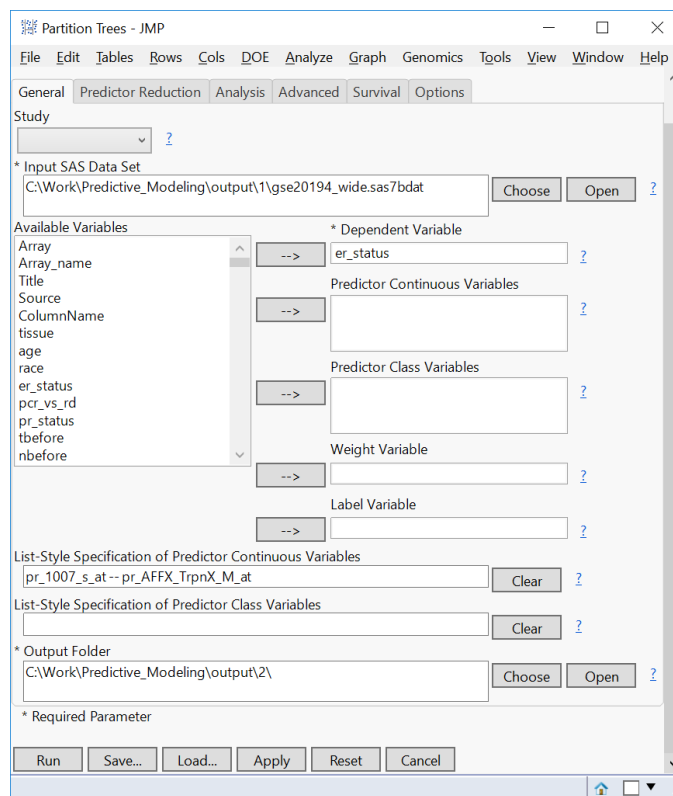


Whole Model Fit

Through previous analyses surrounding the GSE20194 data set downloaded from GEO, we have created models and tested them on **Test Data Sets** to evaluate their fit. In the *Cross Validation Model Comparison Step-Guide*, we created and evaluated a **Partition Trees** and **Radial Basis Machine** model. We then determined our data had an adequate number of samples in the *Learning Curves Step-Guide*. Cross validation and test set model comparison have shown that the **Partition Trees** model is the best fitting model. The last step in the predictive modeling workflow is to fit the whole model to the **Training Data Set** and to assess its effectiveness.

1. From the **Genomics Starter** load the **Partition Tree** model settings. Select **File > Load Life Sciences Setting** from the toolbar. Then find the settings folder created in the *Cross Validation Model Comparison Step-Guide*. The tabs should appear as below.
 - The input data is the ***gse20194_wide.sas7bdat*** set created in the *Importing and Transposing a Data Set Step-Guide*.
 - Note that these settings can be created again by selecting **Predictive Modeling > Main Methods > Partition Trees** from the **Genomics Starter** menu.
 - For further information on the model settings see the *Cross Validation Model Comparison Step-Guide*.



General tab.

Radial Basis Machine - JMP

File Edit Tables Rows Cols DOE Analyze Graph Genomics Tools View Window Help

Process Description

General Lock-In Predictor Reduction Analysis Genetic Algorithm Survival Options

Filter

Transform

Standardize

Test

☒ Use statistical testing to filter predictors ?

Maximum Number of Filtered Predictors
20 ?

Statistical Testing Method for Continuous Predictors
Unequal Variance T-Test ?

☐ Nominalize Continuous Dependent Variables ?

Asymmetric Loss Fitting Proportion [0,1]
? ?

Quantile Level [0,1]
? ?

Multiple Testing Method
FDR ?

$-\log_{10}(\text{p-Value})$ Cutoff
1.3 ?

Absolute Mean Difference Cutoff for Continuous Predictors
0.5 ?

Absolute Proportion Difference Cutoff for Class Predictors
0 ?

Cluster

☒ Use K-Means clustering to reduce predictors ?

Usage of K-Means Clusters ?

