

ALTERNATIVE ASSESSMENT 1

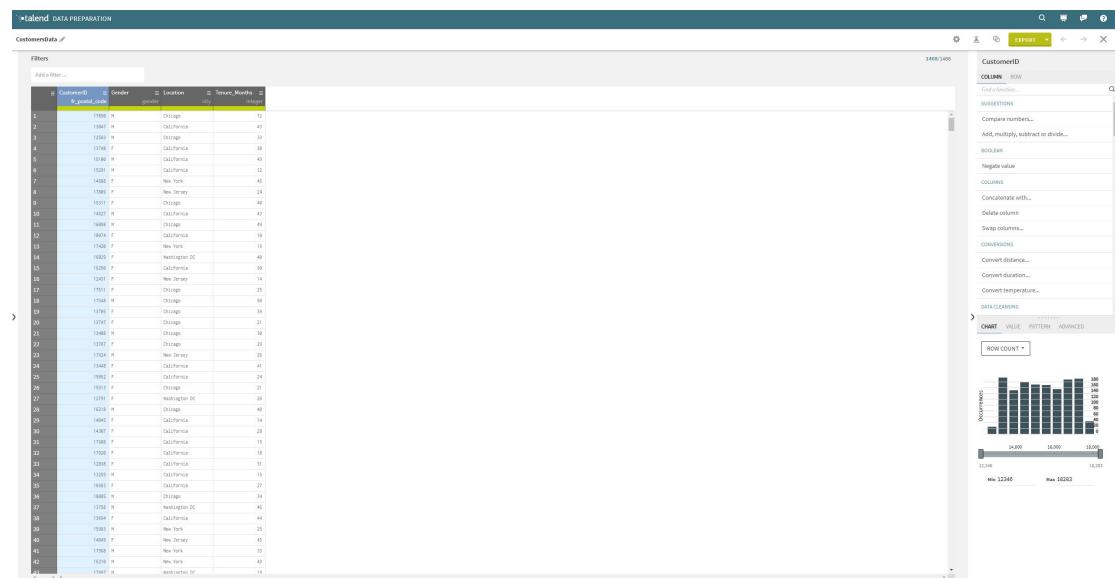
Lau Chyuan Jinn 22096017

The dataset is found from Kaggle, link:

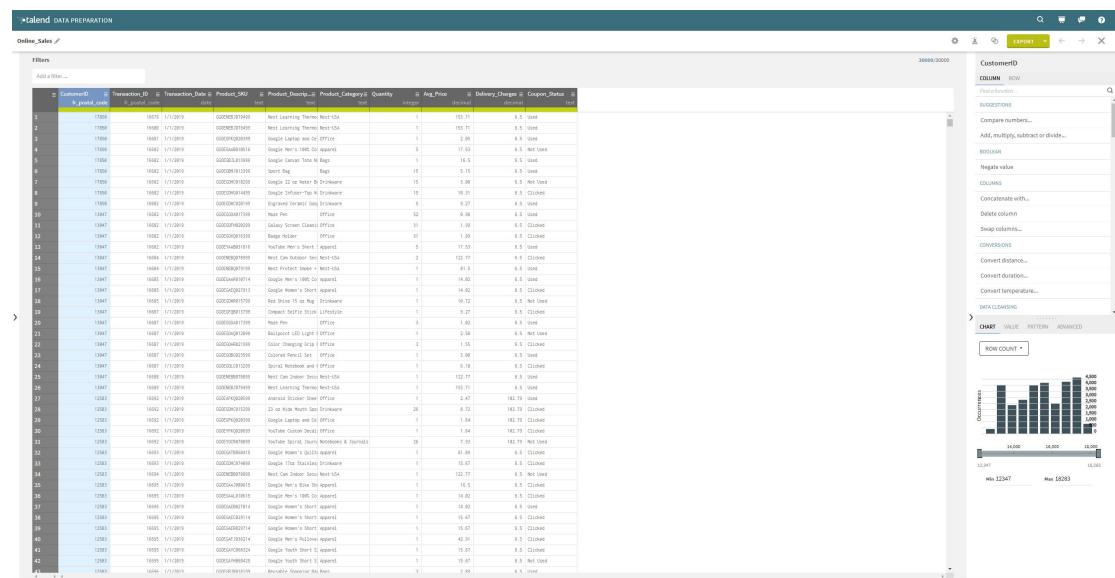
<https://www.kaggle.com/datasets/rishikumarrajvansh/marketing-insights-for-e-commerce-company/data>

Download the CustomersData.csv and Online_Sales.csv from Kaggle.

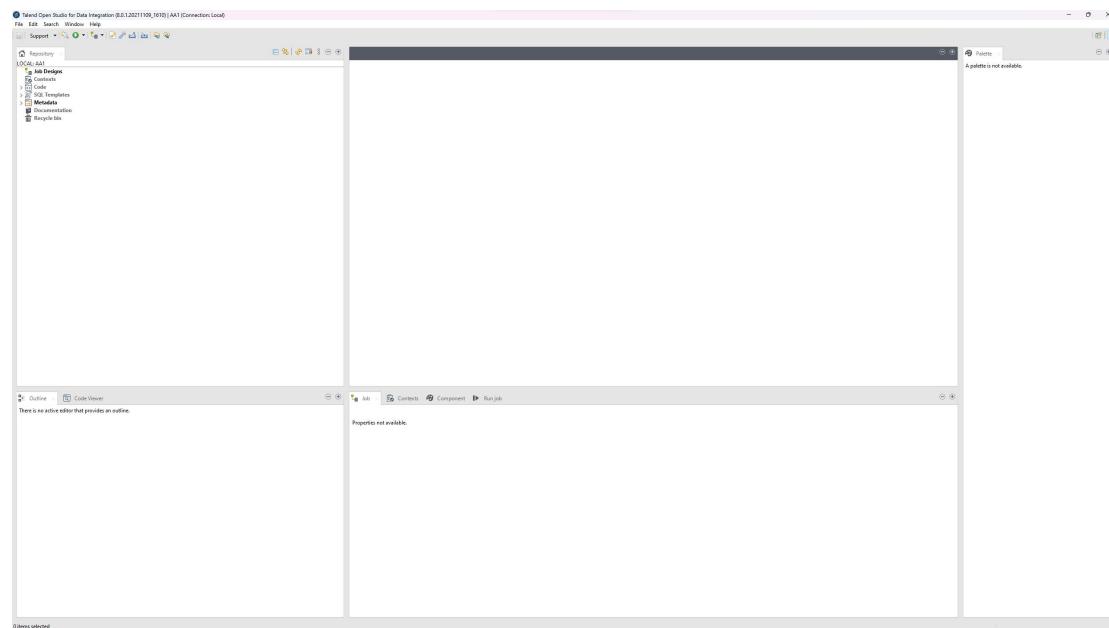
1. Using Talend Preparation to filter the CustomersData.csv, remove the null value row, it is the process for data cleaning.



