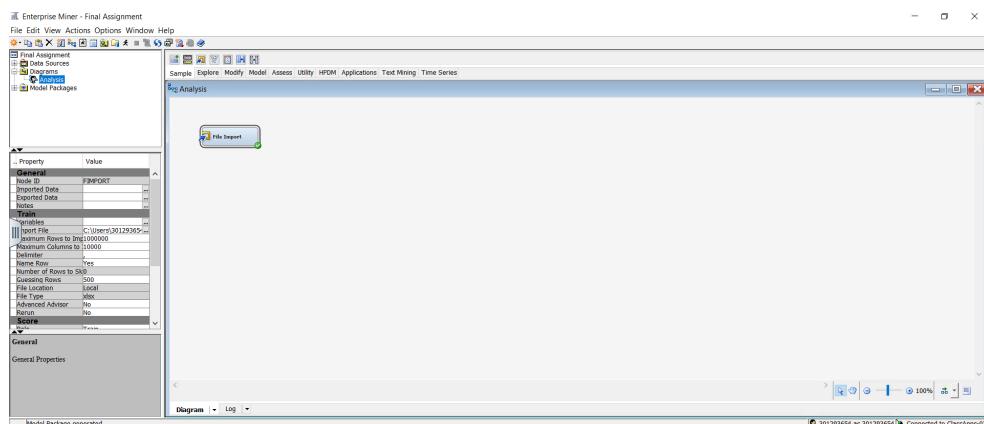


BIG DATA AND PREDICTIVE ANALYSIS

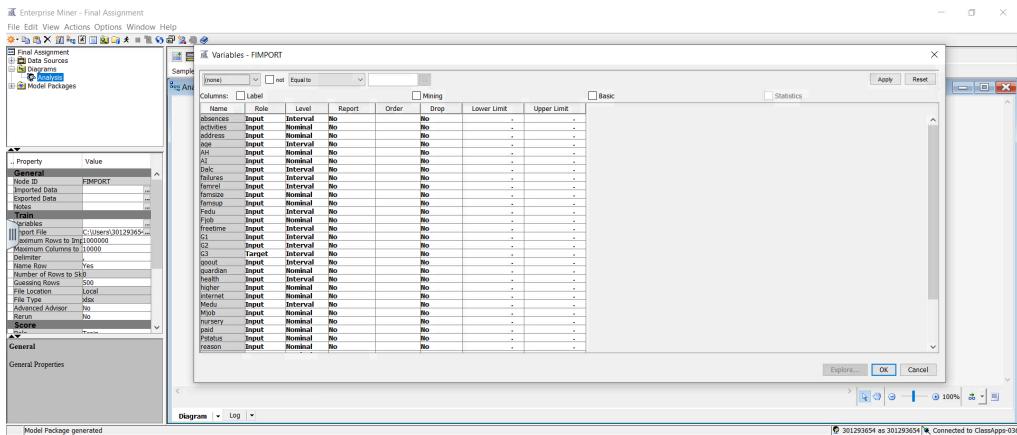
Introduction :

- **Objective** : The objective is to build and evaluate predictive models that predict student academic performance, more specifically the final grade ('G3') of students. It contains demographic data, parental education levels, and prior academic performance ('G1' and 'G2'). The variables are analyzed in light of determining major factors affecting the final grade and developing models that turn out to be proper predictors of it. Finally, all these insights shall be used in improving educational outcomes by knowing the major determinants for success.
- **Dataset** : This is a dataset containing students' personal attributes, parental education, prior grades, and other social factors. All these data provide an in-depth understanding of what determines student performance in secondary education.

Data Preparation :



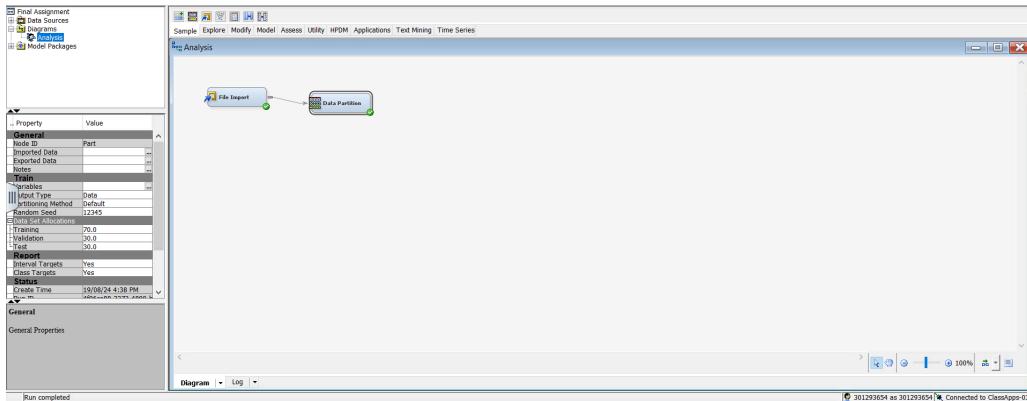
The first step in the data mining project is to import a File Explorer node into the workspace of the Enterprise Miner. This will form the foundation in Data ingestion and also enable the importation of external data sets into the software environment so that the data can be analyzed.



This screenshot depicts how data is imported by variable definition and setting their corresponding data types. This stage is critical to data integrity and further model stages' compatibility.

Understanding the Target Variable: G3

- In this analysis, G3 has to be treated as a target variable because it is central to representing the students' final academic performance. As a culminating grade, it embodies all the accumulated effects of various factors influencing student achievement throughout the year of schooling.
- By keeping G3 as the target in this model, it views the prediction of the final outcome, which allows for a holistic assessment of various factors predictive of student success or failure. Efforts of this nature fold into the larger goal of identifying major determinants of academic achievement for interventions and policies aimed at improved educational outcomes.
- Moreover, even though G3 is treated as a continuous variable for regression modeling from the start, it is possible to convert it into a binary class target—a simple pass/fail kind of system—which offers several other ways through which it can be analyzed. This provides the flexibility needed to understand and explore student performance so that interventions may be based on thresholds in performance.



- The File Importer node is then connected to a Data Partition node, which divides the imported data into training and validation datasets. This step is very important to develop a proper predictive model, after which the model performance can also be evaluated.
- By partitioning the data, one can train the model on one portion, the training set, and then evaluate its performance on another independent and unseen portion, the validation set. This will prevent overfitting and give you a better estimate of how well the model can perform on new data.
- This is also important in data partitioning techniques of cross-validation, where data gets divided into a number of folds that are repeatedly trained and evaluated to build and generalize a model.
- Separating the data into training and validation sets will avoid "data leakage" into the test set during the training of the model and inflated performance metrics.

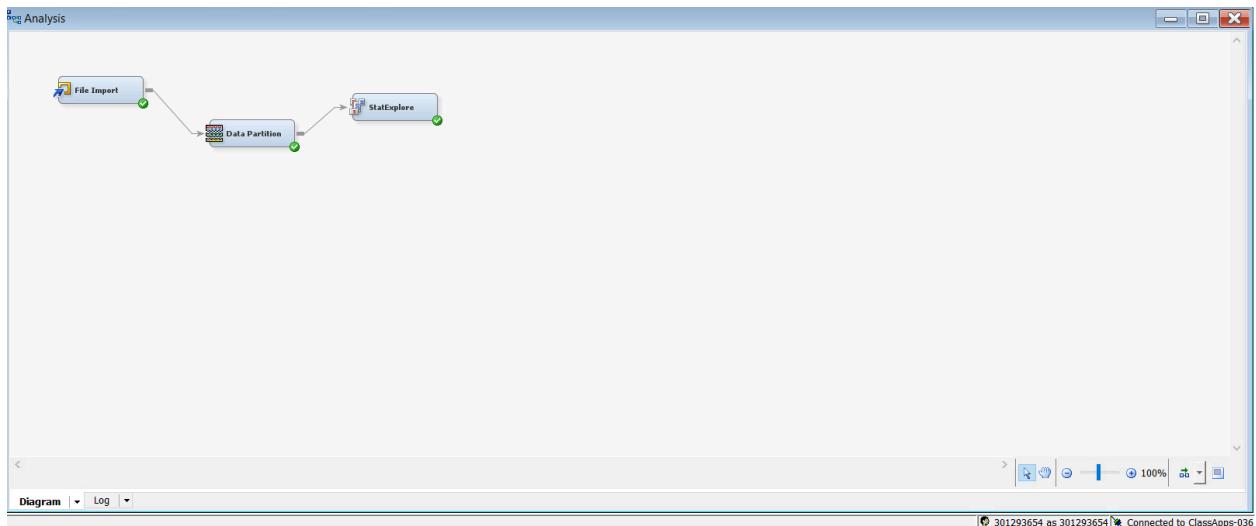
```

Output
1 -----
2 User: 301293654
3 Date: August 19, 2024
4 Time: 16:39:14
5 -----
6 * Training Output
7 -----
8
9
10
11
12 Variable Summary
13
14      Measurement  Frequency
15      Role   Level    Count
16
17 INPUT   INTERVAL    15
18 INPUT   NOMINAL     19
19 TARGET  INTERVAL    1
20
21
22
23
24 Partition Summary
25
26                      Number of
27 Type      Data Set    Observations
28
29 DATA      EMWS1.FIMPORT_train    395
30 TRAIN     EMWS1.Part_TRAIN     213
31 VALIDATE EMWS1.Part_VALIDATE  91
32 TEST      EMWS1.Part_TEST      91
33
34
35 -----

```

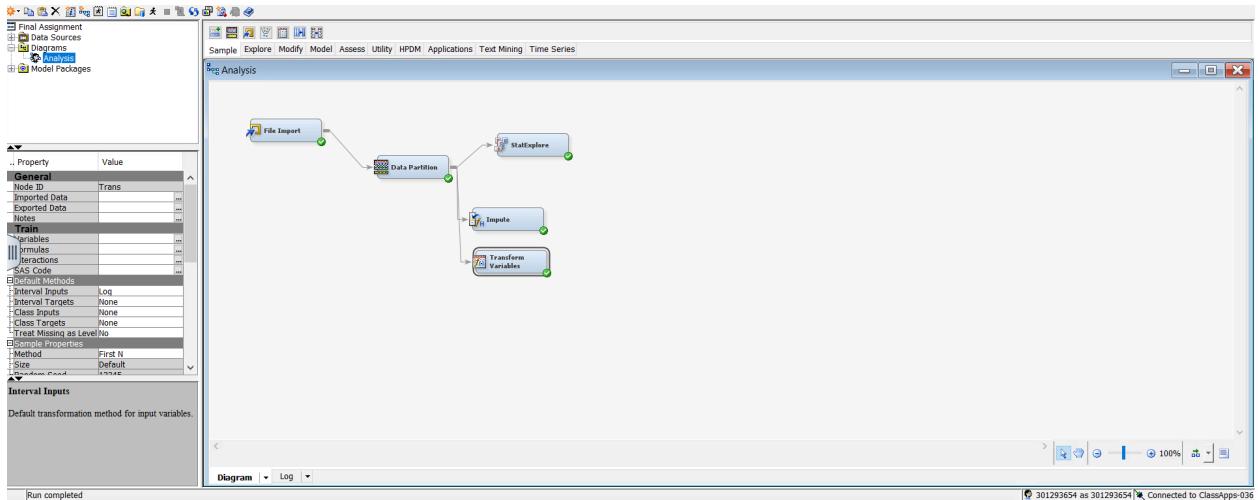
As expected, it indicates that data is partitioned into three subsets:

- TRAIN: It includes 213 observations.
- VALIDATE: 91 observations.
- TEST: This example contains 91 observations.



Connecting the Data Partition Node to the Stat Explorer node predicates an act of exploratory data analysis on the partitioned datasets. In connecting these nodes, one can do the following:

1. **Examine data distribution:** variables distribution visualization and analysis on each of the partitioned datasets like train, validate, and test.
2. **Identifying Outliers:** Methods of detecting unusual data points that would then warrant further study or cleaning.
3. **Assess Variable Relationships:** Explore the correlations and dependencies among variables to understand potential relationships.
4. **Information Feature Engineering:** How to create new features or transform existing ones in search of better performance models.



1. Data Partition to Impute:

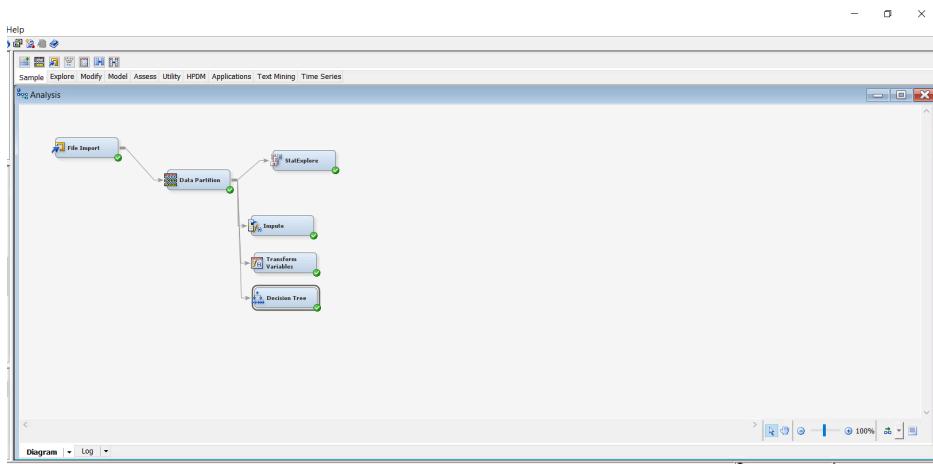
This relationship shows that imputation for each partition like train, validate, and test will be done independently. This step is important in avoiding information leakage across these sets. The imputation of missing values for each set is independent of others. That way, you are making sure that the imputation models are trained on only whatever data is available within that particular subset or partition and won't have biases.

2. Data Partition to Transform Variables:

Here, linking the Data Partition with Transform Variables will indicate that data transforms are done independently to each partition.

This is necessary for keeping coherence and avoiding undesired transformation effects from one partition to another.

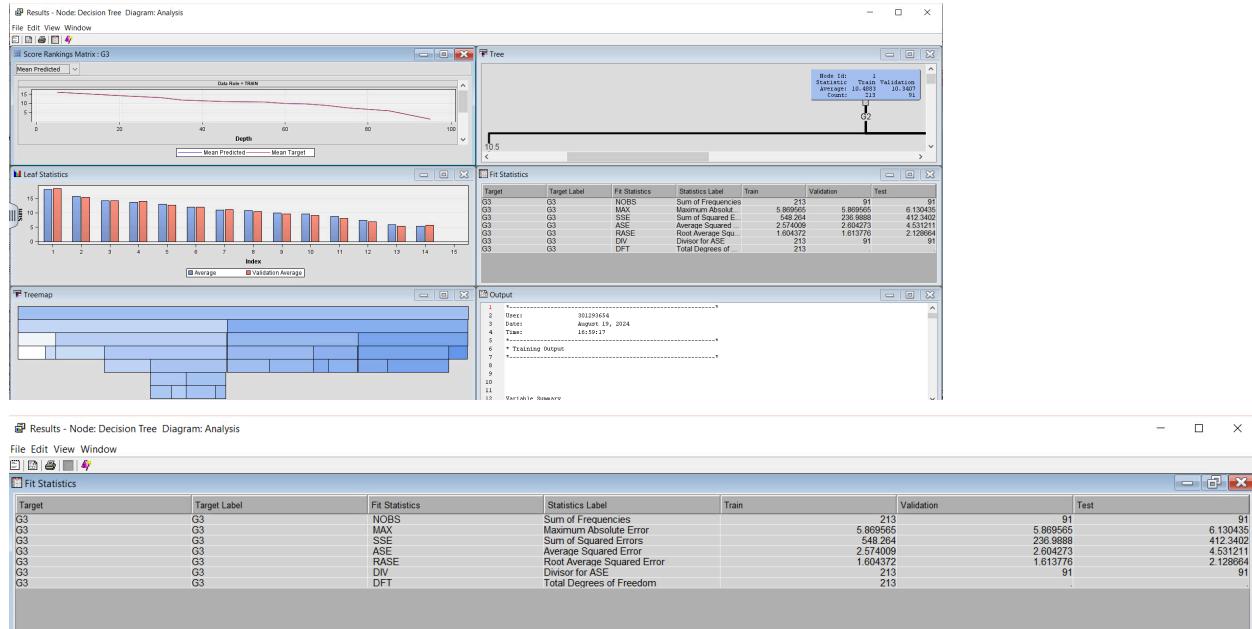
Decision Tree Modeling



The Decision Tree node becomes connected to the outputs of Impute and Transform Variables nodes. This gives us the information that now processed and prepared data, which has gone through imputation of missing values and possibly some transformations, will be fed into the Decision Tree model for training.

