

PharmaSUG 2012 - Paper AD22

Using JMP® Partition to Grow Decision Trees in Base SAS®

Mira Shapiro, Analytic Designers LLC, Bethesda, MD

ABSTRACT

Decision Tree is a popular technique used in data mining and is often used to pare down to a subset of variables for more complex modeling efforts. If your organization has only licensed Base SAS and SAS/STAT you may be surprised to find that there is no procedure for decision trees. However, if you are licensed JMP 9 user, you can build and test a decision tree with JMP. The Modeling→Partition analysis provides an option for creating SAS Data Step scoring code. Once created, the scoring code can be run in Base SAS. This discussion will provide a brief overview of decision trees and illustrate how to create a decision tree with Partition in JMP and then create the SAS Data Step Scoring code.

INTRODUCTION

As analysts and statisticians we are often faced with the question of what relationships are present in our data. JMP provides methods to find an answer. The JMP Partition Platform provides an easy way to create decision trees with numerous options for perfecting and interpreting the results. This discussion will focus on creating a decision tree and will give an introduction to the many options that are available for improving and interpreting the results within JMP. Additionally, the option to create a decision tree in JMP and export the code to be run in Base SAS will be demonstrated.

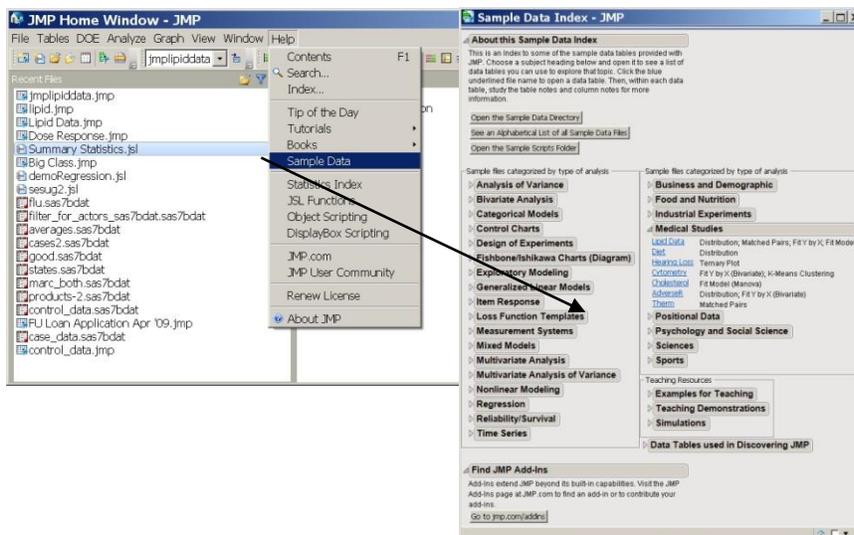
WHAT IS A DECISION TREE?

A decision tree model allows for examination of the relationship between a response variable and multiple possible predictors. The potential predictors are evaluated using statistical methods appropriate to their type and assessed as to their predictive value for the response variable. The data is then split into two groups based on the value of the predictor. As the tree is built by recursive splitting, the predictors are re-evaluated at each stage.

There are numerous resources available for understanding the underlying statistical techniques and algorithms used in the decision tree modeling process. Such a discussion is beyond the scope of this paper. Several recommended books, papers and online resources are listed in the References and Recommended Reading sections at the end of this paper.

JMP SAMPLE DATA USED IN EXAMPLES

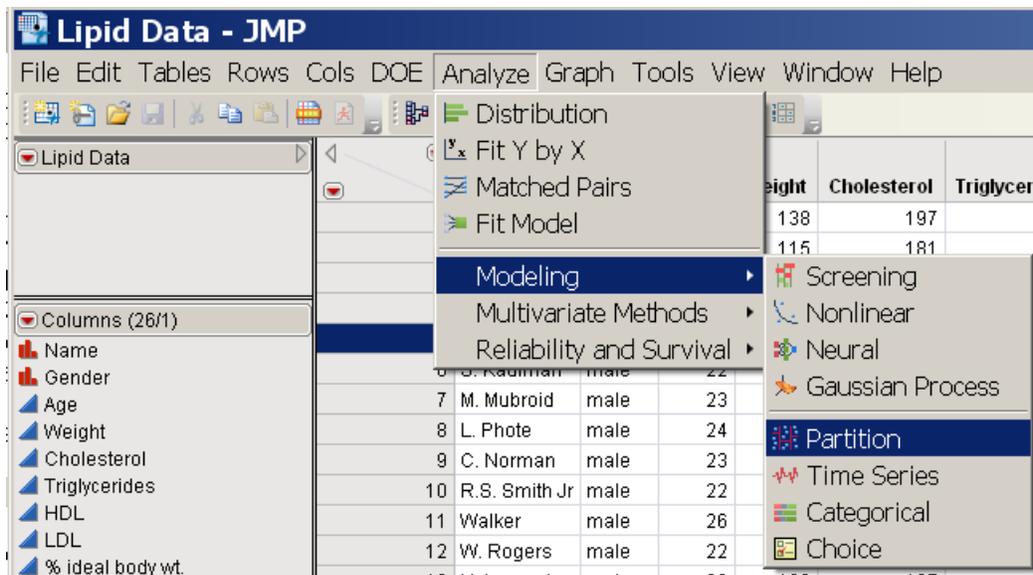
A variety of sample data sets are included with JMP. The index of data sets is helpful to find a sample data set by subject area or analysis type and provides opportunities for new users to explore JMP's capabilities. The examples in this paper use the sample data set "Lipid Data". The sample JMP data sets are found under the *JMP Home Window* → *Help* → *Sample Data*.



