



**UNIVERSITI  
MALAYA**

**WQD7005 DATA MINING  
2023/2024 S1**

**ALTERNATIVE ASSESSMENT 1**

Matric Number:	S2164046
Name:	HU LIANGLIANG

## Case Study: E-Commerce Customer Behaviour Analysis

### Tasks Data Import and Preprocessing:

- Import your dataset into SAS Enterprise Miner, handle missing values, and specify variable roles. [15 marks]
- Decision Tree Analysis: Create a decision tree model in SAS Enterprise Miner to analyse customer behaviour. [20 marks]
- Ensemble Methods: Apply Bagging and Boosting, using the Random Forest algorithm as a Bagging example. [10 marks]

### Deliverables:

- A report detailing each step of the process, including the rationale behind your choices and any challenges faced.
- An analysis of the decision tree and ensemble methods, with insights into customer behavior and suggestions for business strategy. [5 marks]

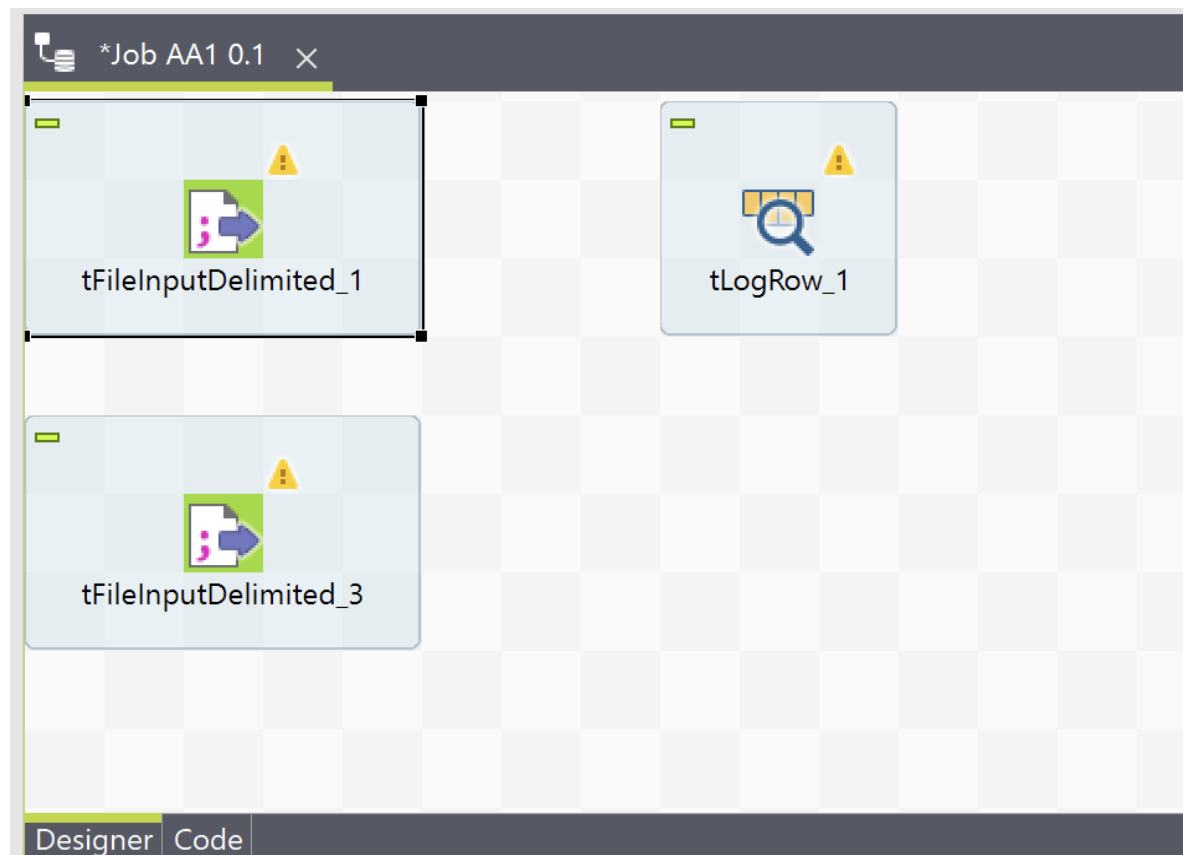
I will sequentially process the dataset in the following order to conduct effective analysis of e-commerce customer behavior. I will use Data Integration (DI) to integrate the entire dataset, employ Data Preparation (DP) for initial transformation and exploration of the data, and finally utilize SAS EM to handle tasks such as imputing missing values and specifying variable roles. This comprehensive approach is aimed at facilitating subsequent decision tree analysis and employing various ensemble methods.

