**Paper 357-2011**

# JMP® 9 and Interactive Statistical Discovery

## John Sall, SAS Institute Inc., Cary, NC, USA

## ABSTRACT

JMP 9 represents a major new revision of statistical discovery software from SAS® for engineers, scientists, and business analysts.

## INTRODUCTION

Though there are many improvements throughout all the components in the JMP system, most of the excitement is focused on a few special features:

- Better scripting and a plugin system to make using scripts easy

- Improved presentation graphics

- More powerful interfaces to Excel, R, and SPSS as well as SAS

- A JMP Pro version that adds more data mining fits

- A new look on Windows

## EXTENDING JMP THROUGH SCRIPTING

Suppose that you need to do an analysis of hyperspectral data. Our eyes and most of our cameras see in only three colors, but hyperspectral cameras can see in many colors, or wavelengths, including infrared and ultraviolet. Even the Landsat satellites send back images in seven wavelengths. Using the extra wavelengths helps discriminate land cover. Let's analyze an extreme example, Aviris data from NASA, which has 240 wavelengths. You can download some public Aviris data from a NASA website containing a number of 614-by-512 pixel images in 240 wavelengths. The script to do this is not hard, using some new features. Here is the housekeeping code to choose a directory and get the names of the rfl files.

```
NamesDefaultToHere(1);
defaultDir = "/C:/Users/sall/Documents/JMP9 Demo/Hyperspectral Data/";
inDir  = PickDirectory("Navigate to directory with AVIRIS rfl files",
        defaultDir,showEdit(1));
outDir = PickDirectory("Navigate to directory for resulting jpg PCA images",
        defaultDir,showEdit(1));
filelist = filesInDirectory(inDir);
ni = nItems(filelist);
rflFiles = {};
for(i=1,i<=ni,i++,if(EndsWith(fileList[i],".rfl"),InsertInto(rflFiles,fileList[i])));
   nRFL = NItems(rflFiles);
```

The first command, NamesDefaultToHere(1), says that all the names in this script will be kept in the local environment so that different scripts won't interfere with each other. This is one of the most important enhancements for JMP 9.

Now we sequence through all the files we found. Each one is brought in as one large binary chunk, a 'blob'—binary large object, using a new option in JMP 9. We figure out how many lines are in the image.

```
for(iFile=1,iFile<=nRFL,iFile++,
  file = rflFiles[iFile];
  fBlob = loadTextFile(inDir||file,blob);
  sz = length(fBlob);
  nchannel = 224;      // number of colors/wavelengths
  nsample = 614;       // number of pixels across each line
  nLines = 512;        // number of lines in the image
  rawData = Blob To Matrix( fBlob, "int", 2, "big", nchannel ); fBlob = 0;
  nLines = (sz/2)/(nsample*nchannel); // sometimes it's not 512
  n = nsample*nLines;  // number a data table rows
```

Next we standardize the data, taking care not to divide by zero for constant columns, and calculate the Correlation matrix. The Correlation function is new, and it is very fast, using multithreading.

JMP 9 and Interactive Statistical Discovery, continued

```
stdv = vstd(rawData); stdv += stdv==0; // don't divide by zero
stdData = (rawData-J(n,1,1)*vmean(rawData)):/(J(n,1,1)*stdv); //standardize
rawData = 0;
corr = Correlation(stdData);
```

Next, we calculate the eigen decomposition of the correlation matrix and use the eigenvectors to score the first four principal components, with the scores reshaped back to the dimensions of the image.

```
{m,e} = eigen(corr); show(m[1::10,0]);
pc1v = shape(stdData*e[0,1],nLines,nsample)`;
pc2v = shape(stdData*e[0,2],nLines,nsample)`;
pc3v = shape(stdData*e[0,3],nLines,nsample)`;
pc4v = shape(stdData*e[0,4],nLines,nsample)`;
```

Now we scale the results so that we can colorize by values from 0 to 1.

```
pc1s = (.5+pc1v/(2*max(pc1v))):*(pc1v>0)+(.5-pc1v/(2*min(pc1v))):*(pc1v<0);
pc2s = (.5+pc2v/(2*max(pc2v))):*(pc2v>0)+(.5-pc2v/(2*min(pc2v))):*(pc2v<0);
pc3s = (.5+pc3v/(2*max(pc3v))):*(pc3v>0)+(.5-pc3v/(2*min(pc3v))):*(pc3v<0);
pc4s = (.5+pc4v/(2*max(pc4v))):*(pc4v>0)+(.5-pc4v/(2*min(pc4v))):*(pc4v<0);
stdData = u = v = 0;  // free memory
```

Now we convert to the principal images using default heatmap colorization: blue for low, to gray for medium, to red for high. Image objects and functions are new in JMP 9.

```
imagepc1 = NewImage(HeatColor(pc1s));
imagepc2 = NewImage(HeatColor(pc2s));
imagepc3 = NewImage(HeatColor(pc3s));
imagepc4 = NewImage(HeatColor(pc4s));
```

