# Salford Predictive Modeler®

# Model Compression via ISLE and RuleLearner®

*This guide elaborates on powerful ways to combine the TreeNet®
and GPS engines to achieve model compression and more.*

Minitab

**Minitab** ➤®

# Background and Motivation

The principal idea behind model compression is that it may be possible to use modern regression techniques to post-process a complex model to create a radically simplified version that can perform almost as well as the original.  It can go even further: if you build the original complex model in a specific way there is an excellent chance that the post-processed simplified model will perform even better than the best possible complex model!

Model compression is important for several reasons:

♦ Compressed models can perform dramatically faster in scoring.  For real time scoring compression may be the only way to feasibly deploy a high performance TreeNet®.

♦ Compressed models may use many fewer predictors; compression may thus offer an alternative way to accomplish feature selection.

♦ Compression may provide models that are easier to explain using the component trees.

In this chapter we will cover model compression basics and try to explain why it can be expected to work. To begin, you need to understand that any model compression is a four stage process.

**First**, we use one of our essential data mining engines to build a predictive model.  The engine could be CART®, MARS®, TreeNet®, or Random Forests®.

**Second**, we use the predictive model obtained to convert the original data (train and test partitions) into a new transformed version of the data.

**Third**, we use a regularized regression (GPS) to build a new model using the transformed data.

**Finally**, we extract a new compressed version of the original model.

✓ In our first release of model compression the first stage model must be a TreeNet, but in later releases the other engines can be used instead.
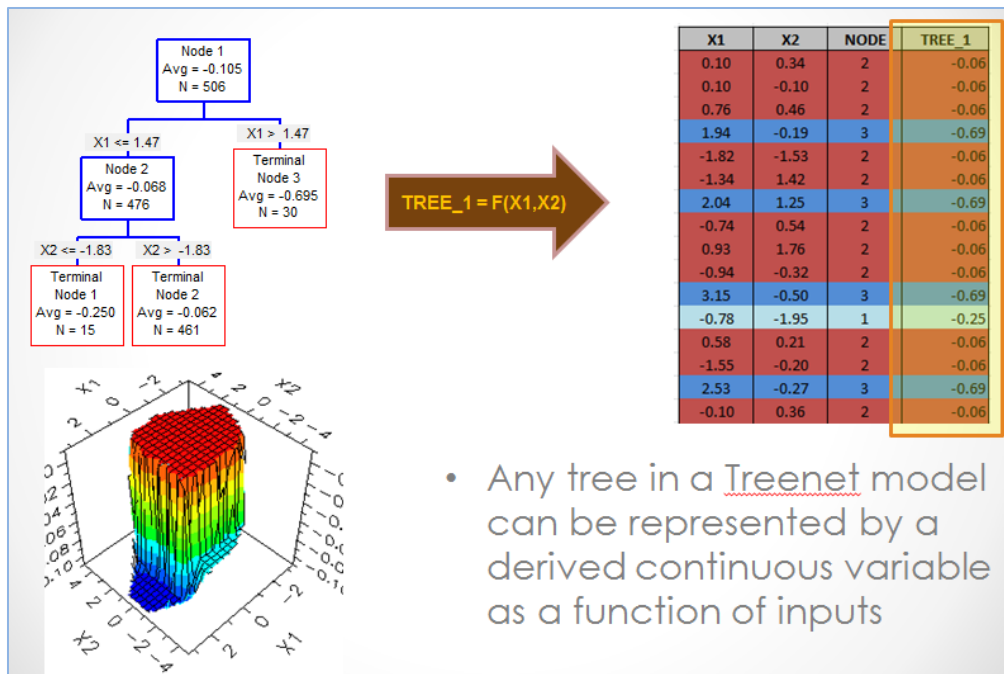
To understand how to work with model compression you need to have a basic grasp of working with the GPS regularized regression engine described earlier.  If you have not already done so, please review those basics as we will assume this understanding in the discussion that follows.

The second step above requires further elaboration as there are several possible ways to re-express a TreeNet model into a linear combination of components later to be fed into the GPS engine for compression purposes.  Two fundamentally different approaches arise here, one is called *ISLE* and the other is called *RuleLearner®*.  We start by explaining the first one and then follow up with the second.

**ISLE** is an acronym for *Importance Sampled Learning Ensembles*, introduced by Jerome Friedman with co-author Bogdan Popescu in 2003.

✓ The technical paper is available on our website at http://www.salford-systems.com

While the ISLE approach has a far reaching theoretical setup, for all practical purposes it can be introduced as decomposing a TreeNet model into a set of variable transformations represented by each individual tree in that model.  This stems from a simple observation that any tree can be viewed as a potentially multivariate transformation of the original predictors.  This is illustrated below.

Minitab ▸®

One can therefore represent any TreeNet model as a linear combination of individual tree variables with coefficients 1.0.

