**Paper 407-2013**

# Statistical Analyses Using SAS® Enterprise Guide®

R. Scott Leslie, MedImpact Healthcare Systems, Inc.

## ABSTRACT

Conducting statistical analyses involves choosing proper methods, understanding model assumptions and displaying clear results.  The latest releases of SAS® Enterprise Guide® offer conveniences, such as point-and-click wizards and integrated syntax help, to ease the burden on users.  This tutorial demonstrates how to perform statistics quickly and easily using some handy features of SAS Enterprise Guide.  Examples of multiple linear regression, logistic regression and survival analysis are covered as well as some hints on how to navigate Enterprise Guide menus. This tutorial is intended for SAS users with beginning to intermediate experience with the above mentioned statistics and basic experience with SAS Enterprise Guide.

## INTRODUCTION

This paper was created for the Tutorial section of 2011 Western Users of SAS Software (WUSS) conference.  The main purpose of this tutorial is to demonstrate ways to conduct statistical analyses using the many available tasks and functions in Enterprise Guide.  It is intended for SAS programmers with less experience with Enterprise Guide but knowledge of commonly used statistical methods.  One common misconception of Enterprise Guide is that users cannot perform all the functions available in Display Manager, or the traditional SAS BASE application.  This paper demonstrates that Enterprise Guide leverages the powerful capabilities of SAS combined with a user friendly interface.

SAS® Enterprise Guide® is described ([EG Fact Sheet](#)) as "an easy-to-use interface" that allows users "to access the power of SAS data access, reporting and analytics, enabling you to quickly create and publish SAS Stored Processes that can be leveraged from Microsoft Office applications or Web browsers."

This figure taken from the SAS website is a visualization of the description above.  It accurately displays how the powerful SAS engine is under the hood of the nice graphical user interface of Enterprise Guide allowing users to perform jobs quickly and easily.
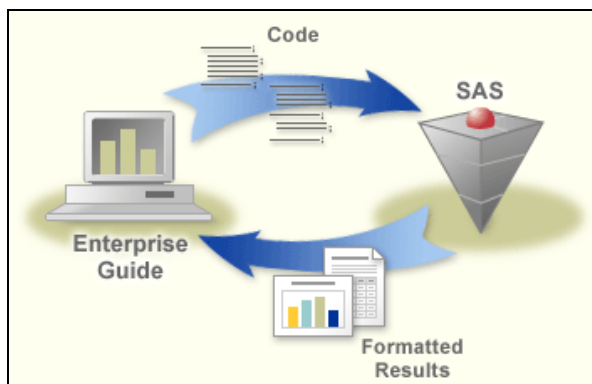


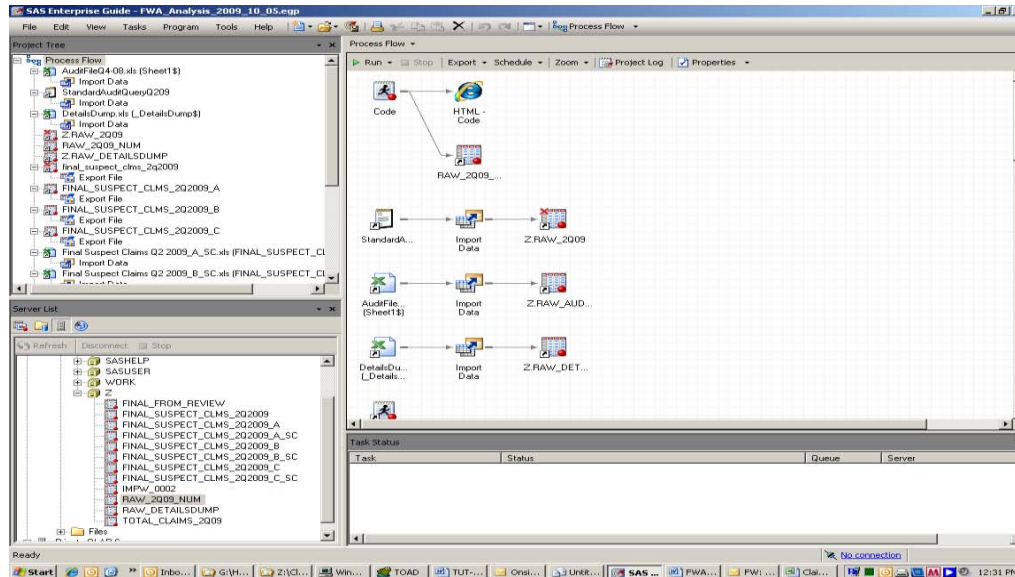**Figure 1. SAS Enterprise Guide**

## GETTING FAMILIAR WITH THE ENTERPRISE GUIDE INTERFACE

This section is a brief overview on the Enterprise Guide interface.  For those familiar with the look of Enterprise Guide, please look forward to the next sections of this paper.