2. Also same to Online_Sales.csv. After that, Export both csv data.



3. Using Talend Open Studio for Data Integration for next step.



4. Import both csv files in Delimited File.

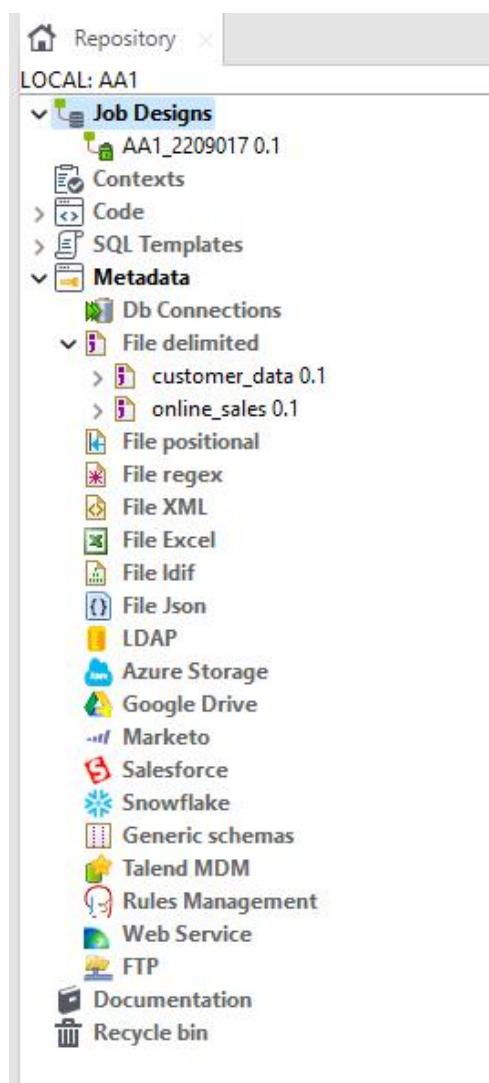
The screenshot shows the "File - Step 3 of 4" configuration dialog for a delimited file import. The dialog is divided into several sections:

- File Settings:** Encoding is set to "US-ASCII", Field Separator is "Comma", Row Separator is "Standard EOL", and Corresponding Character is "\r\n".
- Rows To Skip:** Header is set to 1, Footer is unchecked, and there is an option to "Skip empty row".
- Escape Char Settings:** CSV is selected, and the escape character is set to "Empty".
- Limit Of Rows:** There is an input field for "Limit" which is currently empty.
- Preview:** This section shows a preview of the data with 6 rows. The columns are CustomerID, Gender, Location, and Tenure_Months. The data is as follows:

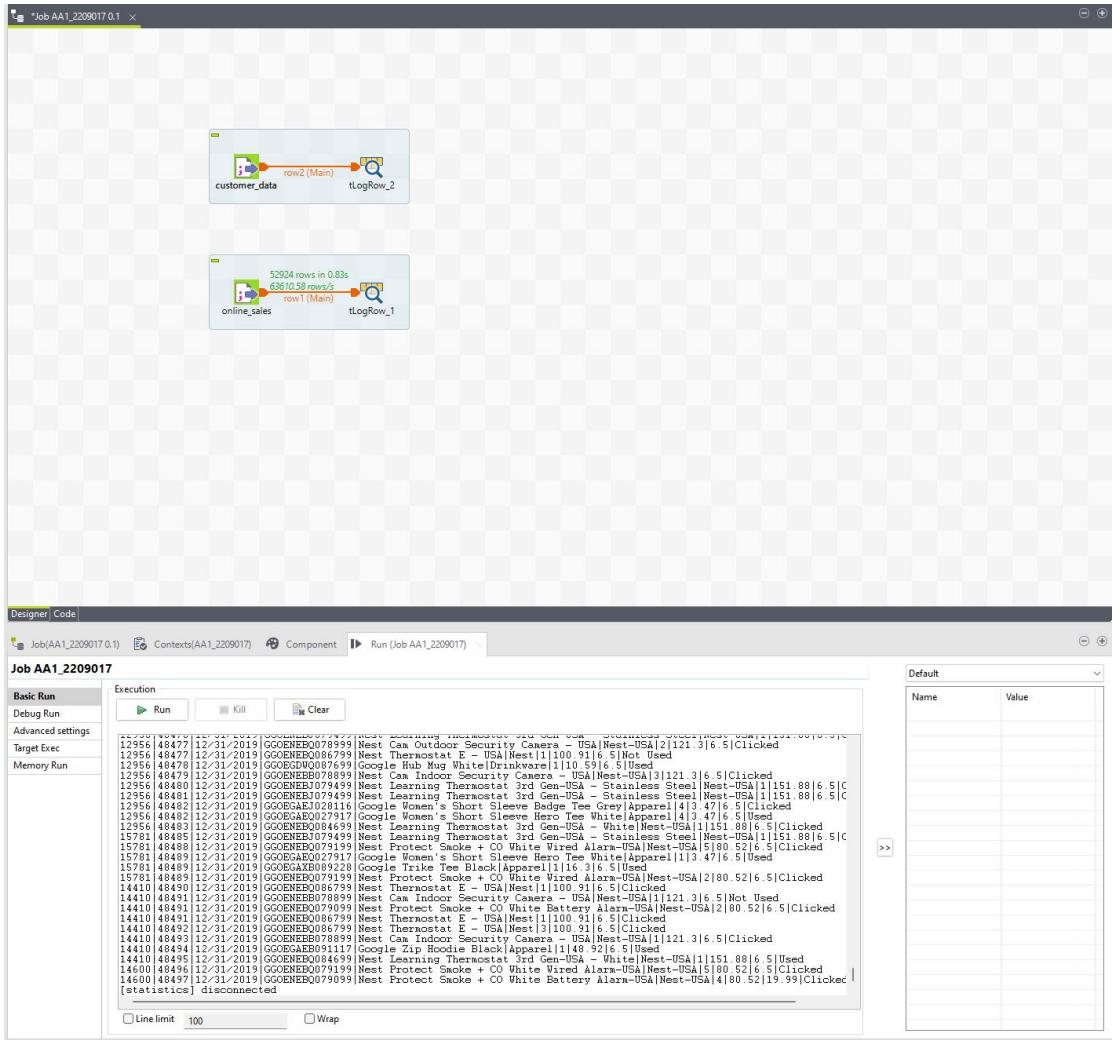
CustomerID	Gender	Location	Tenure_Months
17850	M	Chicago	12
13047	M	California	43
12583	M	Chicago	33
13748	F	California	30
15100	M	California	49
15291	M	California	32

At the bottom of the dialog, there are buttons for "Export as context", "Revert Context", and navigation buttons: "< Back", "Next >", "Finish", and "Cancel".