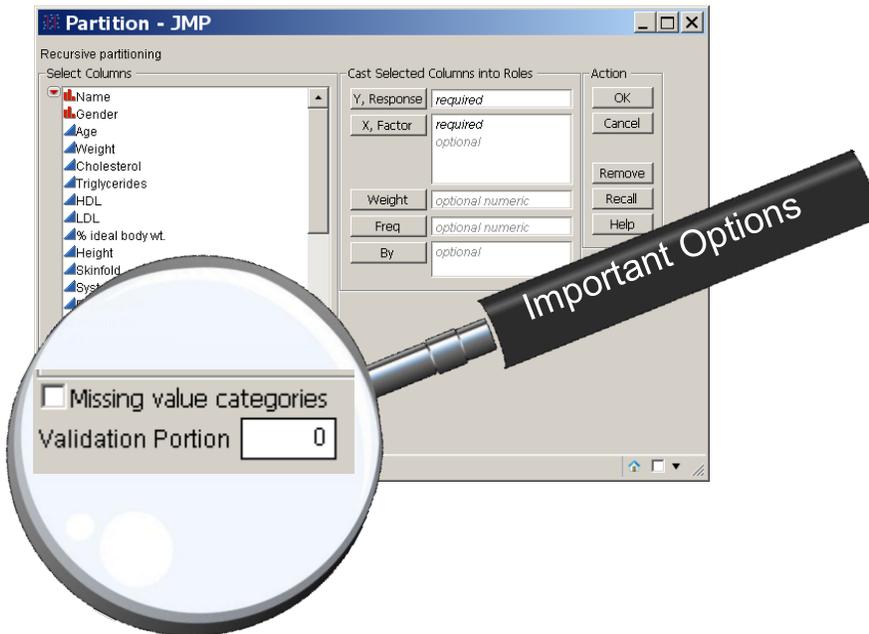
	Name	Gender	Age	Weight	Cholesterol	Triglycerides	HDL
1	J. Suds	male	22	138	197	152	4
2	T. Wilson	female	22	115	181	59	€
3	D.S. Quintent	male	22	190	190	117	4
4	R. Beal	female	22	115	131	54	€
5	R. James	male	25	160	172	93	4
6	S. Kaufman	male	22	150	233	176	4
7	M. Mubroid	male	23	154	194	79	4
8	L. Phote	male	24	185	155	89	4
9	C. Norman	male	23	178	234	307	2
10	R.S. Smith Jr	male	22	158	201	88	€
11	Walker	male	26	188	258	299	3
12	W. Rogers	male	22	150	212	52	€
13	M. Lumpole	male	22	123	137	158	2
14	D. Fineman	female	27	138	285	98	€
15	R. Smith	male	22	143	218	101	4
16	J. Newman	male	24	139	167	71	€
17	D. Smith	male	22	156	170	81	4
18	R. Heckleton	male	22	150	157	86	3
19	T. Plotkus	female	24	135	215	71	€
20	R. Humble	male	25	219	194	71	4
21	B. Beer	male	28	173	207	107	€
22	L. Henry	male	22	151	198	80	4
23	L. Aycock	male	23	182	189	47	€
24	P. Pilgrim	male	24	161	216	95	3
25	G. Lucas	male	22	176	212	140	4
26	D. Whitsel	male	26	177	175	77	4
27	G. Regular	male	23	174	158	57	3
28	B. Hill	male	20	234	115	95	2
29	D. Chappell	male	20	150	130	100	2

GROWING A DECISION TREE

Once you have opened or created your data table in JMP, there are multiple ways to invoke the JMP Partition Platform. From the JMP Home Window, select **Analyze** → **Modeling** → **Partition** to begin.



Next, the JMP Partition selection pane opens, providing the opportunity to select variables for roles in the model and make use of the initial options for the model.



Missing Value Categories:

Y Response

	Checked	Unchecked
Categorical	Additional variable level is created for the missing values	Excluded
Continuous	Excluded	Excluded

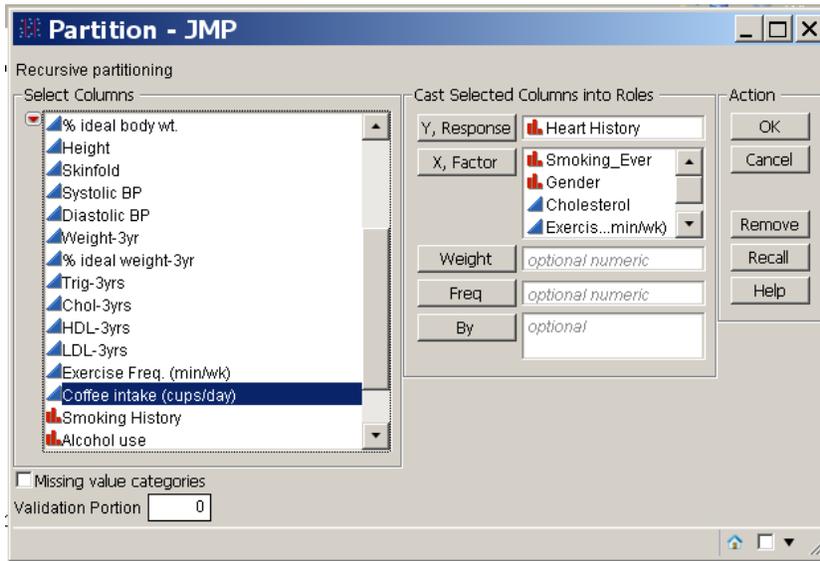
X Predictors (used as splitting variable)

	Checked	Unchecked
Categorical	Additional variable level is created for the missing values	Random assignment of value to one side of the split
Continuous	Random assignment of value to one side of the split	Random assignment of value to one side of the split

Validation Portion:

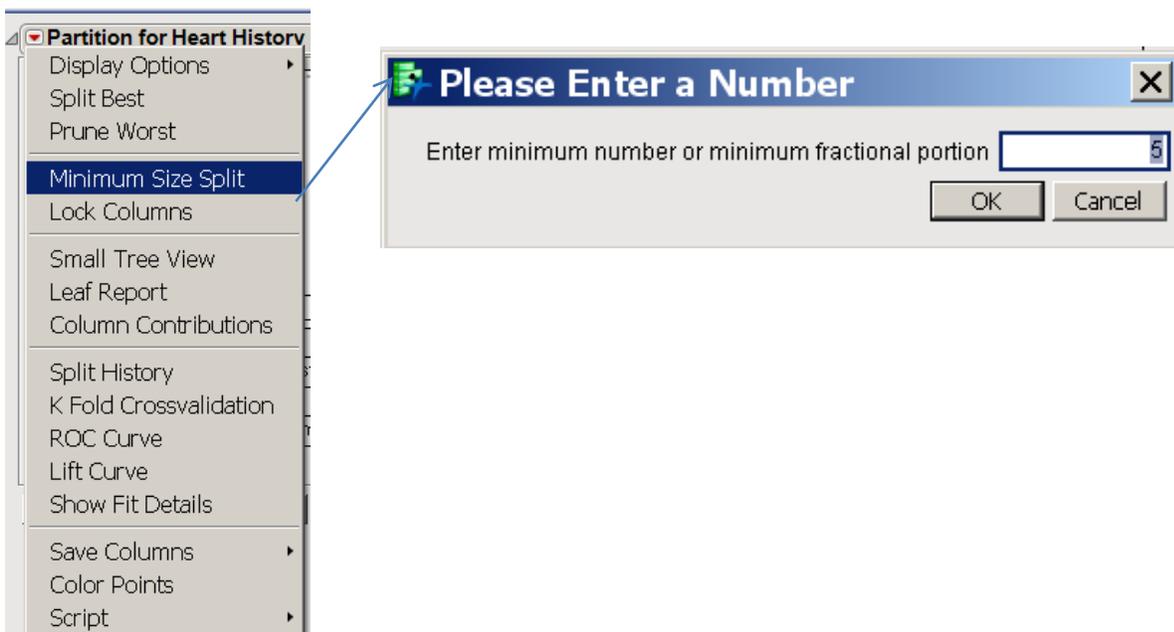
This option allows the user to use a portion of their data to estimate the model parameters, leaving the remaining portion to validate the model.