At first glance, the application looks very menu driven and can be unpleasant to traditional SAS programmers.  It took me a while to reconcile the log and output windows.  After time, I came to realize that the interface has many more features than the Display Manager seen in SAS.  The key to getting accustomed to Enterprise Guide is learning how to work with process flows.

First, all work is organized in one file called a **Project** with the file extension ".epg".  A sample project is shown below.  All code, data, log and output are contained in the project and represented as **Nodes** (the icons) shown in the **Process Flow** window below.  So the process flow represents your desktop that charts all your work.  I find this
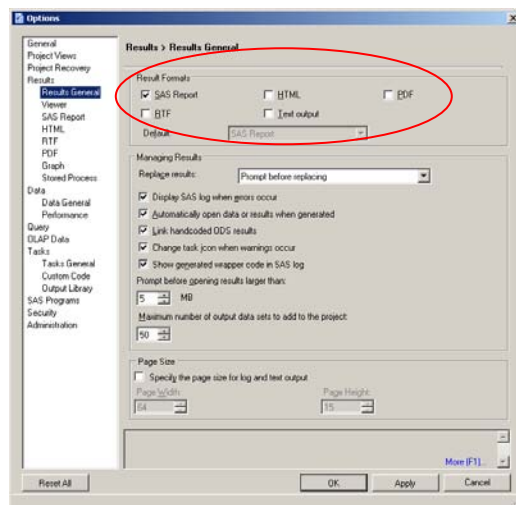
process flow window useful because projects can be broken into steps.  For example, instead of 1000 lines of code in one program that don't have natural breaks you can have 4 or 5 programs as steps within a process flow. This is explained later in the paper. You can always export all code to a single program file by going to **File → Export → Export All Code**.  The **Process Flow** and **Project Tree** windows give programmers an easy view of your entire project.  You can also pull existing SAS programs into your project with shortcuts (non-embedded programs) or place the programs into the project (embedded programs).



**Display 1. SAS Enterprise Guide Interface**

## CHANGING SAS ENTERPRISE GUIDE SETTINGS

If you would like to change the Enterprise Guide environment such as Project Tree views and outputted results, go to the **Tools → Options**.  The most relevant option for this tutorial is the **Results→ Results General** option, where you change the format of your output.  You can choose to generate SAS Report, HTML, PDF, RTF, and/or text files.



**Display 2. Changing your Environment with the Options Menu**

## EXAMPLE 1: LINEAR REGRESSION

The next two examples of this paper use the SASHELP.CARS data set available from the SASHELP library.  This data set has 428 observations and 15 variables.  It basically contains data for make and models of cars, such as
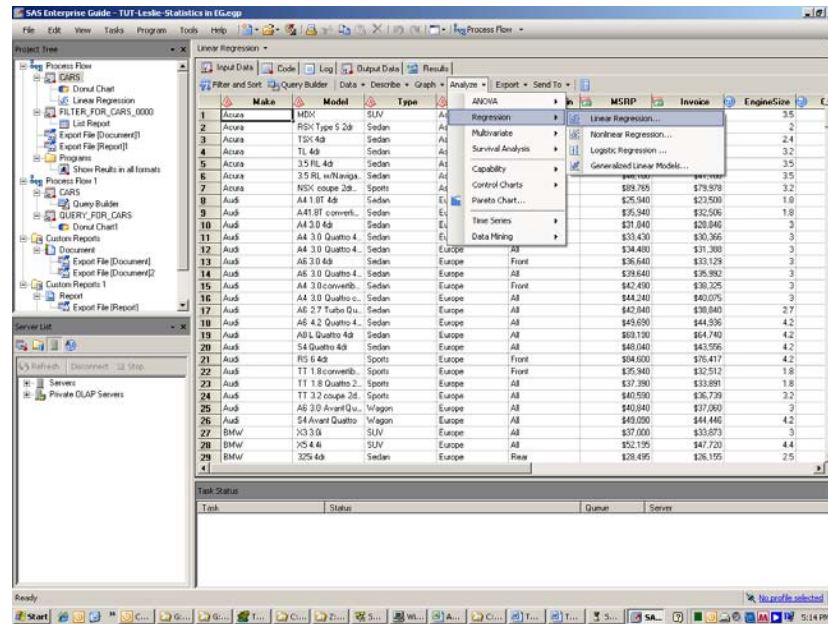
miles per gallon, number of cylinders, cost, etc.  The data set contains 5 character variables, 2 currency variables and 8 numeric variables.  The first step is to import the data set into your project by using the menus **File→ Open→ Data** (browse to SASHELP library)