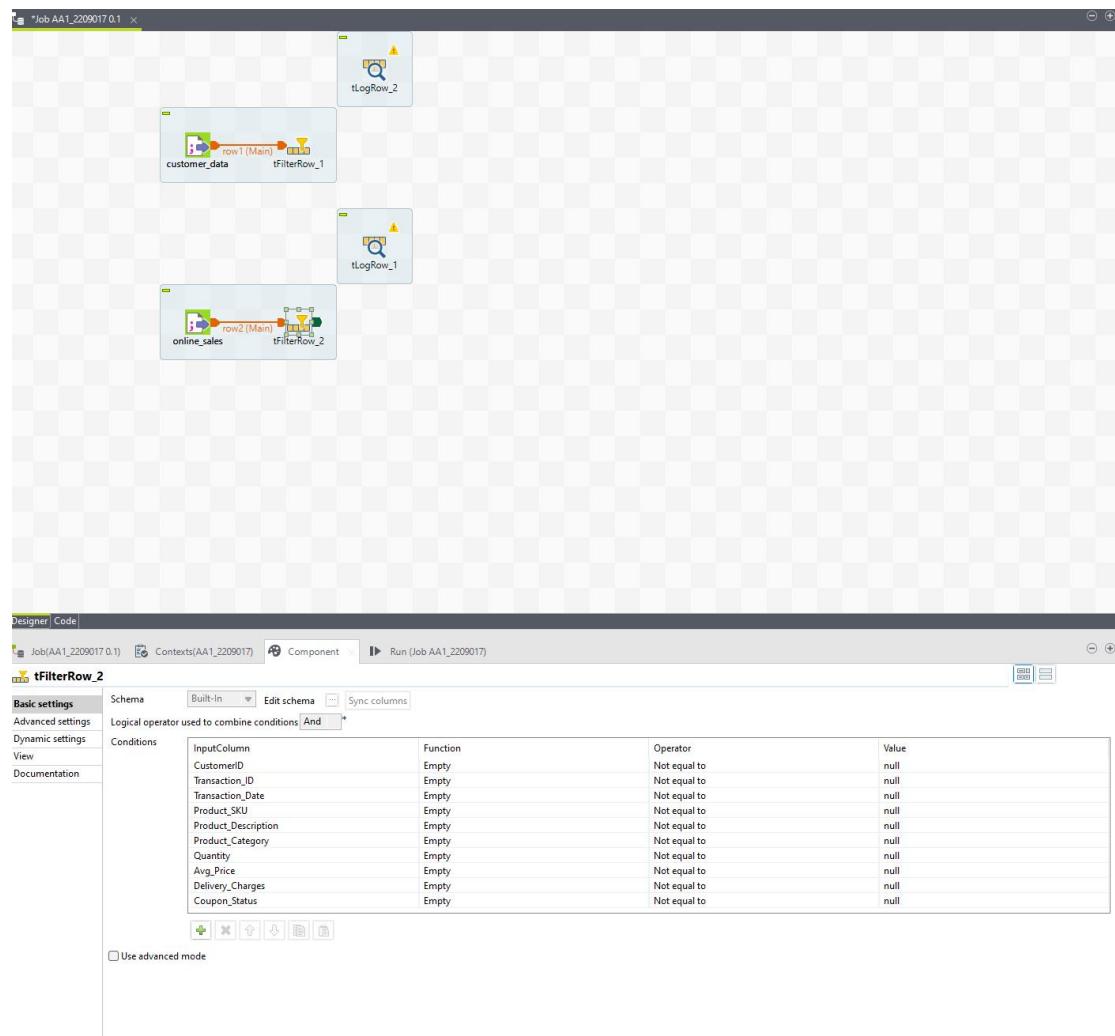
5. The left hand side, Repository will look like that.



6. Then drag customer_data 0.1 and online_sales0.1 to Job diagram. Next connect with tLogRow to view the data whether success display.