✓  This assumes that the learn rate and initial guess adjustments were incorporated into the trees, which is trivial to do.

Thus, after these variables have been fed as inputs to the GPS engine, a new linear combination will have different coefficients.  Given the GPS's inherent ability to do variable extraction by trying different variable selection strategies (elasticities), a fraction of predictors may have zero coefficients which in this context is equivalent to removing the corresponding tree from the model thus achieving model compression.  As an added bonus, it is possible that the new model will have superior performance.

The ISLE compression is visualized on the next diagram.



Note how some trees have been eliminated (zero coefficients) while the others were adjusted to provide the best fit.

Another way to think about the model compression is to recall that a TreeNet model is an ensemble model consisting of many small CART-style decision trees.  Each tree produces an output "score", and the final model is simply the grand total of all the individual scores.  Usually TreeNet approaches its final result slowly, with each tree contributing a small part of the final model result.  When we "compress" the

TreeNet model via GPS regression we select a subset of these trees and reweight them in order to arrive at a model which has fewer trees and possibly better performance as well. The abstract technical mechanism by which this is accomplished is best understood by studying GPS; here we provide an intuitive explanation for what is going on.
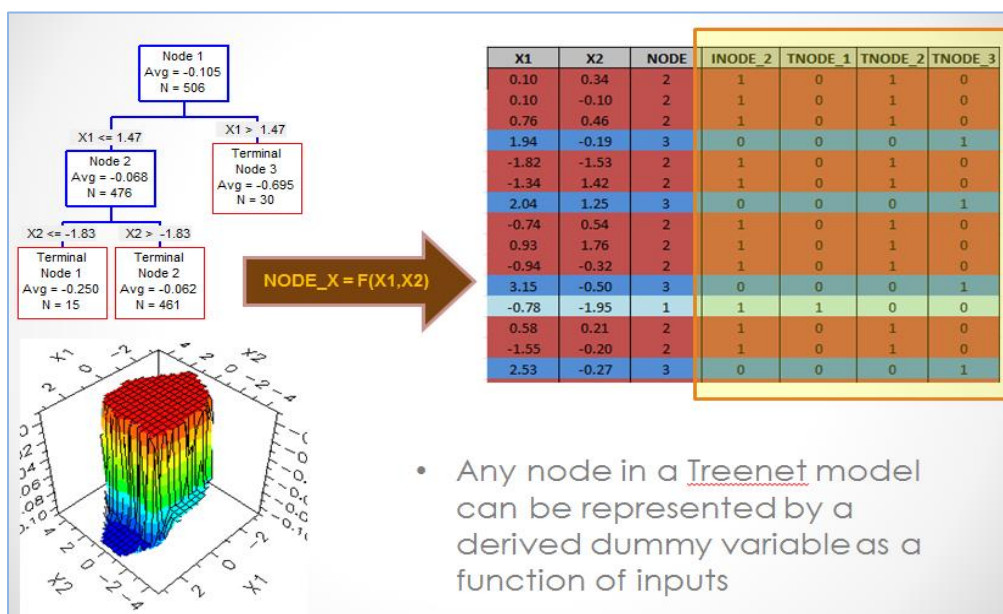
A TreeNet model is built myopically, looking only one step ahead as the model is built tree-by-tree. To obtain good results, each step is small, meaning that each tree can add only a very small amount to the overall TreeNet prediction. When the TreeNet model is complete we typically have very many of these trees and there is an excellent chance that some tree structures are repeated exactly, or almost exactly. The compression phase can then work as follows:

♦ Identical or almost identical trees can be aggregated into just one tree with an appropriate weight. If the same tree was grown 20 times we could replace these trees with just one tree and provide it with a weight of 20. In real world modeling the trees may differ in details which make exact reproduction in this way impossible, but we may be able to produce a very close approximation.

♦ When the TreeNet model is grown we do not know what the next tree will bring and in the process of arriving at final scores there may be some oscillation around a final solution or even backtracking in which some trees are essentially cancelled out of the final result. Once the model has been constructed we can accomplish any "cancelling out" by simply deleting trees and we can skip over any oscillations. The end result will be a smaller collection of trees.

The final result of the ISLE Model compression is a smaller model of the same type we started with. In this case, we always start with a TreeNet and we end up with a TreeNet. But the new TreeNet has fewer trees and each tree has its own weight which could be far different than 1.0!
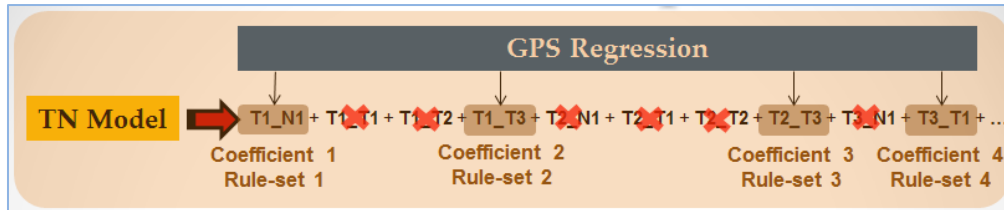
The second approach (here introduced as the **RuleLearner**) of decomposing a TreeNet model into individual derived variables is introduced by pushing the decomposition process further down, all the way to the individual nodes. Observe that any tree spawns a set of dummy variables, one variable per each terminal and internal node.