To demonstrate multiple linear regression, the CARS data set is used to check for possible relationships between number of cylinders of the car, weight, length and gas mileage in the city (MPG_City).  So the model equation would be,

$$MPG\_city = \beta_0 + \beta_1 Cyclinders + \beta_2 Length + \beta_3 Weight + \epsilon$$
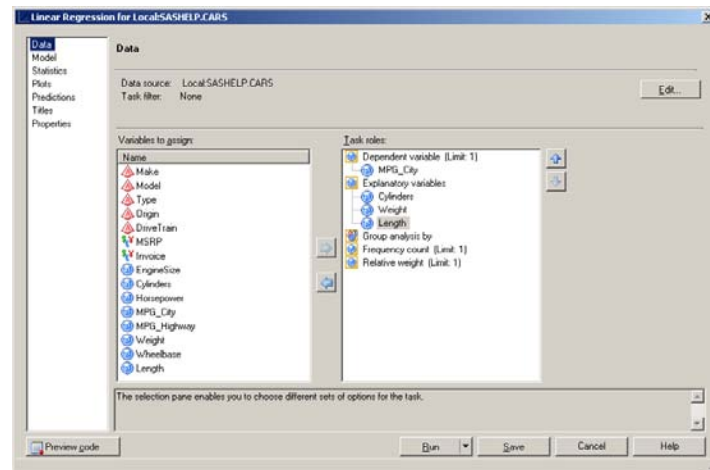
Where MPG_City is a function of the variables cylinders, length and weight plus an error term

Once opening the CARS data set into your project by using the menus **File→ Open→ Data** (browse to SASHELP library), you can use the Context tabs at the top of the data set to run linear regression analysis.   Here's a view of the data set and drop down menu for linear regression.
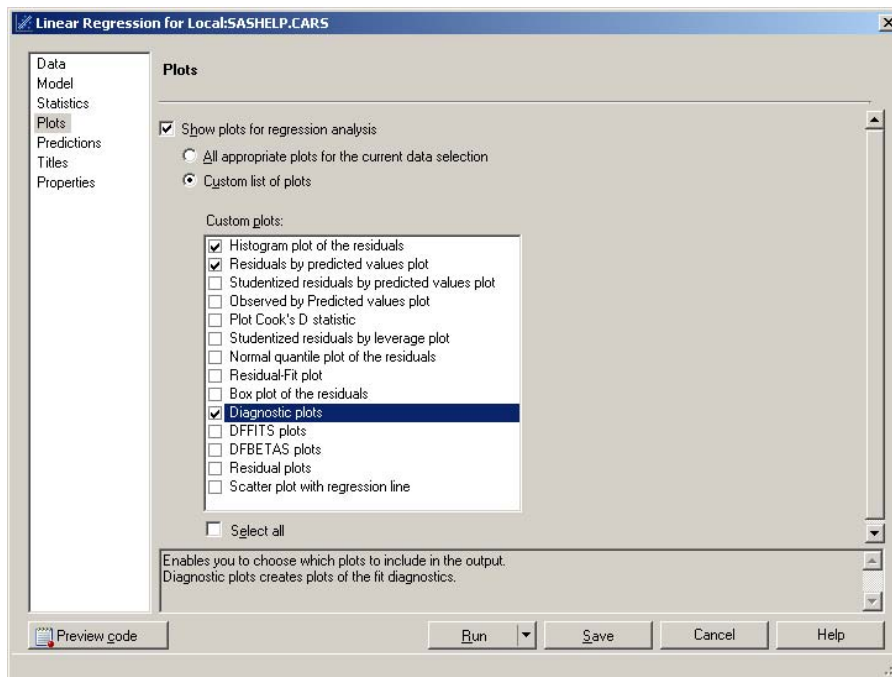


**Display 3. Using Context Menu to Run Linear Regression**

Once launching the Linear Regression menu, the task wizard opens allowing the user to select the dependent and independent variables plus select many other options.
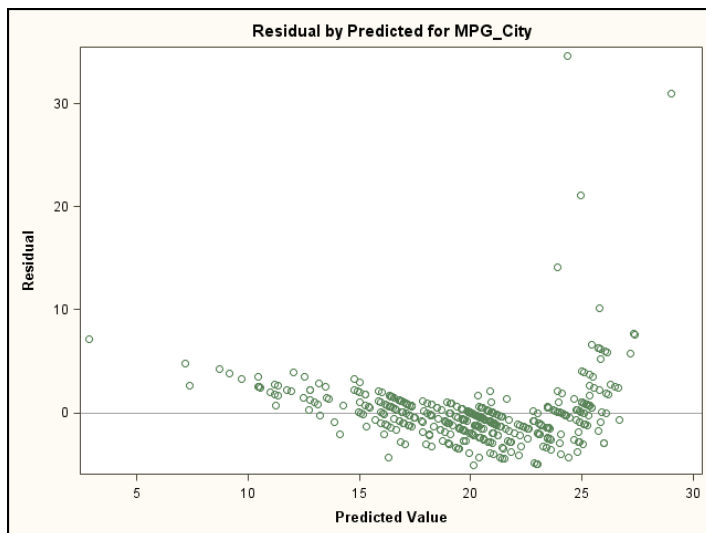