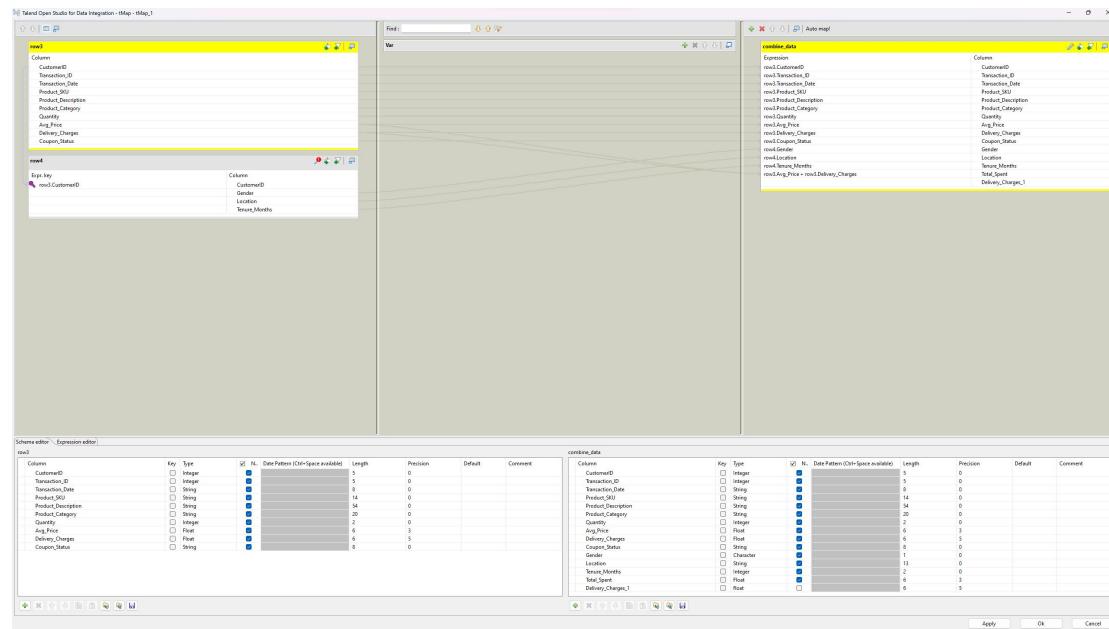


7. Using tFilterRow to filter the null value for both tables.

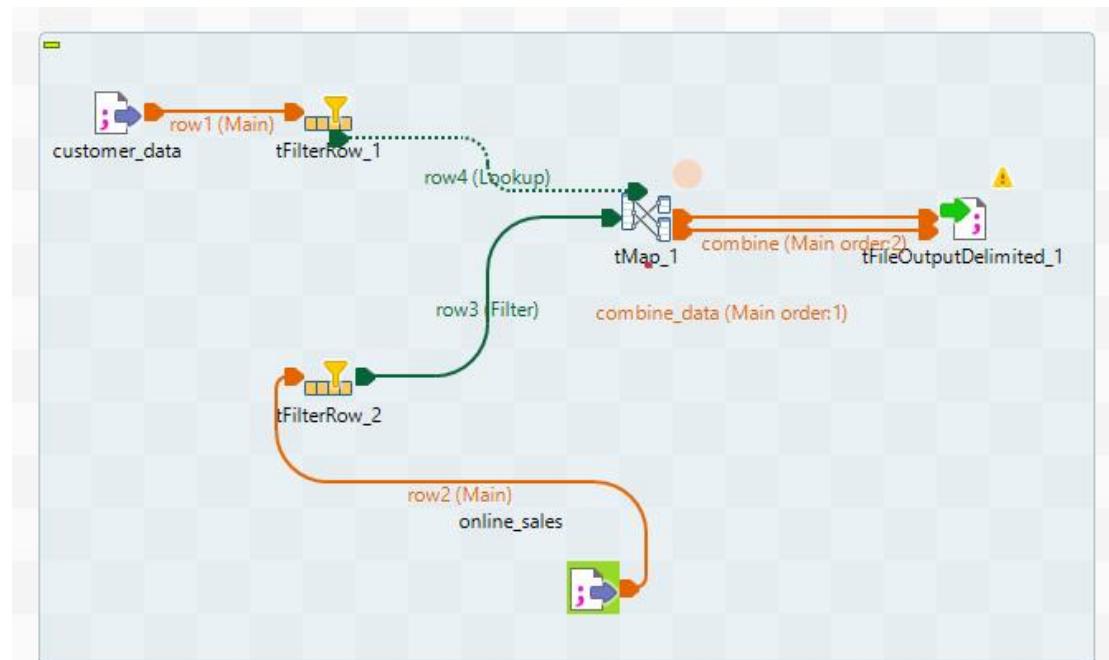


8. Using tMap to join both data, result to combine_data.csv

9. Set CustomerID as expression key.



10. Finally export the data using tFileOutputDelimited. The diagram shown as below.



11. Since combine_data.csv still haven't churn variable, which indicates whether the customer has stopped purchasing. Will process to next step to create churn variable.

12. Import combine_data.csv in Google Cloud, then using tool in Google Cloud called BigQuery to apply SQL.

The screenshot shows the Google Cloud Explorer interface. At the top, it says "Google Cloud" and "My Project 85865". Below that is a "Sandbox" button and a message "Set up billing to upgrade to the full E". The main area is titled "Explorer" with a "+ ADD" button and a "Home" icon. A search bar says "Type to search". Under "Viewing resources", there is a "SHOW STARRED ONLY" link. A tree view shows a project folder "bionic-comfort-410416" containing an "External connections" folder and a "online_sales" dataset. Inside "online_sales" is a table named "combine_data".

13. With the query below, will add a variable named churn which indicates whether the customer has stopped purchasing more than 30 days before the next year.

```
Untitled 3
RUN SAVE DOWNLOAD SHARE SCHEDULE MORE
1 WITH LastPurchaseDate AS (
2     SELECT
3         CustomerID,
4         MAX(Transaction_Date) AS Last_Transaction_Date
5     FROM
6         'bionic-comfort-410416.online_sales.combine_data'
7     GROUP BY
8         CustomerID
9 )
10 , churn AS (
11     SELECT
12         t.CustomerID,
13         t.Transaction_Date,
14         CASE
15             WHEN lpd.Last_Transaction_Date IS NULL THEN 1
16             WHEN lpd.Last_Transaction_Date + 30 < cast('2020-01-01' as date) THEN 1
17             ELSE -- Active
18         END AS Churn
19     END AS Churn
20 )
21 FROM
22     'bionic-comfort-410416.online_sales.combine_data' t
23 LEFT JOIN
24     LastPurchaseDate lpd ON t.CustomerID = lpd.CustomerID
25 )
26
27 SELECT t.CustomerID, t.Transaction_ID, t.Transaction_Date, Product_SKU, Product_Description, Product_Category, Quantity, Avg_Price, Delivery_Charges, Coupon_Status, Total_Spent, churn.Churn
28 FROM `mpg7005.online_sales.online_sales_data` t
29 LEFT JOIN churn
30 ON t.CustomerID = churn.CustomerID AND t.Transaction_Date = churn.Transaction_Date
```

14. The dataset is prepare enough and proceed to next step, which is analyse the dataset in SAS Enterprise Miner.

15. First, import the dataset combine_data.csv into SAS Enterprise Miner.

The screenshot shows the SAS Enterprise Miner interface. On the left, the 'File Import' node is selected in the diagram palette. The properties pane on the right shows the following configuration:

- General**: Node ID FIMPORT, Imported Data, Exported Data, Node Name sales.
- Train**: Variables, Import File D:\Acer\Documents\JMK\combine_data.csv, Maximum Rows to Import 1000000, Maximum Columns to Import 100000, Delimiter , Name Row Yes, Number of Rows to Skip 0, Guessing Rows 500, File Type Local, Advanced Advisor No, Rerun No.
- Score**: Role Train.
- Report**: Summarize No.
- Status**: Create Time 1/7/24 6:16 AM, Run ID c403d5d3c3c2d842afe1, Last Error, Last Status Complete, Last Run Time 1/7/24 6:20 AM, Run Duration 0 Hr. 0 Min. 2.90 Sec.
- Exported Data**: Set of tables exported by this node.

The 'Results - Node File Import Diagram: sales' window displays the contents of the imported dataset:

```

25   The CONTENTS Procedure
26
27   Data Set Name      FIMPORT.FIMPORT_DATA
28   Member Type       DATA
29   Engine            VS
30
31   Created          07/01/2024 06:20:51
32   Last Modified     07/01/2024 06:20:51
33   Protection        0
34   Data Set Type    0
35   Label
36   Data Representation SOLARIS_X86_64, LINUX_X86_64, ALPHA_T9064, LINUX_IA64
37   Encoding         wtf-8 Unicode (UTF-8)
38
39
40   Engine/Host Dependent Information
41
42   Data Set Page Size 131072
43   Number of Data Set Pages 52
44   First Data Page 1
45   Max Obs per Page 1022
46   Obs in First Data Page 992
47   Max Number of Data Set Pages 9
48   Filenames          /home/u03650100/AII/Workspaces/EMW31/fimport_data.sas7bdat
49   Release Created 9.40LIM7
50   Host Created Linux
51   Inode Number 222038702
52   Access Permission r-w-r--r-
53   Owner Name 065650100
54   Path /
55   File Size (bytes) 646616
56
57
58   Alphabetic List of Variables and Attributes
59
60   # Variable      Type Len Format Informat
61
62   8 Avg_Price     Num 8 BEST12. BEST12.
63   15 Churn        Num 8 BEST12. BEST12.
64   10 Coupon_Status Char 4 $4. $4.
65   1 CustomerID   Num 8 BEST12. BEST12.
66   9 Delivery_Charges Num 8 BEST12. BEST12.
67   11 Gender       Char 1 $1. $1.
68   12 Location     Char 13 $13. $13.
69   6 Product_Catogory Char 9 $9. $9.
70   5 Product_Description Char 14 $14. $14.
71   4 Product_SKU   Char 14 $14. $14.
72   7 Quantity     Num 8 BEST12. BEST12.
73   13 Tenure_Months Num 8 BEST12. BEST12.
74   14 Total_Spent  Num 8 BEST12. BEST12.
75   3 Transaction_Date Num 8 MMDDYY10. MMDDYY10.
76   2 Transaction_ID Num 8 BEST12. BEST12.
77
78
79  ****
80  * Score Output *
81

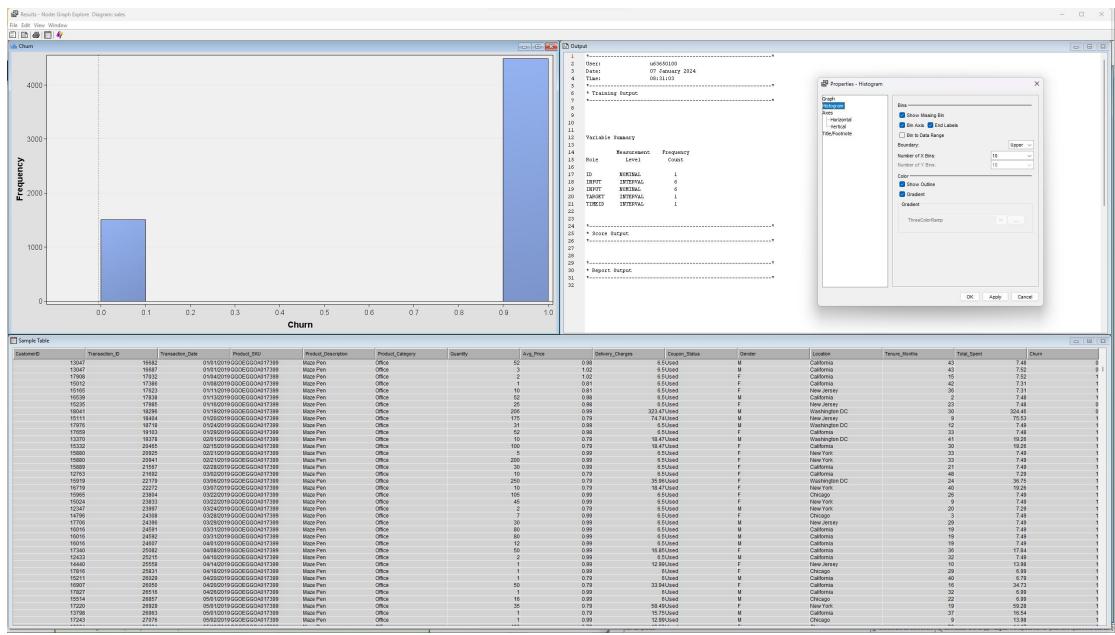
```