The images are saved on disk.

```
imagepc1<<SaveImage(outDir||rflFiles[iFile]||"_PC1.jpg");
imagepc2<<SaveImage(outDir||rflFiles[iFile]||"_PC2.jpg");
imagepc3<<SaveImage(outDir||rflFiles[iFile]||"_PC3.jpg");
imagepc4<<SaveImage(outDir||rflFiles[iFile]||"_PC4.jpg");
```

In addition to the principal images, we also create images of the eigenvector loadings so that we can make a legend across the top showing which wavelengths are contributing to each principal image.

```
imageEig1 = NewImage(DirectProduct(HeatColor(e[0,1]`/2+.5),J(12,2,1)));
imageEig2 = NewImage(DirectProduct(HeatColor(e[0,2]`/2+.5),J(12,2,1)));
imageEig3 = NewImage(DirectProduct(HeatColor(e[0,3]`/2+.5),J(12,2,1)));
imageEig4 = NewImage(DirectProduct(HeatColor(e[0,4]`/2+.5),J(12,2,1)));
```

Now we arrange the images into a presentation and show it:

```
        NewWindow(rflFiles[iFile]||" PCA",
             LineupBox(ncol(2),spacing(3),
                    imageEig1,imageEig2,
                    imagepc1,imagepc2,
                    imageEig3,imageEig4,
                    imagepc3,imagepc4));
```

This is a loop, which we need to close with a parenthesis.

```
);
```

If you run this script, you will get a window for each hyperspectral 'rfl' file, looking like the images shown in Figure 1.

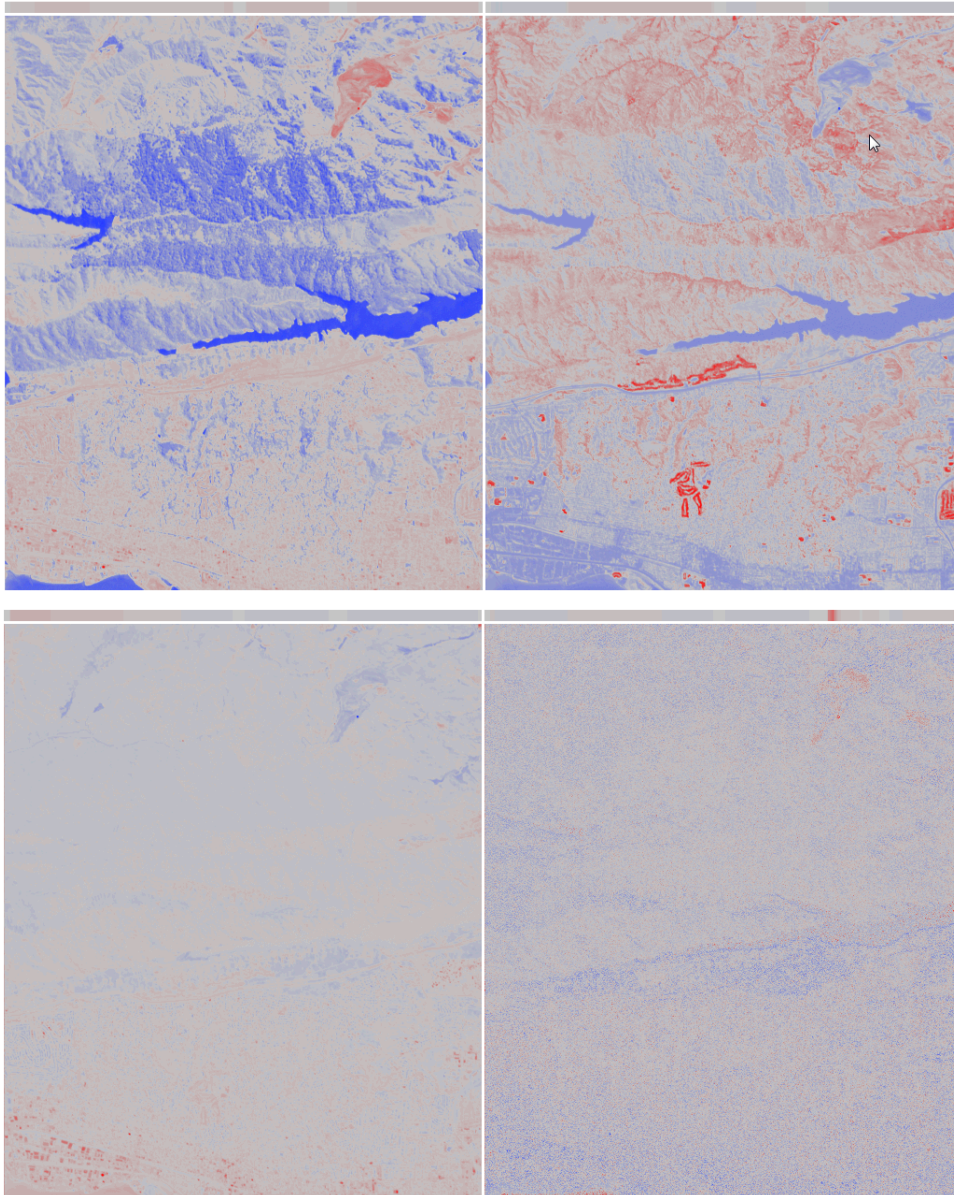JMP 9 and Interactive Statistical Discovery, continued