**Display 4. Linear Regression Options**

3

The user is also able to perform a variety of linear regression diagnostics to check for assumptions of multiple linear regression (e.g., linearity, normality, independence, equal variance), identify influential data points and check for correlations between variables (e.g., weight and length of car) to check for possible multicollinearity problems.  You can also output predicted values and create diagnostic plots.
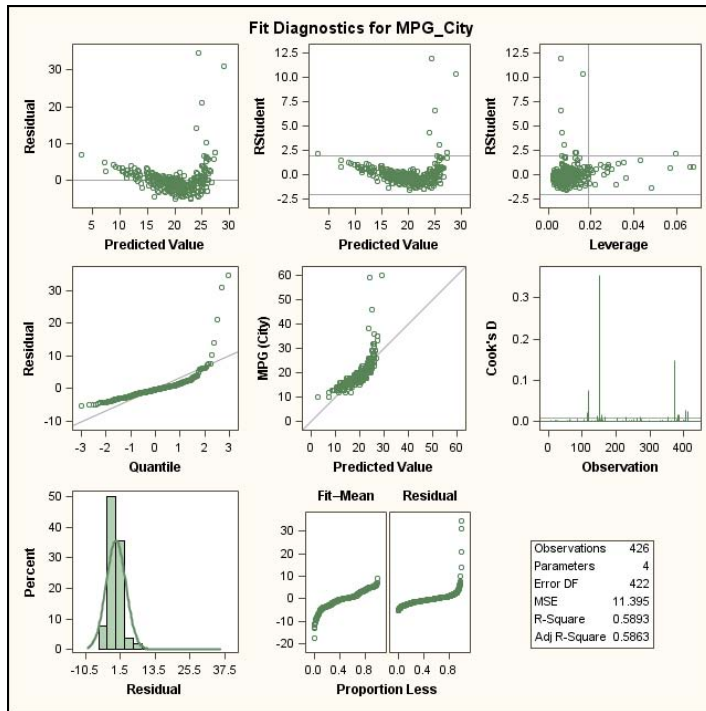


**Display 5. Plot Menu in the Linear Regression Task**

For example the Residual by Predicted Value plot can be used to assess the linearity assumption.  You can also get a panel of Fit Diagnostics.
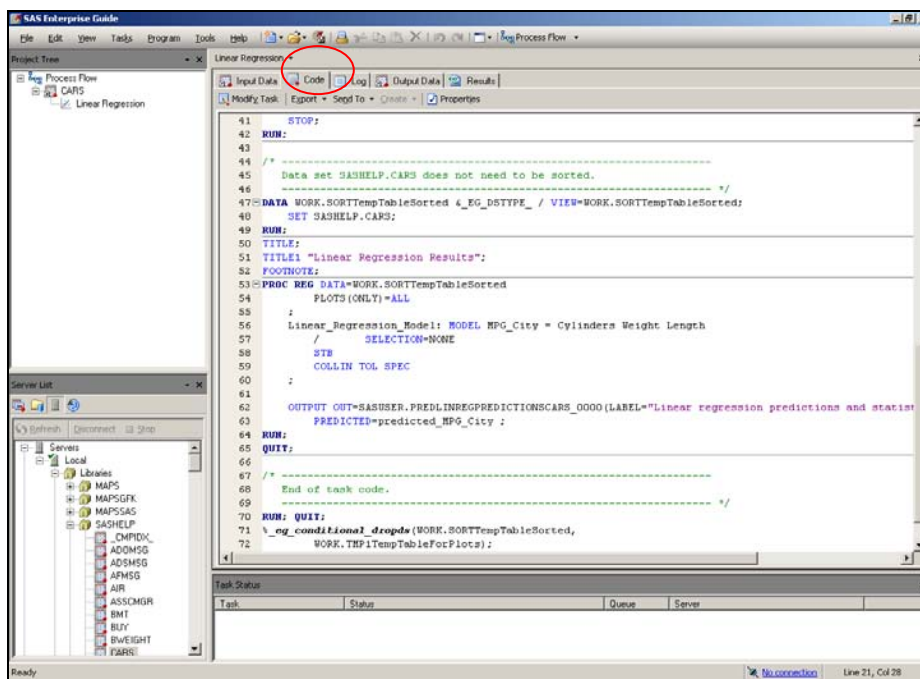