At this point, drag and drop columns to cast them into roles. For this example, *Heart History* is chosen as the Y response, and several others are cast into the role as potential predictors.

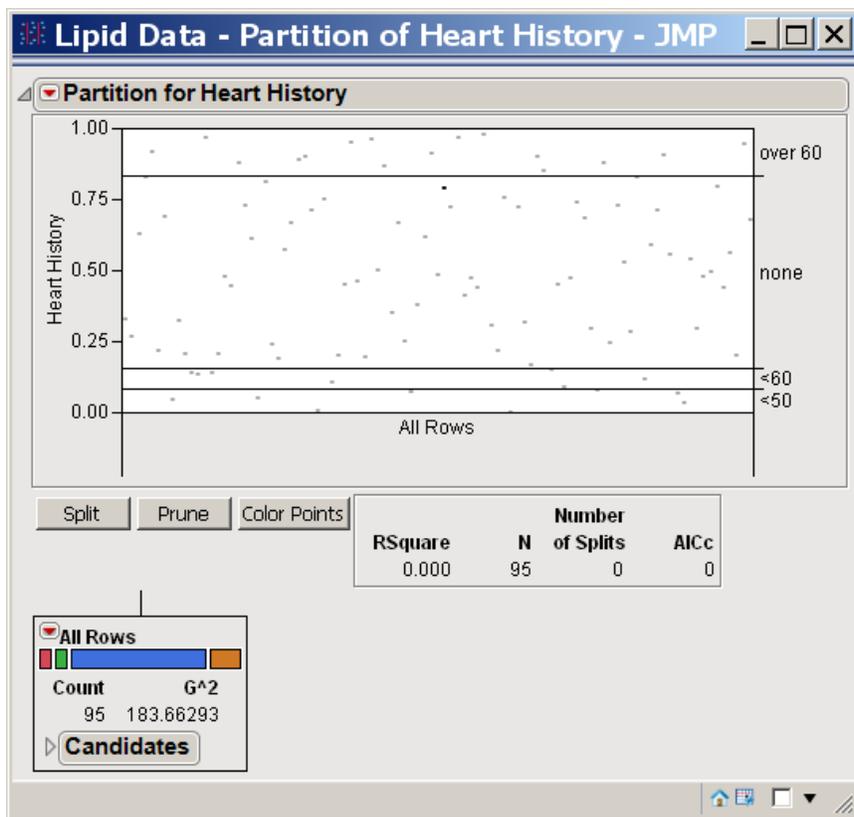


Minimum Size Split

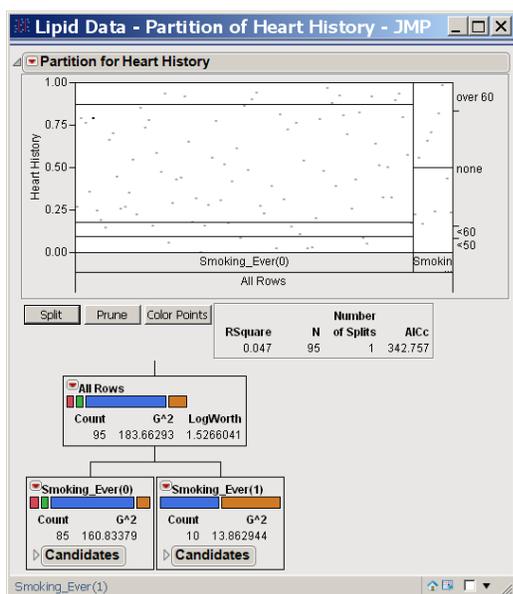
The user has the option to define the minimum size of a group in a split. To access this option, click the hotspot next to the main partition window, and mouse down to **Minimum Size Split**. The size can be entered as a number or a fractional portion. The Partition process will not create a group that violates this criterion.



Once the **OK** button is clicked, the Partition Platform is launched. The points are shown on the initial pane, with the values of **Heart History** shown on the right. It is now up to the user to click the **Split** button to begin the process.

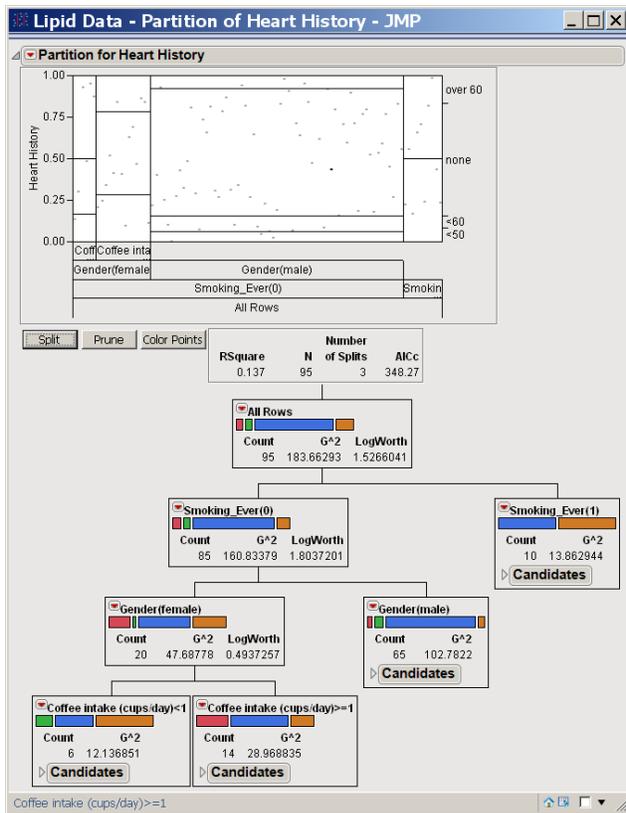


After the **Split** button is clicked below the main Partition window, JMP begins to create the decision tree. The first split selected is the **Smoking_Ever** variable that was created from the smoking history variable and coded as 0/1.