**Figure 1.  Images Created from Principal Components Computed from .rfl Hyperspectral Data**
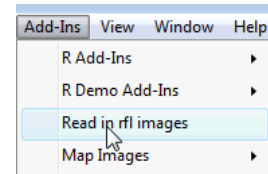
JMP 9 and Interactive Statistical Discovery, continued

Notice that the first principal image concentrates on general intensity (the loading legend shows reddish values across the spectrum). The second principal image compares red wavelengths and blue (the legend loads positively in the lower middle wavelengths and negatively on the higher ones). The third principal image concentrates on infrared (legend positive at the low wavelengths). The fouth principal image loads in a very narrow band in the blue range of the spectrum, and the image looks rougher than the others. Each image emphasizes different features in the scene.

Considering that this is a 240 column by 314,368 row problem, it is amazing that it only takes about 30 seconds on my 8-core laptop to do this and everything else in this script per image.

Now consider that you want to enable this script to be easy to use for yourself. You install a menu item to invoke it.

Also, you want it easy to install for many users. JMP 9 now supports 'add-ins' which zip together a script with some information about how to install it in the menus. This add-in file can now be put on a website. Any user can download it and drag it into JMP, and it automatically installs with a menu item to invoke it.

## EXPLORING DATA

Exploring data is JMP's special talent.

### COLOR JMP DATA TABLE CELLS

One of the graphs in JMP is a cellplot, sometimes called a heatmap as shown on the left in Figure 2, which shows the data with colored rectangles. With JMP 9, the spreadsheet itself can be colorized, either by value across a column or individually per cell.



**Figure 2.  Color using a Cell Plot (left) and JMP Data Table Cells (right)**

### FEATURES IN GRAPH BUILDER

Graph Builder has many new skills. In the Graph Builder example on crime data (Figure 3 and Figure 4), I can click on a column name  'State' which identifies U.S. States, drag it to the new Shape drop zone, and it brings up a map and sets the graph coordinates to latitude by longitude. Graph Builder has looked through all the maps it has to find out which supports the values in the column.

Then I drag the column name  'Violent Rate' to the Color drop zone and the states are colorized, with red for high and blue for low.

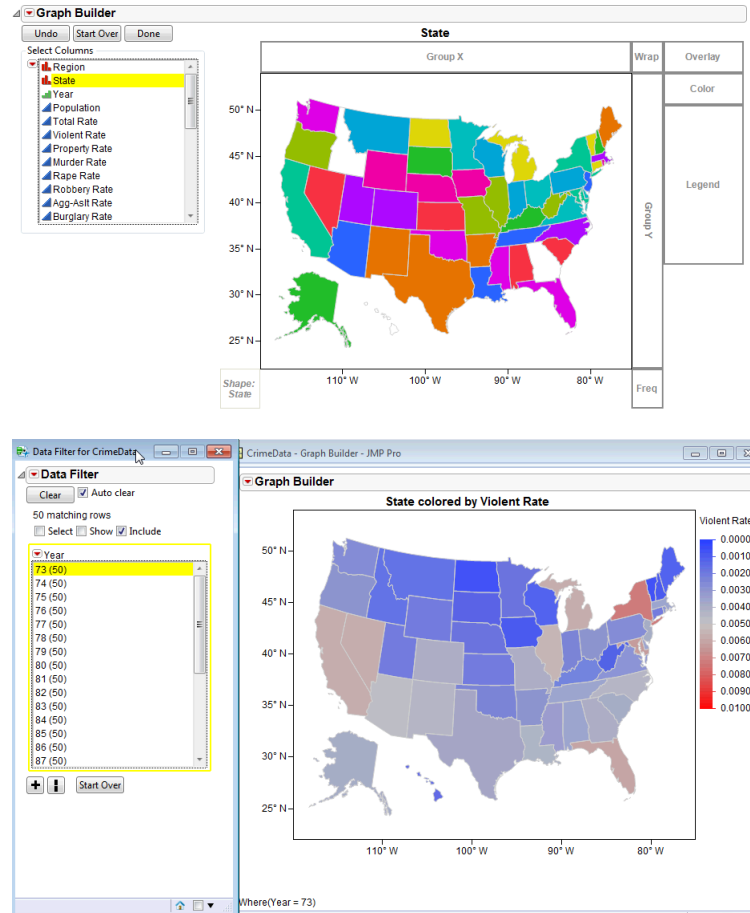JMP 9 and Interactive Statistical Discovery, continued