I randomly generated a dataset according to different rules, aiming to produce results that meet the task requirements and closely resemble real-life scenarios.

**First, let's preview the datasets. There are two distinct datasets available—one containing detailed customer personal information and the other containing customer purchase information. To proceed, we need to merge these two disparate datasets.**

1	CustomerID	Age	Gender	Location	Membersh	85751	41	3561.19	Home Goc	Profession	9	0	#####	78
2	86667	49	Female	Shenzhen	Bronze	85705	46	3842.98	Books	Profession	6	0	#####	87
3	86533	21	Male	Guangzho	Gold	86155	10	440.58	Clothing	Unemploy	7	1	2023/2/4	331
4	85972	24	Female	Hangzhou	Bronze	85984	23	1854.61	Sports	Unemploy	2	0	2023/9/2	121
5	86227	64	Male	Chongqing	Bronze	85985	42	5119.77	Home Goc	Self-Empc	9	0	2023/9/4	119
6	86346	23	Male	Chongqing	Bronze	86557	87	35488.83	Home Goc	Self-Empc	18	0	#####	127
7	86645	32	Male	Wuhan	Silver	86462	26	2751.18	Electronics	Self-Empc	29	0	2023/7/2	183
8	85701	65	Female	Beijing	Gold	85719	38	1680.13	Electronics	Unemploy	4	0	#####	87
9	85905	38	Female	Wuhan	Bronze	85777	93	17380.37	Electronics	Self-Empc	9	1	#####	277
10	85980	45	Male	Shenzhen	Bronze	86434	52	7977.06	Clothing	Self-Empc	16	0	#####	133