✓ The only exception is the root node which includes the entire sample and would result to a trivial constant dummy variable.



• Any node in a Treenet model can be represented by a derived dummy variable as a function of inputs

Thus the complete set of node dummies represents all structural information contained in any TreeNet model. Note that the set is overrepresented, internal nodes largely overlap with their parents and completely include their children. This opens up even greater opportunities for possible model simplification by identifying the subset of influential nodes and finding coefficient adjustments associated with each node such that the original model fit is essentially recovered. Again, just as in the case of ISLE approach, the GPS regularized regression engine is perfectly suited for this task. Furthermore, it is possible that the resulting model may achieve higher performance compared to the original TreeNet model while maintaining an excellent compression ratio.

This process is illustrated on the following diagram.



Observe that each node dummy represents a collection of rules in terms of the original predictors leading into this node. Therefore, another way to look at this type of compression is in terms of extracting a collection of informative rule-sets (segments). Each rule-set operates on a subset of data or, to paraphrase, each rule-set represents a segment of the data. The regression coefficient assigned to this segment by the GPS is simply the segment-based adjustment to the overall model prediction. Note that segments can overlap with each other and may have either positive or negative adjustments of varying magnitude thus making the model prediction extremely versatile.

Both types of TreeNet model compression (ISLE and RuleLearner) are implemented in SPM®. In the following sections we describe the basic steps needed to setup model compression as well as explain various output windows produced along with the possible usage modes.

**Minitab >®**

# Model Compression in Action – ISLE EXAMPLE (coupled with RULELEARNER sub-example)

At this point it will suffice to say that setting up a model compression run from a user standpoint is very straightforward.  We are essentially creating a "pipeline", a sequence of two already introduced engines, first a TreeNet engine using the original predictors then followed up by a GPS engine that uses decomposed TreeNet model as its own inputs.  The decomposition itself is done automatically by SPM based on the user selected method, which could be either ISLE or RuleLearner (as well as additional trivial modifications introduced below).  In addition, it is assumed that each engine component of the pipeline (TreeNet and GPS) is configured using its own set of control parameters as usual.

The end result of the compression pipeline is the overall compression summary plus the usual output of the TreeNet engine (as if it were run by itself) as well as the standard GPS output of the follow up compression phase.

The output displays are deliberately designed for ease of comparison between the original uncompressed model and the resulting compressed model.  The user can choose any desired model for viewing, scoring, and translation using the "point and click" approach thus balancing the degree of compression versus model size.

## Setting up Model Compression

We start with the GOODBAD.CSV data set which is rather on the small side but runs very quickly and illustrates our main points well.  The dataset comes from a credit default study and has 664 observations and 14 variables including a binary response variable named TARGET, with 1 indicating a bad credit standing account and 0 otherwise.  The remaining predictors provide a nice mix of continuous and categorical measures which represent a typical modeling environment.  The goal of the study is to build a predictive model for bad credit using the TreeNet approach and then possibly compress the model using either ISLE or RuleLearner approaches.

Use the **File/Open/Data File…** menu item to open the GOODBAD.CSV dataset, then click the **[Model…]** button in the resulting activity window to open the **Model Setup** window.

In the Model tab:

♦ Set **Analysis Engine** to **ISLE Model Compression** – this activates the pipeline part of the SPM – the **ISLE** tab. This exercise is useful for elaborating on the **RuleLearner** component as well.

✓ Note that this mode activates the **TreeNet** tab and the **GPS** and **GPS Advanced** tabs, thus giving you full control over other parts of the pipeline.

♦ Set **Sort:** to **File Order** – this sorts all variables according to the position in the dataset.

♦ Check TARGET variable as the **Target**.

♦ Check all remaining variables as predictors.