**Figure 3.  Color Map Data in Graph Builder by Crime Rates**

Here is the crime picture for 1991, the height of the 'crack' epidemic, then in 1999 when violent crime had calmed down. This can be played like a movie, letting you watch violent-crime rate patterns both geographically and across time.
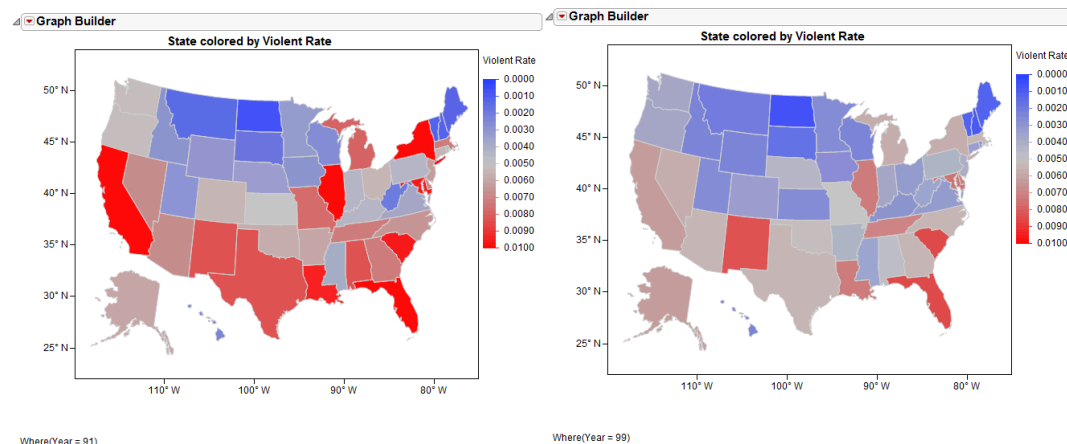


**Figure 4.  Look at Changing Crime Rates over Time**

JMP 9 and Interactive Statistical Discovery, continued

## EXPLORING MULTIVARIATE DATA

Suppose you need to characterize the multivariate distribution of several variables, in this case 30,000 rows of FCS data. The KMeans platform has many features for doing this. The data may be noisy, so as a first step you find the distance of each point to the k[th] nearest neighbor, as shown in Figure 5. Selecting the 900 or so points that are far from the 15[th] nearest neighbor can be done by dragging a rectangle in the 15[th] plot below. Notice that in JMP 9 the points that are not selected are faded, where in earlier versions of JMP the selected points were bigger. The fade selection mode is much better for when there are many points and you need to spot where the selected points are.



**Figure 5.  Nearest Neighbor Plots**

Once these points are selected, you can see where they are in the space, using the 3D scatterplot, and the scatterplot matrix.



**Figure 6.  Noisy Points Showing in 3D Plot and Scatterplot Matrix**

Now that we have concentrated the sample to points that are near other points, it is ready to find the structure of the distribution, using normal mixtures. JMP 9 has a much faster Normal Mixtures facility. There are a variety of ways to see the resulting mixture. There is the principal component score plot (left-hand plot in Figure 7) with shaded ellipses showing the normal contours and with circles indicating the portions in each cluster. The same idea can be done in 3D with semi-transparent ellipsoids (middle plot in Figure 7).

JMP 9 and Interactive Statistical Discovery, continued

Parallel coordinates plots show the structure of each cluster, assigning points to the closest cluster in the Mahalanobis sense (right-hand plot in Figure 7).



**Figure 7.  PCA Score Plot (left), 3D Plot with Ellipsoids (middle), Parallel Coordinate Plot (right)**

Also, there is a feature to output the scatterplot matrix with ellipsoids projected into each component pair's space. Here we faded the points by changing the transparency setting.



**Figure 8.  Scatterplot Matrix of Normal Mixture Clusters using Transparency Features**


## ALIAS-OPTIMAL EXPERIMENTAL DESIGN

JMP has a tradition of strong support for experimental design. Suppose that you have to make a screening experiment for six factors in 12 runs. The plain d-optimal design is efficient for the main effects, but if there is a big two-factor interaction lurking that is not estimable here, it will contaminate the main effects estimates due to the large number of 1/3 and -1/3 correlations confounding main effects with two-factor interactions (left picture in Figure 9). But if I change to an 'alias optimal' criterion, all these will be cleaned up—to zero in this situation, as seen in the all-zero (all-blue) section in the upper right and lower left (right picture in Figure 9).

JMP 9 and Interactive Statistical Discovery, continued



**Figure 9.  Correlation Maps for D-Optimal Design (left) and Alias Optimal Design (Right)**

The alias-optimal design is trading some efficiency in the main effects with robustness to two-factor interactions. This is not always a trade you will want to make, but if you suspect strong two-factor interactions, but don't have the runs to make them estimable, this kind of design is worth considering.