### 1. Data Integration (DI)

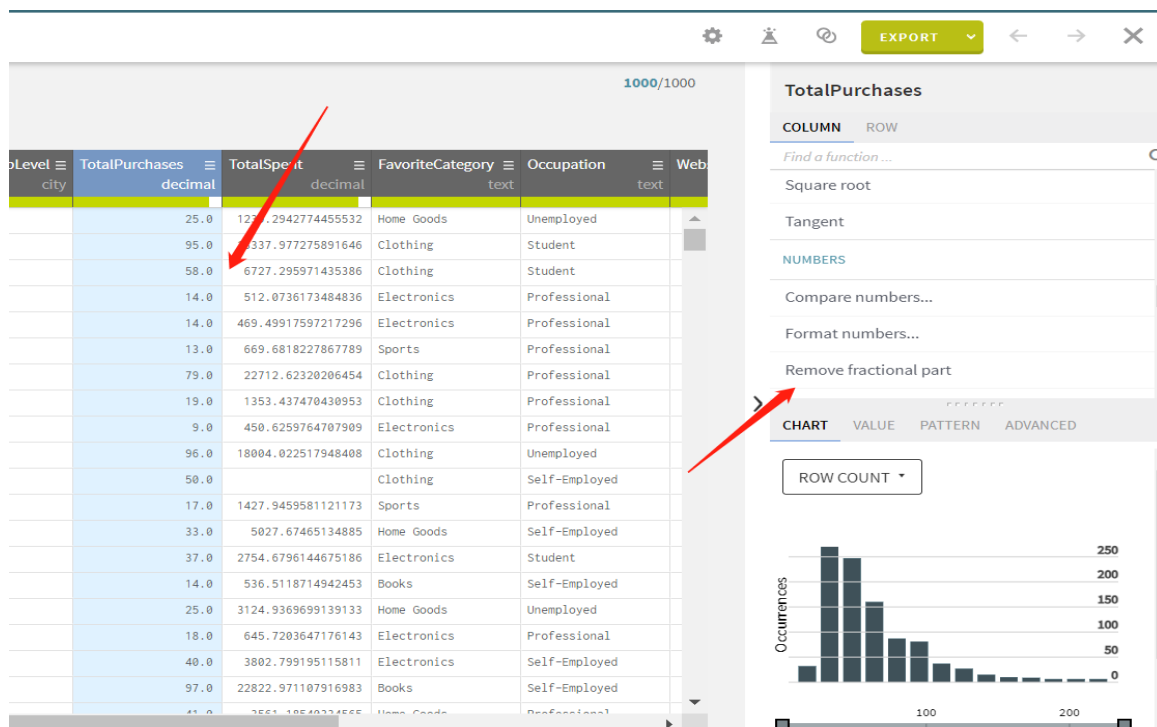


Connecting two different datasets

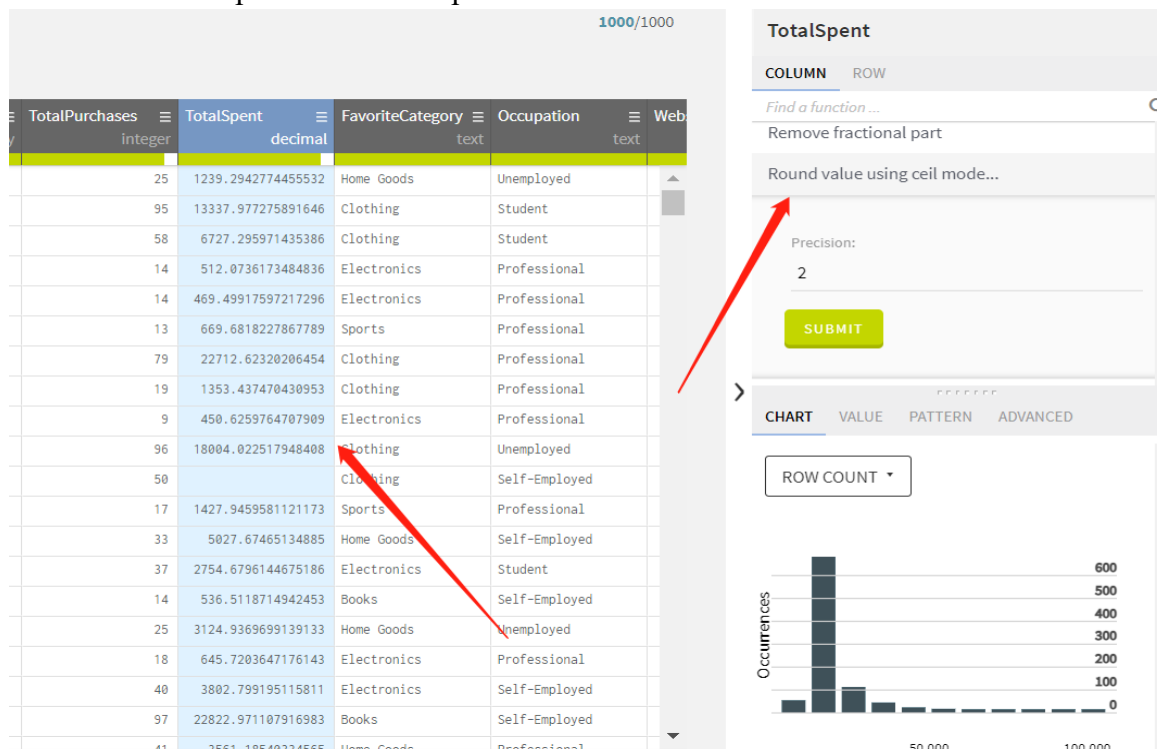
## 2. Data Preparation (DP)

After merging the distinct datasets, we need to conduct initial exploration and processing of the data.

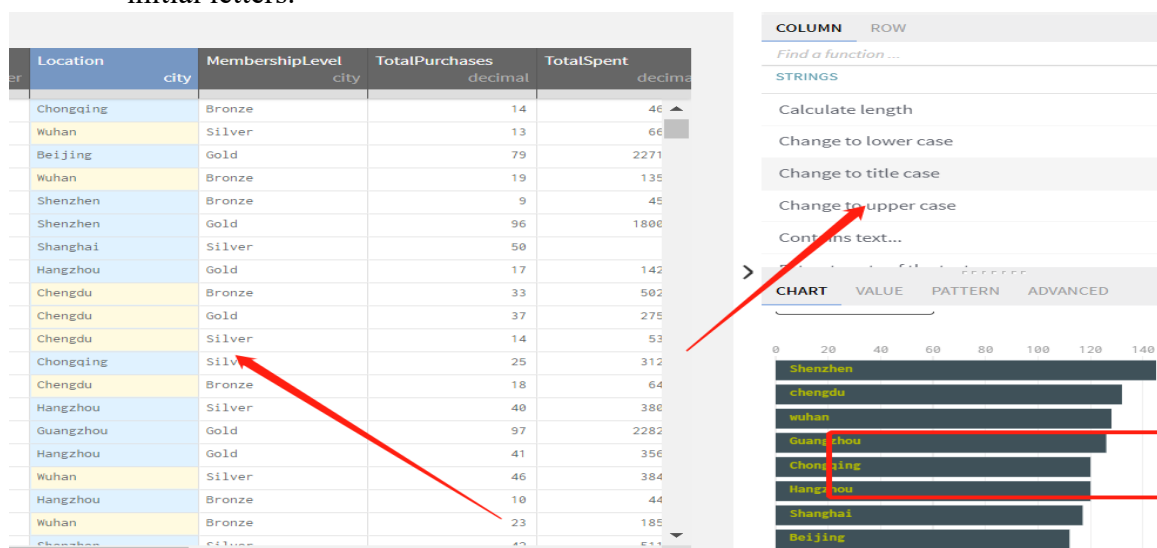
- 2.1 It is not reasonable for TotalPurchases to have decimal points, so I removed the decimal points from TotalPurchases.