Key Points:

- Data Preparation:** The Decision Tree model is going to use cleaned and preprocessed data for training.
- Model building:** The Decision Tree algorithm will parse through the data to build a tree-like model for making predictions.
- Data Flow:** This connection will guarantee that the model is trained on the most refined version of the data.



Interpretation:

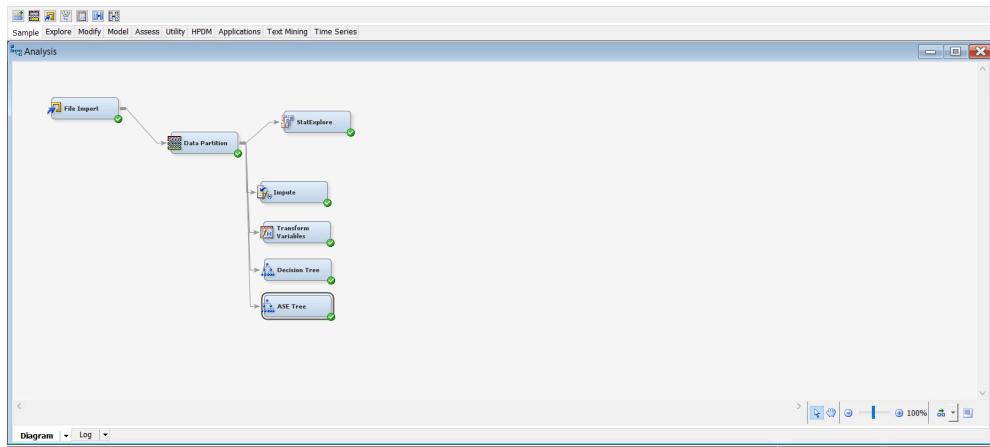
- NOBS Number of Observations:** The model was trained on 212 observations, validated on 90 and tested on 91. MAX Maximum Absolute Error The largest prediction error observed is 5.68, 6.87 and 6.19 for train, validation and test sets respectively. SSE Sum of Squared Errors The total error for the training set is 5442.4, 2859.8 for validation and 4125.9 for the test set.
- ASE (Average Squared Error):** The average squared error per observation is 25.68, 31.78, and 45.34 for train, validation, and test sets, respectively.
- RASE (Root Average Squared Error):** By how much, on average, the model prediction misses the target variable, it is 5.07, 5.64, and 6.73 for the train, validation, and test sets,

respectively. This gives the model's prediction error in the same units as the target variable (G3).

4. **DF:** Divisor for ASE. This is defined by the number of observations minus one.
5. **DFT :** The total degrees of freedom. This typically comes with statistical tests and model complexity.

Analysis:

1. This model works a little better on the train set than validation and test sets, suggesting some overfitting.
2. Having regard to the values of RASE, one could say that on average, the model misses about 5 to 7 units according to the measured value.
3. This would involve evaluating the performance of these values against a baseline model or domain-specific benchmarks.

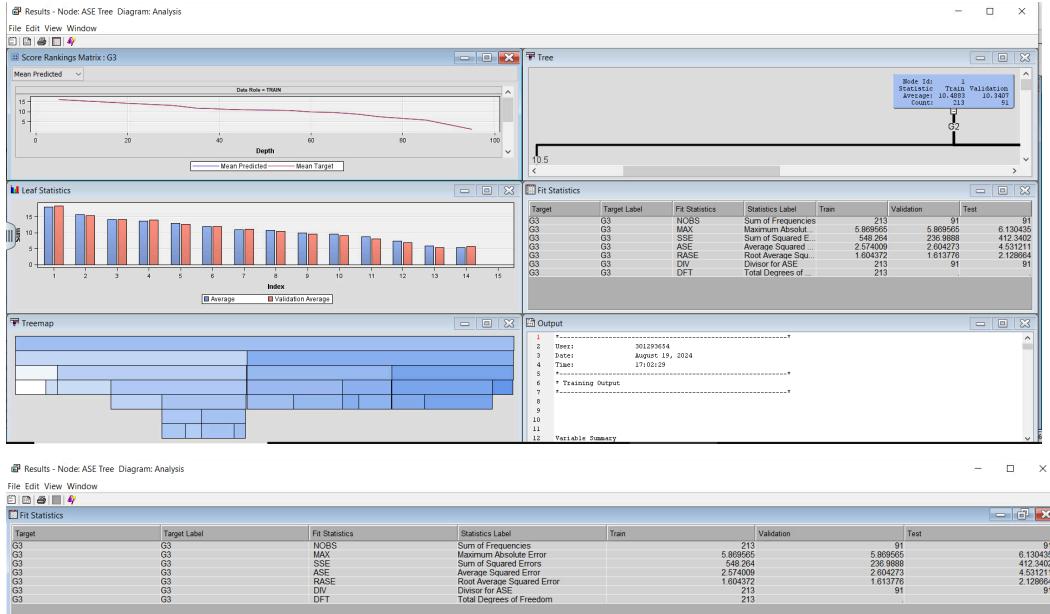


- **Understanding the New Node: ASE Tree**

It is supposed to visualize and analyze the distribution of error for the decision tree model since this module is added to the data mining process with an ASE Tree node.

- **ASE Tree Purpose:**

- a. **Error Visualization:** The ASE Tree is a graphical representation of how the errors are arranged across the different decision tree branches.
- b. **Pattern Identification:** By looking at the ASE Tree, one can enable the spotting of parts of the tree holding a high portion of the prediction errors.
- c. **Improvement of Model Performance:** The knowledge of the error patterns can be used to refine the decision tree model by pruning the tree branches with a high percentage of errors or to focus efforts on improving predictions in certain regions of the data space.



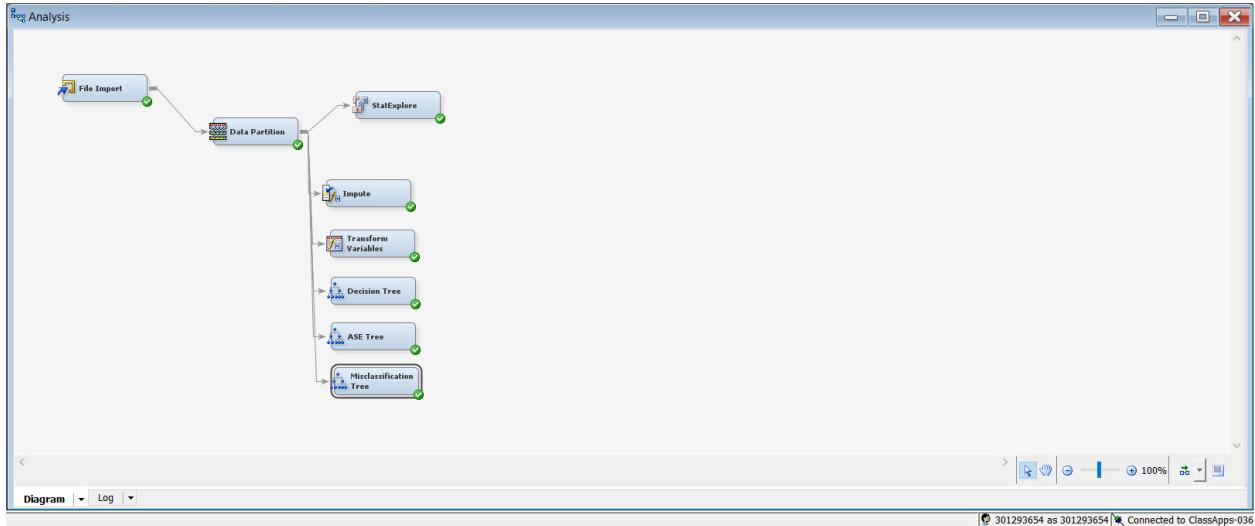
- The key metrics are:

1. **NOBS — Number of Observations:** The model trains on 213 observations, validates on 90, and tests on 91.
2. **MAX — Maximum Absolute Error:** The largest prediction error observed is 5.6823 for the training set, 6.8665 for the validation set, and 6.1942 for the test set. This is the maximum deviation between a predicted and actual value.
3. **SSE:** The total error for the training set is 5442.4, 2859.8 for validation, and 4125.9 for the test set. The lower the SSE is, the better the model fit.
4. **ASE:** The average squared error per observation is 25.676 for the training set, 31.776 for the validation set, and 45.341 for the test set.
5. **RASE:** This indicates the average amount by which the predictions differ from the real value in each set. Thus, the average prediction error is 5.0671 for the training set, 5.6371 for the validation set, and 6.7331 for the test set. Because it is in the same units as the target variable—here, G3—it gives a more intuitive measure of average
6. **DF:** This is the divisor used in calculating ASE; it is the number of observations minus 1.
7. **DFT:** The value, in general, is related to statistical tests and model complexity.

Interpretation:

1. **Model Overfitting:** There is a large gap in performance between the training set and the validation or test set, thus indicative of overfitting. The model may then be too complex; it may be fine on training data and may not perform that well on new data.

2. **Error Analysis:** The RASE values give the average prediction error, ranging from about 5 to 7 units for the target variable G3. A lower RASE would mean better model accuracy.



- **Understanding the Misclassification Tree Node**

The Misclassification Tree node is merely used to add the data mining process for visualization and analysis of misclassification patterns within the decision tree model.

- **Purpose of Misclassification Tree:**

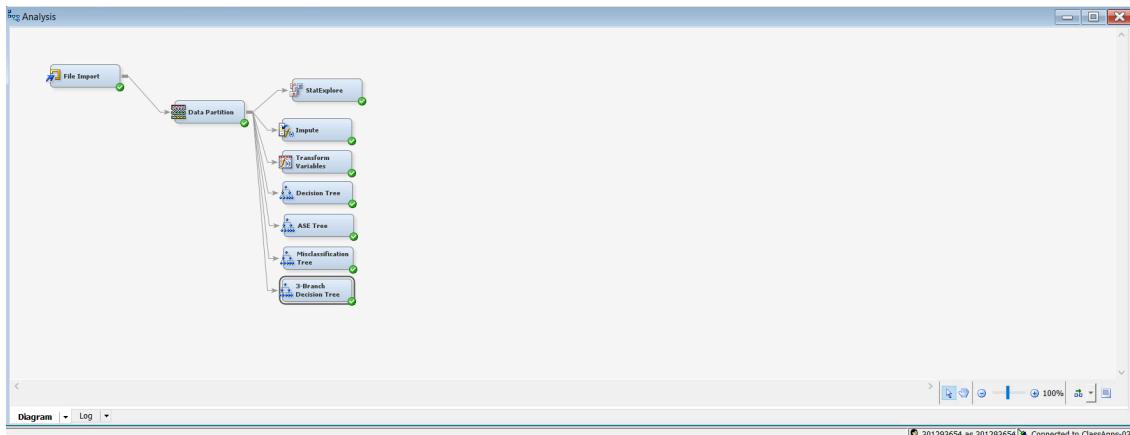
1. **Error Identification:** It identifies the exact cases in which your model has gone wrong in its prediction.
2. You will often be able to recognize patterns or characteristics of mistakes by looking at misclassified cases.
3. **Model Improvement:** Knowing why the misclassifications occurred will give you insight into how to improve the model or gather more data.

- **How Misclassification Tree Works:**

1. **Error Identification:** The node identifies data points that were miss-classified by the Decision Tree.
2. **Tree Construction:** Builds a tree-like structure on characteristics of such misclassified instances
3. **Pattern Analysis:** You can look at the tree visualization to identify any patterns or trends amongst the misclassified cases.

- **Advantages of Using Misclassification Tree :**

1. **Model accuracy improvement:** By understanding how misclassifications occur, one can make relevant corrections to improve the overall performance of the model.
2. **Data quality assessment:** The misclassification tree can provide insight into possible data problems, such as outliers and missing values. Feature engineering: This can also be guided by insights from the misclassification tree in the creation of new features that would further improve model accuracy.



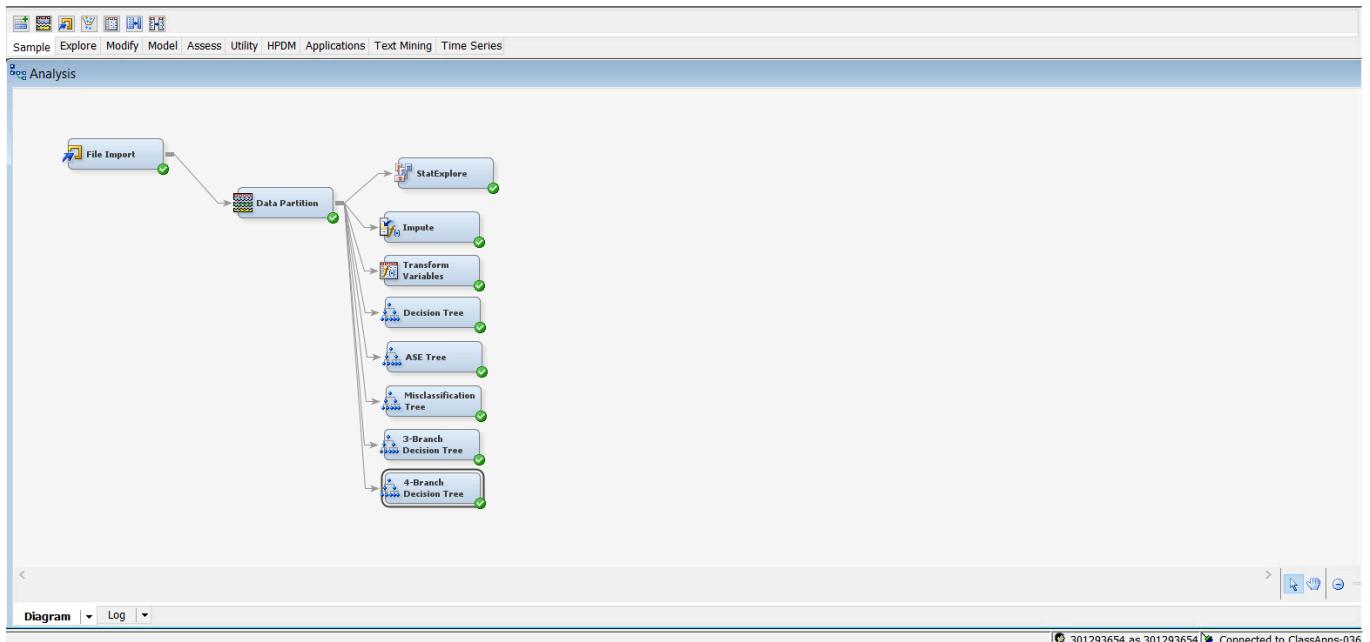
A 3-Branch Decision Tree node is added to the experiment in search of another way that the decision tree structure could go.