**Output 1. Predicted Value Plot from a Linear Regression Task**

**Output 2. Fit Diagnostics Panel from a Linear Regression Task**

Task wizards run code very quickly without writing the code. But you can add and change code when needed. After running a task wizard, the code that Enterprise Guide creates to generate the analysis is viewable and modifiable. I find this easy to copy this code into a new program and then make any necessary modifications. This allows you to leverage Enterprise Guide to write the code and learn additional SAS syntax. The window below shows the PROC REG code generated when running the Wizard. The code is created and accessible to you from the Code tab in the Context Menu. This feature (new starting with version 4.2) allows users to modify and re-run tasks quickly and easily. In the project explorer window, selecting the tab opens the code for viewing. You can now edit the code and rerun.
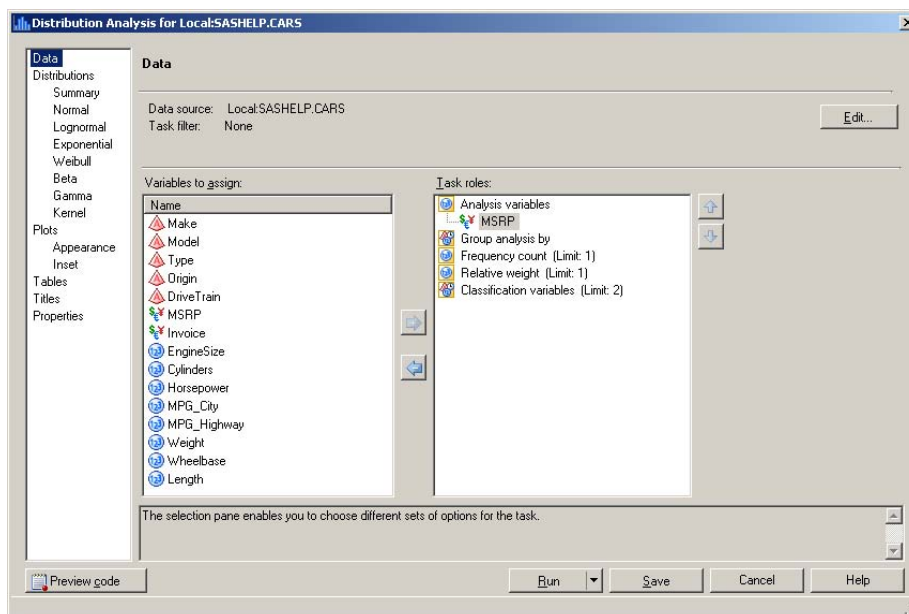


**Display 6. Modifying Task-Generated Code to Re-Run a Linear Regression Task**

## EXAMPLE 2: LOGISTIC REGRESSION

Logistic regression is a statistical technique to describe the relationship of several independent variables to a dichotomous dependent variable.  Using the SASHELP.CARS data set we could check for relationships between car characteristics and manufacturer's suggested retail price of the car (MSRP) if we categorize the MSRP price by high or not high using an arbitrary threshold.  Selecting type of car, origin and horsepower as possible explanatory variables the model equation would be,
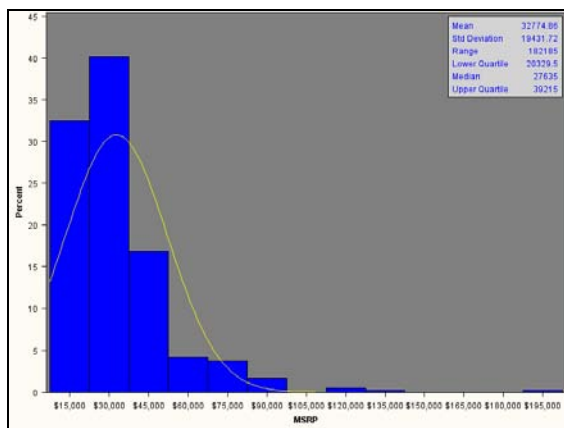
$$P \text{ (High priced car)} = 1/\ 1 + e^{-(\acute{a} + ß1Origin + ß2Horsepower + \epsilon)}$$

To make the dependent variable dichotomous you could create a new variable, high priced car indicator variable, to flag those observations that are considered high priced.  A look at the distribution of the MSRP variable using the Distribution Analysis Task wizard can give ideas on where to create a threshold for "high price".  After opening the CARS data set into your project by using the menus **File→ Open→ Data** (browse to SASHELP library), you can use the Context tabs at the top of the data set to run a distribution analysis.  Here's a view of the data set and drop down menu for distribution analysis.