16. Next, set the churn as target variable.

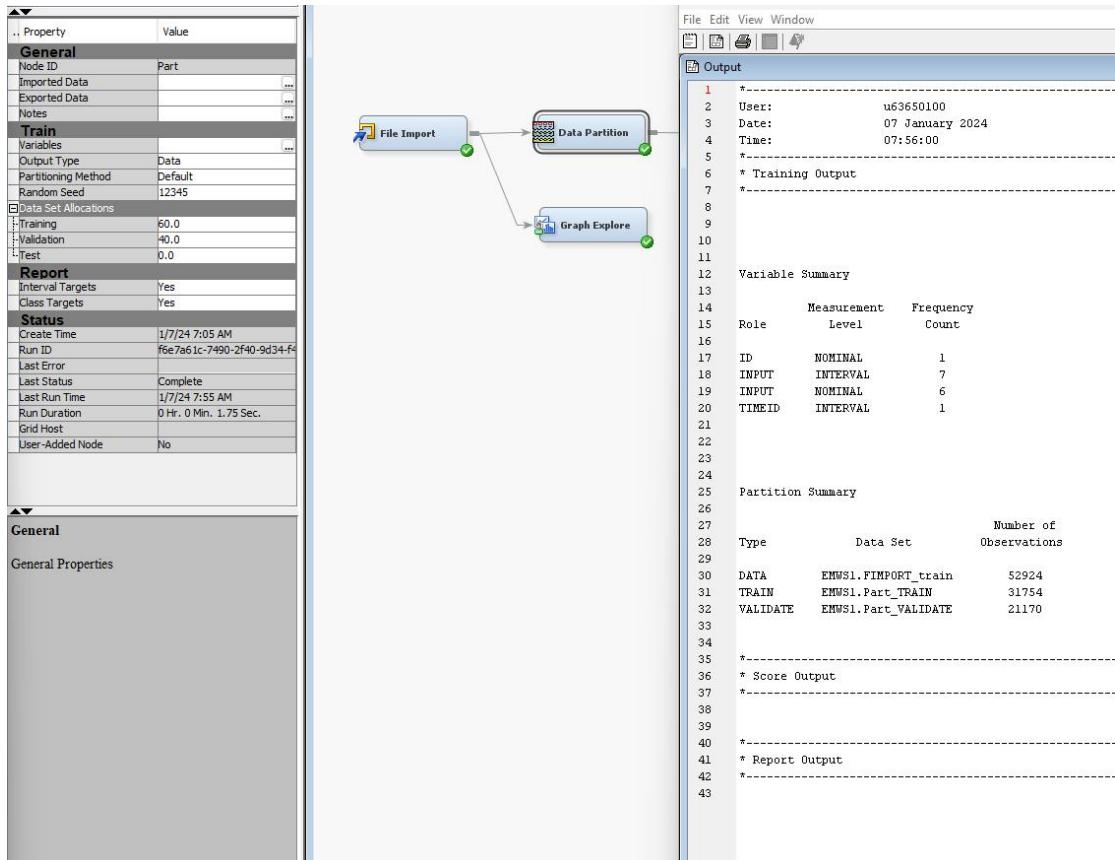
The screenshot shows the 'Variables - FIMPORT' table. The 'Churn' variable is highlighted in blue, indicating it is the target variable. The table includes columns for Name, Role, Level, Report, Order, Drop, Lower Limit, and Upper Limit.

Name	Role	Level	Report	Order	Drop	Lower Limit	Upper Limit
Avg_Price	Input	Interval	No		No	*	*
Churn	Target	Interval	No		No	*	*
Coupon_Status	Input	Nominal	No		No	*	*
CustomerID	Input	Interval	No		No	*	*
Delivery_Charge	Input	Interval	No		No	*	*
Gender	Input	Nominal	No		No	*	*
Location	Input	Nominal	No		No	*	*
Product_Catogory	Input	Nominal	No		No	*	*
Product_Descrip	Input	Nominal	No		No	*	*
Product_SKU	Input	Nominal	No		No	*	*
Quantity	Input	Interval	No		No	*	*
Tenure_Months	Input	Interval	No		No	*	*
Total_Spent	Input	Interval	No		No	*	*
Transaction_Dat	Time ID	Interval	No		No	*	*
Transaction_ID	ID	Nominal	No		No	*	*

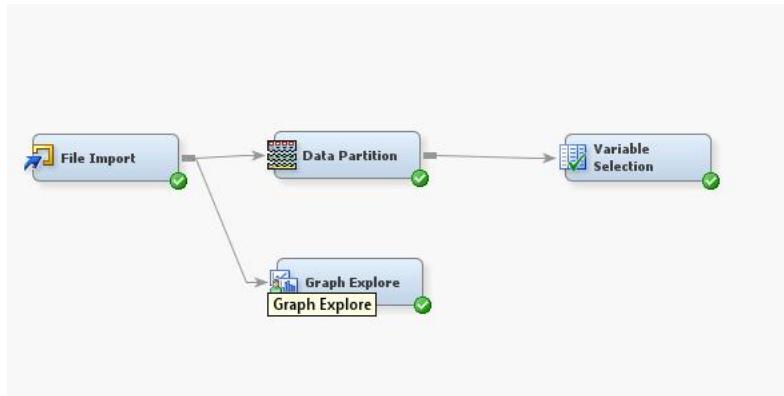
17. With histogram, check the dataset whether have null value. (From diagram shown is 0).