☒ Select one representative predictor from each cluster
☐ Use the cluster means themselves as predictors

Maximum Number of K-Means Clusters / Predictors
300 ?

Correlation Radius for Clustering [0,1]
0.8 ?

Forest

Regress

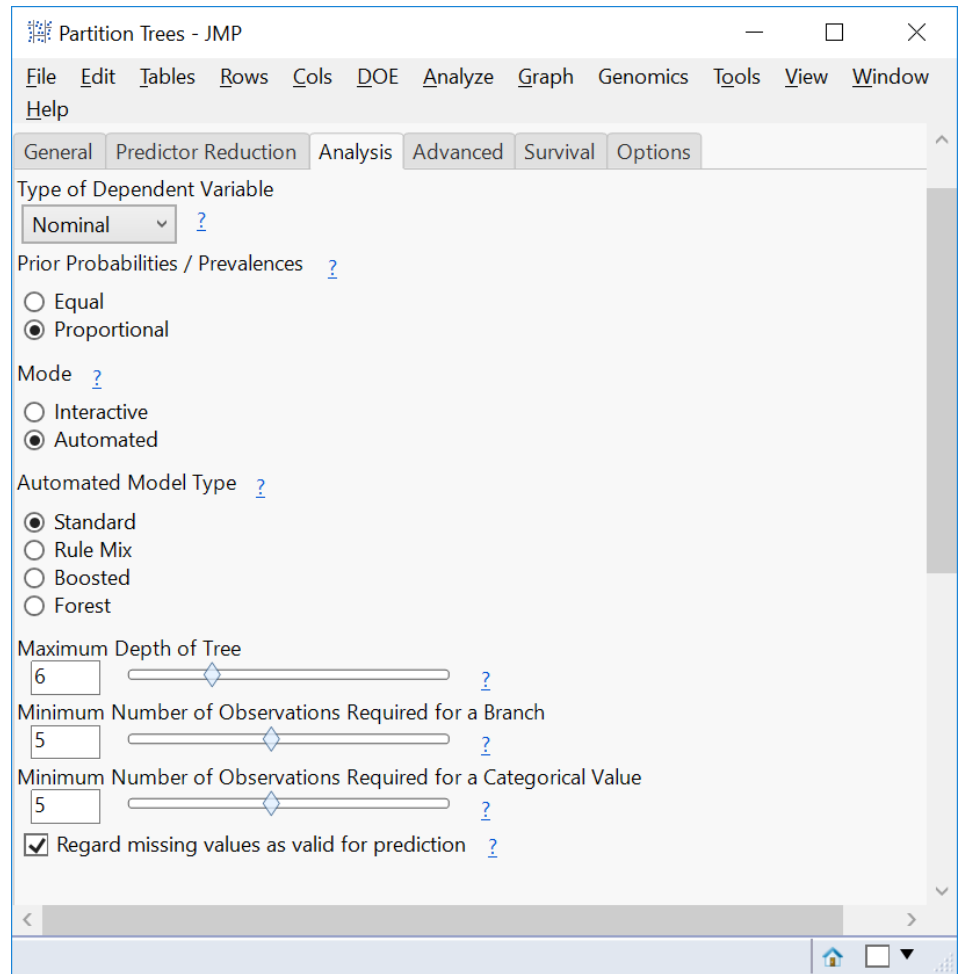
Optimize

☐ Standardize Predictors Row-Wise ?

* Required Parameter

Run Save... Load... Apply Reset Cancel

Predictor Reduction tab.



Analysis tab.