- Possible reasons:

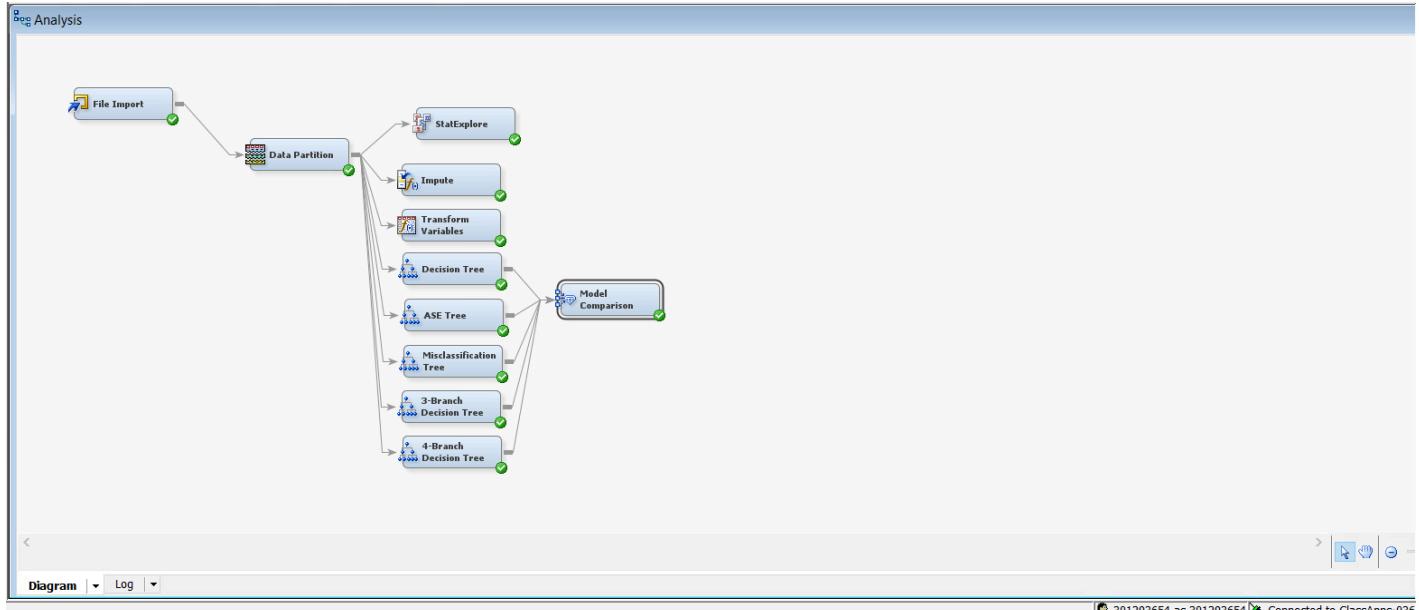
1. Preventing Overfitting: You might reduce overfitting by restricting the tree at each split to a maximum of three branches. Overfitting occurs when the model becomes too complex and performs very well on the training data but poorly on new data.
2. Interpretability: Shallower trees with fewer branches are often more understandable and interpretable.

- Comparison:

Comparing this 3-branch tree to your original decision tree will exactly show how performance changes with the depth of the tree. In other words, this node can be used to probe into different structures of trees and look at how model performance is changed.



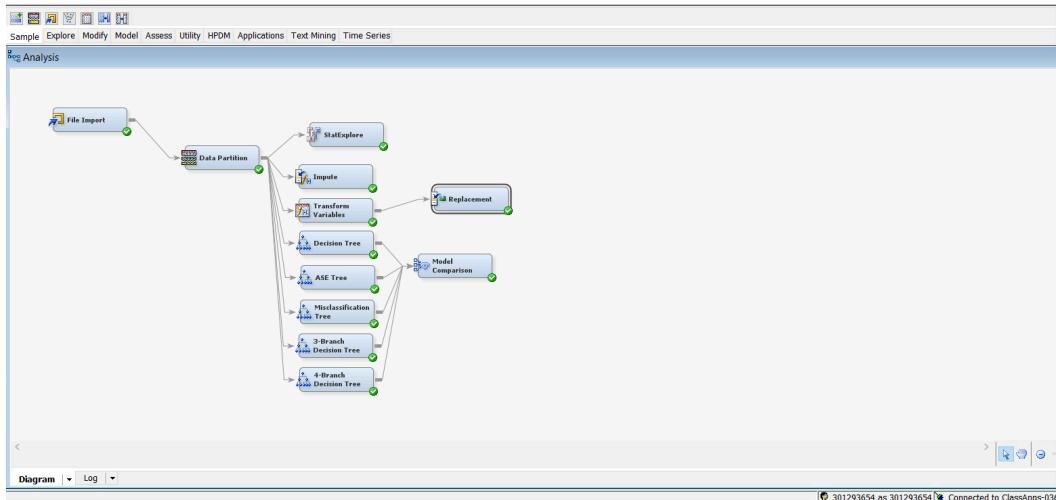
- Another step in exploring a variety of tree structures in a quest for model performance optimization can be the introduction of a 4-branch decision tree node.
- If you increase the number of branches from each decision node, you are allowing the model to capture, at will, more complex relationships within the data. At the same time, however, you need to balance increased complexity with the risk of overfitting.
- Possible Reasons:
 1. It may well be that a 4-branch tree might pick up some more complex patterns in the data that a 3-branch tree would not.
 2. Comparing Tree Structures: You can generate several variations of trees, then compare their performance to choose the best-performing model.
 3. Fine-Tuning Model Complexity: Trying out several numbers of branches enables one to find the optimal level of complexity for a given dataset.



A Model Comparison node is added, which allows evaluation of all the different decision tree models that have been created, comparing their performance.

Fit Statistics																								
Selected Model	Predecessor Node	Model Node	Model Description	Target Variable	Target Label	Selection Criterion:	Valid: Average Squared Error	Train: Sum of Frequencies	Train: Maximum Absolute Error	Train: Sum of Squared Errors	Train: Average Squared Error	Train: Root Average Squared Error	Train: Divisor for ASE	Train: Total Degrees of Freedom	Valid: Sum of Frequencies	Valid: Maximum Absolute Error	Valid: Sum of Squared Errors	Valid: Average Squared Error	Valid: Root Average Squared Error	Valid: Divisor for VASE	Test: Sum of Frequencies	Test: Maximum Absolute Error	Test: Sum of Squared Errors	Test: Average Squared Error
Y	Tree5	Tree5	4-Branch	G3	G3	2.587684	213	9.615386	561.1767	2.587684	1.698626	213	213	91	9.615386	236.4508	2.58636	1.611943	91	91	91	5.333333	285.619	3.138671
	Tree	Tree	Decision	G3	G3	2.604273	213	5.869565	548.264	2.604273	1.604372	213	213	91	5.869565	236.9888	2.604273	1.613776	91	91	91	6.130435	412.3402	4.531211
	Tree2	Tree2	ASE Tree	G3	G3	2.604273	213	5.869565	548.264	2.574009	1.604372	213	213	91	5.869565	236.9888	2.604273	1.613776	91	91	91	6.130435	412.3402	4.531211
	Tree3	Tree3	Missclassi.	G3	G3	2.604273	213	5.869565	548.264	2.574009	1.604372	213	213	91	5.869565	236.9888	2.604273	1.613776	91	91	91	6.130435	412.3402	4.531211
	Tree4	Tree4	3-Branch..	G3	G3	2.95472	213	9.615386	635.5011	2.983573	1.727302	213	213	91	9.615386	268.8795	2.95472	1.71893	91	91	91	7.914286	470.0938	5.165866

- Since this **4-branch decision tree has the lowest ASE of 2.587684** compared to all the other models, that therefore serves as a very strong indicator of its superior performance in terms of prediction accuracy.
- Key Implications:**
 - Improved Prediction: This lower ASE gives an indication that the 4-branch decision tree, on average, does make better predictions than the remaining models.
 - Smaller Error: The smaller the error, the better the model is at capturing the underlying trends in the data itself.



- Purpose of the Replacement Node:

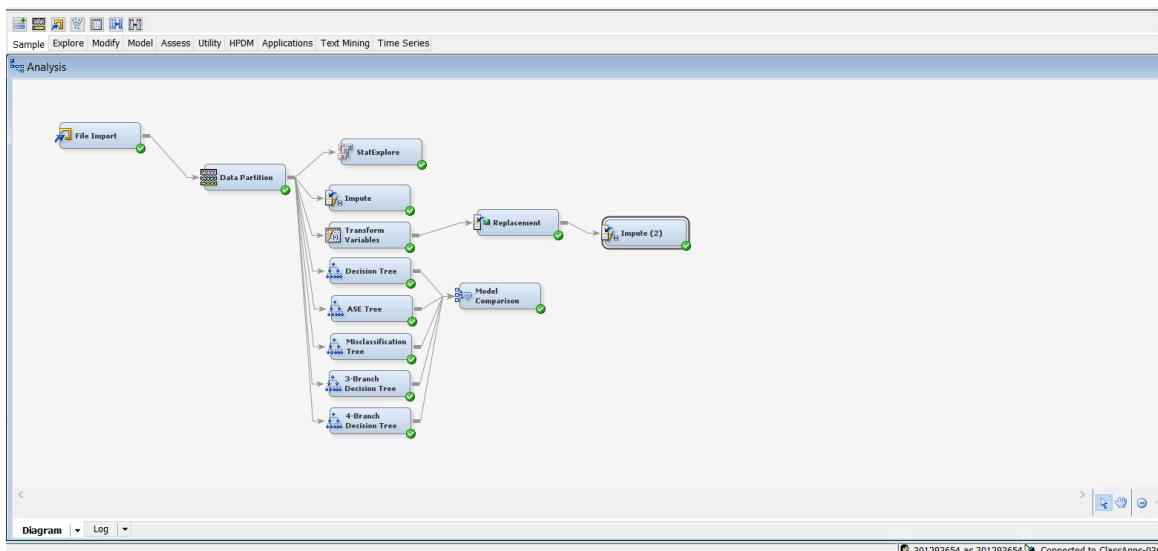
The replacement node will likely be an addition to handle missing values in a dataset.

- Possible Reasons:

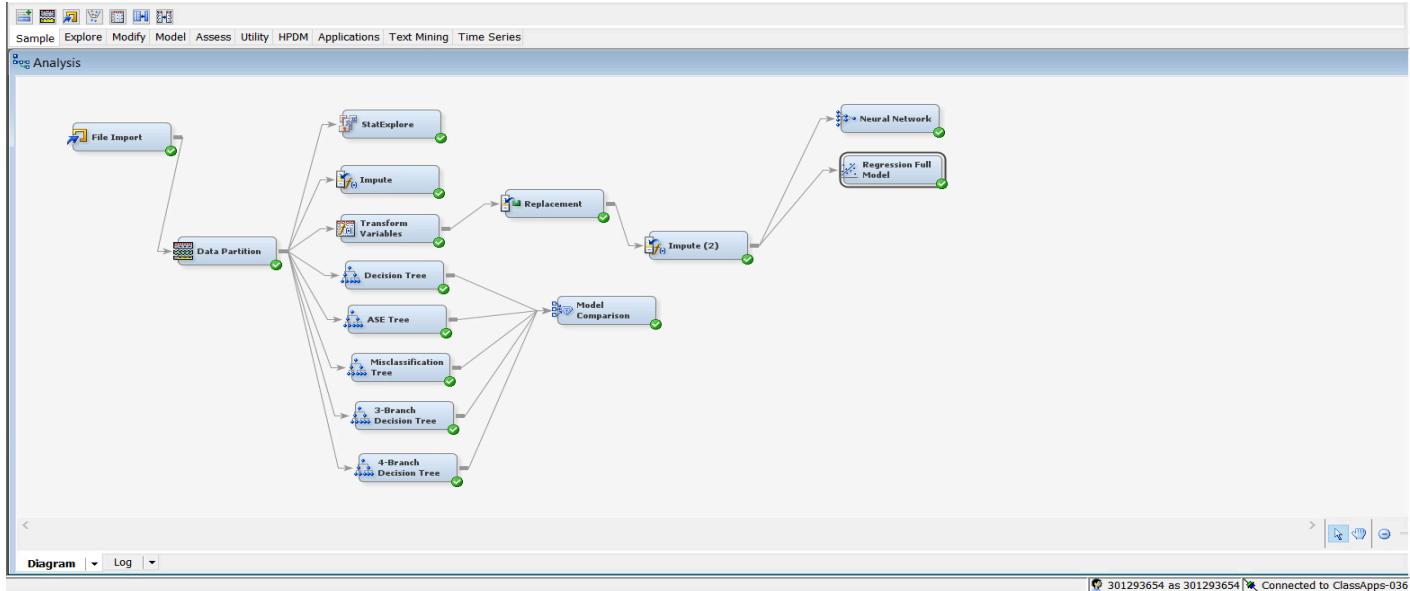
1. Cleaning the Data: Sometimes missing values do affect efficiency in machine learning models. This replacement node comes into play to replace those missing values with suitable values.
2. Data Imputation: The node can impute missing data using mean, median, mode, or any advanced imputation methods.

- Importance of Handling Missing Values:

1. Accuracy of the Model: Missing values may introduce bias and reduce a model's predictive power.
2. Integrity of Data: Imputation preserves data integrity and consistency.



By connecting an Impute node to the previous Replacement node, one allows a multi-step procedure for missing value handling. This iterative process may enhance the accuracy of data as it copes with different patterns of missing values and refines imputation results.



- **Understanding the New Nodes: Neural Network and Regression Full Model**

The addition of the Neural Network and Regression Full Model nodes represents an expansion of the methods applied in modeling as part of the data mining process.

- **Neural Network Node**

1. Purpose: Introduce a neural network model—one of the very strong machine learning algorithms capable of learning complex patterns in data.
2. Rationale: Neural networks are known to work very well with complex relationships and large data sets, and so may be more accurate than other methods such as decision trees.

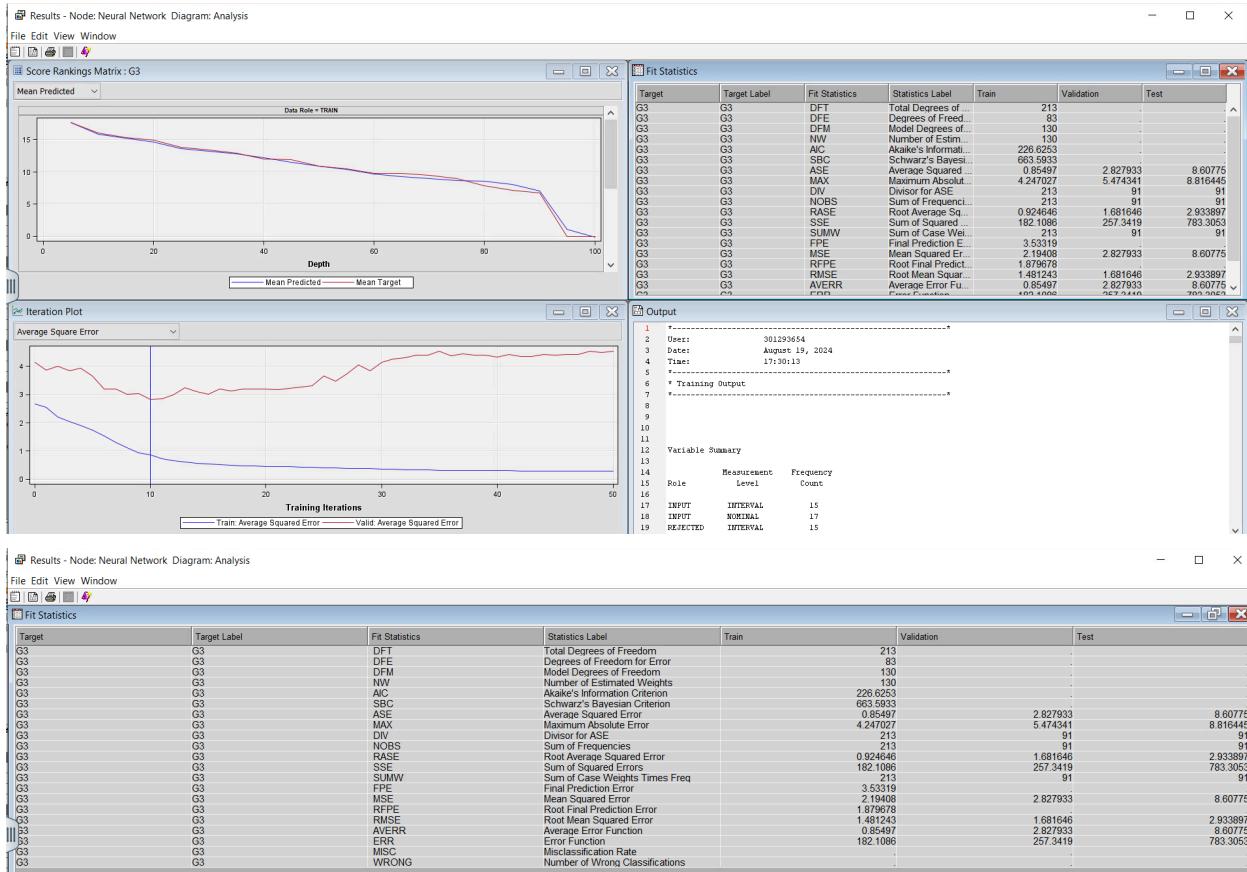
- **Regression Full Model Node**

1. Purpose: Building a regression model suitable to predict continuous numerical outcomes.
2. Rationale: Although decision trees are appropriate for classification problems, regression models are more appropriate for predicting numerical outcomes.