♦ Select **Target Type** to **Classification/Logistic Binary** to take advantage of the natural capability of both TreeNet and GPS engine to handle binary targets effectively.
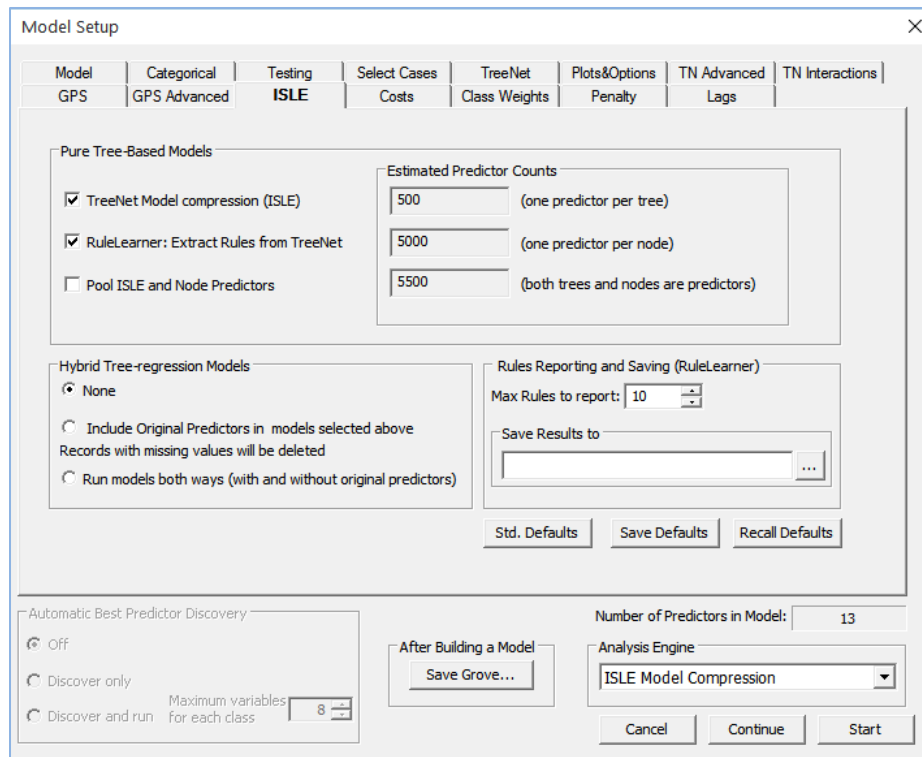
Minitab ⬐®

Switch to the Testing tab, and make sure that **Fraction of cases selected at random** for testing is set to 0.2.

We now need to set up the TreeNet component of the pipeline, as if you were setting up a stand-alone TreeNet model. Using the TreeNet tab, we set the following controls.

♦ Set the **Learn Rate** to 0.05.

♦ Set the **Number of trees** to 500.

♦ Set the **Criterion Determining Number of Trees Optimal for Logistic Model** to **ROC area**.

♦ Leave the remaining controls at their defaults (consult with the screenshot below).

✓ See the guide on TreeNet for detailed description of all of the available controls.

Next you may want to set up the GPS component of the pipeline using the **GPS** and **GPS Advanced** tabs. The default settings in those tabs are good enough for our purposes, you may however, review the GPS guide and later experiment with different number of steps, points, as well as elasticities and learn rates. Finally, switch to the ISLE tab to setup the actual pipeline "linkage" parameters, which could be as simple as specifying whether you want ISLE or RuleLearner compression or both (Pooled ISLE and Node Predictors).

Note that we have requested both ISLE and RuleLearner which will trigger two independent compression runs after the initial TreeNet model has been built: one run to try the ISLE compression and the other to try the RuleLearner compression. This gives you the convenience of not building the TreeNet model (the first stage of the pipeline) twice and thus saving the overall run time.

The control options on the **ISLE** tab under the Pure Tree-Based Models section are:

♦ **TreeNet Model Compression (ISLE)** – uses tree-based variables as intermediate inputs to the GPS compression stage.

♦ **RuleLearner: Extract Rules from TreeNet** – uses node-based variables as intermediate inputs to the GPS compression stage.

♦ **Pool ISLE and Node Predictors** – technically a hybrid approach, uses the combined set of both tree-based and node-based variables as intermediate inputs to the GPS compression stage. Note that the end result of this could be a compressed model that has both types of components.

♦ **Estimated Predictor Counts** – reports the total number of intermediate variables that will be created during the pipeline process. Note that the RuleLearner approach dramatically increase the memory and computational burden of the process, use it judiciously!

♦ **Hybrid Tree-regression Models** – allows to include the original predictors into the intermediate set of inputs to the GPS compression stage. This sometimes helps to compress the model even further in the presence of strong linear effects.

💣 Keep in mind that each pipeline mode, especially when it involves the node-based variables (RuleLearner) may take significant amount of run time. It is best to start experimenting with a single ISLE or RuleLearner model first and then proceed with more elaborate setups.

♦ **Max Rules to Report** – sets the default number of the top important rules that will be reported and saved.

♦ **[Save Results to]** – click this button if you want to save the top important rules into a file.

**Minitab** ▶®

Press the **[Start]** button to initiate the modeling pipeline.

Depending on your prior system settings, the following window may occur:



This indicates that some of the categorical predictors available in the dataset had under-represented categories that could not be assigned to both learn and test partitions. You may ignore this warning for now and press the **[OK]** button to proceed with the modeling run.

As the run develops, you will first see the TreeNet engine progress report followed by the GPS engine progress report for the ISLE run and then again for the RuleLearner run. This indicates that the two separate pipelines are being properly executed.