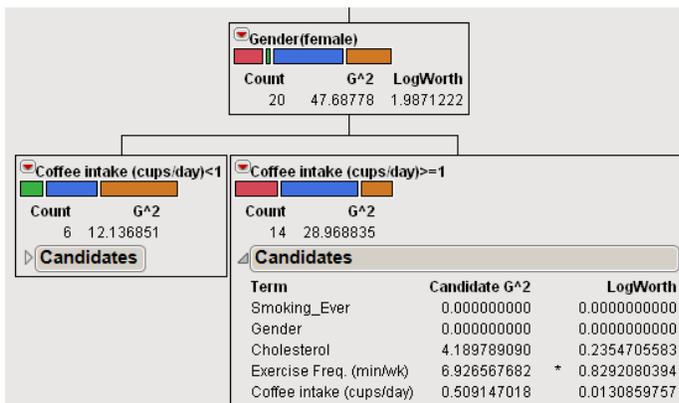
SPLITTING AND PRUNING YOUR DECISION TREE

Partition is a recursive process, and the user has multiple options for continuing and controlling the process. The **Split** button was clicked again; the “best” split selected by JMP was **Gender** for non-smokers (**Smoking_Ever=0**), and the next was **Coffee Intake** for females.

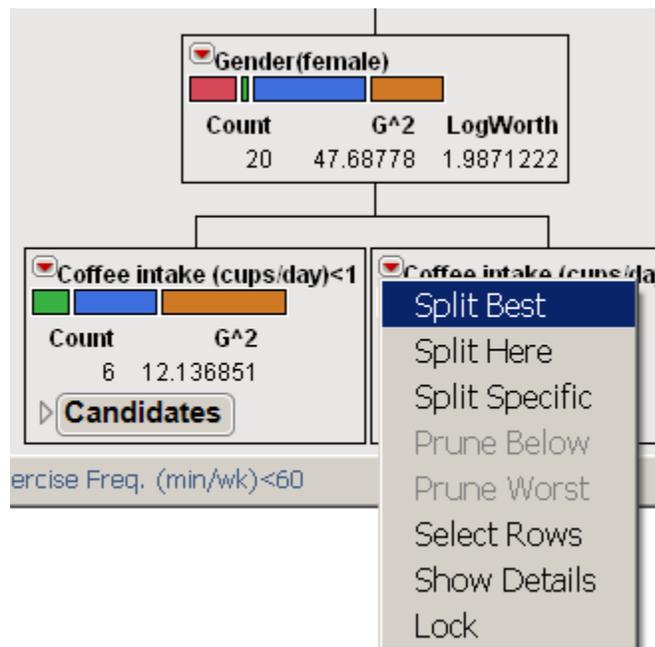


Picking a Candidate

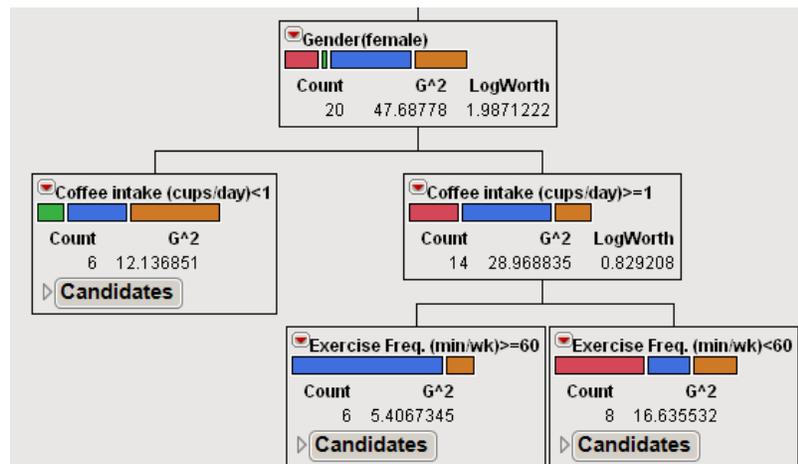
In many cases the user has some subject area knowledge that may direct them to pick a particular split that may not show up as the “best” in the JMP evaluation criteria. By clicking the gray triangle to the left of **Candidates** in any node, the user can examine all of the variables that can be used for a split and their associated statistics.



Alternatively, clicking under the hotspot at any node will reveal available options. Here we click **Split Best** to allow JMP to select the variable for the split.