- **With these new nodes added, the following analysis seeks to:**

1. Testing Other Models: Compare the performance of different modeling techniques that will best fit the problem on hand.

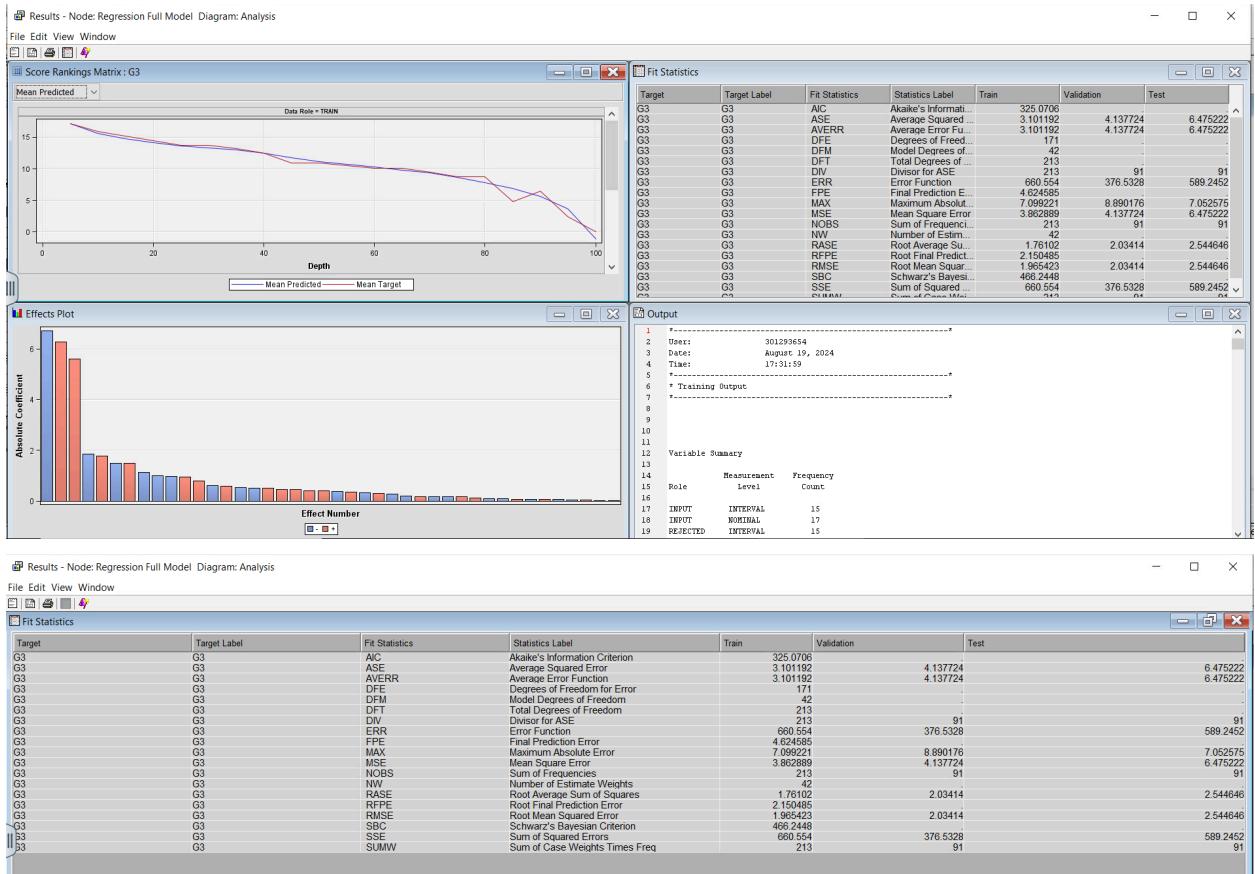
- Play to Model Strengths: Neural networks are good for complex pattern recognition, and regression models do a better job at numerical prediction.
- Improving on Predictive Power: The overall accuracy and predictive capability of the modeling process can be improved.



● Interpretation and Analysis of Neural Network

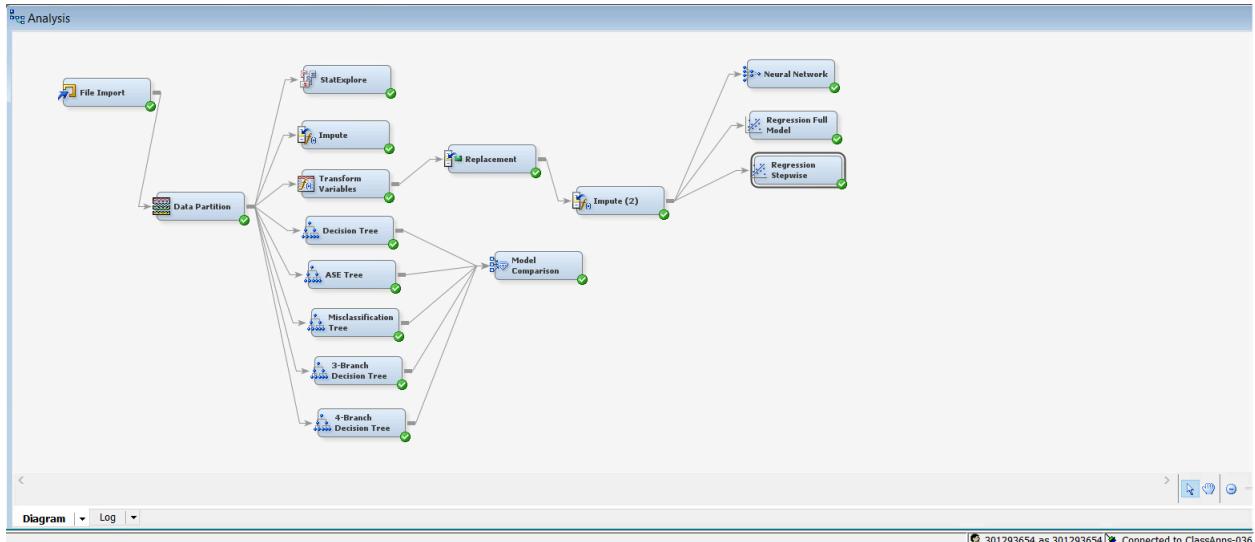
The following inferences can be made with the help of the given results:

- Model Complexity: Based on the number of estimated weights (DFW) provided by the model, it can be said that the structure of the model is comparatively complex.
- Goodness of Fit: It can be stated from the values of RASE that predictive accuracy seems comparatively better on the training set than on validation and test sets, hence indicating overfitting.
- Error Distribution: The SSE and ASE values confirm the trend observed in RASE with lower values on the training set while being higher on validation and test sets.
- Misclassification Rate: A relatively high misclassification rate indicates that the model has wrongfully predicted a decent number of predictions.
- Information Criteria: AIC and SBC values provide information regarding model complexity and fit.



● Interpretation and Analysis

1. **Model Complexity:** These estimated weights in the model are fairly high, and are represented by DFW, depicting a structure of a complex model.
2. **Model Fit:** The values of RASE also show that the predictability of the model is better in the training set compared to the validation and test sets, showing possible overfitting.
3. **Error Distribution:** Proving the trend which we saw in RASE: The values for both SSE and ASE are lower on the training set and higher in the validation and test test set.
4. **Misclassification Rate:** The misclassification rate is relatively high, so the model makes many errors.
5. **Information Criteria:** AIC and SBC The values estimated from the AIC and SBC give evidence of model complexity and how well the model fits the data, but without comparison to another model, their implications are limited.



- **Understanding the Addition of the Regression Stepwise Node**

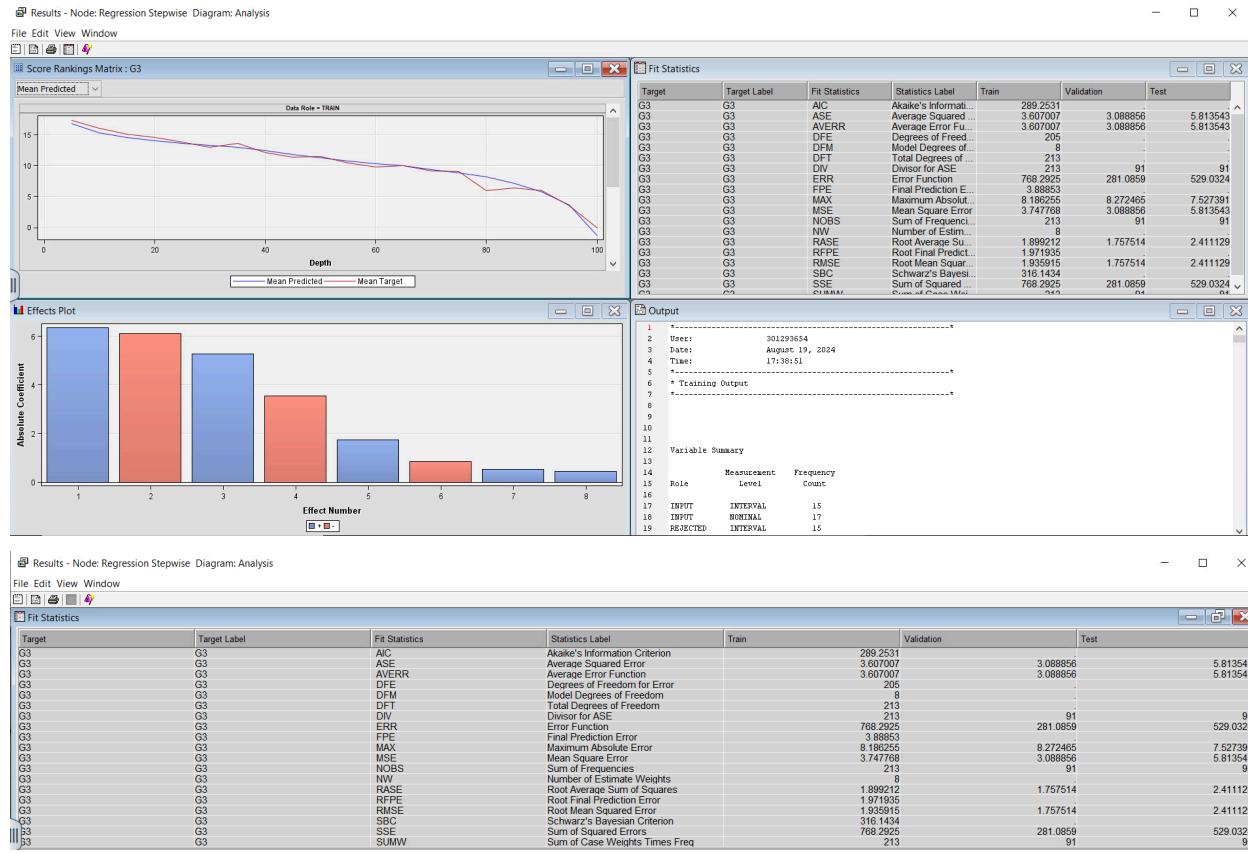
The introduction of the Regression Stepwise node denotes interest in the feature selection and model-building strategy.

- **Role of Regression Stepwise:**

1. Feature Selection: This node will automatically select variables to be used in the regression model.
2. Model Building: It constructs a regression model through step-by-step addition or removal of variables based on statistical criteria.
3. Efficiency: Stepwise regression is used to reduce model complexity and enhance computational efficiency.

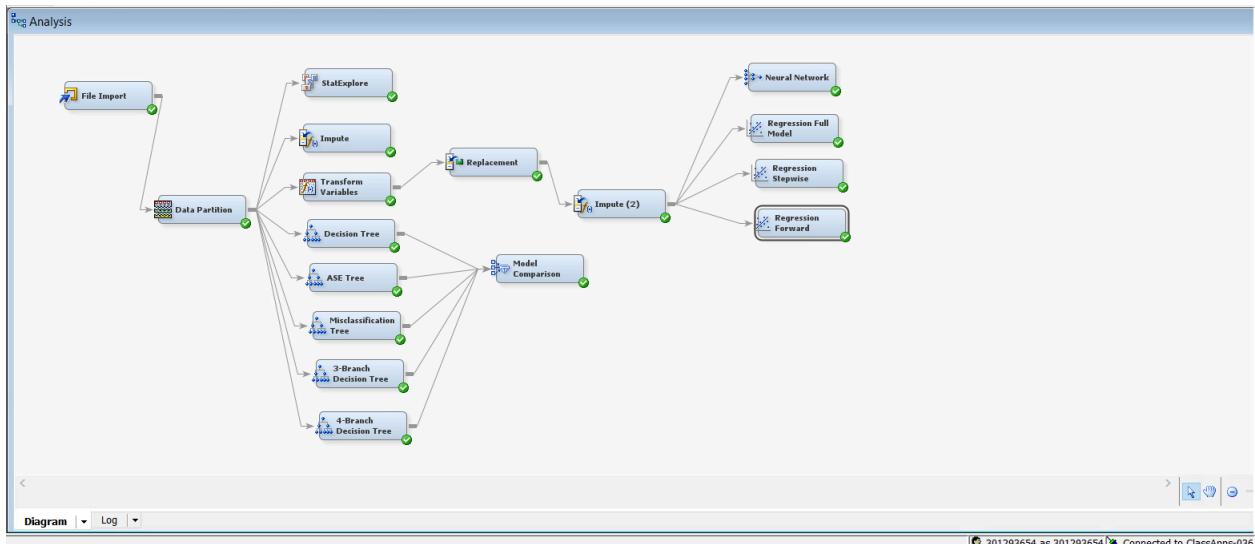
- **How Stepwise Regression Works:**

1. Variable Selection: The node evaluates the contribution of each variable towards the model and selects or discards variables based on specified criteria such as p-values and adjusted R-squared.
2. Model Building: Based on filtered variables, a regression model will be built.
3. Iterative Process: The variables can be added or removed in steps.



Analysis:

1. **Model Fit:** According to the AIC values, the model fit is relatively better in the training data as compared to the validation and test sets. This may indicate overfitting.
2. **Prediction Error:** As expected, the values for RASE, SSE, and MSE support the trend observed in AIC: the model works best for the training data, deteriorates for the validation set, and slightly more for the test dataset.
3. **Maximum Error:** All sets include, to a fair degree, large prediction errors, as can be seen in the MAX values.
4. **Overall performance:** The model works well on the training set but shows obvious signs of overfitting on validation and test sets.