## MODELING

The Concrete Slump Test data, contributed to the UCI Machine Learning lab, is a very large-scale response surface experiment with three responses, 7 factors, and 103 runs. When expanded to all the quadratic terms, the model has 36 parameters for each response. To validate the predictive ability of fits, I created validation and test sets, identified by a Validation column in the data table. Using the Stepwise personality, if you fit all the terms for the first response (Slump), you get an Rsquare of 0.79 for the training data. However, the Rsquare for the validation data is very terrible ( -1.415)—i.e. the validation data is fit far worse than just a simple mean (left-hand report in Figure 10). If you clear the terms and then click Go to step in terms until the crossvalidation Rsquare is maximized, then the resulting model validates much better with a validation Rsquare of 0.38. All these cross validation features are new to JMP 9.



**Figure 10.  Stepwise Example of Cross Validation**

In order to fit all the responses in one click on the Go button, I hold down the control key to broadcast that Go. Then I can control-click the Run Model button to run all the resulting models in the Standard Least Squares platform. The fits are assembled in the new Fit Group, allowing different models to support a single combined Profiler to explore the response surface across all the factors and responses (see Figure 11). In earlier versions of JMP, this would have involved many more steps and saving lots of prediction formulas as data table columns, and the result would not have included confidence intervals.

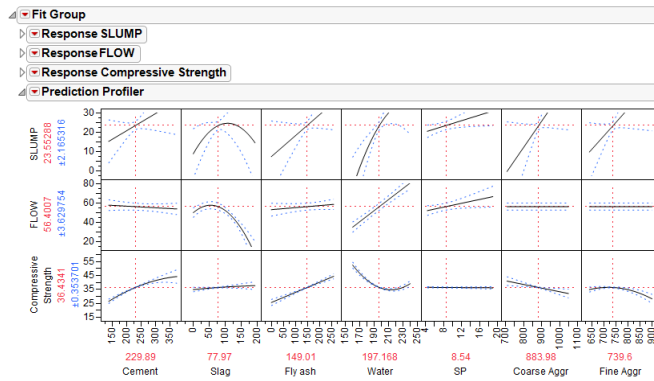JMP 9 and Interactive Statistical Discovery, continued



**Figure 11.  Profiler for Concrete Slump Test Data**

## SURVEY ANALYSIS

To see JMP at work analyzing surveys, I went to the website of a very large survey on Civic Engagement in America (http://www.bowlingalone.com/data.htm), the 'Bowling Alone' data. A convenient form of the data is in an SPSS save file. You can now easily download that and open it in JMP, and it preserves all the value labels.



**Figure 12.  Partial Listing of the Bowling Alone Data**

Notice that I opted to use the long names, which are very descriptive. However, these long names have a common phrase "(freq last 12 months)" repeated that will distort the compactness of the results (see Figure 12). It is easy to change them. I select all 104 variables for these frequency value and bring up the Search dialog to change them, clicking 'Replace All'. While all these columns are still highlighted, I can click the analysis type and change them to Ordinal so that their values will be charted with an ordinal theme. Also I check the column properties to verify the value labels from SPSS (see Figure 13).

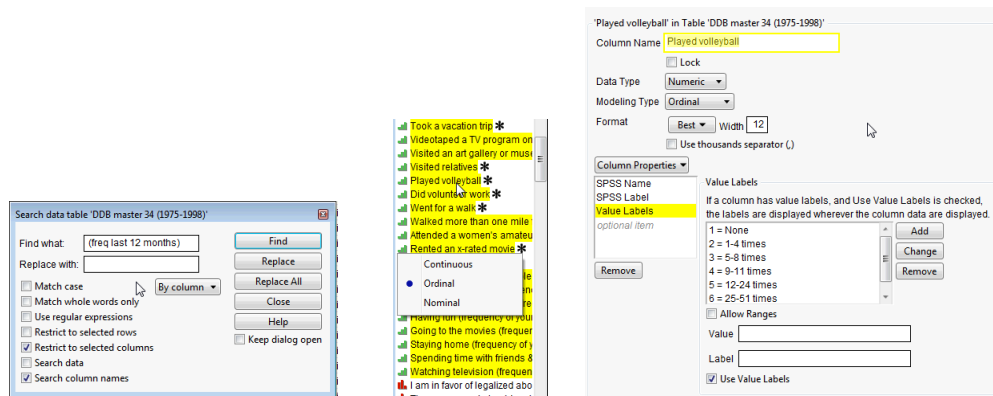JMP 9 and Interactive Statistical Discovery, continued

**Figure 13.  Search Dialog (left), Change Modeling Type (middle), Assign Value Labels (right).**

Now I invoke the Categorical platform and click Separate Responses for the selected 104 columns, and then Year of Survey for my X grouping column.
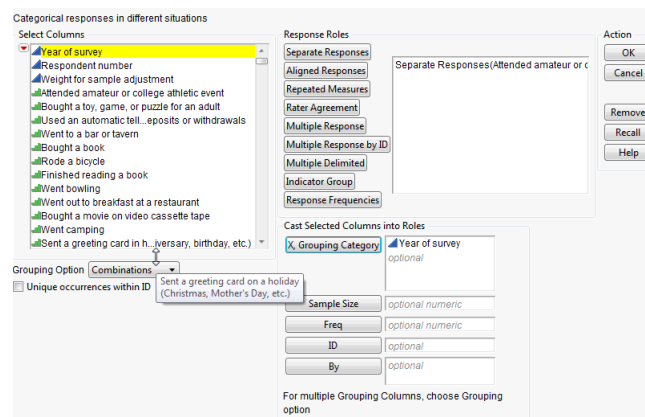
**Figure 14.  Categorical Analysis Launch Dialog**