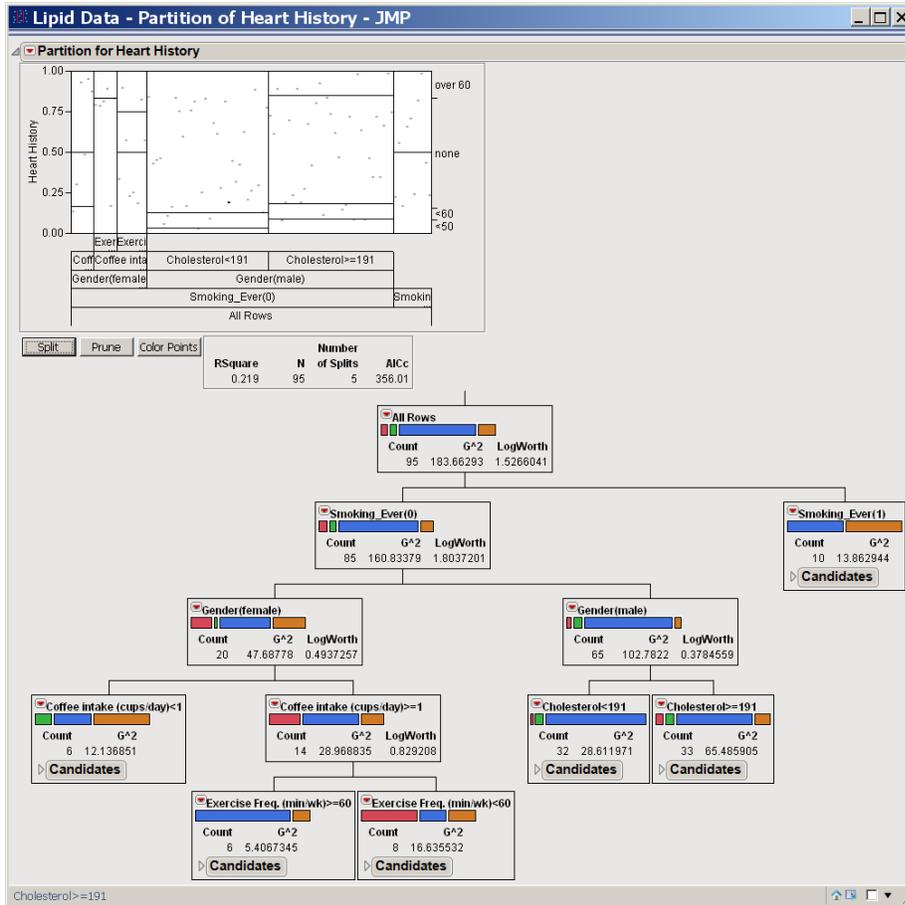


The best split under **Coffee Intake** of 1 cup a coffee or greater per day is **Exercise Freq.**



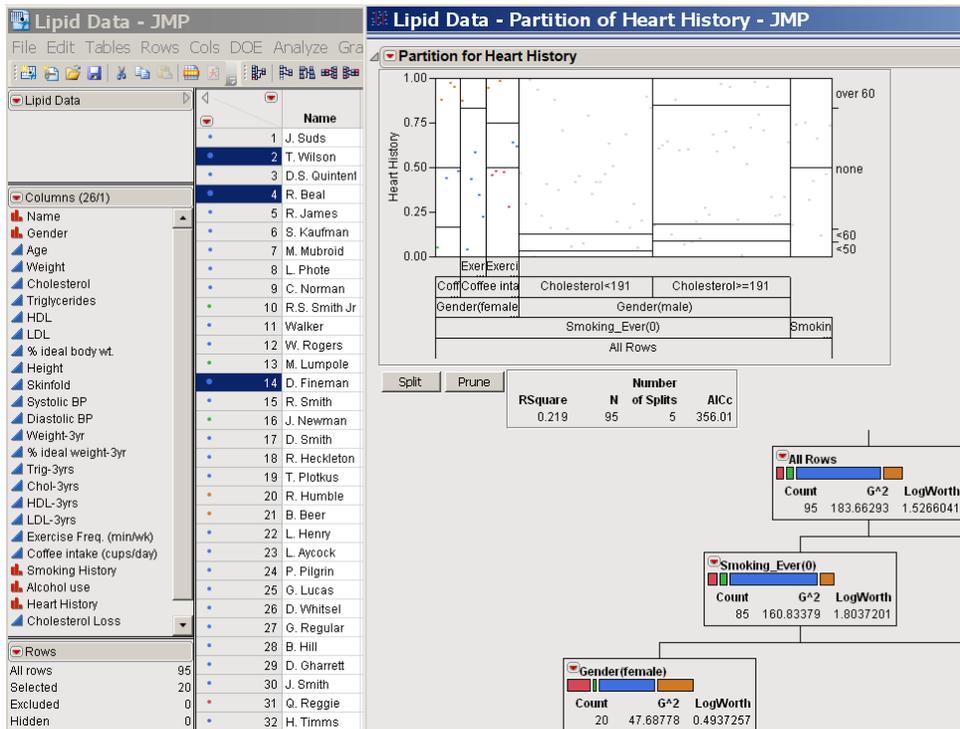
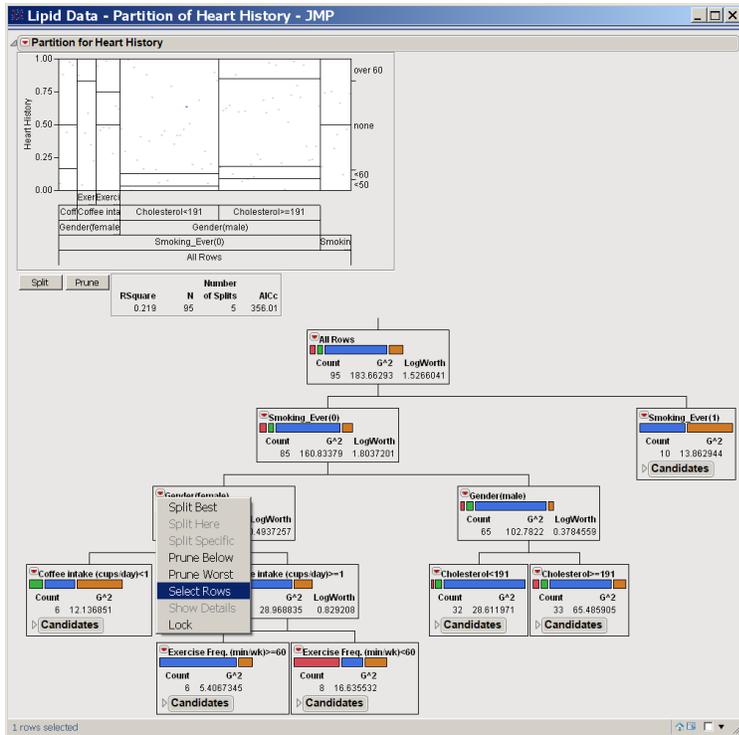
Recursive Splitting

The user can use a combination of techniques to build their decision tree model. No node will be created that violates the minimum split criterion specified in the Partition Platform Launch Window.