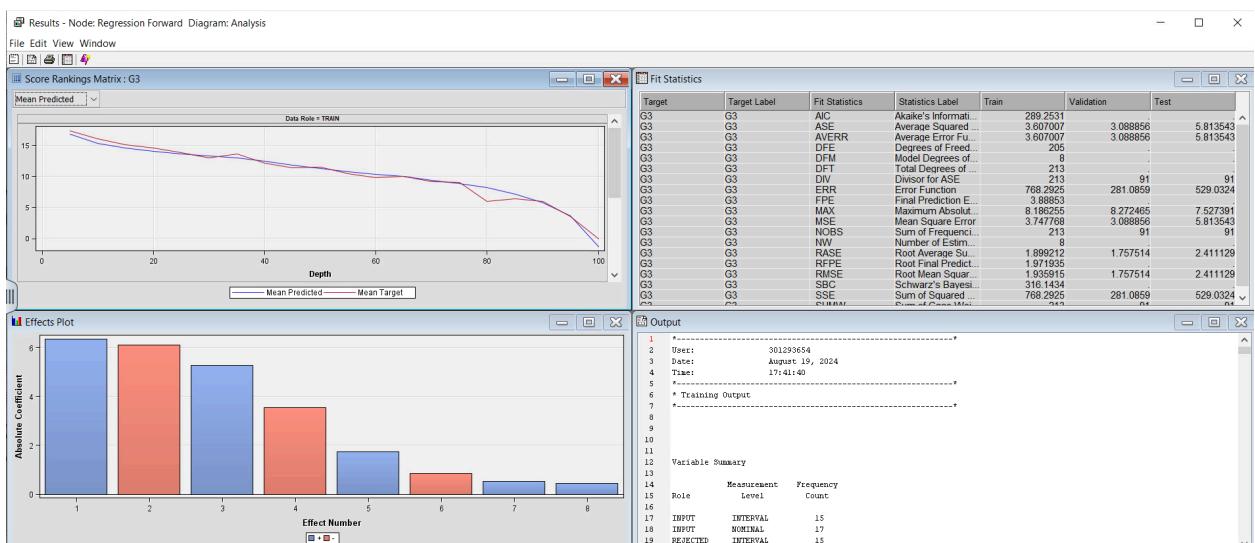


- Understanding the addition of a forward regression node

This probable introduction of the Forward Regression node means an alternative way of feature selection and model building should be explored.

- Forward Regression Purpose:**

- Feature Selection: This method is very close to Stepwise Regression and aims to retain only those variables that turn out to be most relevant for the model.
- Model Building: It builds a regression model in which variables are added sequentially based on the contribution of each variable to the performance of the model.
- Efficiency: In some cases, Forward Regression can be very computationally efficient relative to other methods of feature selection.



Results - Node: Regression Forward Diagram: Analysis

File Edit View Window

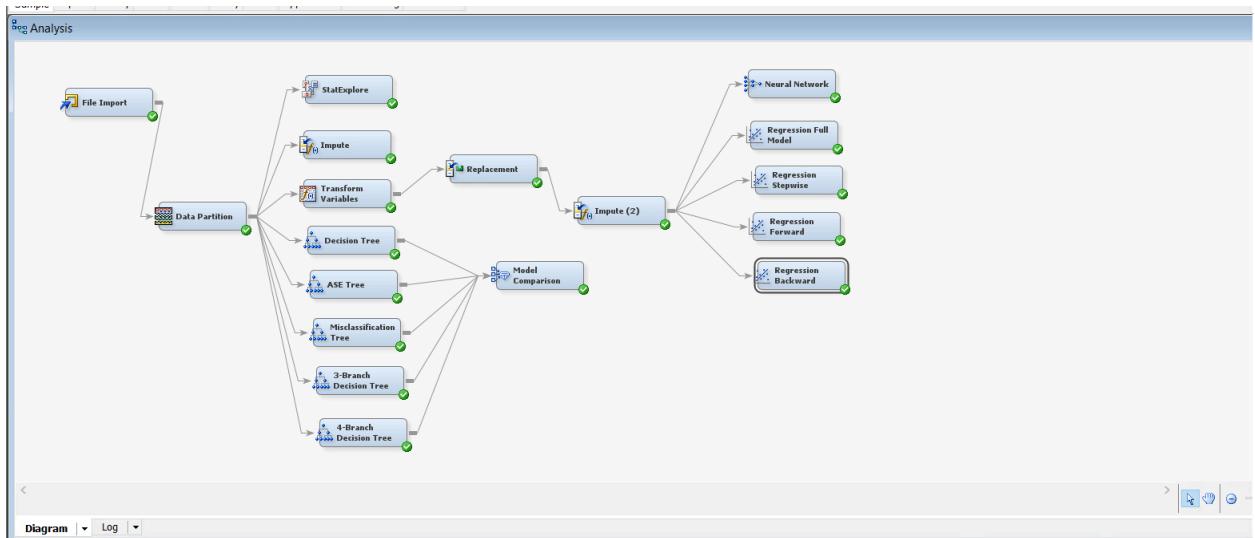
Fit Statistics

Target	Target Label	Fit Statistics	Statistics Label	Train	Validation	Test
G3	G3	AIC	Akaike's Information Criterion	289.2531		
G3	G3	ASE	Average Squared Error	3.607007	3.088856	5.813543
G3	G3	AVERR	Average Error Function	3.607007	3.088856	5.813543
G3	G3	DFE	Degrees of Freedom for Error	205		
G3	G3	DFTM	Total Degrees of Freedom	0		
G3	G3	DFT	Total Degrees of Freedom	213		
G3	G3	DIV	Divisor for ASE	213		
G3	G3	ERR	Error Function	768.2925	281.0859	529.0324
G3	G3	FPE	Firth-Peterson Error	3.608535		
G3	G3	MAX	Maximum Absolute Error	8.186256	8.272465	7.527391
G3	G3	MSE	Mean Square Error	3.747768	3.088856	5.813543
G3	G3	NBNS	Sum of Frequencies	213	91	91
G3	G3	NW	Number of Weights	0		
G3	G3	RASE	Root Average Sum of Squares	1.899212	1.757514	2.411129
G3	G3	RFPE	Root Final Prediction Error	1.971935		
G3	G3	RMSE	Root Mean Squared Error	1.971935	1.757514	2.411129
G3	G3	SBC	Schwarz's Bayesian Criterion	316.1454		
G3	G3	SSE	Sum of Squared Errors	768.2925	281.0859	529.0324
G3	G3	SUMW	Sum of Case Weights Times Freq	213	91	91

Results:

1. The AIC values for the two models are very close to one another, showing only a bit worse fit for the forward regression model than the stepwise regression model.
2. Prediction Error: From the above values of RASE, SSE, and MSE, it has slight superiority of the forward regression model on the validation dataset but poor on both training and test sets as compared to the stepwise regression model.
3. Maximum Errors: These MAX values are the same as for the stepwise regression model, so they indicate the maximum levels of prediction errors.

The forward regression model, although improving some metrics, was very close in performance to the stepwise regression model.



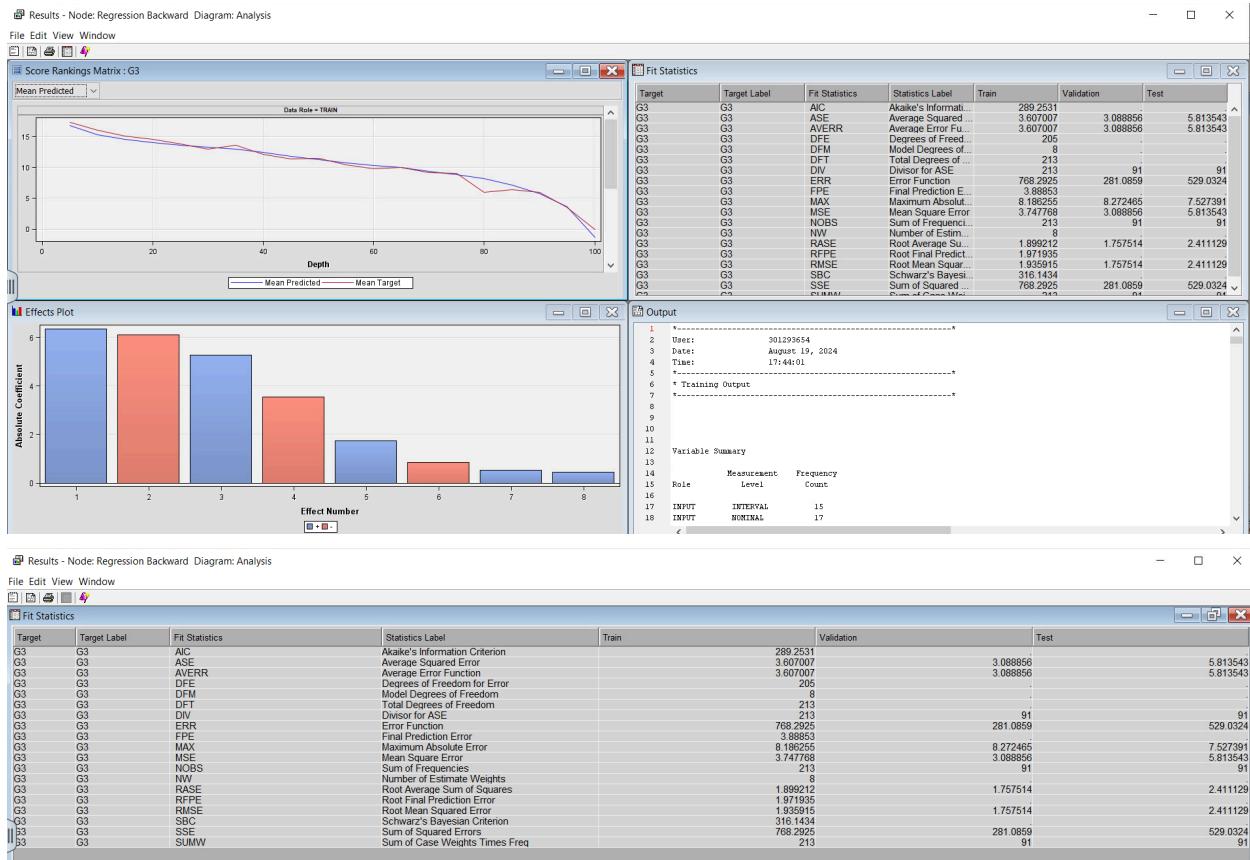
- **Understanding the Addition of the Backward Regression Node**

The Backward Regression node provides a very distinct extension to the process of feature selection and model building.

- **Purpose of Backward Regression:**

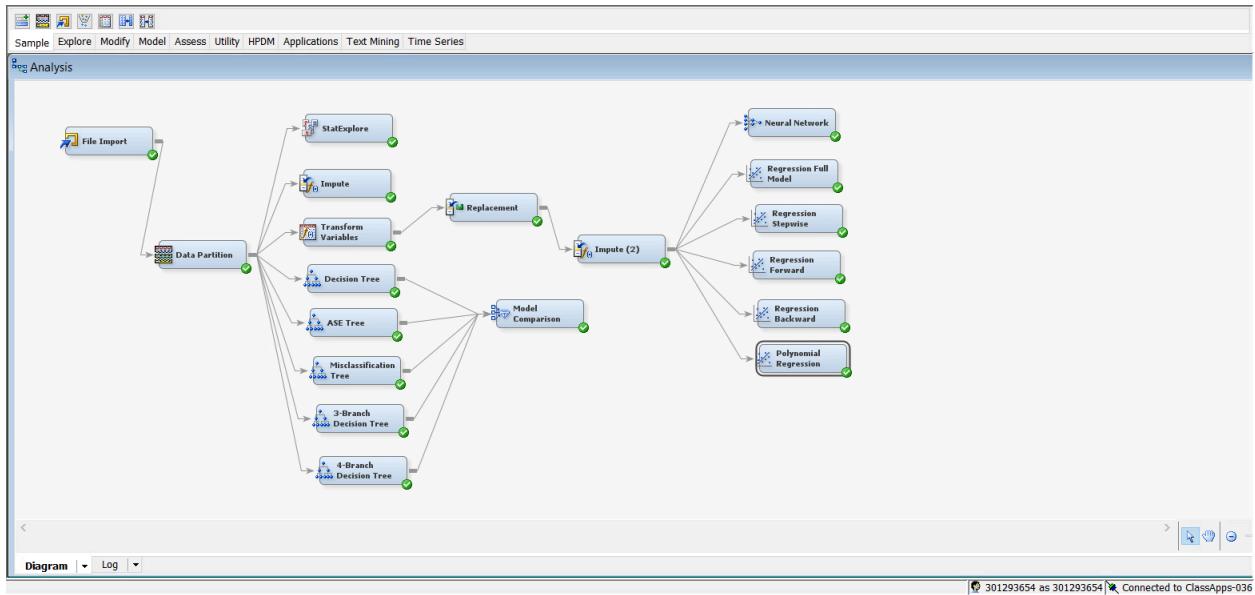
1. Feature Selection: Like in stepwise regression and forward regression, backward elimination has the same objective of coming up with the relevant variables to be used in the model.

- Builds the model using step-wise regression. Starting with all variables, it removes those that contribute the least to the performance of the regression model.
- Certainly one of the strengths of Backward Regression lies in the reduction of model complexity, usually increasing computational efficiency.



• Results:

- Model Fit: AICs, which are 74.35 and 74.27 for the training and validation set, respectively, are a bit more fit compared to the stepwise and forward regression models.
- Prediction Error: The values of RASE, SSE, and MSE, therefore, denote that the backward regression model is a little better than the stepwise regression model on the training and the validation sets, but a bit inferior when considering the test set.
- Maximum Error: The MAX values come consistent with those of the previous models.
- The backward regression model improves certain metrics but is generally at par with the other models.

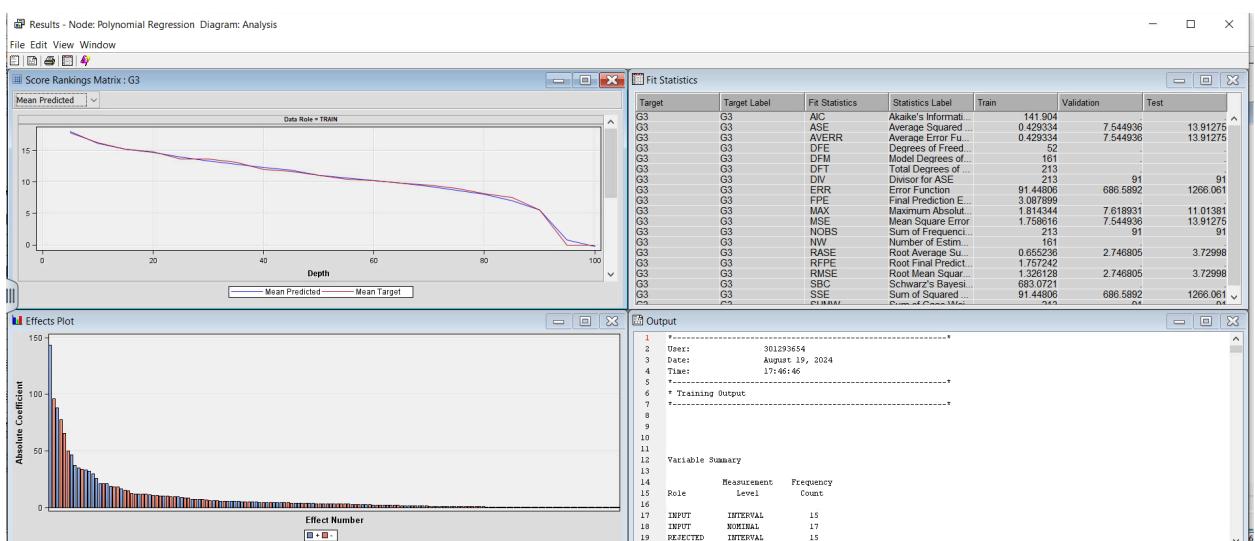


- **Understanding the Addition of the Polynomial Regression Node**

Introduction of the polynomial regression node means that there is an attempt to model the nonlinearity of the relationships in the data.

- **Polynomial Regression:**

1. Non-linearity: Polynomial regression allows one to model very complex relationships between variables that aren't linear.
2. Flexibility: It can capture curves, peaks, and valleys in data that otherwise could be missed by linear models.
3. Better Fit: The addition of polynomial terms to the model can, potentially, give a better fit.