The default report is to show the frequency tables and response rate tables. To make the result more compact, I can uncheck the Frequencies and Share of Responses and Legend to get only the Share Charts, as shown by the right-hand charts in Figure 15. On the right, you see the change over time in two of the variables, riding a bicycle and reading a whole book. There are 102 more like this.

JMP 9 and Interactive Statistical Discovery, continued



**Figure 15. Frequency Table and Share Charts from Categorical Platform**

Now I want to see how the response 'Finished reading a book' breaks down by region, sex, and educational level. To find these columns among the 389 columns in the table, I use the new search control in the column selection list. Then I select 'Each Individually' to get a separate breakdown for each of the three X variables.
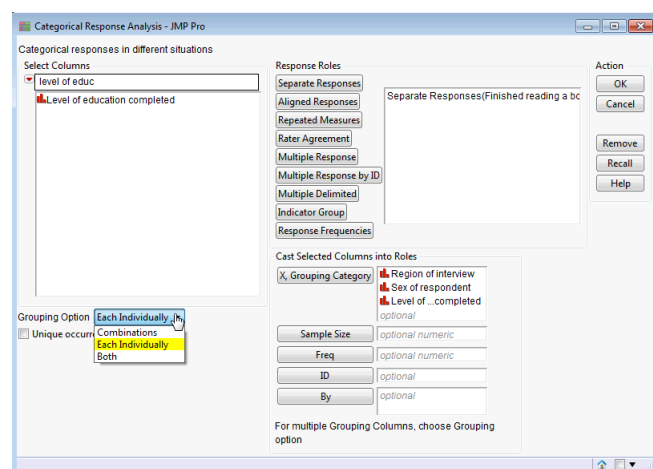


**Figure 16. Categorical Launch Window Specifying Grouping Variables and Grouping Option**

To the resulting analysis, I make similar option selections to get just the charts, and now I see which regions, which genders, and which educational levels are associated with more frequent reading (Figure 17).

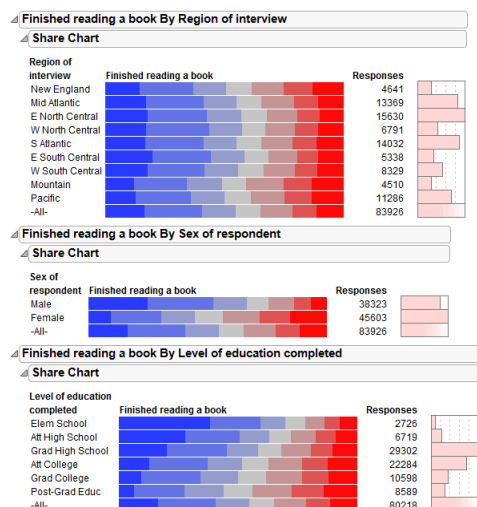JMP 9 and Interactive Statistical Discovery, continued