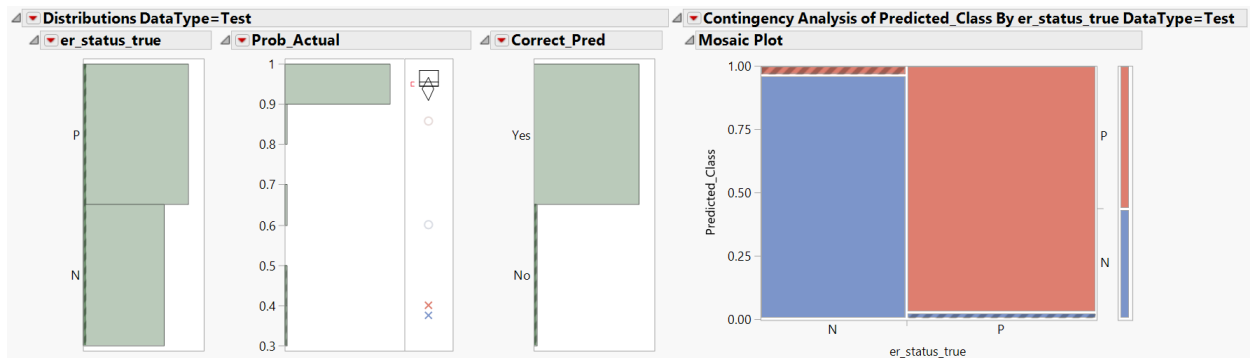
2. On the **Options** tab, select one of the **Test Data Sets** created in the *Subsetting a Test Data Set for Test Set Model Comparison*. We will use ***gse20194_wide_test_2.sas7bdat***.
 - An optional **Validation Data Set** can be added on the **Options** tab if desired.
3. Click **Run** to begin the analysis.

Results

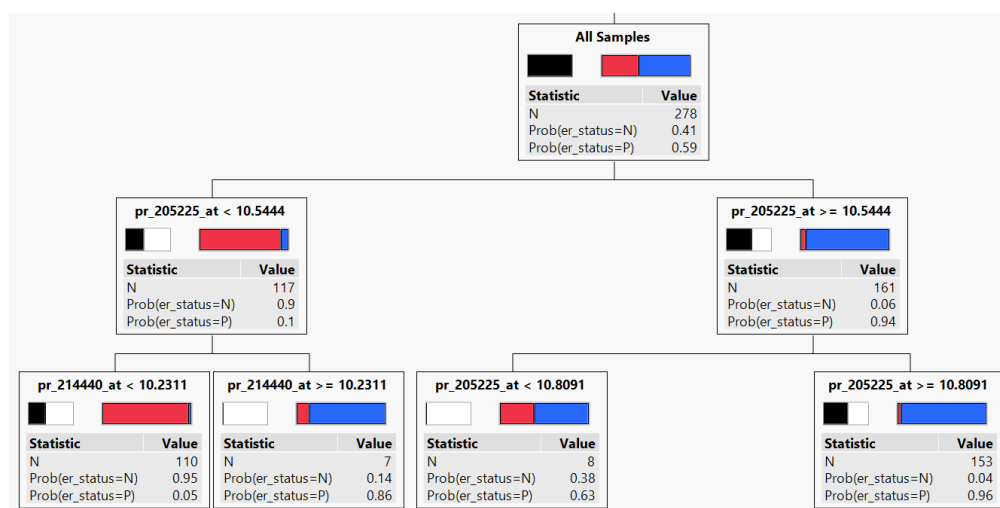
4. When the analysis finishes running, the results will appear in a new window. At the very top of the **Model Results** tab are summary results for the model from the test set and training set. This section also contains the **Final Selected Variables** which shows us that there are 3 final probes included in the model.

Model Results	Tree Diagram	Subtree Assessment	Variable Importance	ROC	SAS Output
Predictor Reduction Settings: Stat Test = Unequal Variance T-Test, Multiple Testing Method = FDR, -log10(p-value) Cutoff = 1.3, Mean Difference Cutoff = 0.5, K-Means = 300 Analysis Settings: Tree Model Type = Standard, Priors = Proportional Final Selected Variables: pr_205225_at, pr_210735_s_at, pr_214440_at Test Set Criteria: Root Mean Square Error = 0.1288, Mean Absolute Error = 0.0628, Area Under ROC Curve = 0.9984, Accuracy = 0.9677, Accuracy_N = 0.9630, Accuracy_P = 0.9714 Training Set Criteria: Root Mean Square Error = 0.2123, Mean Absolute Error = 0.0901, Area Under ROC Curve = 0.9672, Accuracy = 0.9460, Accuracy_N = 0.9386, Accuracy_P = 0.9512					

- Beneath the summary of fit statistics there are distribution plots for the test set and then the training set. Beneath each set of plots are summary statistics.