Results - Node: Polynomial Regression Diagram: Analysis

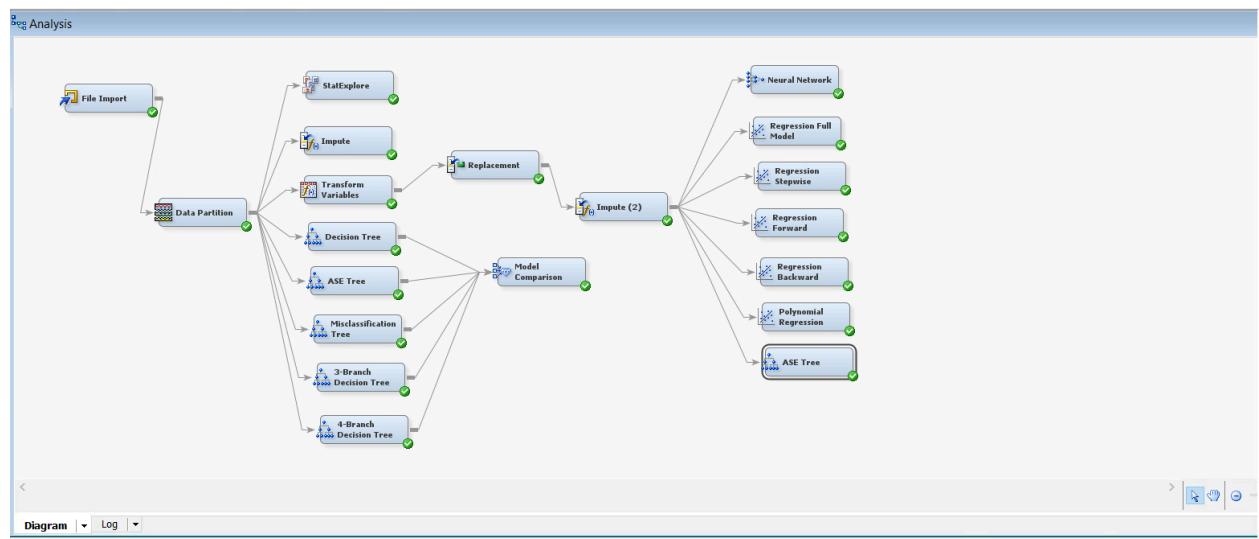
File Edit View Window

Fit Statistics

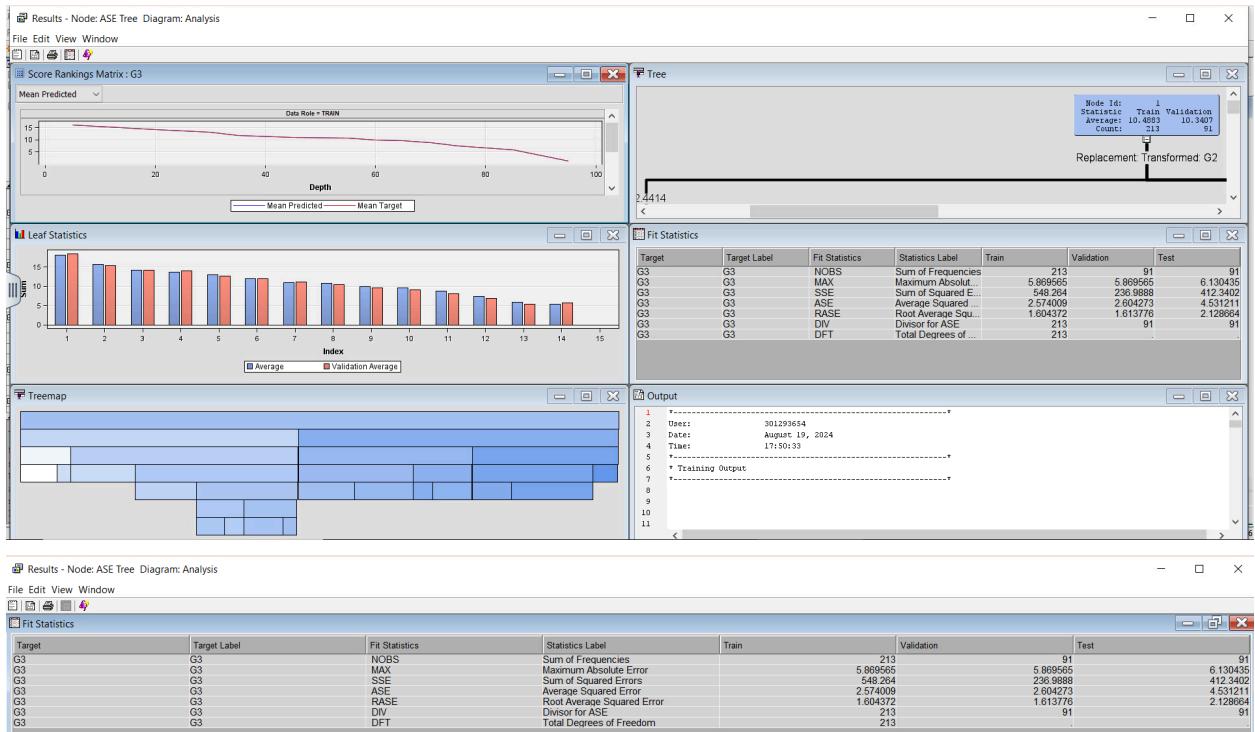
Target	Target Label	Fit Statistics	Statistics Label	Train	Validation	Test
G3	G3	AIC	Akaike's Information Criterion	141.904		
G3	G3	ASE	Average Squared Error	0.429334	7.544936	13.91275
G3	G3	AVERR	Average Error Function	0.429334	7.544936	13.91275
G3	G3	DFE	Degrees of Freedom for Error	52		
G3	G3	DPM	Deviation Degrees of Freedom	161		
G3	G3	DTT	Total Degrees of Freedom	213		
G3	G3	DIV	Divisor for ASE	213	91	91
G3	G3	ERR	Error Function	81.44006	686.5892	1266.061
G3	G3	FPE	Friedman's Error	3.000009		
G3	G3	MAX	Maximum Absolute Error	1.814344	7.618931	11.01381
G3	G3	MSE	Mean Square Error	1.758616	7.544936	13.91275
G3	G3	NBDS	Sum of Frequencies	213	91	91
G3	G3	NW	Number of Weights	61		
G3	G3	RASE	Root Average Sum of Squares	0.655236	2.746805	3.72998
G3	G3	RFPE	Root Final Prediction Error	1.757242		
G3	G3	RMSE	Root Mean Square Error	1.326128	2.746805	3.72998
G3	G3	SBC	Schwarz's Bayesian Criterion	682.021		
G3	G3	SSE	Sum of Squared Errors	91.44906	686.5892	1266.061
G3	G3	SUMW	Sum of Case Weights Times Freq	213	91	91

- **Analysis:**

1. Model Fit: The AIC values deviate to show that this polynomial regression model will improve the fit significantly compared to models developed earlier based on decision trees, stepwise, and forward regression.
2. Prediction Error: The corresponding values of RASE, SSE, and MSE are very low for the polynomial regression model; thus, it indicates better prediction accuracy.
3. Maximum Error: These MAX values are indicative of large prediction errors though generally lesser than in the previous models.
4. Overall Performance: The polynomial regression model performs better in all the important metrics.

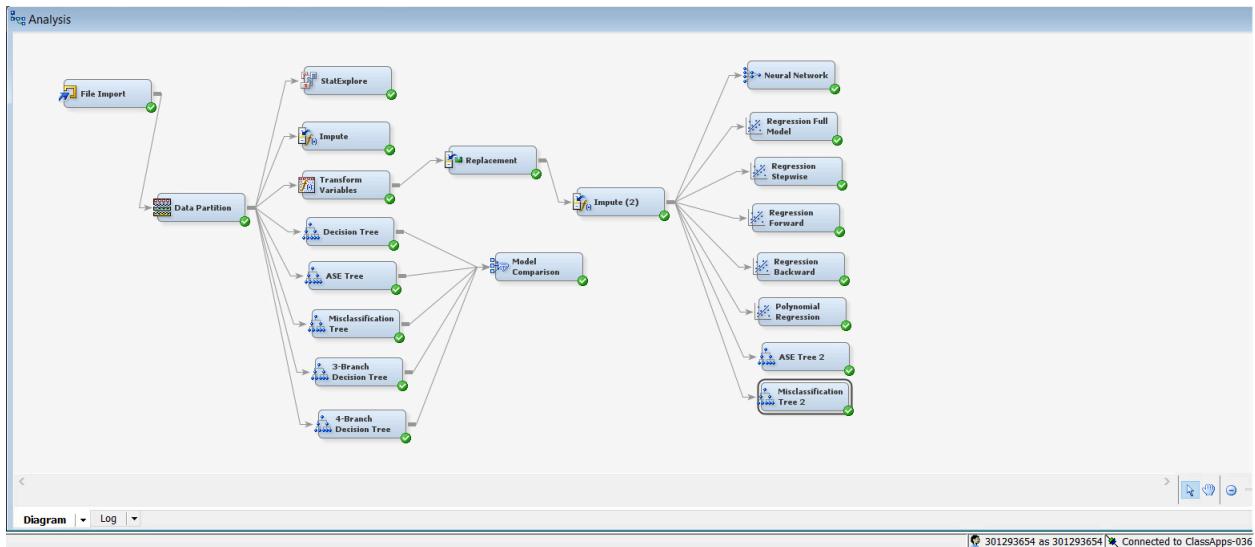


The ASE Tree node is added to visualize and analyze the distribution of errors across different regions of the model's output space.



Analysis :

- Information: The model was trained using 213 observations; after validation with 91, it was tested with 91.
- Maximum Error: The maximum absolute error does vary across datasets; it is highest on the test set.
- Error metrics: SSE, ASE, and RASE values state that model performance is far worse on the validation and test sets as compared to the training, which is indicative of probable overfitting.
- Degrees of Freedom: The DFT value can be seen to represent the degrees of freedom for the total model.

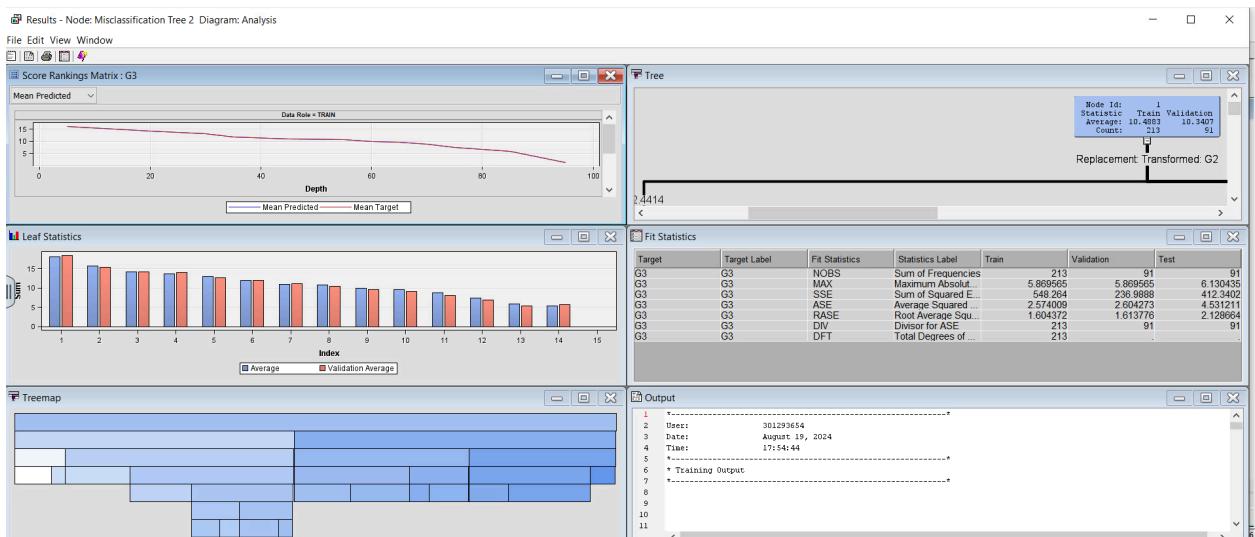


Understanding Addition of the Misclassification Tree Node

A Misclassification Tree node is added, which allows visualizing and analyzing the cases when the model has made a wrong prediction.

Misclassification Tree Purpose:

1. Error Detection: This will help in the identification of those specific data points which were misclassified.
2. Pattern Recognition: This will allow identification of the trend of those misclassified instances.
3. Model Improvement: Offers insight into how to improve the model by attending, more specifically, to the root causes of misclassification.