**Figure 17.  Share Charts from Categorical Analysis Results**

## EXCEL MODELING

When you estimate models in JMP, the Profiler is a great way to explore the response surface. But what if your model is a set of business calculations in an Excel spreadsheet? Moving all those calculations to JMP data table formulas is too laborious to be practical in most cases. With JMP 9, you can profile a model that is still in Excel. Here is a financial model for the Airbus 380 in Excel. Note that additional command items have been added for JMP. With the 'Edit/New Model' button, I can specify which cells are dependent responses, and which cells correspond to inputs that I want to profile against (Figure 18).
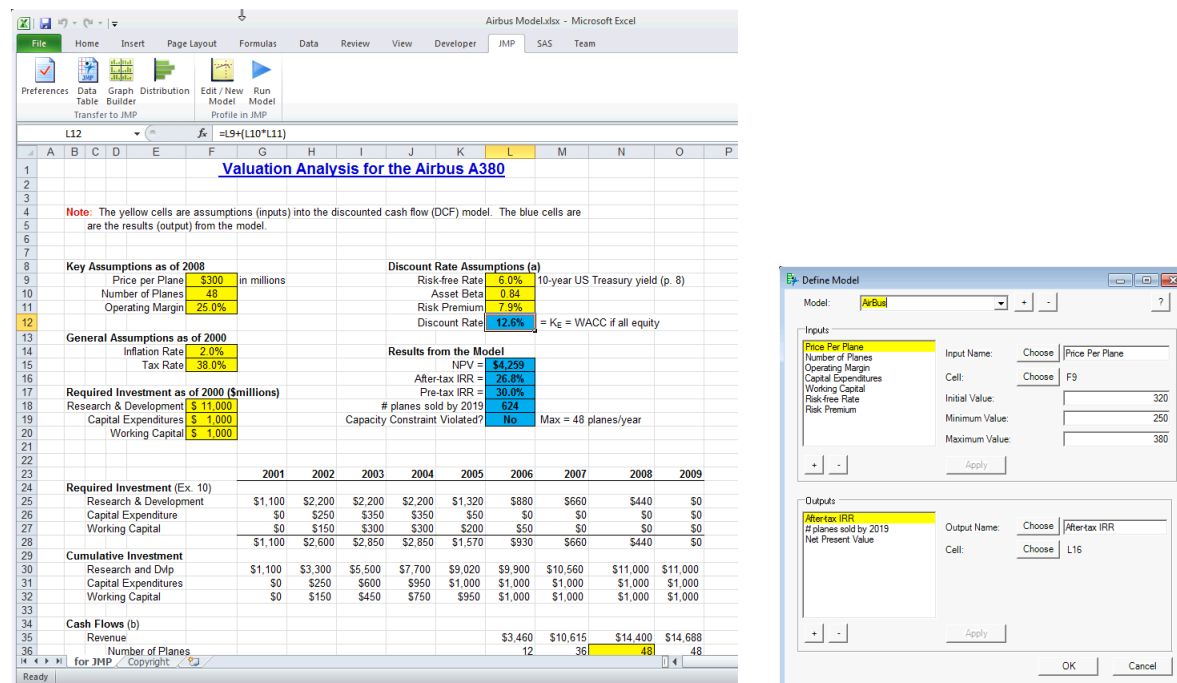


**Figure 18.  Excel Spreadsheet with Formulas and Define Model Dialog for JMP**

When I click Run Model, it connects the model with JMP and JMP repeatedly has Excel re-evaluate the model under many different factors settings, resulting in the Profiler in Figure 19. You can drag the current values around to see how the other profile traces are affected. This gives you a great view of how changes in any variable would affect all the responses.

JMP 9 and Interactive Statistical Discovery, continued

You can even bring up a simulator to explore distributions of the responses with respect to random distributions in one or more factors.
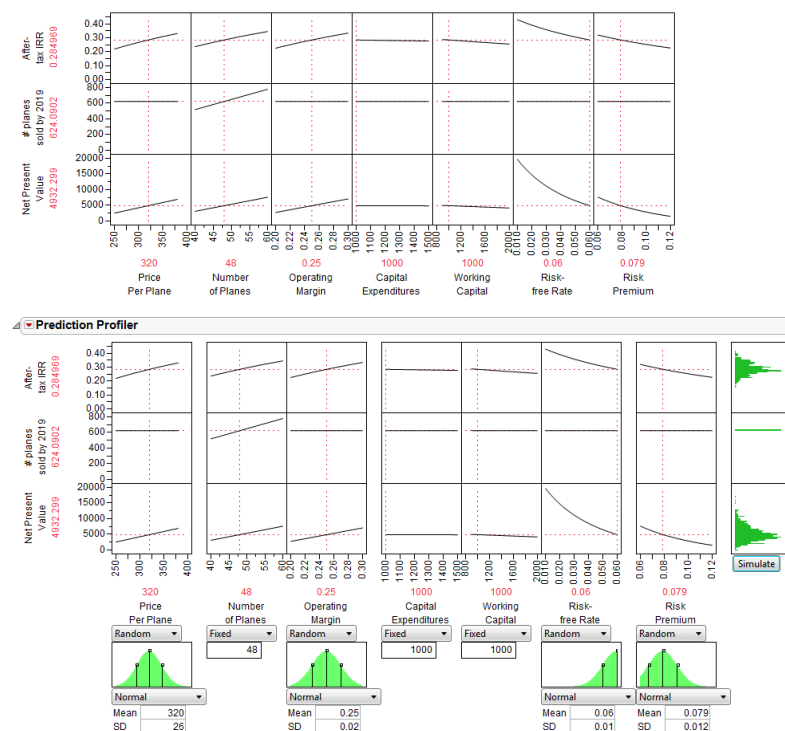


**Figure 19.  Profiler Display and Simulator Using Excel Formula**

## DATA MINING

In JMP Pro, Version 9, the Partition platform, which creates decision trees, includes new crossvalidation features—the splitting stops when the validation measures of fit stop improving the fit.

In addition, there are two new measures implemented. In bootstrap forest,

- a number of decision trees are built based on different bootstrap samples
- candidate set columns are selected randomly

These features usually perform better than a straight decision tree. For example, lets use several methods to fit the Wisconsin breast cancer data from the UCI repository (Figure 20 and Figure 21).