## Viewing Model Compression Results

A successfully finished compression run with more than one pipeline will produce the **Pipeline Results: Models** summary window shown below.

In our case, three different models are reported:

♦ **Original** – the original best TreeNet model before compression.

♦ **ISLE** – best model using ISLE (tree-based variables) compression.

♦ **RuleLearner** – best model using RuleLearner (node-based variables) approach.

The graph area displays the currently selected performance measures for each model.

Note that the original TreeNet model had 498 trees (coefficients) while the ISLE compressed model has only 247 trees (50.40% compression rate) while improving the test sample ROC performance.

Likewise, the original TreeNet model had 4980 nodes (6-node default regression trees used by the TreeNet engine actually have 10 nodes each, counting both terminal and internal nodes and excluding the root node).  The RuleLearner compressed model has 2789 selected nodes (55.77% compression rate) while improving the test sample ROC performance as well.

Thus, we have successfully accomplished the original goal of reducing the complexity of the TreeNet model by a significant amount while improving the test sample performance.

The following key controls are located on the results window:

♦ **Charts** – one can switch between the **[Performance]** and **[Coefficients]** display mode of the graph.

♦ **Chart Type** – one can switch between the **[Bar]** and **[Line]** style of the graph.

♦ **[Show Optimal Performance]** – press this button to highlight the best performing model using green background in the table.

♦ **Performance Measure:** - choose among a variety of different performance measures.  Once a new measure is selected, all displays will be updated.

♦ **Sample –** these buttons select among the learn and test data partitions.

♦ **[Score…]**, **[Translate…]**, **[Commands…]**, **[Save Grove…]** – these are standard action buttons available in any engine.

## Scoring Compressed Model

Press the **[Score…]** button to activate the **Score Data** control window.

Use the **Select model to score:** selection box to pick the model of interest (the Original, ISLE, or RuleLearner).

Press the **[Select…]** button in the **Select Specific Model** section if you want to pick a different model on the GPS path for the current compression approach.

The rest of the scoring controls have been described earlier in the manual.

Press the **[Score]** button to initiate the scoring process.

Minitab ▶®

## Translating Compressed Model

Press the **[Translate…]** button to activate the **Model Translation** control window.



Use the **Select Specific Model** selection box to pick the model of interest (the Original, ISLE, or RuleLearner).

Press the **[Select…]** button in the **Select Specific Model** section if you want to pick a different model on the GPS path for the current compression approach.

The rest of the translation controls have been described earlier in the manual.

Press the **[OK]** button to initiate the model translation process.

We used the above procedure to translate the 2-tree ISLE model into SAS. An example output code would have the following layout:

```
/* Tree 2 of 276 */

tnscore = -0.844412208; /* GPS ISLE Tree Coefficients */

<Block of code representing the tree logic to generate "response" update for the
current tree>

tnscore = tnscore + 0.126800358 * response;

/* Tree 5 of 276 */

<Block of code representing the tree logic to generate "response" update for the
current tree>
tnscore = tnscore + 1.25642133 * response;
return;
```

ISLE model translation only includes trees selected by the compression method (non-zero coefficients). Observe how the new updated ISLE coefficients are introduced directly into the model code.

Minitab ►®

✓ A similar code is produced when translating a RuleLearner model. However, because the RuleLearner model decomposition takes place at the node level, there is no simple way to reduce the model code by eliminating the redundant nodes. As the interim implementation, we now produce the full original TreeNet model code with the node-level adjustments augmented for all nodes selected by the RuleLearner process. Thus the code itself produces the correct RuleLearner scores but it also includes large parts that are unnecessary. This inconvenience has now been fully resolved in the latest release of the RuleLearner component.

## Viewing ISLE Model Compression Details

Double click on the ISLE row in the **Pipeline Results: Models** window. This will open the **ISLE** window reporting details of the ISLE compression run.

✓ This window automatically opens up instead of the **Pipeline Results: Models** window when ISLE is the only compression method requested in the **ISLE** tab of the **Model Setup** window.



This window shows detailed ISLE compression results indexed by model size in trees. The two test sample performance curves are aligned such that one can easily compare the performance of two models with identical size. Note that the 247-tree TreeNet model has test ROC 0.885 while the 247-tree ISLE compressed model has ROC 0.906. Even though both models have identical number of trees, they differ in the actual trees. The TreeNet model takes the first 247consecutive trees from the TreeNet sequence (all combined with the unit coefficients) while the ISLE model takes some specific subset of 247 trees from the original 500 trees available in the sequence, combined with ISLE specific coefficients.

✓ The actual subset of trees can be determined by looking into the ISLE Details or translating the model.

The top part of the window reports basic stats on the optimal TreeNet model versus optimal ISLE model in terms of ROC performance.

**Minitab** ▶®

✓ You can switch to alternative performance measures using the Performance Measure: selection box.