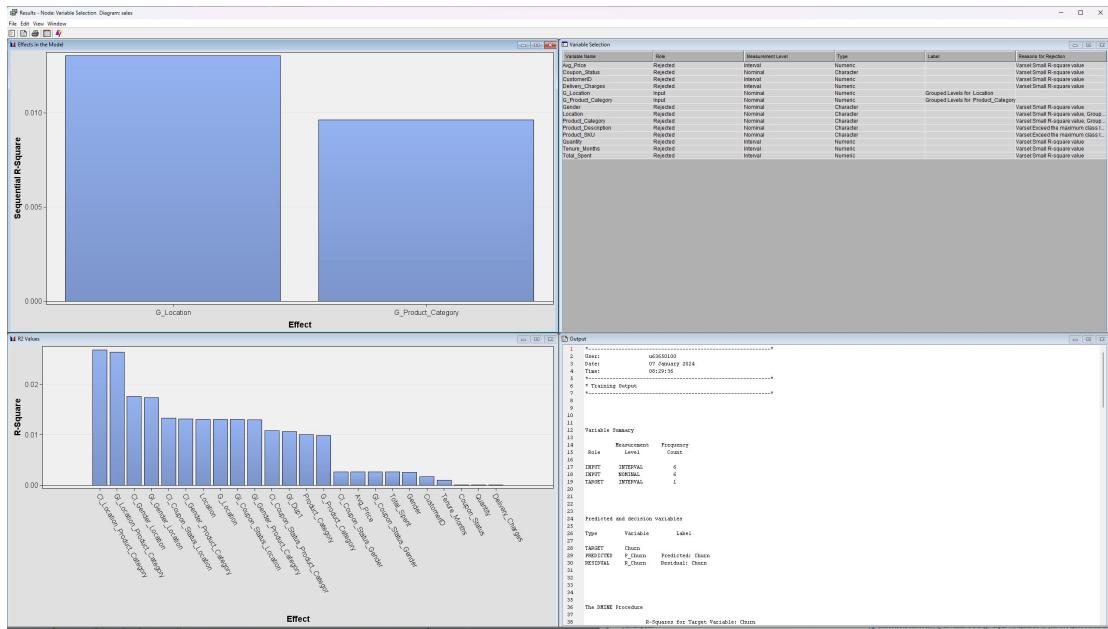
18. Add a Data Partition node, in SAS Enterprise Miner, the Data Partition node is used to divide a dataset into training, validation, and testing sets. It helps in building, fine-tuning, and evaluating predictive models by ensuring proper separation of data for these purposes. The node supports randomization, stratified sampling, and may allow for k-fold cross-validation. Its primary role is in assessing model performance on unseen data and is a crucial step in the model development process.



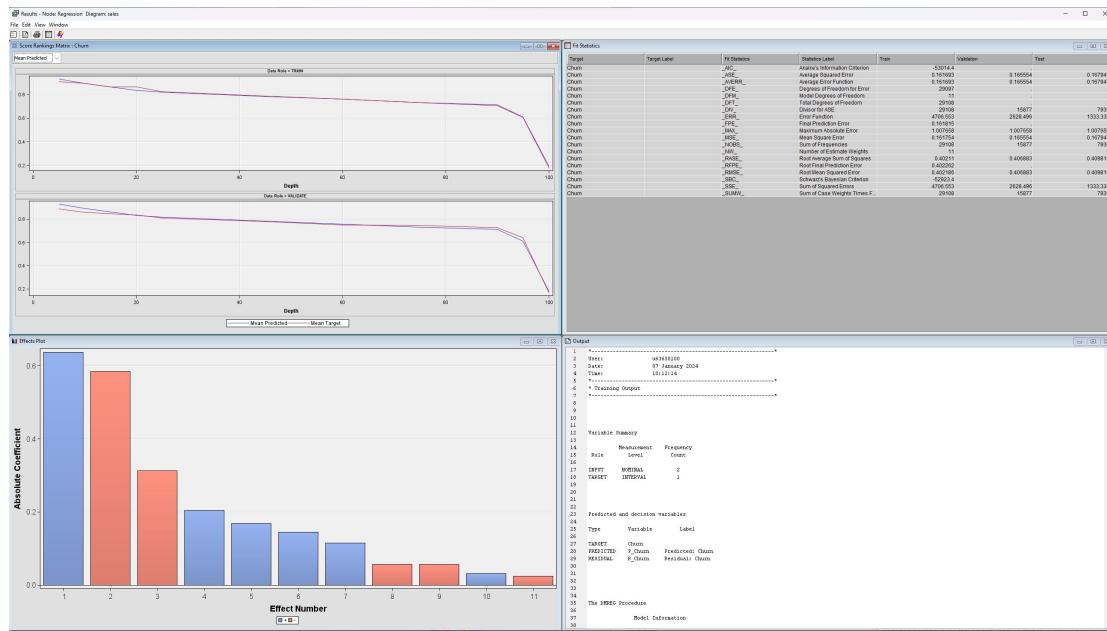
19. Then add a Variable Selection node, Variable selection is a crucial step in building predictive models as it helps identify the most relevant features or variables that contribute to the model's performance.