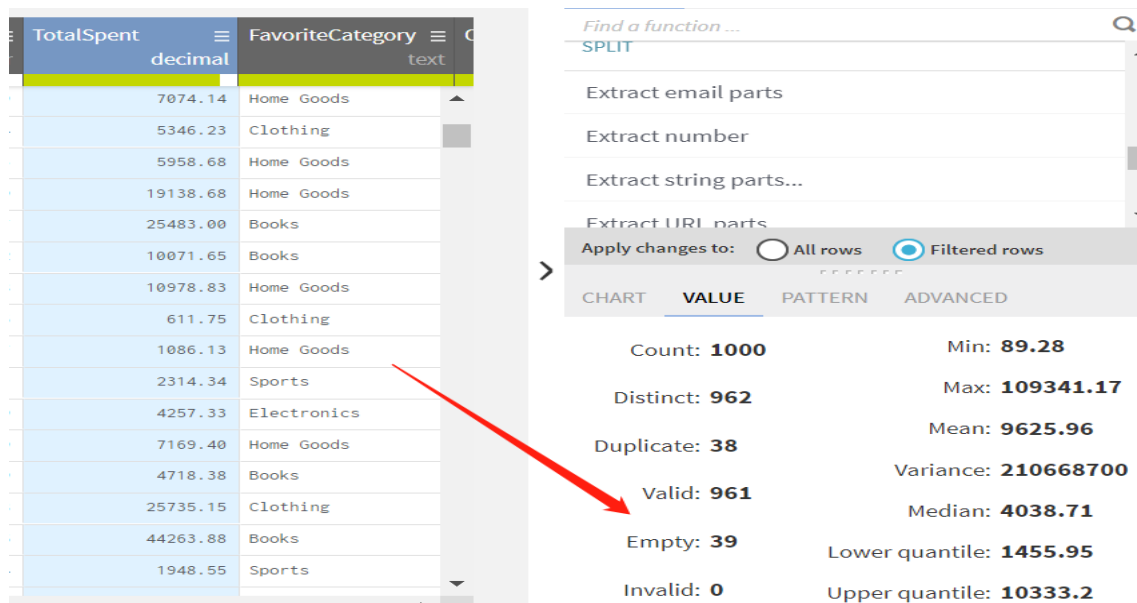
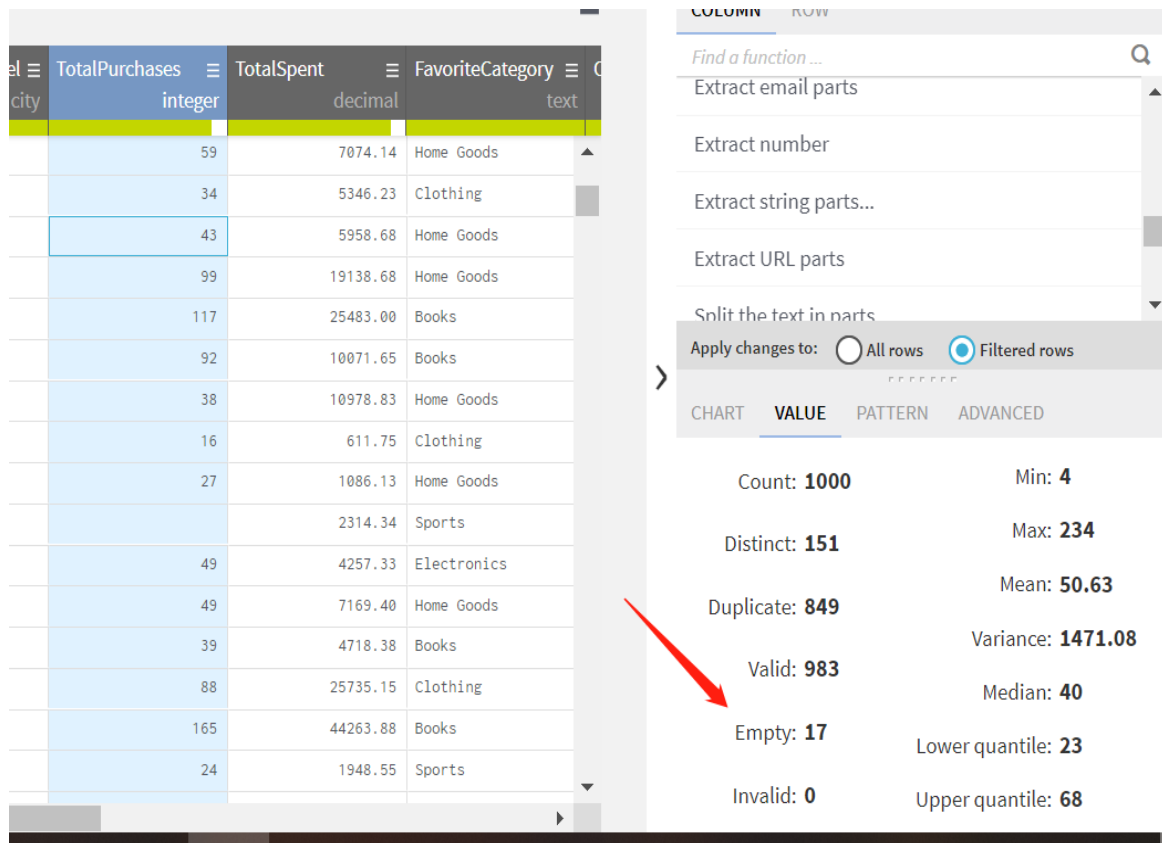
- 2.2 Many decimal places in TotalSpent are not reasonable, so I retained two decimal places for TotalSpent.



