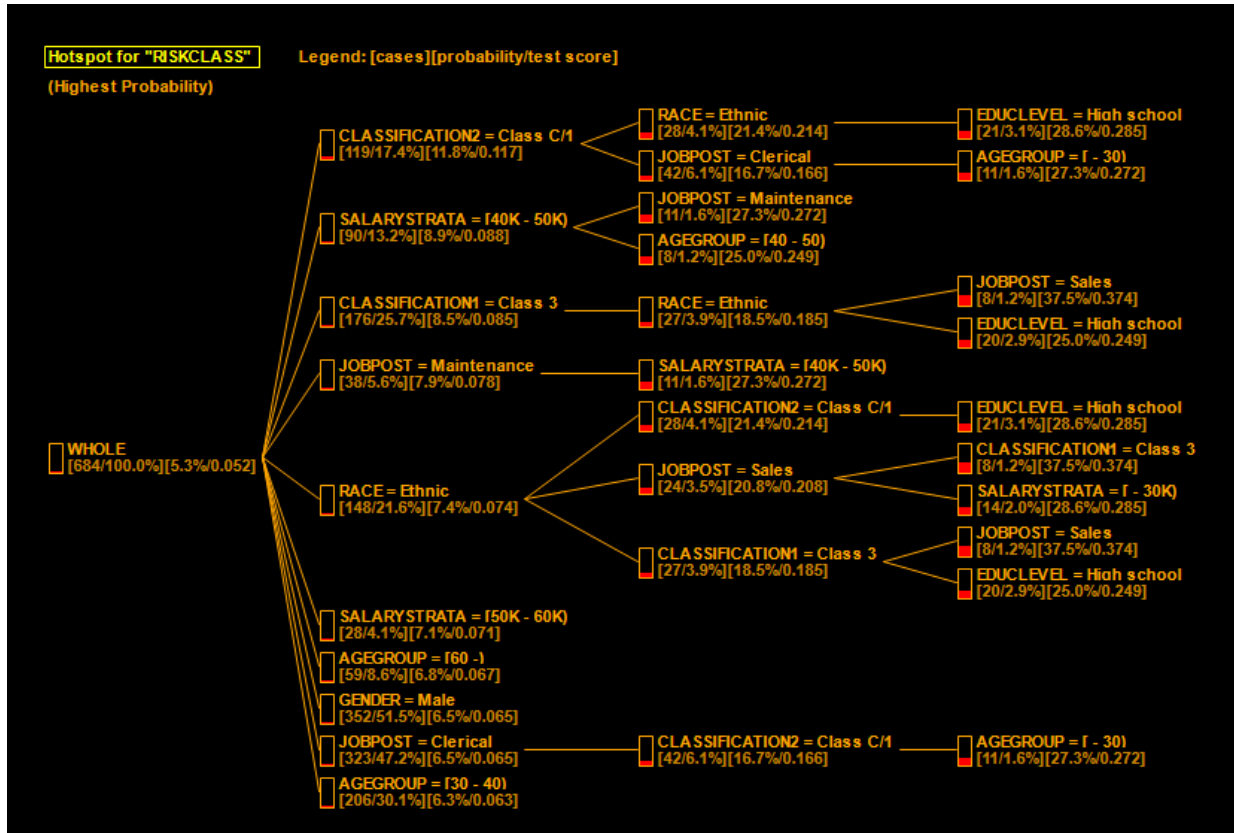# Variable Relevancy Analysis in Predictive Modeling

## Variable Relevancy Analysis

The purpose of relevancy analysis is to find variables which contain *predictive information* to be used in predictive models. Low relevancy variables do not provide information that can be used in predictive models. Rather they can cause *overfitting*. Selection of modeling variables through relevancy analysis is very important to improve *predictive accuracy* and avoid overfitting. Most of predictive modeling work is spent on data preparation and variable relevancy analysis. The rest of modeling work is rather straightforward. Three relevancy analysis methods are described in the subsequent sections.