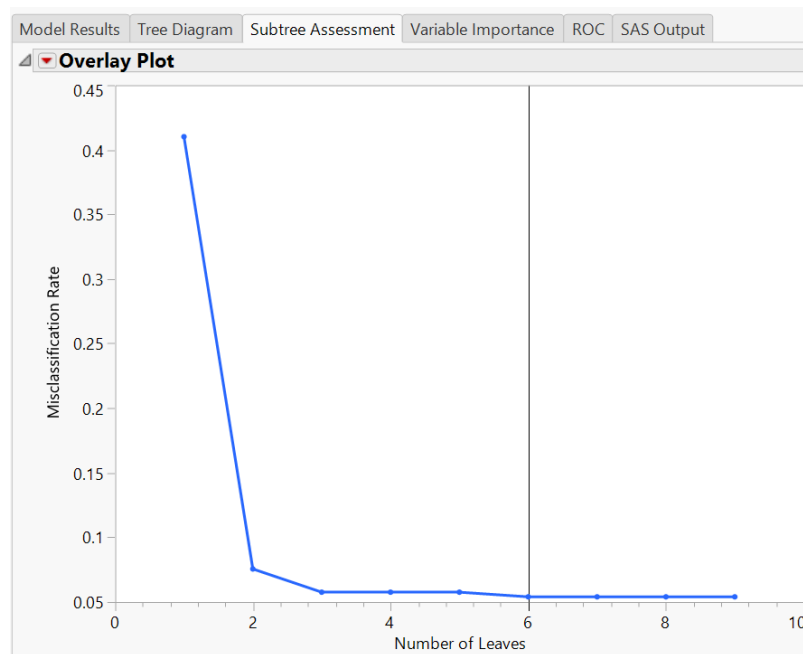
- The first plot shows the distribution of the dependent variable. Recall that we designated estrogen receptor status as the dependent variable when creating the model. This shows the distribution of both positive (P) and negative (N) status.
 - The middle plots both show statistics regarding the accuracy of the model. The **Prob_Actual** plot is the distribution of the posterior probability. Values below 0.5 here are a misclassification. The **Correct_Pred** is the simple distribution of correct vs. incorrect calls that the model made. You can see below the plot under the **Frequencies** tab, that the model was 96.7% correct for its predictions in the test set.
 - The **Mosaic Plot** sorts all correct and incorrect predictions into cells for both classifications of er_status.
- The **Tree Diagram** tab gives statistics for each of the three probes designated by the model including the number of data points for each level of each probe and the distribution of er_status for each level.



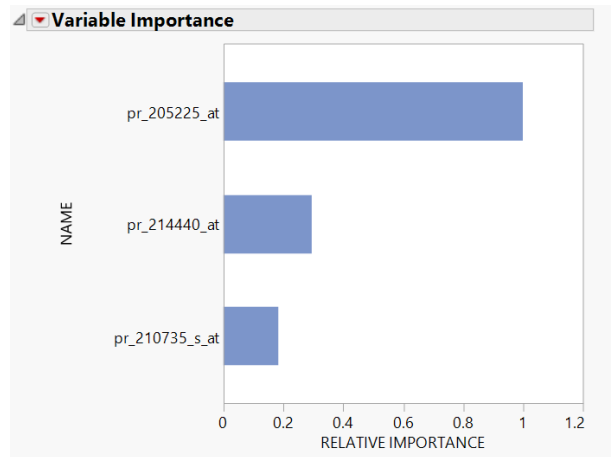
- The top of the **Tree Diagram** shows the probability of er_status for both N and P for the entire data set (278 samples). As the diagram moves further

down, each branch is the statistically optimal way to segment the data for accurate prediction.