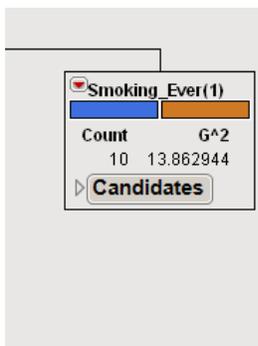
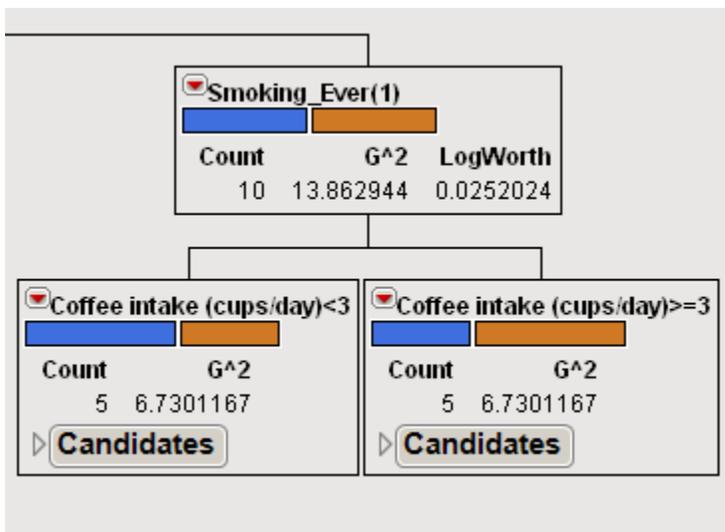
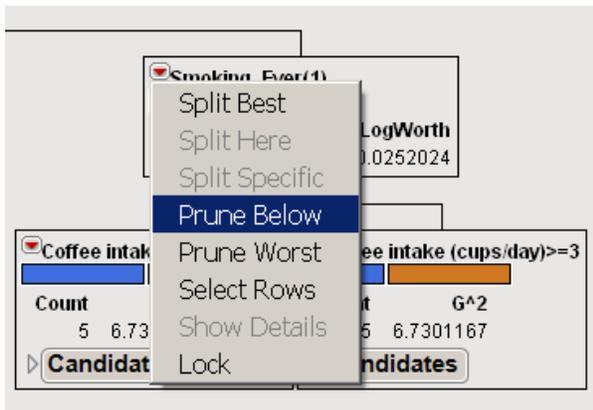
Dynamic Linking

As in all other JMP operations, Dynamic Linking of the data table and results allows the user to visualize the relationships. Clicking on the hotspot on any node and selecting **Select Rows** causes the data table to become visible, and the associated rows will be highlighted.



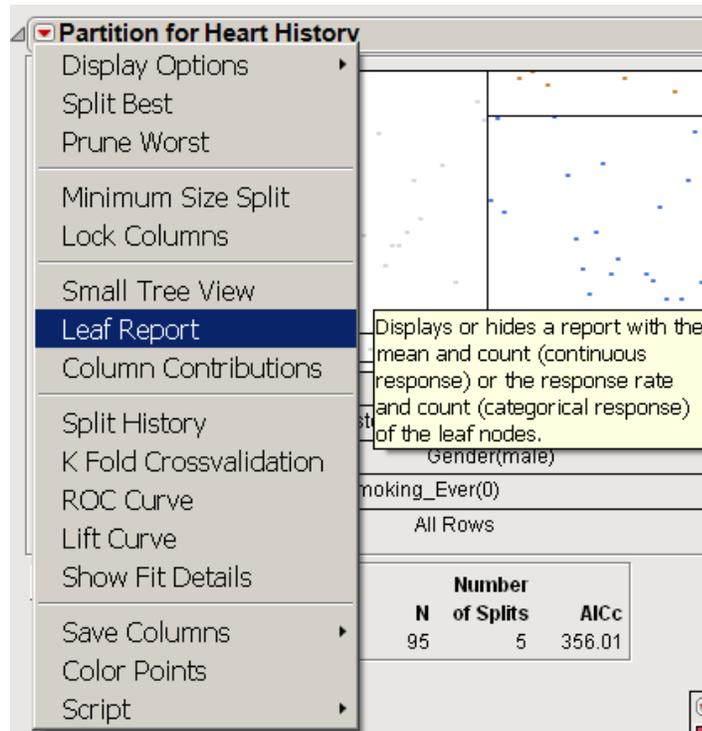
Pruning A Decision Tree

Just as the tree can be built with JMP, nodes can be removed by pruning areas of the tree. Looking at the level that we added in the previous step, we've decided to remove the **Coffee Intake** split. By clicking the hot spot next to **Smoking_Ever**, we reveal the available options, including **Prune Below**. Selecting **Prune Below** will result in the removal of all nodes below the selected node.

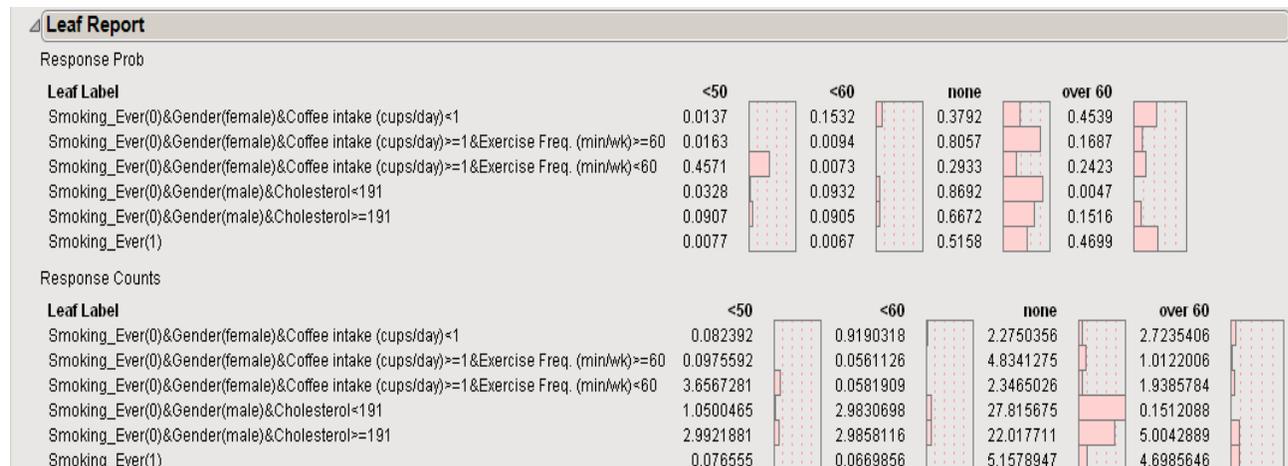


OPTIONS FOR UNDERSTANDING AND EVALUATING YOUR DECISION TREE

There are numerous statistics and reports that are available for evaluating your decision tree model. Generating and showing or hiding each of these options is accomplished by clicking the hotspots and making use of the hide/unhide icon. As shown in the following examples, most of the options can be moused over to reveal a box with a description of the particular option. The JMP Help facility and online resources have detailed information about the statistics generated.