Press the **[TreeNet Details]** button to open the usual **TreeNet Output** window described in the TreeNet engine chapter of this manual.

Press the **[ISLE Details]** button to open the usual **GPS Results** window described in the GPS engine chapter of this manual.

✓ Note that the GPS component of the ISLE pipeline uses trees as input variables. Therefore, all parts of the GPS output windows referencing the predictors (Coefficients, etc.) will now report names like Tree 1, Tree 2, Tree 3, etc.

The remaining control items are similar to what we have already described earlier.

## Viewing RuleLearner Model Compression Details

Double click on the RuleLearner row in the **Pipeline Results: Models** window. This will open the **RuleLearner Results: RuleLearner** window reporting details of the RuleLearner compression run.

✓ This window automatically opens up instead of the **Pipeline Results: Models** window when RuleLearner is the only compression method requested in the **RuleLearner** tab of the **Model Setup** window.



This window shows detailed RuleLearner compression results indexed by model size in nodes. The two test sample performance curves are aligned such that one can easily compare the performance of two models with identical size.

Note that the 2789-node TreeNet model has test ROC 0.875 while the 2789-node RuleLearner compressed model has ROC 0.90584. Even though both models have identical number of nodes, they differ in the actual trees selected. The TreeNet model takes the first consecutive trees from the TreeNet sequence (all combined with the unit coefficients and the original node updates) while the RuleLearner model takes some specific subset of 2789 nodes from the original complete set available in the sequence. Furthermore, each selected node has a new associated update.

✓ The actual subset of nodes can be determined by looking into the RuleLearner Details or translating the model.

The top part of the window reports basic stats on the optimal TreeNet model versus optimal RuleLearner model in terms of ROC performance.

| Original TreeNet | | Best Rule Extraction | |
|---|---|---|---|
| Tree Size: | 6 | % Compression: | 55.77 |
| Trees Grown: | 500 | % Gain/Loss: | 0.009 |
| Rules Optimal: | 6,306 | Rules Extracted: | 2,789 |
| Performance: | 0.89735 | Performance: | 0.90584 |

✓ You can switch to alternative performance measures using the Performance Measure: selection box.

**Minitab** ▶®

Press the **[TreeNet Details]** button to open the usual **TreeNet Output** window described in the TreeNet engine chapter of this manual.

Press the **[RuleLearner Details]** button to open the usual **GPS Results** window described in the GPS engine chapter of this manual.

✓ Note that the GPS component of the RuleLearner pipeline uses nodes as input variables. Therefore, all parts of the GPS output windows referencing the predictors (Coefficients, etc.) will now report names like Tree 1 Node 1, Tree 1 TNode 1, Tree 2 Node 1, etc.

The remaining control items are similar to what we have already described earlier except for the **[Rules…]** button which opens up a RuleLearner specific **Explore** window.

## Exploring the Rules

Press the **[Rules…]** button in the **Pipeline Results: RuleLearner** window to open the **Explore RuleLearner** window shown below.



Recall that the RuleLearner compression method decomposes a TreeNet model into a collection of individual nodes (internal and terminal) and then harvests individual nodes using GPS. Each node represents a segment of the original data described by a collection of rules leading into this node. Hence, analyzing the rules associated with the top selected nodes might provide additional modeling insights.

The nodes are listed on the left part of the **Explore** window according to the order of their importance (as reported by the GPS component). The actual rules associated with the currently selected node (segment) are shown on the top right part while the basic stats about the node (segment) are shown on the bottom right part of the window.

## Important Rule Metrics

Note the use of the following terms generally accepted in the associative rules research community:

♦ **Support** – fraction of the sample (either learn or test) covered by the segment.

♦ **Lift** – segment mean of the target divided by the overall sample mean of the target.

✓ In the case of binary target, this is the same as segment fraction of 1s divided by the sample fraction of 1s.

♦ **Mean, Std. Dev.** – the usual statistical definitions of the mean and standard deviation in the segment of interest.

Click on the [Rules Table] button to open up even more detailed information about each node. The resulting table can be sorted by clicking on any of its columns. This allows quick and easy identification of the potential rules of interest. For example, you may focus on finding the segments having the highest test sample lift or segments having the highest support and so on.

For example, internal node 4 of tree 241 is at the top of the importance list, it covers almost 50% of the available data (very large support) and is described by a joint collection of rules involving POST_BIN, OCCUP_BLANK, and GENDER.



Finally, click on the **[Support vs. Lift]** button to see the graphical representation of the usual tradeoff involved. Note that the larger support segments tend to have smaller lift.

This concludes our discussion of the ISLE model compression approach now available in SPM.