## 1. Hotspot Drill-down Profiling

CMSR Hotspot Drill-down tools can identify *profiles* of major risky customer segments. The following figure shows a hotspot drill-down example;



Nodes of the figure show customer segments with highest ratios of "risk" occurrences. Red portions of boxes indicate ratios (or probabilities) of risk. *Variables and categorical items appearing in this output can be a good candidate for predictive modeling.*

## 2. Correlation Analysis

CMSR correlation analysis can be used to identify relevant variables that show significant (linear) correlation to risk. The following figure shows a CMSR correlation analysis output;
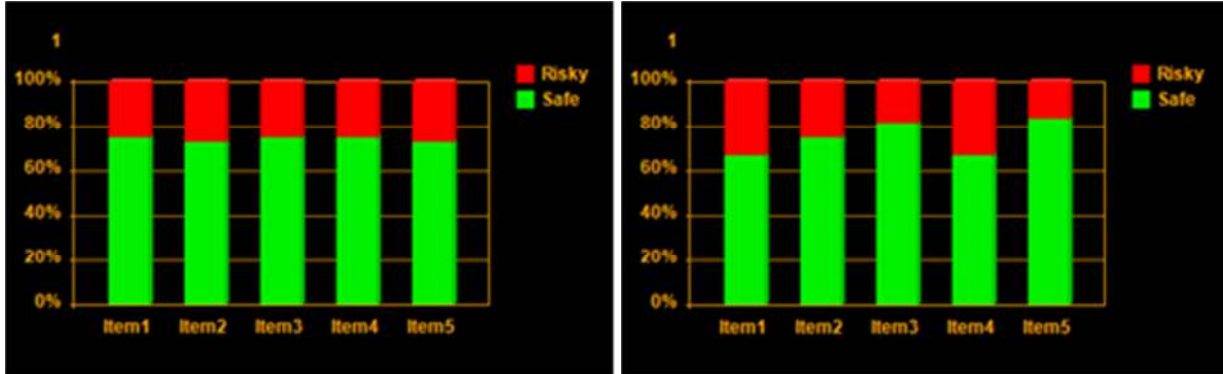


"r-value" is the coefficients for linear correlation. Having 0 means no correlation. 1.0 is the perfect *positive* linear correlation and -1.0 is the perfect *negative* linear correlation. "r-square" is the squared values of "r-value". *Variables with high absolute "r-value" are good candidates for predictive modeling.*

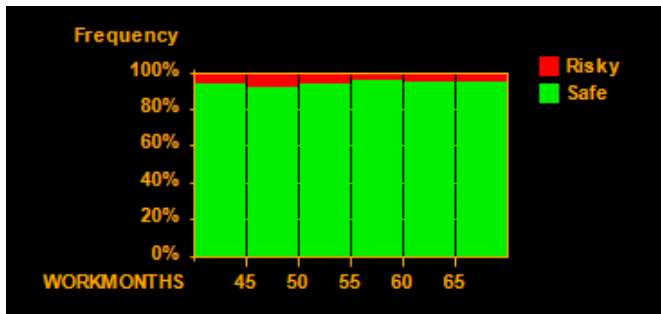# Variable Relevancy Analysis in Predictive Modeling

## 3. Categorical Bar/Histogram Charts

Relevant variables can be identified using CMSR hotspot drill-down and correlation analysis tools, as explained in the previous sections. In addition, CMSR categorical bar and histogram charts can be used to *identify* and to *verify* relevant variables. The following two bar charts show risk distribution by categorical items of two different categorical variables.



The left bar chart shows little variation in risk distribution amongst categorical items. This type of variables is not useful for predictive modeling. Note that predictive modeling relies on variations amongst different items. *The right bar chart, on the other hand, shows significant variations in risk distribution. This type of variables can be useful in predictive modeling and can be used in predictive models.*

Similar analysis can be performed on linear (= numerical) variables using categorical histogram charts. The following figure shows categorical distribution of risk over a numerical variable. Red parts show only some variations amongst different numerical value ranges. But close examination shows relative size differences of red proportions, which represent risk, are significantly large. So this variable *may* be used in predictive models.



Categorical bar and histogram charts can be used to *verify* suitability of variables identified using hotspot drill-down analysis and correlation analysis.

(It is important to note that your local anti-discrimination laws may bar to use certain demographic variables if your models are used in discriminatory purposes. Check your local anti-discrimination laws.)

**Rosella**

**Rosella Software**
**www.roselladb.com**