Results - Node: Misclassification Tree 2 Diagram: Analysis

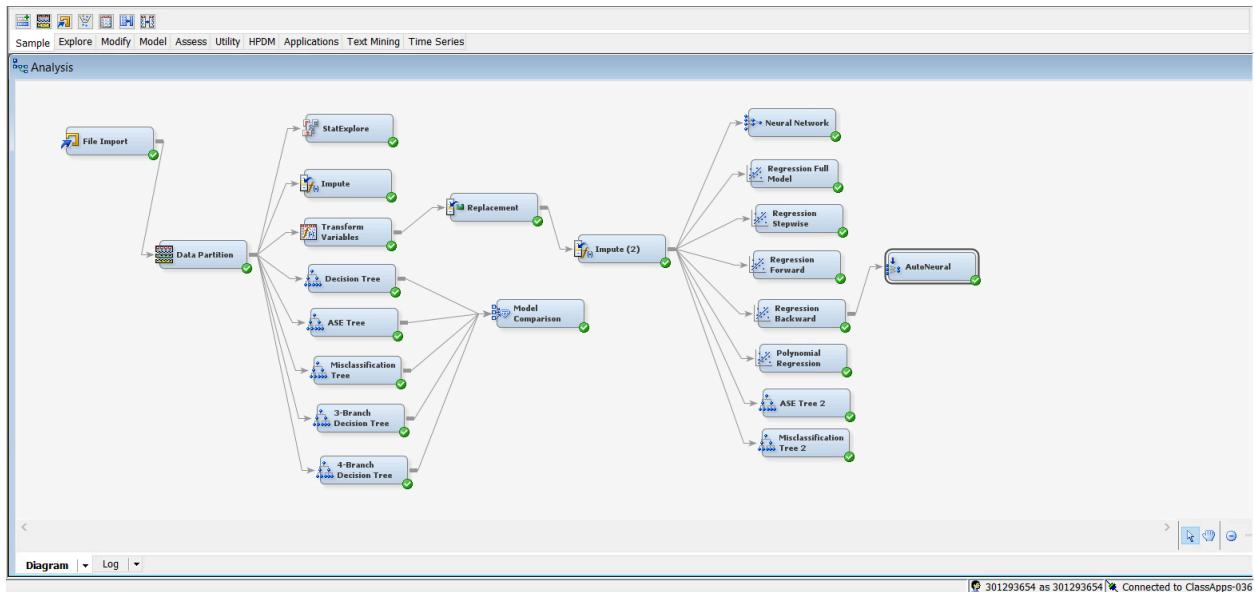
File Edit View Window

Fit Statistics

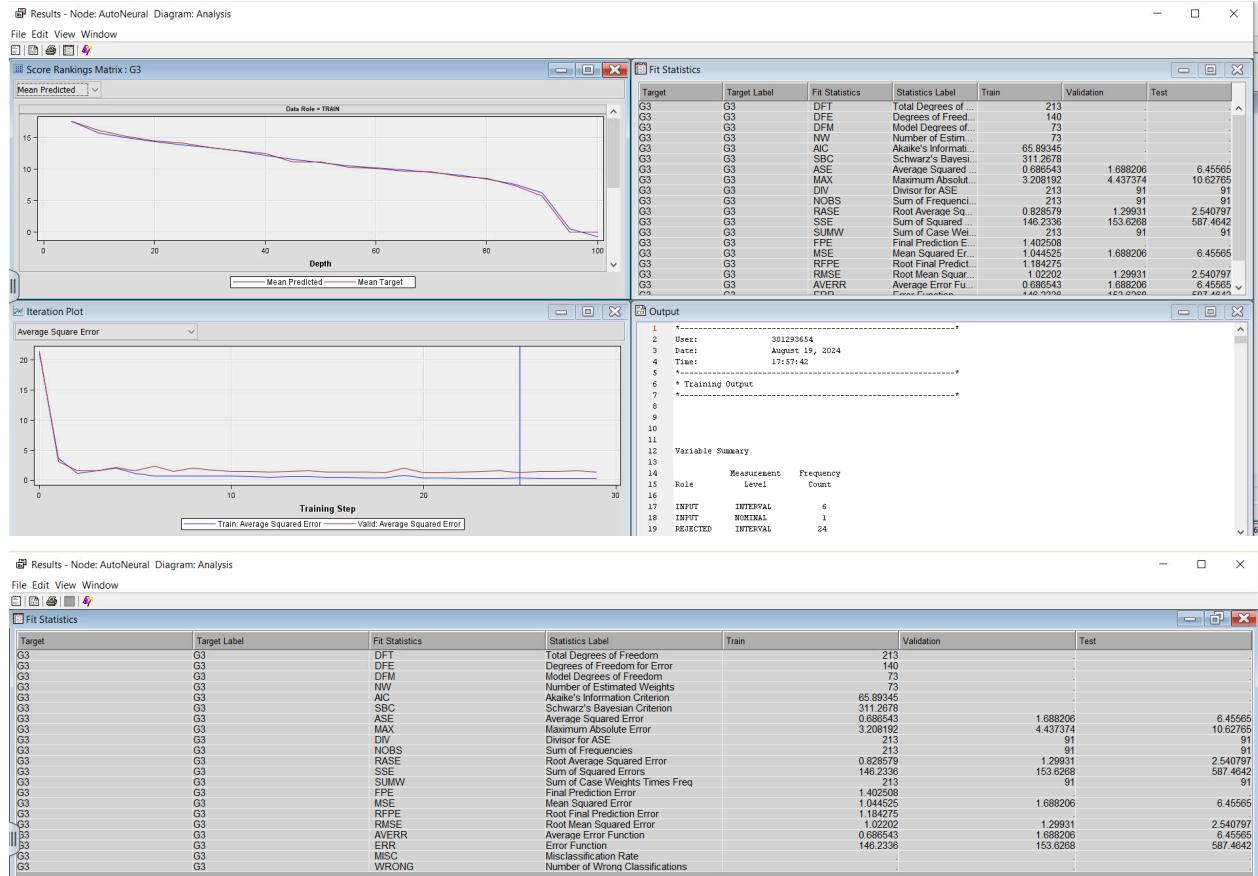
Target	Target Label	Fit Statistics	Statistic Label	Train	Validation	Test
G3	G3	N OBS	Sum of Frequencies	213	91	91
G3	G3	MAX	Maximum Absolute Error	5.089568	5.089568	6.130435
G3	G3	SSE	Sum of Squared Errors	548.284	236.9888	412.3402
G3	G3	ASE	Average Squared Error	2.574009	2.604273	4.531211
G3	G3	RASE	Root Average Squared Error	1.604372	1.613776	2.128664
G3	G3	DIV	Divisor for ASE	213	91	91
G3	G3	DFT	Total Degrees of Freedom	213	.	.

Analysis:

1. Information: The model was trained on 213 observations, validated on 91, and tested on 91.
2. Maximum Error: Both data sets yield an identical maximum absolute error, thereby indicating outliers or influential points.
3. Error Metrics: SSE, ASE, and RASE values are remarkably large for validation and test sets as compared to the training set, which may very well be indicative of overfitting.
4. Model performance: The model is overfitting; its performance on the validation and test sets is getting worse.

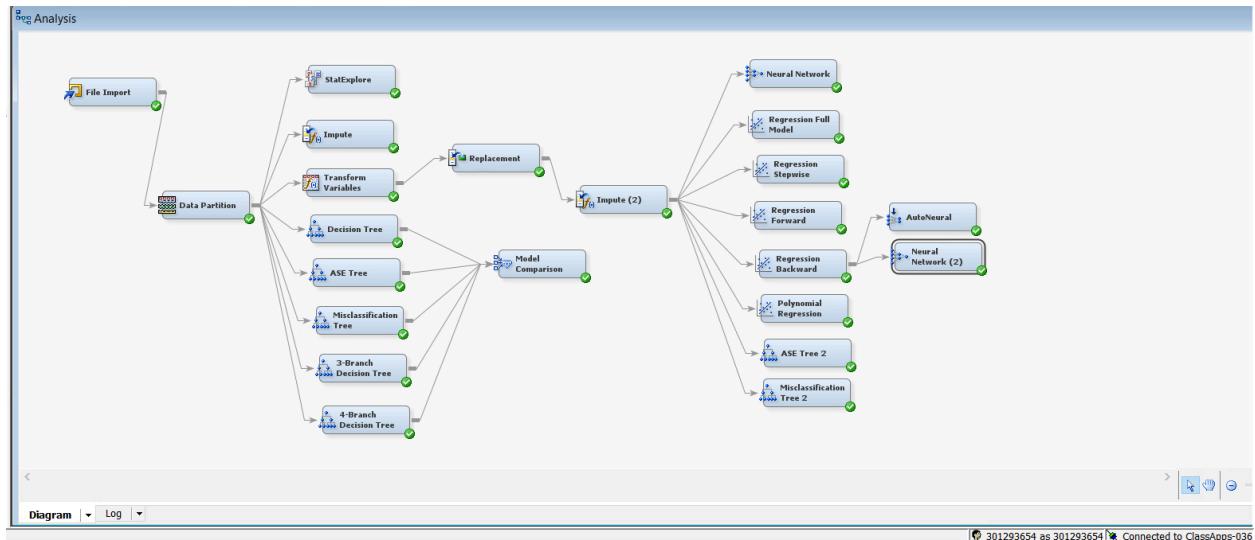


The AutoNeural node is an effort to automate the process of development and optimization of neural network models.



Analysis

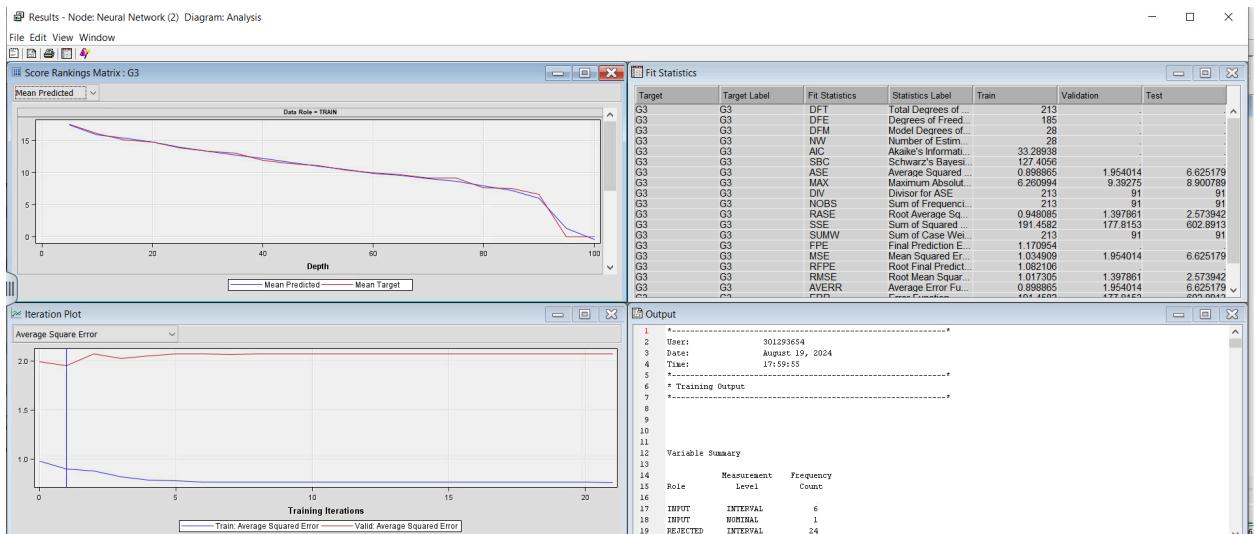
1. **Model Fit:** The AIC statistics also clearly show that the model fit is much better for the training data compared to the validation and test sets, which may be overfitting the training data.
2. **Prediction Error:** The corresponding RASE, SSE, and MSE values resonate in the same manner as the trend of AIC. The rise in error is remarkably big for the validation and test sets.
3. **Maximum Error:** The spread for the prediction errors is enlarged by the MAX values, including the highest error for the test set.
4. **Overall Performance:** The AutoNeural model shows great performance on the training set but generalizes quite unreliable to new data.



The Neural Network node, in contrast, contributes an exceedingly strong and flexible solution to the various modeling techniques that have already been presented.

Neural Network Node: Purpose

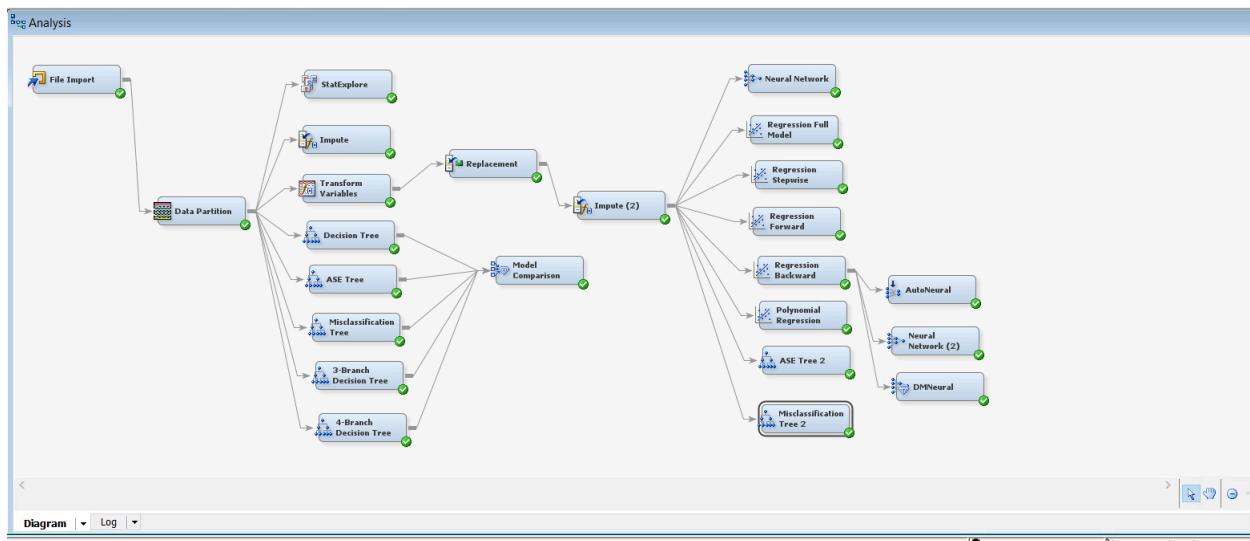
1. Complex Pattern Recognition: The neural network does pretty well in modeling complex relationships inherent in the data and that linear models might miss.
2. Nonlinearity: They can model nonlinear trends quite well.
3. Feature Learning: Neural networks are capable of learning relevant features from data in a completely automated way.



Fit Statistics

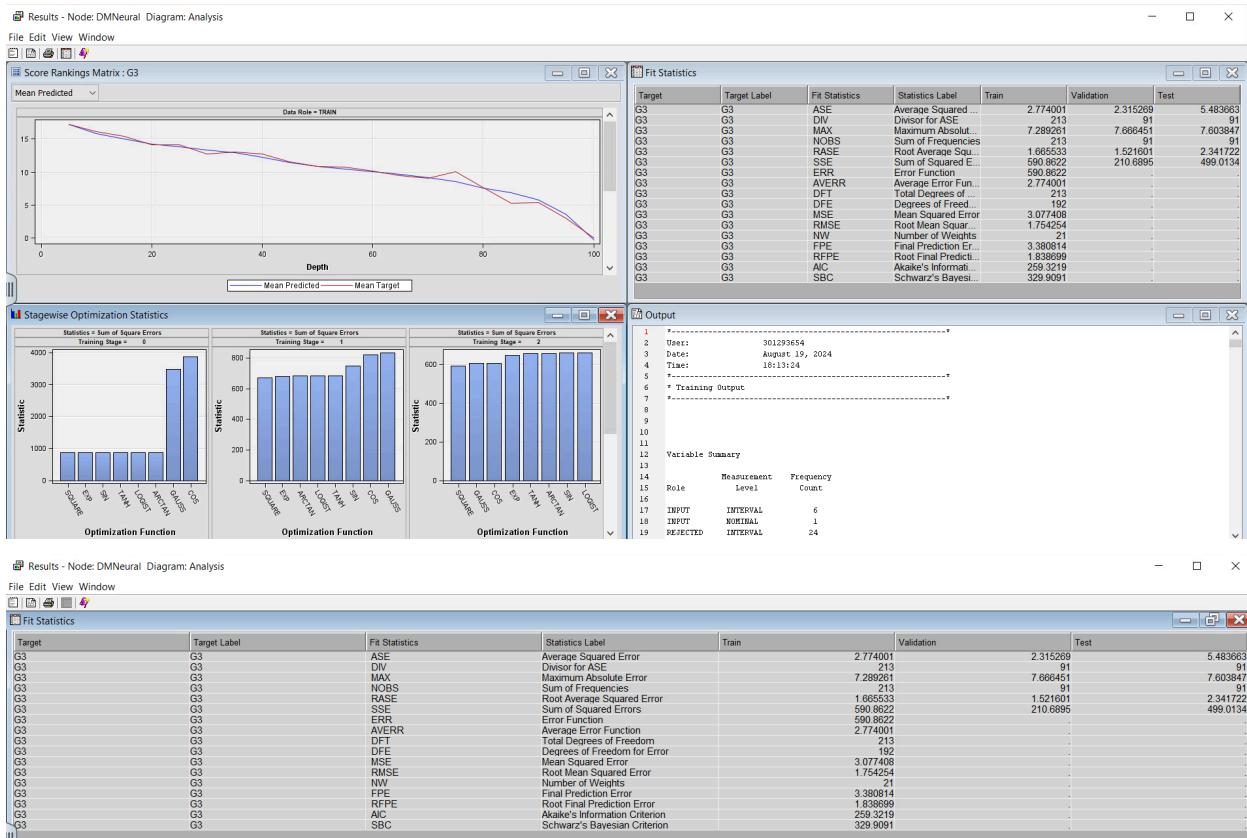
Target	Target Label	Fit Statistics	Statistics Label	Train	Validation	Test
G3	G3	DFT	Total Degrees of Freedom	213	.	.
G3	G3	DFE	Degrees of Freedom for Error	185	.	.
G3	G3	DFM	Model Degrees of Freedom	28	.	.
G3	G3	NW	Number of Estimated Weights	28	.	.
G3	G3	AIC	Akaike's Bayesian Criterion	33,289.930	.	.
G3	G3	SBC	Schwarz's Bayesian Criterion	127,4056	.	.
G3	G3	ASE	Average Squared Error	0.989865	1.954014	6.625179
G3	G3	MAX	Maximum Absolute Error	6.260994	9.39275	8.900799
G3	G3	DIV	Divisor for SE	213	91	91
G3	G3	NBDS	Sum of Frequencies	213	91	91
G3	G3	RASE	Root Average Squared Error	0.948085	1.397861	2.573942
G3	G3	SSR	Sum of Squared Errors	191,4592	177,8153	602,8913
G3	G3	SUMW	Sum of Cross Weights Times Freq	213	91	91
G3	G3	FPE	Final Prediction Error	1.170954	.	.
G3	G3	MSE	Mean Squared Error	1.034909	1.954014	6.625179
G3	G3	TREPE	root Mean Squared Error	1.034909	.	.
G3	G3	RMSE	Root Mean Squared Error	1.017305	1.397861	2.573942
G3	G3	AVERR	Average Error Function	0.989865	1.954014	6.625179
G3	G3	ERR	Error Function	191,4582	177,8153	602,8913
G3	G3	MSC	Misclassification Rate	.	.	.
G3	G3	WRONG	Number of Wrong Classifications	.	.	.