## RuleLearner Example

At this point it will suffice to say that setting up a model compression run using ISLE from a user standpoint is very straightforward.  We are essentially creating a "pipeline", a sequence of two already introduced engines, first a TreeNet engine using the original predictors followed by a GPS engine that uses a decomposed TreeNet model as its own inputs.  The decomposition itself is done automatically by SPM based on the user-selected method, which could be either ISLE or RuleLearner (as well as additional trivial modifications introduced below).  In addition, it is assumed that each engine component of the pipeline (TreeNet and GPS) is configured using its own set of control parameters as usual.

The end result of the compression pipeline is the overall compression summary plus the usual output of the TreeNet engine (as if it were run by itself) as well as the standard GPS output of the follow-up compression phase.

The output displays are deliberately designed for ease of comparison between the original uncompressed model and the resulting compressed model.  The user can choose any desired model for viewing, scoring, and translation using the "point and click" approach thus balancing the degree of compression versus model size.

## Setting up RuleLearner

Again, we start with the GOODBAD.CSV data set which is rather on the small side but runs very quickly and illustrates our main points well.  The dataset comes from a credit default study and has 664 observations and 14 variables including a binary response variable named TARGET, with 1 indicating a bad credit standing account and 0 otherwise.  The remaining predictors provide a nice mix of continuous and categorical measures which represent a typical modeling environment.  The goal of the study is to build a predictive model for bad credit using the TreeNet approach and then possibly compress the model

**Minitab**

using either ISLE or RuleLearner approaches.

Use the **File/Open/Data File…** menu item to open the GOODBAD.CSV dataset, then click the **[Model…]** button in the resulting activity window to open the **Model Setup** window.

In the Model tab:

♦ Set **Analysis Engine** to **RuleLearner** – this activates the pipeline part of the SPM – the **RuleLearner** tab. Note that this mode activates the **TreeNet** tab and the **GPS** and **GPS Advanced** tabs thus giving you full control over other parts of the pipeline.

♦ Set **Sort:** to **File Order** – this sorts all variables according to the position in the dataset.

♦ Check TARGET variable as the **Target**.

♦ Check all remaining variables as predictors.

♦ Select **Target Type** to **Classification/Logistic Binary** to take advantage of the natural capability of both TreeNet and GPS engine to handle binary targets effectively.



Switch to the Testing tab, and make sure that **Fraction of cases selected at random for testing** is set to 0.2.

We now need to set up the TreeNet component of the pipeline, as if you were setting up a stand-alone TreeNet model.

♦ Set the **Learn Rate** to 0.05.

♦ Set the **Number of trees** to 500.

♦ Set **the Optimal Logistic Model Selection Criterion** to **ROC area**.

♦ Leave the remaining controls at their defaults (consult with the screenshot below).

✓ See the TreeNet chapter for detailed description of all of the available controls.

**Minitab ›**

Next you may want to set up the GPS component of the pipeline using the **GPS** and **GPS Advanced** tabs. The default settings in those tabs are good enough for our purposes, you may however, review the GPS part of the manual and later experiment with different number of steps, points, as well as elasticities and learn rates. Finally, switch to the ISLE tab to setup the actual pipeline "linkage" parameters, which could be as simple as specifying whether you want ISLE or RuleLearner compression or both (Pooled ISLE and Node Predictors).

Note that we have requested only the RuleLearner rule extraction feature which will trigger only one independent compression runs after the initial TreeNet model has been built: you could always try one run for the ISLE compression and the other to try the RuleLearner compression like what we showed earlier for the ISLE algorithm which would give you the convenience of not building the TreeNet model (the first stage of the pipeline) twice and thus saving the overall run time.

The control options on the **RuleLearner** tab under the Pure Tree-Based Models section are:

♦ **TreeNet Model Compression (ISLE)** – uses tree-based variables as intermediate inputs to the GPS compression stage.

♦ **RuleLearner: Extract Rules from TreeNet** – uses node-based variables as intermediate inputs to the GPS compression stage.

♦ **Pool ISLE and Node Predictors** – technically a hybrid approach, uses the combined set of both tree-based and node-based variables as intermediate inputs to the GPS compression stage.  Note that the end result of this could be a compressed model that has both types of components.

♦ **Estimated Predictor Counts** – reports the total number of intermediate variables that will be created during the pipeline process.  Note that the RuleLearner approach dramatically increase the memory and computational burden of the process, use it judiciously!

♦ **Hybrid Tree-regression Models** – allows to include the original predictors into the intermediate set of inputs to the GPS compression stage.  This sometimes helps to compress the model even further in the presence of strong linear effects.

♦ Keep in mind that each pipeline mode, especially when it involves the node-based variables (RuleLearner) may take significant amount of run time.  It is best to start experimenting with a single ISLE or RuleLearner model first and then proceed with more elaborate setups.

♦ **Max Rules to Report** – sets the default number of the top important rules that will be reported and saved.

♦ **[Save Results to]** – click this button if you want to save the top important rules into a file.

**Minitab** ▷®

Press the **[Start]** button to initiate the modeling pipeline.