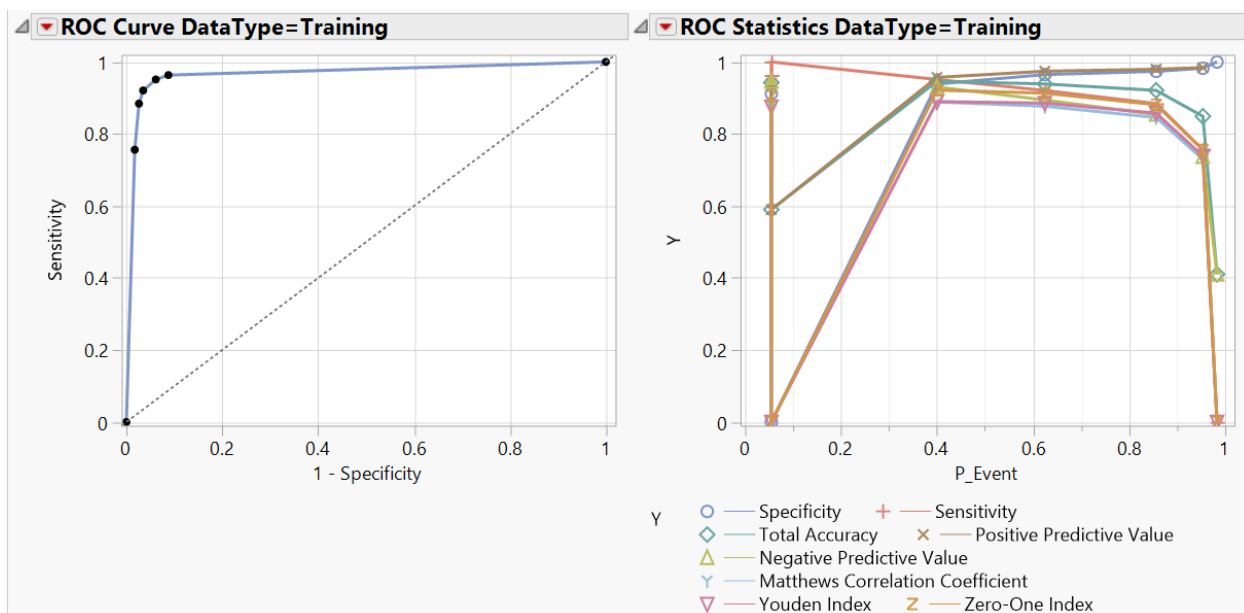
- The first level of branches has the probe “pr_205225_at”. This shows that the 117 samples with a value less than 10.54 for the probe have 0.9 probability of being negative for er_status, and for values greater than or equal to 10.54, samples have a 0.94 probability of being positive.
 - Each subsequent branch of the tree further breaks up the samples from the branch above it.
7. The **Subtree Assessment** tab plots the rate of misclassification vs the number of leaves in the tree. The vertical line is drawn at the actual number of fitted leaves in the model.
- This model has 6 leaves, or segments of the tree that are not further segmented. The model classifies each sample into one of the 6 leaves before making a prediction.



8. The **Variable Importance** tab gives the relative importance for each of the three probes chosen as branches in the tree.
- The probe “pr_205225_at” has the most influence on this model.



9. The **ROC** tab plots the plots the correctly identified proportion of true positives vs correctly identified true negatives for both the test and training set.
- The area under the curve on the first plot is a measure of the sorting efficiency of this model. A value of 1 here would represent perfect sorting.
 - The second plot show several performance statistics for different posterior probability cutoffs on the x-axis.



Summary

This document served as a guide to fitting a **Partition Trees** model to a set of markers to predict estrogen receptor status in human breast cancer data. The model was chosen after comparison to another model through cross-validation (See *Cross Validation Model Comparison Step-Guide*). The **Partition Trees** model makes predictions by choosing significant markers and then segmenting the data based on

values for those markers. This specific model split the data into 6 groups before making a prediction, and was over 90% accurate.