Leaf Report



ROC Curve

Partition for Heart History

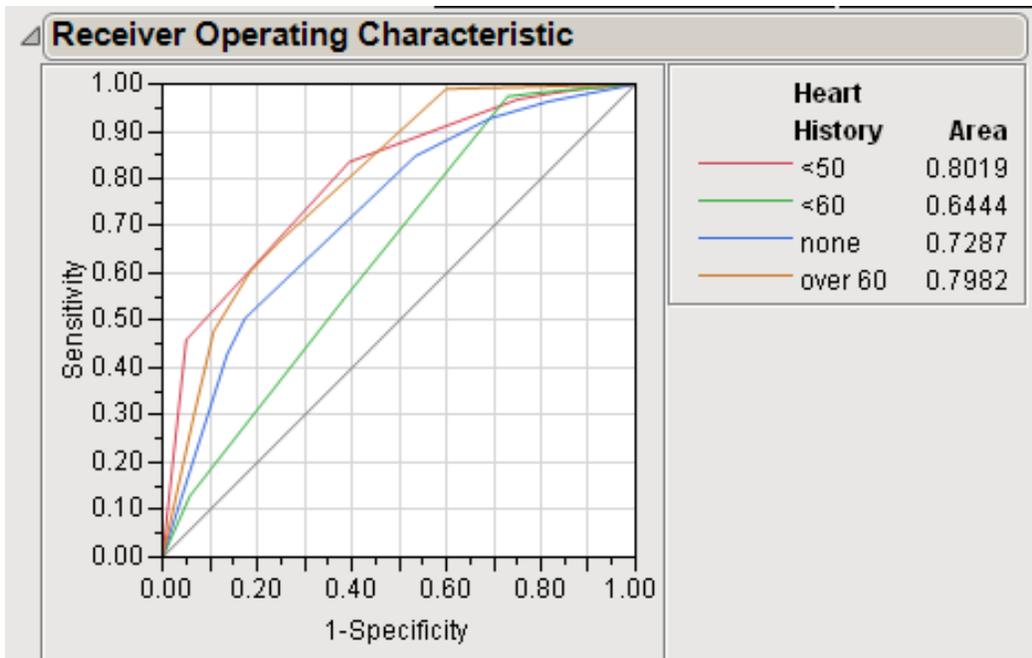
- Display Options
- Split Best
- Prune Worst
- Minimum Size Split
- Lock Columns
- Small Tree View
- Leaf Report
- Column Contributions
- Split History
- K Fold Crossvalidation
- ROC Curve**
- Lift Curve
- Show Fit Details
- Save Columns
- Color Points
- Script

Cholesterol<191 Cholesterol>=191

Gender(male)

Plots the response-category sorting efficiency of the model predictions.

	Number		
	N	of Splits	AICc
	95	5	356.01

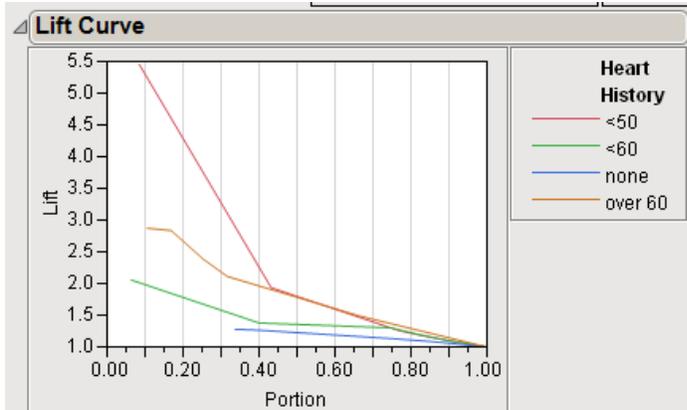


Lift Curve

Lift Curve
 Show Fit Details
 Save Columns
 Color Points
 Script

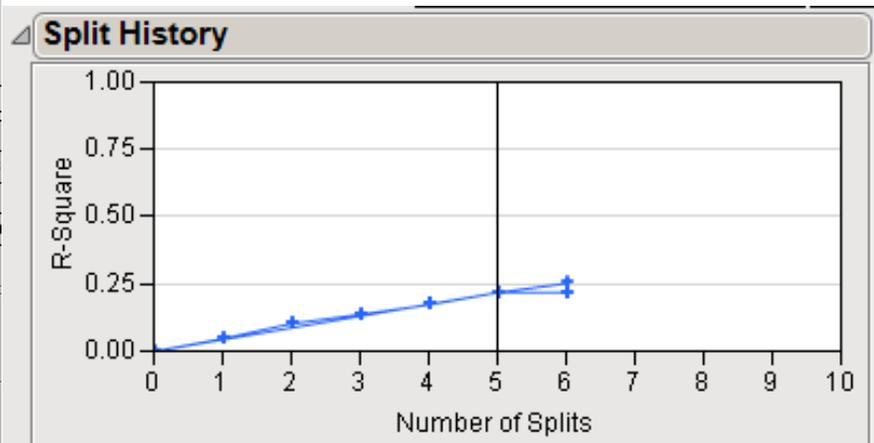
Plots how much more saturated the top x-percent of predicted values are compared to the whole population.

Count	85	160.0
-------	----	-------



Split History

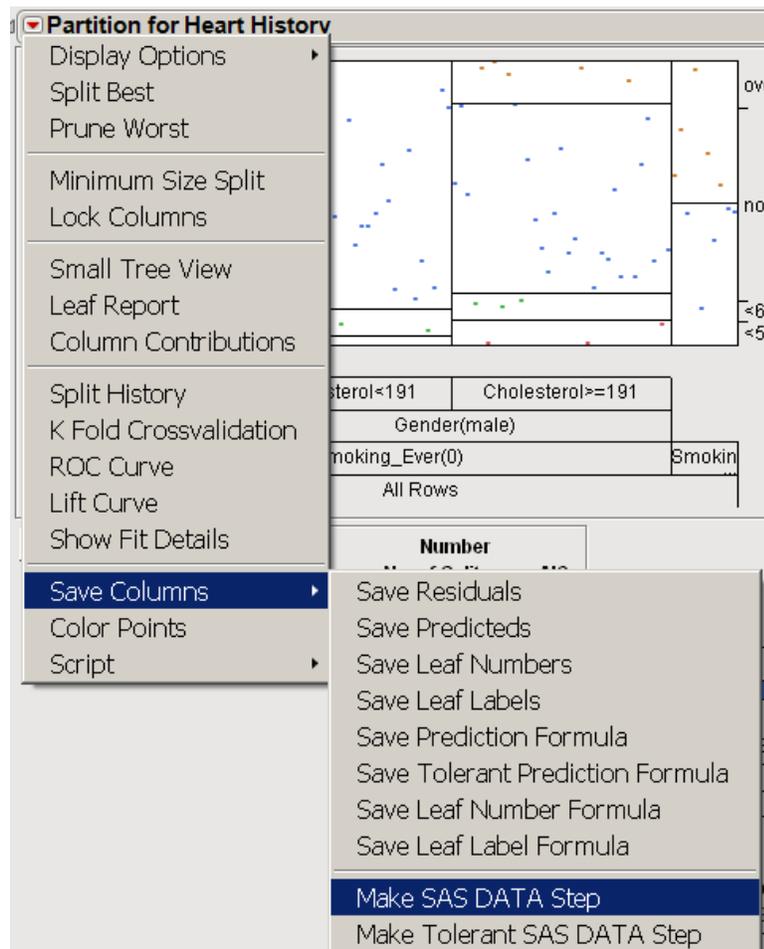
Partition for Heart History
 Display Options
 Split Best
 Prune Worst
 Minimum Size Split
 Lock Columns
 Small Tree View
 Leaf Report
 Column Contributions
 Split History
 K Fold Crossvalidation
 ROC Curve
 Lift Curve
 Show Fit Details
 Save Columns
 Color Points
 Script