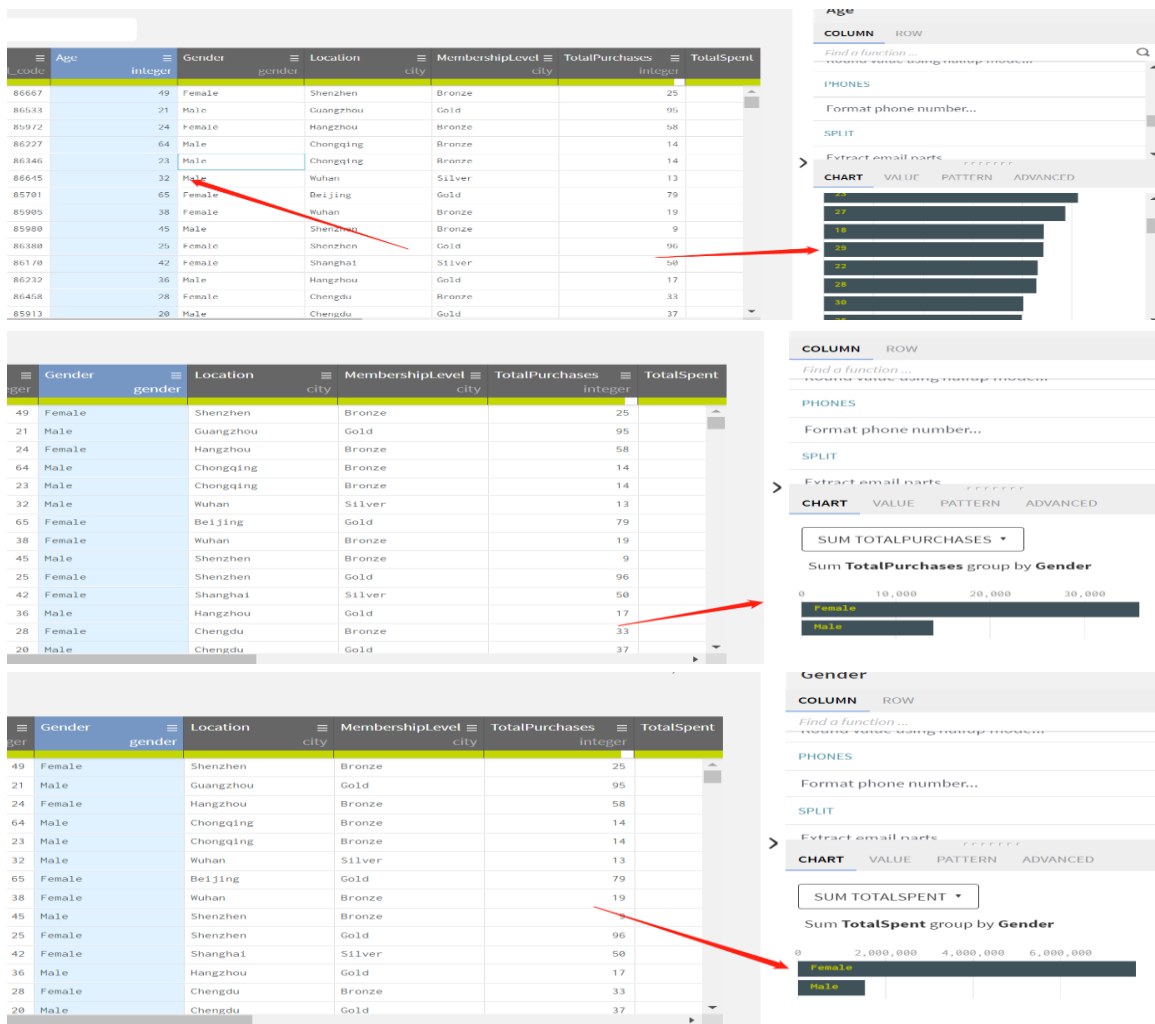
- 2.3 The capitalization of city names in the location data is inconsistent, with some having uppercase initial letters and others having lowercase initial letters. This is not reasonable, so I have standardized all city names to have uppercase initial letters.