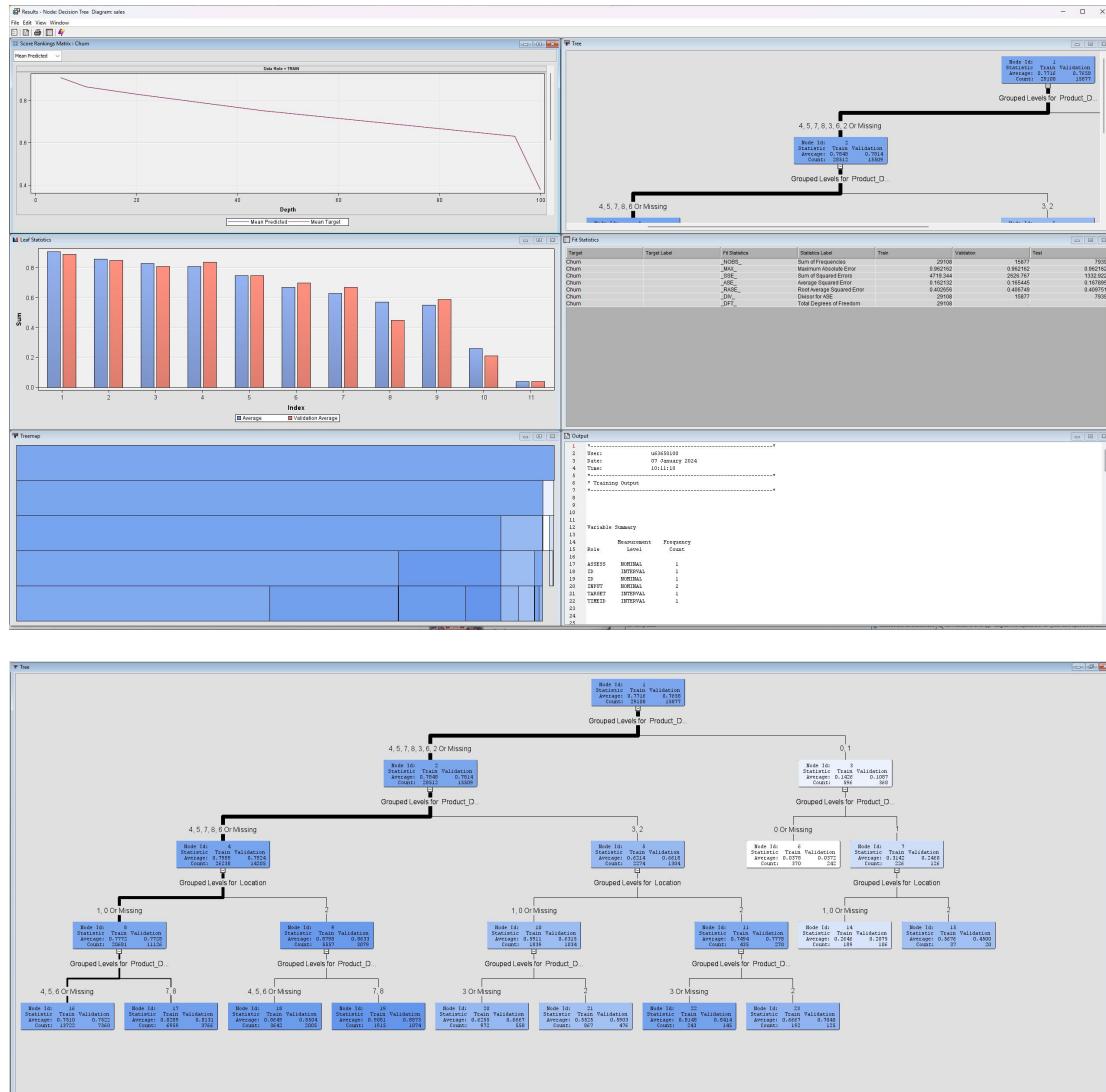
20. The result of Variable Selection shown as below.



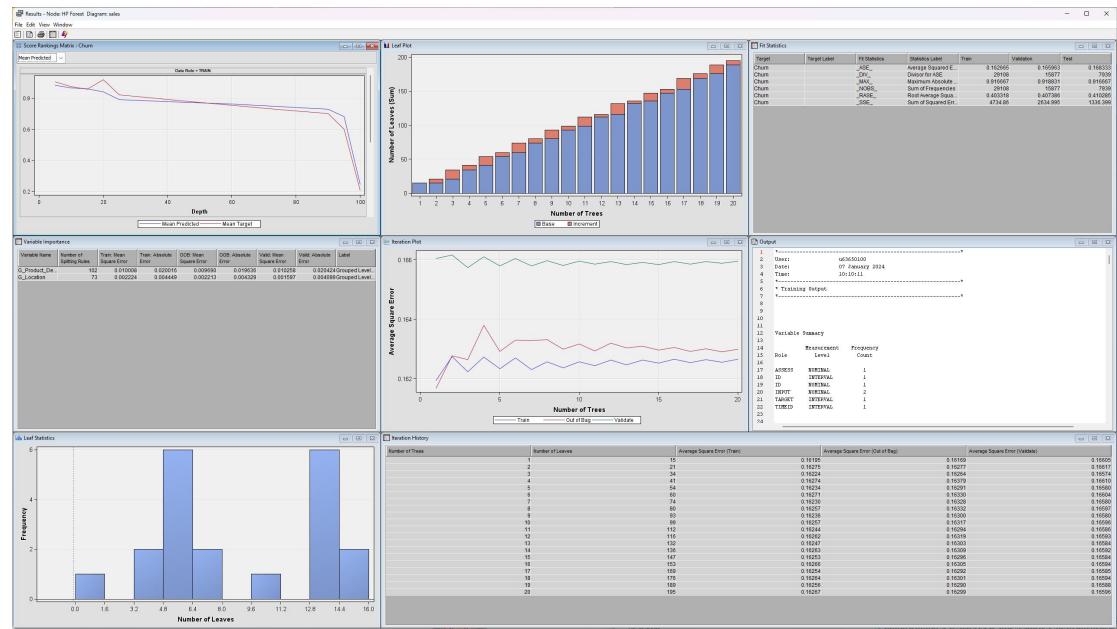
21. Add a Regression node, the result will shown as below. The Regression node is used for building regression models. The regression analysis helps predict a continuous target variable based on one or more predictor variables. The top left section displays two line graphs side by side, showing some trends or data over time. The X-axis is labeled "Days" and the Y-axis has numerical values, indicating it's a plot of some variables over days. In the top right section, there's a table with multiple columns including "Name," "Mean," "Std Dev," "Minimum," etc., suggesting statistical data of different variables. In the bottom left section, there's a bar graph labeled "Effect Number" on the X-axis and "Value" on the Y-axis. It displays five bars with varying heights. The bottom right section contains text information about model summary including factors like 'A', 'B', 'AB', etc., and their levels. There are also details about random seed and total runs.



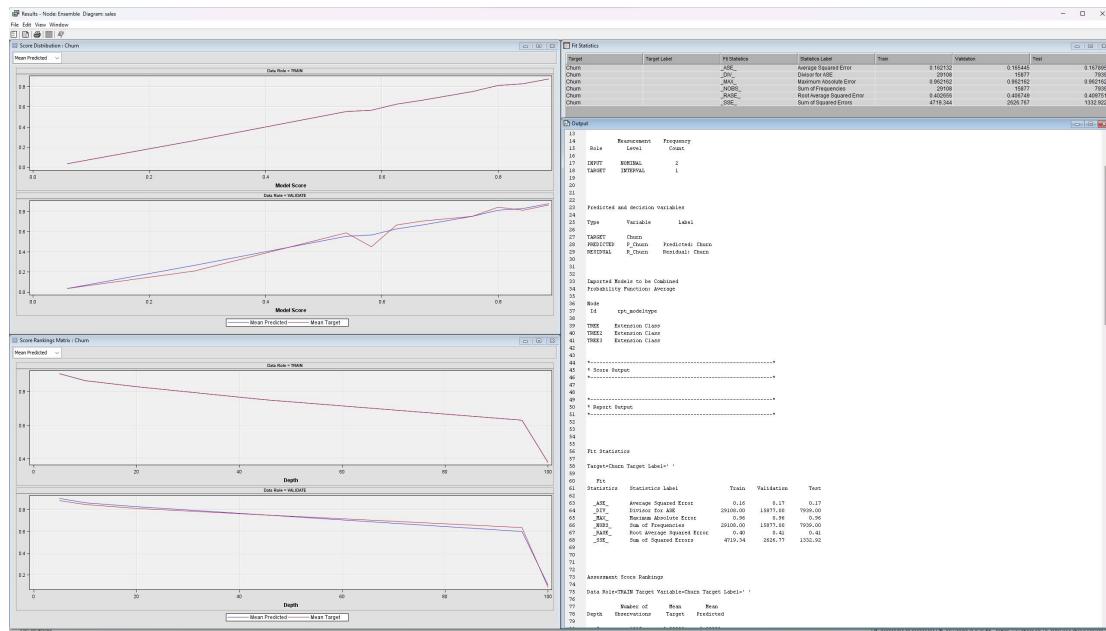
22. Next, add a Decision Tree node, the result will shown as below. The Decision Tree node is used to build decision tree models. Decision trees are a type of predictive modeling technique that recursively splits data into subsets based on the values of predictor variables, ultimately creating a tree-like structure. The image shows a complex flowchart with multiple blue boxes connected by arrows. Each box contains text and numbers, indicating steps or processes in data handling. There are different sections labeled as "0 Missing", "0-3 Missing", "4,5,7,8 < 0 Missing" etc., suggesting categories based on missing data. Inside the boxes are parameters like "Show Size", "Hide Size", "Validation" and numerical values indicating some form of measurement or evaluation criteria. Arrows connect these boxes indicating the flow of processes or steps from one to another. The background is plain white making the blue boxes and black text prominent.