**Display 7. Distribution Analysis Task**
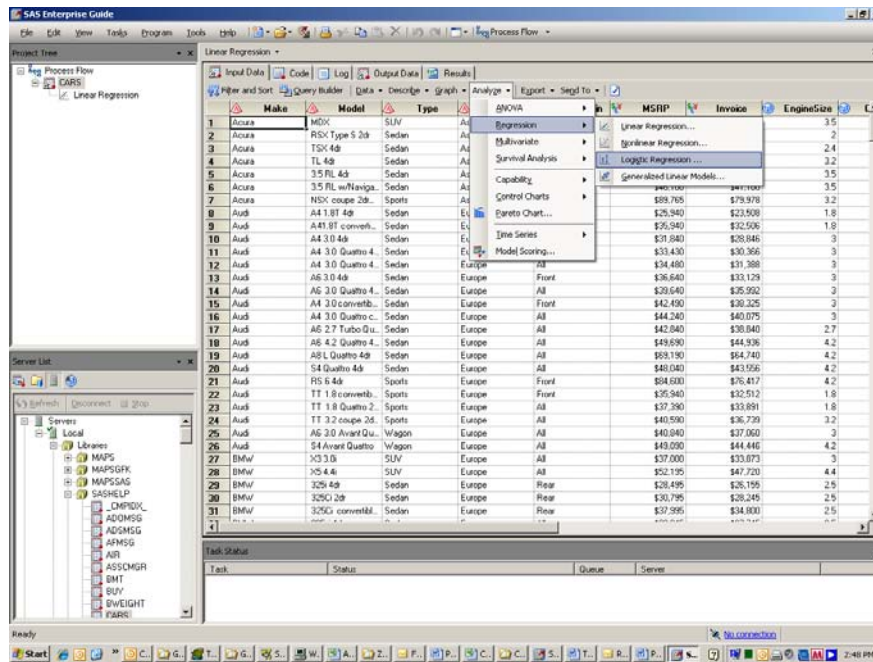
Here's the resulting histogram.



**Output 3. Histogram from a Distribution Task**

Using the upper quartile as the threshold to indicate a high-priced car, a data step could be used to create a new variable, "MSRP_HIGH" with values of 1=yes and 0=no.

```
DATA CARS;
FORMAT MSRP_HIGH 8.;
SET SASHELP.CARS;
IF MSRP >= 39215 THEN MSRP_HIGH=1;
ELSE MSRP_HIGH=0;
LABEL MSRP_HIGH='HIGH-PRICED CAR';
RUN;
```

Now from this new data set, the same analyze context menu can launch the logistic regression task wizard. This is shown below.



**Display 8. Modifying Task-Generated Code to Re-Run a Linear Regression Task**

Select results of the output are shown below. Odds ratio estimates using this data set find that cars from Europe are over 14 times more likely to be high-priced than those cars from the USA and cars from Asia are almost 43% less likely to be high-priced than cars from the USA. Also, there are increased odds of the car being high-priced with increasing horsepower. The wizard has an option in to specify the units for quantitative variables. Horsepower in units of 10 is shown in the second odds ratio estimates table.

| Analysis of Maximum Likelihood Estimates | | | | | | |
|---|---|---|---|---|---|---|
| Parameter | | DF | Estimate | Standard Error | Wald Chi-Square | Pr > ChiSq |
| Intercept | | 1 | -11.8989 | 1.3434 | 78.4466 | <.0001 |
| Horsepower | | 1 | 0.0441 | 0.00524 | 70.8717 | <.0001 |
| Origin | Asia | 1 | -1.2661 | 0.3035 | 17.4020 | <.0001 |
| Origin | Europe | 1 | 1.9644 | 0.2974 | 43.6392 | <.0001 |

| Odds Ratio Estimates | | |
|---|---|---|
| Effect | Point Estimate | 95% Wald Confidence Limits |
| Horsepower | 1.045 | 1.034 1.056 |
| Origin Asia vs USA | 0.567 | 0.213 1.508 |

7

| Odds Ratio Estimates | | |
|---|---|---|
| **Effect** | **Point Estimate** | **95% Wald Confidence Limits** |
| **Origin Europe vs USA** | 14.337 | 5.508 | 37.321 |

| Odds Ratios | | |
|---|---|---|
| **Effect** | **Unit** | **Estimate** |
| **Horsepower** | 10.0000 | 1.555 |

**Output 4. Output from a Logistic Regression Task**

Options in the wizard also allow you to run regression diagnostics to identify influential data points and check for the fit of the specified model.  You can also output predicted values and create ROC plots.  As shown in the table below, Hosmer-Lemeshow Goodness-of-fit test was unable to reject the null hypothesis (good fit).