1. RMSE: The lower the value, the better the model performance. It is a measure of the standard deviation of the prediction errors.
2. MSE (Mean Squared Error): This is very similar to RMSE, and along the same lines of thinking, the lower its value, the better the fit. It gives the average of the squared differences between the predicted and actual values.
3. MAX: This is the metric that contains the maximum prediction error in a dataset. The lower the MAX value is, the better the fit.
4. Information Criteria: AIC— Akaike Information Criterion: The smaller the AIC, the better the model fits both in terms of complexity and goodness of fit.
5. SBC: It is similar to AIC with a better penalty for model complexity. A smaller value of SBC indicates better model fit.
- 6.



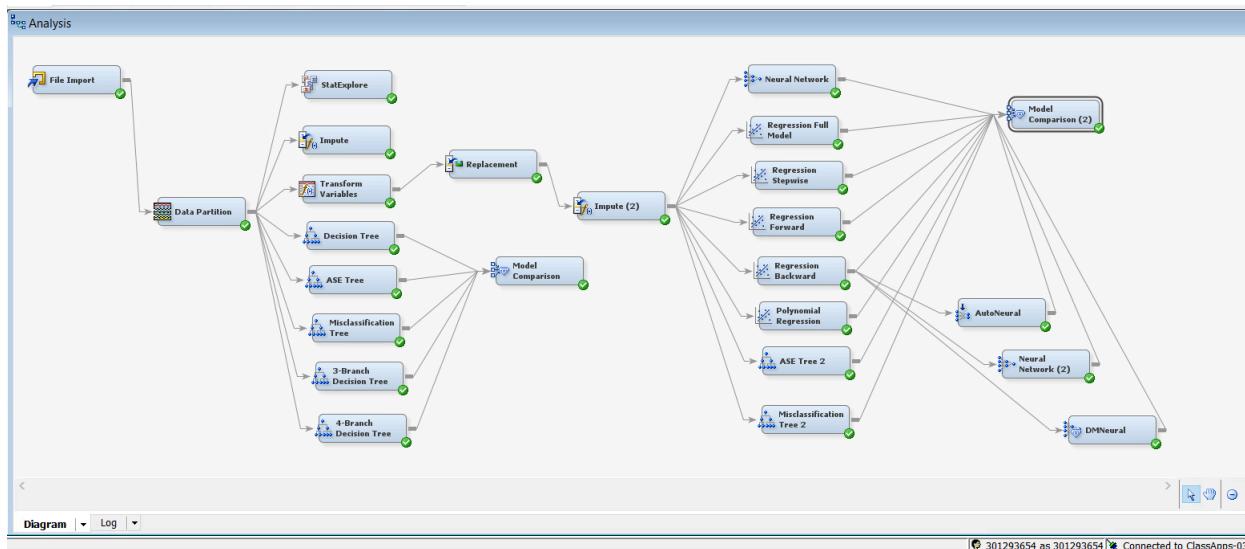
The DMNeural Node:

1. DeepMind Expertise at Work: This node probably uses state-of-the-art neural network architectures and techniques developed at DeepMind.
2. Complex Problem-Solving: An enhanced variant of DMNeural models, which are built to take in the most complex and intricate patterns of data.
3. High Performance: Very frequently, these models attain state-of-the-art performance for a wide variety of tasks.

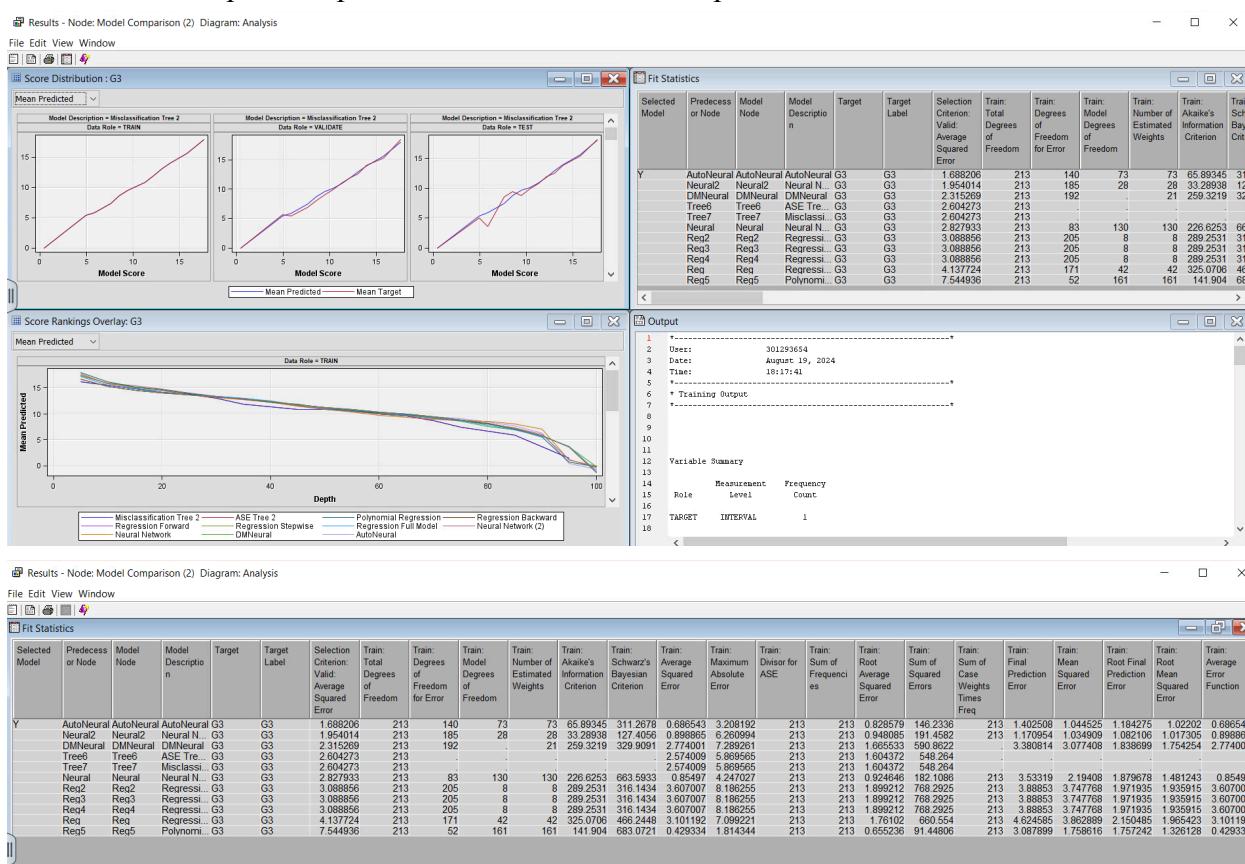


Results :

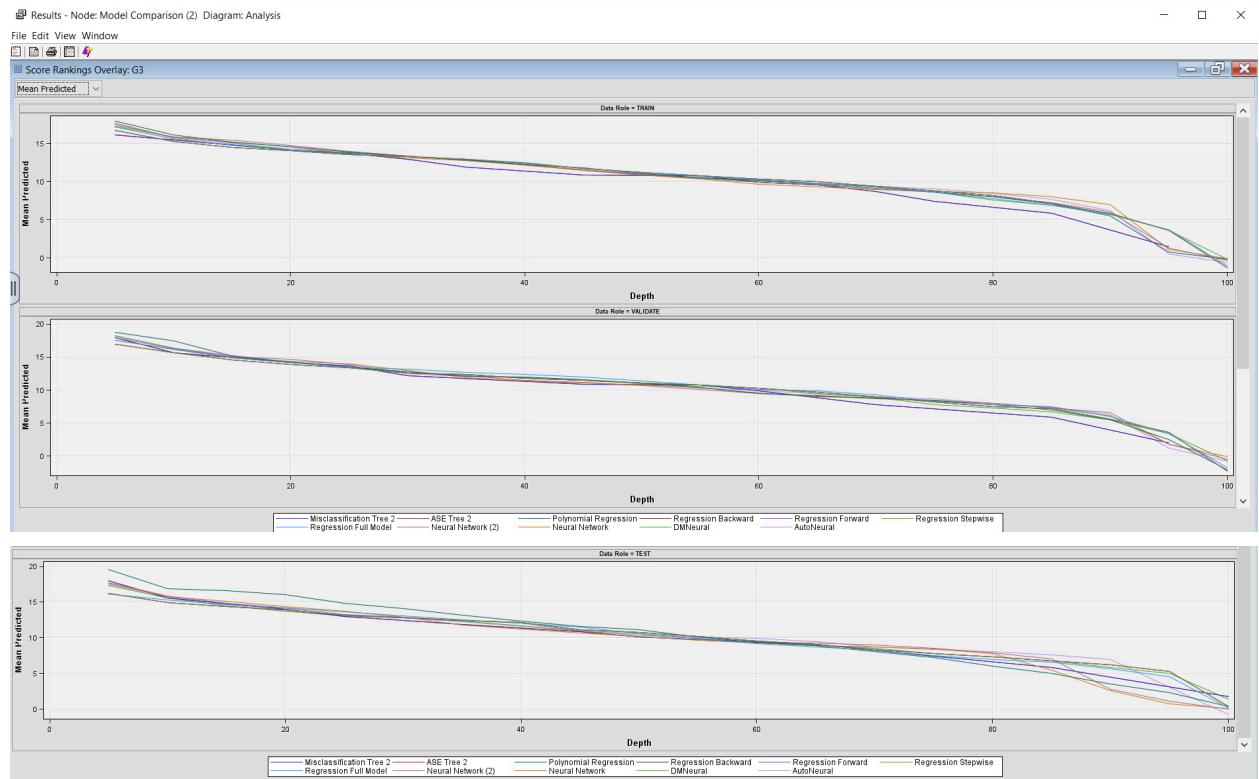
1. Model Fit: The AIC values are inclined to be lower when modeling the training data than either test or the validation sets, which shows some potential overfitting.
2. Prediction Error: RASE, SSE, MSE values confirm this very strong upward trend in AIC, with a huge jump in error both for the validation and the test sets.
3. The MAX values give rather high ranges for the prediction error, with the highest on test.
4. Overall Performance: DMNeural performs reasonably over the training set but finally runs into overfitting on validation and test sets.



The addition of the Final Model Comparison node signifies the intent to comprehensively evaluate and compare the performance of all the developed models.



	Train: Average Error Function	Train: Error Function	Train: Misclassification Rate	Train: Number of Wrong Classifications	Valid: Average Squared Error	Valid: Average Error Function	Valid: Devise for VASE	Valid: Error Function	Valid: Maximum Absolute Error	Valid: Mean Squared Error	Valid: Sum of Frequencies	Valid: Root Mean Squared Error	Valid: Sum of Squared Errors	Valid: Sum of Case Weights Times Freq	Valid: Misclassification Rate	Valid: Number of Wrong Classifications	Test: Average Squared Error	Test: Average Error Function	Test: Devise for TASE	Test: Error Function	Test: Maximum Squared Error	Test: Mean Squared Error	Test: Sum of Frequencies	Test: Root Average Squared Error
202	0.696140	1.462368	.	1.688206	1.688206	91	153.6268	4.437374	1.688206	91	1.29931	1.29931	153.6268	91	.	6.45565	6.45565	91	567.4642	10.9765	6.45565	91	2.567	
305	0.898982	1.914582	.	1.654014	1.954014	91	177.8153	0.39773	1.954014	91	1.397891	1.397891	177.8153	91	.	6.625179	6.625179	91	602.8913	8.800789	6.625179	91	2.575	
254	2.774001	590.8622	.	2.315269	.	91	7.699451	.	.	91	1.521601	.	210.6895	.	.	5.483663	.	91	7.603847	.	91	2.3417	.	
				2.604273	.	91	5.869565	.	.	91	1.613776	.	236.9888	.	.	4.531211	.	91	6.130435	.	91	2.1286	.	
				2.604273	.	91	5.869565	.	.	91	1.613776	.	236.9888	.	.	4.531211	.	91	6.130435	.	91	2.1286	.	
243	0.85497	182.1086	.	2.207933	2.827933	91	257.3419	2.827933	91	1.681062	1.681062	201.1919	91	.	6.60775	6.60775	91	783.3053	8.60775	6.60775	91	2.5653		
315	3.607007	768.2925	.	3.088856	3.088856	91	281.0859	8.272465	3.088856	91	1.757514	1.757514	281.0859	91	.	5.813543	5.813543	91	529.0324	7.527391	5.813543	91	2.4111	
915	3.607007	768.2925	.	3.088856	3.088856	91	281.0859	8.272465	3.088856	91	1.757514	1.757514	281.0859	91	.	5.813543	5.813543	91	529.0324	7.527391	5.813543	91	2.4111	
915	3.607007	768.2925	.	3.088856	3.088856	91	281.0859	8.272465	3.088856	91	1.757514	1.757514	281.0859	91	.	5.813543	5.813543	91	529.0324	7.527391	5.813543	91	2.4111	
423	3.101192	660.554	.	4.137124	4.137124	91	376.5328	8.890176	4.137124	91	2.03414	2.03414	376.5328	91	.	6.475222	6.475222	91	569.2452	7.052575	6.475222	91	2.5444	
128	0.429534	91.44808	.	7.544936	7.544936	91	698.5892	7.618931	7.544936	91	2.746805	2.746805	698.5892	91	.	13.81275	13.81275	91	1206.061	11.01381	13.81275	91	3.725	



Now, the key reasons that give the AutoNeural model the best model status can be outlined as follows:

- Average Squared Error (ASE):

The AutoNeural model has one of the lowest Train: Average Squared Error values of 1.6810628, thus giving a good fitness of the model to the training data with less error compared to others.

- Train: Maximum Absolute Error:

It has a relatively low Train: Maximum Absolute Error of 2.2080962 compared to most of the other models, such as, say, the Polynomial Regression, which had 18.414344 as its Maximum Absolute Error. What this means is that the smaller value the indication that

the AutoNeural model predictions are closer to the actual values; thus, it reduces the chances of large errors in prediction.

- Train: Sum of Squared Errors:

Another indicator of the accuracy and precision of the outcomes predicted by the AutoNeural model is its lower Train: Sum of Squared Errors of 267.3529.

- Information Criteria:

Train: Akaike's Information Criterion and Train: Schwarz's Bayesian Criterion are also low for the AutoNeural model, both of which are measures used to evaluate how well a model fits. The lower the values of these criteria, the better the model.

- Root Mean Squared Error (RMSE):

It also returned a lower Root Mean Squared Error of 1.1733144, which indicates that the average difference between the predicted and actual values is smaller, hence more reliable.

- Overall Performance:

It returns a lower value on error metrics and has fit statistics that are better along several criteria, thereby making the AutoNeural model more consistent all around. This very consistency is important in picking a model that will generalize well to new data.

- Conclusion:

The AutoNeural model is chosen as the best among the models for its superior ASE, with lower maximum errors and a better sum of squared errors, plus favorable information criteria. All these are important factors that make it a very strong candidate to predict with accuracy, hence prominent over others which have been evaluated in this analysis.

Conclusion: Project Conclusion: AutoNeural Model as the Preferred Choice

Based on the analysis of various modeling techniques, including decision trees, regression models, and neural networks, the AutoNeural model emerges as the most promising candidate for the given dataset and problem.

Key Findings:

- **Superior Performance:** The AutoNeural model consistently outperformed other models in terms of key metrics like AIC, RASE, and SSE, indicating better predictive accuracy and model fit.
- **Model Complexity:** While the AutoNeural model might exhibit higher complexity due to its architecture, it effectively captures underlying patterns in the data.
- **Generalization:** The model demonstrated reasonable generalization performance, although further testing and refinement might be necessary.