TRANSPLANTING YOUR DECISION TREE IN SAS

There are times when a user wants to build their decision tree model in JMP and then use the model in Base SAS. In order to create the SAS Data Step code, click on the hotspot at the left of the main partition window, and mouse down to **Make SAS Data Step**. As shown below, the code will appear in a window entitled **Partition SAS Scoring Code-JMP**. This code can be run within base SAS with suitable data sources.

Note that the variable names are created using SAS naming conventions, thereby avoiding the need to surround a variable name with special characters with quotes and following it with an 'n'. For example, **Heart History** becomes **Heart_History**. If you are using another data set, it will be necessary to adjust either the program code or your variables' names accordingly.



Partition SAS Scoring Code - JMP

File Edit Tables DOE Analyze Graph Tools View Window Help

```

/*!PRODUCER: JMP - Partition - Decision Tree */
/*!DATA: Lipid Data */
/*!TARGET: Heart_History */
/*!OUTPUT: Prob__50, Heart_History, "<50" */
/*!OUTPUT: Prob__60, Heart_History, "<60" */
/*!OUTPUT: Prob_none, Heart_History, "none" */
/*!OUTPUT: Prob_over_60, Heart_History, "over 60" */
/*!INPUT: Smoking_Ever */
/*!INPUT: Gender */
/*!INPUT: Cholesterol */
/*!INPUT: Exercise_Freq_min_wk_ */
/*!INPUT: Coffee_intake_cups_day_ */

LABEL Prob__50= 'Predicted: Heart_History=<50';
LABEL Prob__60= 'Predicted: Heart_History=<60';
LABEL Prob_none= 'Predicted: Heart_History=none';
LABEL Prob_over_60= 'Predicted: Heart_History=over 60';

Prob__50=0;
Prob__60=0;
Prob_none=0;
Prob_over_60=0;
IF Smoking_Ever=0 THEN DO;
  IF Gender='female' THEN DO;
    IF Coffee_intake_cups_day_<1 THEN DO;
      Prob__50=Prob__50+0.0137320044296789;
      Prob__60=Prob__60+0.153171966461003;
      Prob_none=Prob_none+0.379172599272267;
      Prob_over_60=Prob_over_60+0.453923429837051;
    END;
  ELSE DO;
    IF Exercise_Freq_min_wk_>=60 THEN DO;
      Prob__50=Prob__50+0.0162598744924326;
      Prob__60=Prob__60+0.00935210673416653;
      Prob_none=Prob_none+0.805687918578284;
      Prob_over_60=Prob_over_60+0.168700100195117;
    END;
  ELSE DO;
    Prob__50=Prob__50+0.457091013494114;
    Prob__60=Prob__60+0.00727386079324064;
    Prob_none=Prob_none+0.293312825560888;
    Prob_over_60=Prob_over_60+0.242322300151758;
  END;
END;
END;
ELSE DO;
  IF Cholesterol<191 THEN DO;
    Prob__50=Prob__50+0.0328139523644141;
    Prob__60=Prob__60+0.0932209319184951;
    Prob_none=Prob_none+0.869239841949033;
    Prob_over_60=Prob_over_60+0.00472527376805779;
  END;
  ELSE DO;
    Prob__50=Prob__50+0.0906723655301667;
    Prob__60=Prob__60+0.0904791398032452;
    Prob_none=Prob_none+0.667203376009356;
    Prob_over_60=Prob_over_60+0.151645118657233;
  END;
END;
END;
ELSE DO;
  Prob__50=Prob__50+0.0076555023923445;
  Prob__60=Prob__60+0.00669856459330144;
  Prob_none=Prob_none+0.515789473684211;
  Prob_over_60=Prob_over_60+0.469856459330144;
END;
END;

```

CONCLUSION

This discussion has provided an overview to the JMP Partition Platform. This facility provides an approach to creating decision trees that is interactive and does not require programming skills. The user has the ability to control the split size, variables and decision tree shape through multiple options within JMP as well as to export Data Step code to SAS. A comprehensive set of statistics and results are available to evaluate the resulting JMP decision tree model.

REFERENCES

<http://support.sas.com/resources/papers/proceedings09/TOC.html>.

Gaudard, Martha, Ph.D. et. al., "Interactive Data Mining and Design of Experiments: the JMP® Partition and Custom Design Platforms". March 2006, Available at

http://www.jmp.com/software/whitepapers/pdfs/372455_interactive_datamining.pdf.

http://www.jmp.com/software/whitepapers/partition_platform/index.shtml

<HTTP://JMP.COM/SUPPORT/>

THE JMP 9.02 HELP FACILITY

http://www.jmp.com/applications/data_mining/

MASTERING JMP®: AN OVERVIEW OF DATA EXPLORATION AND THE JMP® SCRIPTING LANGUAGE, COURSE NOTES, SAS INSTITUTE INC.

RECOMMENDED READING

Data Mining Techniques: For Marketing, Sales, and Customer Relationship Management, Third Edition, by Michael Berry and Gordon Linoff, John Wiley Sons, Inc., April 2011

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

Mira Shapiro
Analytic Designers LLC, Bethesda, MD
mira.shapiro@analyticdesigners.com

SAS and all other SAS Institute Inc. product or service names are registered trademarks or trademarks of SAS Institute Inc. in the USA and other countries. ® indicates USA registration.

Other brand and product names are trademarks of their respective companies.