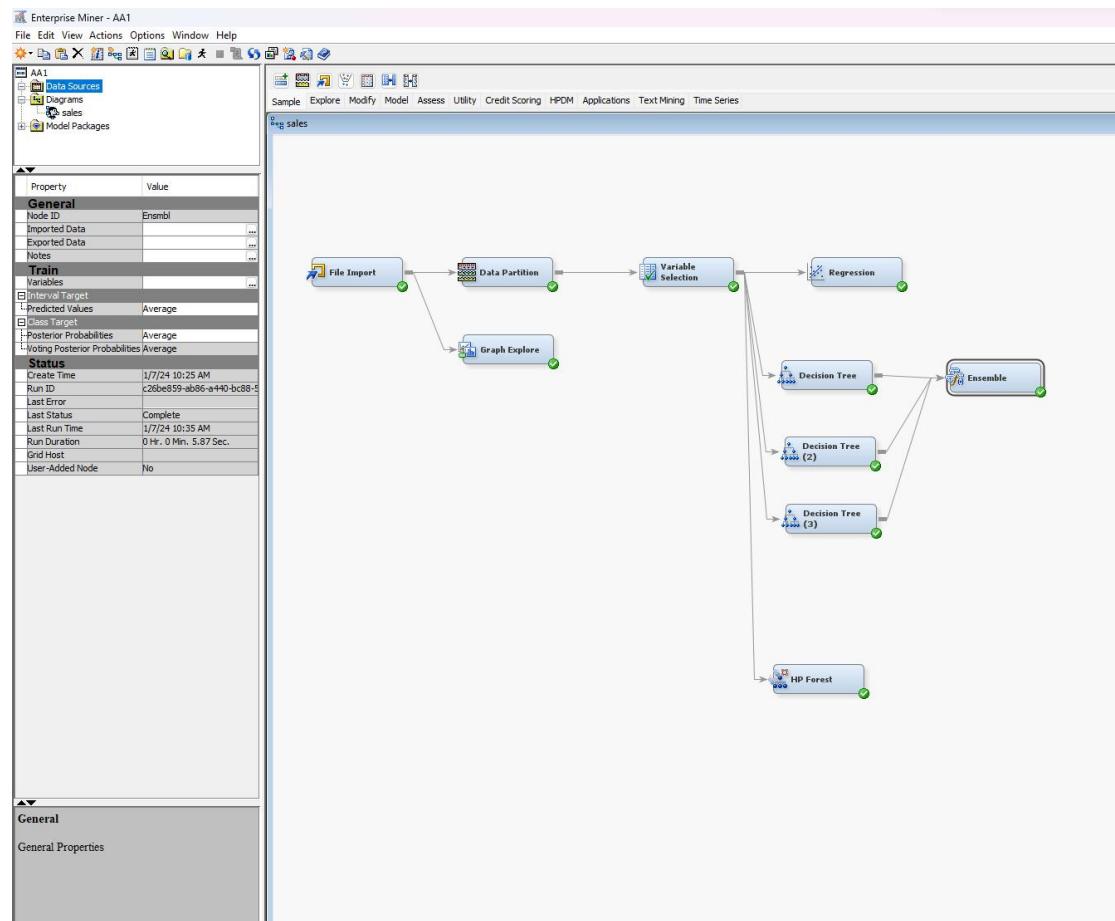
23. Next, add HP Forest, thee result as shown as below. The HP Forest Node is a predictive modeling tool in SAS Enterprise Miner that creates a forest predictive model, which is an ensemble of many decision trees, using the HPFOREST procedure. The top left panel displays a line graph labeled “Run Time vs Type Spectrum” showing two lines that appear to represent different datasets. Adjacent to this, on the top middle panel, there is a bar graph labeled “Number of Types” with bars in two shades representing different categories or datasets. On the top right, there are two small panels; one appears blank and the other displays textual information including values labeled “TOTAL TYPES”, “TOTAL METHODS”, etc. Below these panels in the middle row left side, there’s another line graph labeled “Number of Types”, showing three lines representing different datasets over time. In the middle row right side, there’s a text box containing specific numerical data and labels such as “TOTAL: RANGE” and individual numerical values associated with specific labels like ‘min’, ‘max’, etc. On the bottom left panel, there’s another bar graph labeled “Number of Classes” displaying blue bars at various heights representing different quantities or values. Adjacent to this on the bottom right are two panels displaying tabulated data with rows and columns containing numerical values and text.



24. Next, add 2 more Decision Tree nodes with different leaf (3 and 8). Then connect 3 Decision Tree nodes to a Ensemble node, result is shown as below. Ensemble models combine predictions from multiple individual models to improve overall predictive performance. The left side of the interface shows three graphs plotted, each in its own panel. The top graph is labeled "Big Arm Position (Degrees)" with X-axis labeled as "Time" and Y-axis ranging from 0 to 270. The middle graph is not clearly labeled but has the same X-axis label "Time" and Y-axis labels similar to the top graph. The bottom graph is labeled "Small Arm Velocity (Deg/Sec)" with X-axis labeled as "Time" and Y-axis ranging from 0 to 1000. Each graph displays two lines plotted over time, one in red and another in purple, representing different sets of data or conditions. On the right side, there's a panel displaying numerical data and parameters under different categories like "Parameters," "Display," "Model & Physical Signals," etc. It includes values for parameters like mass, length, inertia along with other settings and configurations for the displayed graphs or related computations. There are also buttons like "Load Config" indicating that users can load different configurations for analysis.



25. The SAS Enterprise Miner diagram will shown as below.



Conclusion:

In conclusion, the comprehensive analytical process undertaken, from data cleaning and integration using Talend to advanced modeling in SAS Enterprise Miner, provides a robust framework for understanding and predicting customer behavior. The decision tree and ensemble models offer valuable insights, particularly in identifying potential churn patterns. The strategic integration of Google Cloud for SQL queries adds depth to the dataset, enhancing the accuracy of predictive models. Moving forward, businesses can leverage these insights to implement targeted customer retention strategies and personalized marketing approaches, thereby improving overall customer satisfaction and optimizing marketing efforts.