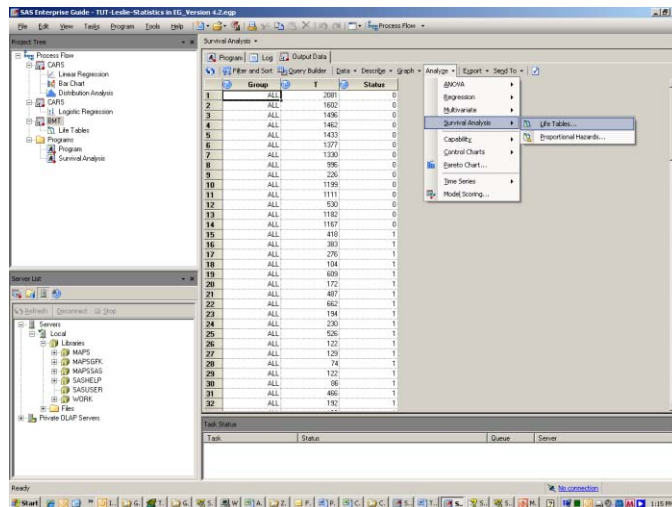| Hosmer and Lemeshow Goodness-of-Fit Test | | |
|---|---|---|
| **Chi-Square** | **DF** | **Pr > ChiSq** |
| 5.3671 | 8 | 0.7177 |



**Output 5. Output from a Linear Regression Task**
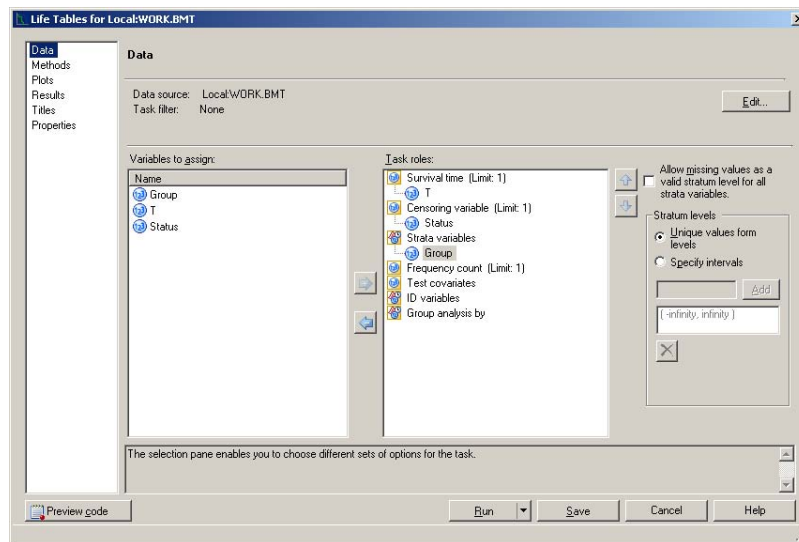
## EXAMPLE 3: SURVIVAL ANALYSIS

Survival analysis is a collection of statistical methods used to analyze time to event data.  To show how Enterprise Guide can perform survival analyses, Example 49.2 Enhanced Survival Plot and Multiple-Comparison Adjustments from the SAS Help and Documentation is used.  This example uses data from Klein and Moeschberger (Klein, J.P. and Moeschberger, M.L. (1997), Survival Analysis: Techniques for Censored and Truncated Data, New York: Springer-Verlag) that contains data on 137 bone marrow transplant patients.  Variables in the data set includes a group variable (3 categories of risk), disease-free survival time (time to death or relapse or the end of the study), and a status variable that indicates censored observations.

Once the data is in the project, you can run the survival analysis using the menus as seen in the window below.

**Display 9. Using Context Menu to Run Survival Analysis**
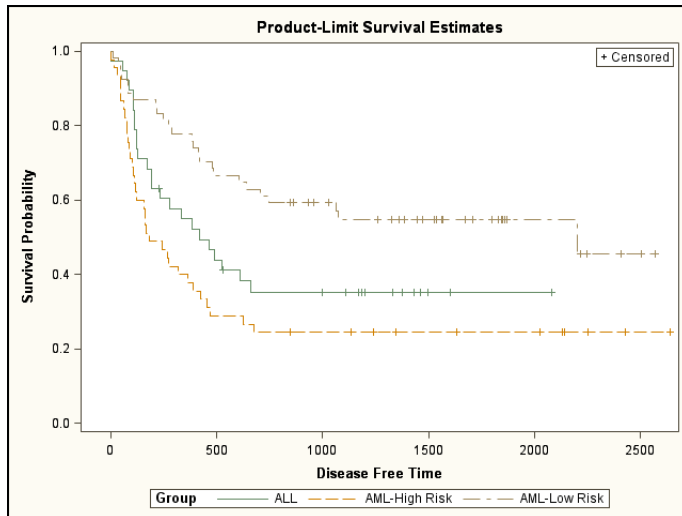
Using the wizard you can assign the Strata (group), Survival Time (time) and Censoring (status) variables.