**Figure 20.  Partial Listing of Wisconsin Breast Cancer Data**

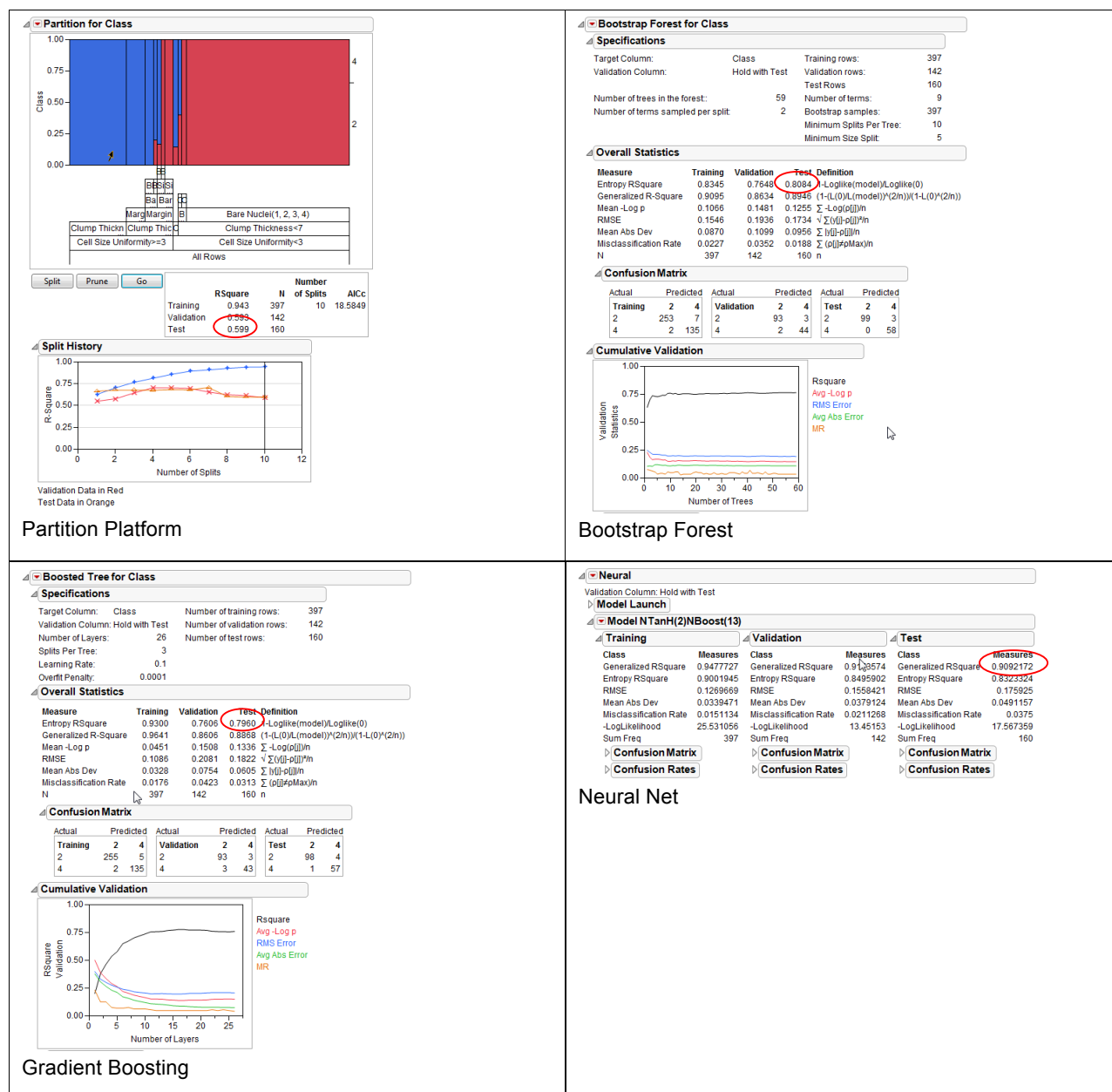JMP 9 and Interactive Statistical Discovery, continued



**Figure 21. Several Methods of Wisconsin Breast Cancer Data**

The best performing method differs for different data sets, and even for how you randomly choose the training, validation, and test sets, and differs each time due to randomness in the methods. But, usually Neural Net is the best and a single decision tree is the worst. Table 1 summarizes the fitting results as reflected by R Square values.

| Method | Test RSquare (entropy) |
|---|---|
| Decision Tree | 0.599 |
| Bootstrap Forest | 0.808 |
| Gradient Boosting | 0.796 |
| New Neural Net | 0.909 |

**Table 1. Summary of Fitting Results**

JMP 9 and Interactive Statistical Discovery, continued

## CONCLUSION

JMP 9 has many new features, and only a few of them are described in this paper. With renewed features in scripting, we expect JMP to be a good implementation language for many new applications.  In addition to JMP's traditional user base of engineers and scientists, we expect other researchers and analysts to discover JMP, both as a stand-alone application, and as a front-end and explorer part of the SAS system.

## CONTACT INFORMATION

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

John Sall
SAS Institute Inc.
SAS Campus Drive
Cary, NC  27513
(919) 677-8000
john.sall@jmp.com
www.jmp.com