Depending on your prior system settings, the following window may occur:



This indicates that some of the categorical predictors available in the dataset had under-represented categories that could not be assigned to both learn and test partitions.  You may ignore this warming for now and press the **[OK]** button to proceed with the modeling run.

As the run develops, you will first see the TreeNet engine progress report followed by the GPS engine progress report for the ISLE run and then again for the RuleLearner run.  This indicates that the two separate pipelines are being properly executed.

## Viewing RuleLearner Results

A successfully finished compression run with one pipeline will produce the **Pipeline Results: TreeNet Rule Extraction with RuleLearner** Results**:** Model summary window shown below:





In our case, two different models are reported:

- ♦ **Original** – the original best TreeNet model before compression.
- ♦ **RuleLearner** – best model using RuleLearner (node-based variables) approach.

## Translating RuleLearner

Press the **[Translate…]** button to activate the **RuleLearner Model Translation** control window.



Use the **Select Specific Model** selection box to pick the model of interest (the Original, ISLE, or RuleLearner).

Press the **[Select…]** button in the **Select Specific Model** section if you want to pick a different model on the GPS path for the current compression approach.

The rest of the translation controls have been described earlier in the manual.

Press the **[OK]** button to initiate the model translation process.

## Viewing RuleLearner Details



This window shows detailed RuleLearner compression results indexed by model size in nodes. The two test sample performance curves are aligned such that one can easily compare the performance of two models with identical size.

Note that the 2789-node TreeNet model has test ROC 0.875 while the 2789-node RuleLearner compressed model has ROC 0.90584. Even though both models have identical number of nodes, they differ in the actual trees selected. The TreeNet model takes the first consecutive trees from the TreeNet sequence (all combined with the unit coefficients and the original node updates) while the RuleLearner model takes some specific subset of 2789 nodes from the original complete set available in the sequence. Furthermore, each selected node has a new associated update.

✓ The actual subset of nodes can be determined by looking into the RuleLearner Details or translating the model.

The top part of the window reports basic stats on the optimal TreeNet model versus optimal RuleLearner model in terms of ROC performance.



✓ You can switch to alternative performance measures using the Performance Measure: selection box.

Press the **[TreeNet Details]** button to open the usual **TreeNet Output** window described in the TreeNet

engine chapter of this manual.

Press the **[RuleLearner Details]** button to open the usual **GPS Results** window described in the GPS engine chapter of this manual.

✓   Note that the GPS component of the RuleLearner pipeline uses nodes as input variables.  Therefore, all parts of the GPS output windows referencing the predictors (Coefficients, etc.) will now report names like Tree 1 Node 1, Tree 1 TNode 1, Tree 2 Node 1, etc.

The remaining control items are similar to what we have already described earlier except for the **[Rules…]** button which opens up a RuleLearner specific **Explore** window.

## Exploring the Rules

Press the **[Rules…]** button in the **Pipeline Results: RuleLearner** window to open the **Explore RuleLearner** window shown below.



Recall that the RuleLearner compression method decomposes a TreeNet model into a collection of individual nodes (internal and terminal) and then harvests individual nodes using GPS.  Each node represents a segment of the original data described by a collection of rules leading into this node.  Hence, analizing the rules asociated with the top selected nodes might provide additional modeling insights.

The nodes are listed on the left part of the **Explore** window according to the order of their importance (as reported by the GPS component).   The actual rules associated with the currently selected node (segment) are shown on the top right part while the basic stats about the node (segment) are shown on the bottom right part of the window.

Note the use of the following terms generally accepted in the associative rules research community:
♦   **Support** – fraction of the sample (either learn or test) covered by the segment.
♦   **Lift** – segment mean of the target divided by the overall sample mean of the target.
✓   In the case of binary target, this is the same as segment fraction of 1s divided by the sample fraction of 1s.

**Minitab** ⊳®

- **Mean-** the average of the target variable for records that satisfy the rule of interest

- **Target Std. Dev.** – the usual standard deviation of the target variable for records that satisfy the rule of interest

- **Rule Std. Dev.:** also referred to as the "scale" of a rule and is given by $\sqrt{support * (1 - support)}$ where support is the support for the rule of interest

- **Importance = Rule Std. Dev. * |Coefficient|**

Click on the [Rules Table] button to open up even more detailed information about each node. The resulting table can be sorted by clicking on the any of its columns. This allows quick and easy identification of the potential rules of interest. For example, you may focus on finding the segments having the highest test sample lift or segments having the highest support and so on.

For example, internal node 4 of tree 241 is at the top of the importance list, it covers more than 50% of the available data (very large support) and is described by a joint collection of rules involving POST_BIN, OCCUP_BLANK, and GENDER.



Finally, click on the **[Support vs. Lift]** button to see the graphical representation of the usual tradeoff involved. Note that the larger support segments tend to have smaller lift.

This concludes our discussion of the new model compression approaches now available in SPM.