**Display 10. Survival Analysis Task Options**

In the Methods window, the use can choose the method of estimation and select a plot of the estimated survival curves in the Plots window.  In this case the product-limit estimation (Kaplan Meier) method is used.  Parts of the results are shown below.

| Test of Equality over Strata | | | |
|---|---|---|---|
| **Test** | **Chi-Square** | **DF** | **Pr > Chi-Square** |
| **Log-Rank** | 13.8037 | 2 | 0.0010 |
| **Wilcoxon** | 16.2407 | 2 | 0.0003 |
| **-2Log(LR)** | 19.5313 | 2 | <.0001 |

**Output 6. Output from a Survival Analysis Task**

As shown earlier in the paper, you can add and change code as needed. After running a task wizard, the code that Enterprise Guide creates is viewable and modifiable. The window below shows a preview of the PROC LIFETEST code to be used when running the wizard. You can insert code using the Insert Code tab.



**Display 11. Modifying Task-Generated Code to Re-Run a Survival Analysis Task**

For instance, you can add the following code to the STRATA statement to get pairwise comparisons.

```
/ test=logrank adjust=sidak diff=control('AML-Low Risk');
```

The result of rerunning this code is below.

| Adjustment for Multiple Comparisons for the Logrank Test | | | | |
|---|---|---|---|---|
| Strata Comparison | | | p-Values | |
| Group | Group | Chi-Square | Raw | Sidak |
| ALL | AML-Low Risk | 5.1400 | 0.0234 | 0.0462 |
| AML-High Risk | AML-Low Risk | 13.8011 | 0.0002 | 0.0004 |

**Output 7. Output from a Survival Analysis Task**

## SAS® ENTERPRISE GUIDE® VERSIONS 4.3 AND 5.1

The analyses and work shown in this paper was done in Enterprise Guide version 4.2.  More recent versions, 4.3 and 5.1, include additional features.  The most noteworthy new features are a new program editor with autocomplete and integrated syntax help.  These features complete and show syntax while you are coding (e.g., typing PROC in the program editor will create a pop-up window that lists all available procedures for you to choose from).  Visit the What's New in Enterprise Guide 4.3 website for details.  Also, the paper SAS® Programmer's Paradise: New Goodies in SAS® Enterprise Guide® 4.3 by SAS Institute® describes many of the new features of EG 4.3.

## CONCLUSION

Enterprise Guide allows programmers to conduct statistical analyses.  Adding a few lines of code in combination with the many drop-down menus allow users to run programs efficiently and even pass on to analysts or persons less familiar with programming.

## REFERENCES

- SAS Fact Sheet, Accessed 3/09/13 http://www.sas.com/technologies/bi/query_reporting/guide/factsheet.pdf

- Getting Started with SAS Enterprise Guide – free learning tutorial available at http://support.sas.com/documentation/onlinedoc/guide/tut42/en/menu.htm

- Fernandez, George. "Quick and Complete Statistical Analyses Using SAS Enterprise Guide". http://www.wuss.org/proceedings08/08WUSS%20Proceedings/papers/tut/tut04.pdf

- Hallahan, Charles and Atkinson, Linda. "Introduction to SAS® Enterprise Guide® 4.1 for Statistical Analysis". *Proceedings of the SUGI 31 Conference*. Paper 109-31. http://www2.sas.com/proceedings/sugi31/109-31.pdf

- SAS tutorial http://support.sas.com/learn/statlibrary/statlib_eg4.1/top_learn.htm

## RECOMMENDED READING

- Little SAS Book for Enterprise Guide 4.2 – By Susan J. Slaughter, Lora D. Delwiche, http://www.sas.com/apps/pubscat/bookdetails.jsp?catid=1&pc=61861

- Data Analysis Using SAS Enterprise Guide- By Lawrence S. Meyers, Glenn Gamst and A. J. Guarino Cambridge University Press, 2009. http://www.cambridge.org/gb/knowledge/isbn/item6026601/?site_locale=en_GB

## CONTACT INFORMATION

Your comments and questions are valued and encouraged. Contact the author at:

Name: R. Scott Leslie
Enterprise: MedImpact Healthcare Systems, Inc.
Address: 10680 Treena Street
City, State ZIP: San Diego, CA 92131
Work Phone: 858-790-6685
E-mail: scottleslie@san.rr.com