for Forecast Means*– the upper confidence limits for the mean value of Y at each of the n values of X.
8. *Residuals* – the n residuals.
9. *Studentized Residuals* – the n Studentized residuals.
10. *Leverages* – the leverage values corresponding to the n values of X.
11. *DFITS Statistics* – the value of the DFITS statistic corresponding to the n values of X.
12. *Mahalanobis Distances* – the Mahalanobis distance corresponding to the n values of X.
13. *Cook's Distances* – Cook's distance corresponding to the n values of X.
14. *Coefficients* – the estimated model coefficients.

Calculations**Regression Model**

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_{p-1} X_{p-1} \quad (20)$$

Error Sum of Squares

$$\text{Unweighted: } SSE = \sum_{i=1}^n \left(y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i - \hat{\beta}_2 x_2 - \dots - \hat{\beta}_{p-1} x_{p-1} \right)^2 \quad (21)$$

$$\text{Weighted: } SSE = \sum_{i=1}^n w_i \left(y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i - \hat{\beta}_2 x_2 - \dots - \hat{\beta}_{p-1} x_{p-1} \right)^2 \quad (22)$$

Coefficient Estimates

$$\hat{\beta} = (X'WX)^{-1}(X'WY) \quad (23)$$

$$s^2 \{\hat{\beta}\} = MSE(X'WX)^{-1} \quad (24)$$

$$MSE = \frac{SSE}{n - p} \quad (25)$$

where $\hat{\beta}$ is a column vector containing the estimated regression coefficients, X is an (n, p) matrix containing a 1 in the first column (if the model contains a constant term) and the settings of the predictor variables in the other columns, Y is a column vector with the values of the dependent variable, and W is an (n, n) diagonal matrix containing the weights w_i on the diagonal for a weighted regression or 1's on the diagonal if weights are not specified. A modified sweep algorithm is used to solve the equations after centering and rescaling of the independent variables.

Analysis of Variance

With constant term:

Source	Sum of Squares	Df	Mean Square	F-Ratio
Model	$SSR = b'X'WY - \frac{\left(\sum_{i=1}^n w_i y_i\right)^2}{\sum_{i=1}^n w_i}$	p-1	$MSR = \frac{SSR}{p-1}$	$F = \frac{MSR}{MSE}$
Residual	$SSE = Y'WY - b'X'WY$	n-p	$MSE = \frac{SSE}{n-p}$	
Total (corr.)	$SSTO = \sum_{i=1}^n w_i (y_i - \bar{y})^2$	n-1		

Without constant term:

Source	Sum of Squares	Df	Mean Square	F-Ratio
Model	$SSR = b'X'WY$	p	$MSR = \frac{SSR}{p}$	$F = \frac{MSR}{MSE}$
Residual	$SSE = Y'WY - b'X'WY$	n-p	$MSE = \frac{SSE}{n-p}$	
Total	$SSTO = Y'WY$	n		

R-Squared

$$R^2 = 100 \left(\frac{SSR}{SSR + SSE} \right) \% \tag{26}$$

Adjusted R-Squared

$$R_{adj}^2 = 100 \left[1 - \left(\frac{n-1}{n-p} \right) \frac{SSE}{SSR + SSE} \right] \% \tag{27}$$

Std.Error of Est.

$$\hat{\sigma} = \sqrt{MSE} \quad (28)$$

Residuals

$$e_i = y_i - \hat{\beta}_0 - \hat{\beta}_1 x_1 - \dots - \hat{\beta}_{p-1} x_{p-1} \quad (29)$$

Mean Absolute Error

$$MAE = \frac{\sum_{i=1}^n w_i |e_i|}{\sum_{i=1}^n w_i} \quad (30)$$

Durbin-Watson Statistic

$$D = \frac{\sum_{i=2}^n (e_i - e_{i-1})^2}{\sum_{i=1}^n e_i^2} \quad (31)$$

If $n > 500$, then

$$D^* = \frac{|D - 2|}{\sqrt{4/n}} \quad (32)$$

is compared to a standard normal distribution. For $100 < n \leq 500$, $D/4$ is compared to a beta distribution with parameters

$$\alpha = \beta = \frac{n-1}{2} \quad (33)$$

For smaller sample sizes, $D/4$ is compared to a beta distribution with parameters which are based on the trace of certain matrices related to the X matrix, as described by Durbin and Watson (1951) in section 4 of their classic paper.

Leverage

$$h_i = \text{diag} \{ X_i' (X'WX)^{-1} X_i \} w_i \quad (34)$$

$$\bar{h} = \frac{p}{n} \quad (35)$$

Studentized Residuals

$$d_i = \frac{e_i \sqrt{w_i}}{\sqrt{MSE_i(1-h_i)}} \quad (36)$$

Mahalanobis Distance

$$MD_i = \left(\frac{h_i - w_i / \sum_{i=1}^n w_i}{1 - h_i} \right) \frac{n(n-2)}{n-1} \quad (37)$$

DFITS

$$DFITS_i = \frac{d_i}{\sqrt{w_i}} \sqrt{\left(\frac{h_i}{1-h_i} \right)} \quad (38)$$

Cook's Distance

$$CD_i = \frac{e_i^2}{pMSE} \left[\frac{h_i}{(1-h_i)^2} \right] \quad (39)$$

Standard Error for Forecast

$$s\{Y_{h(new)}\} = \sqrt{MSE \left(1 + X_h' (X'WX)^{-1} X_h \right)} \quad (40)$$

Confidence Limit for Forecast

$$\hat{Y}_h \pm t_{\alpha/2, n-p} s\{Y_{h(new)}\} \quad (41)$$

Confidence Limit for Mean

$$\hat{Y}_h \pm t_{\alpha/2, n-p} \sqrt{MSE(X_h'(X'WX)^{-1}X_h)} \quad (42)$$