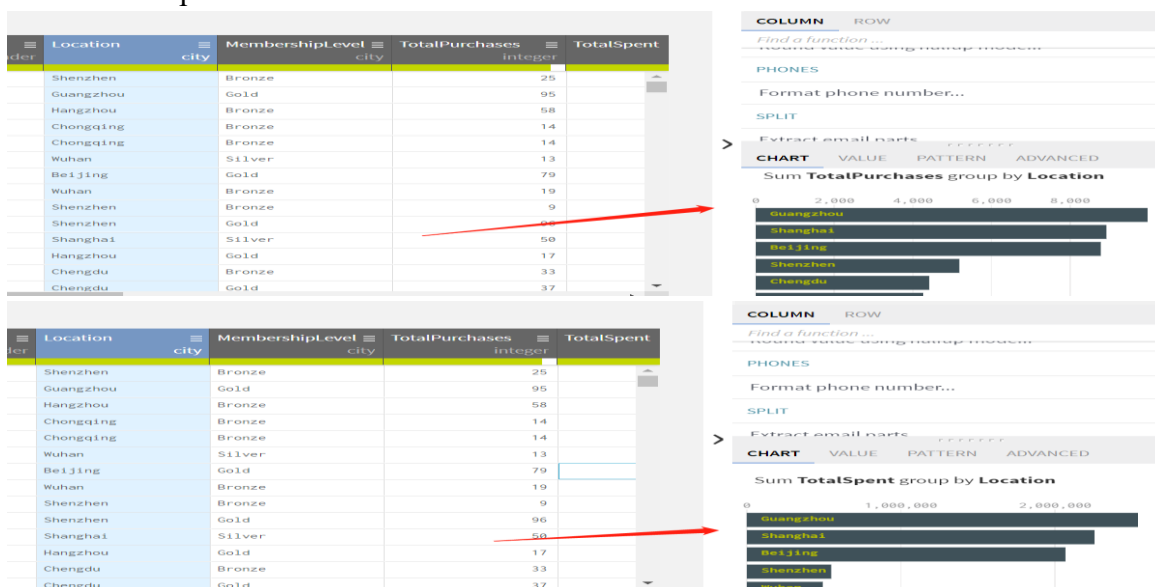
- 2.4 I observed that both TotalPurchases and TotalSpent have missing values, with 17 and 39 missing values respectively. I won't handle them at this point; instead, I will examine the overall structure of the data inside SAS EM before proceeding with any imputation or treatment.



- 2.5 Observing this dataset, it appears that the majority of individuals are young, and the spending capacity of females is significantly higher than that of males. This aligns well with real-world scenarios.

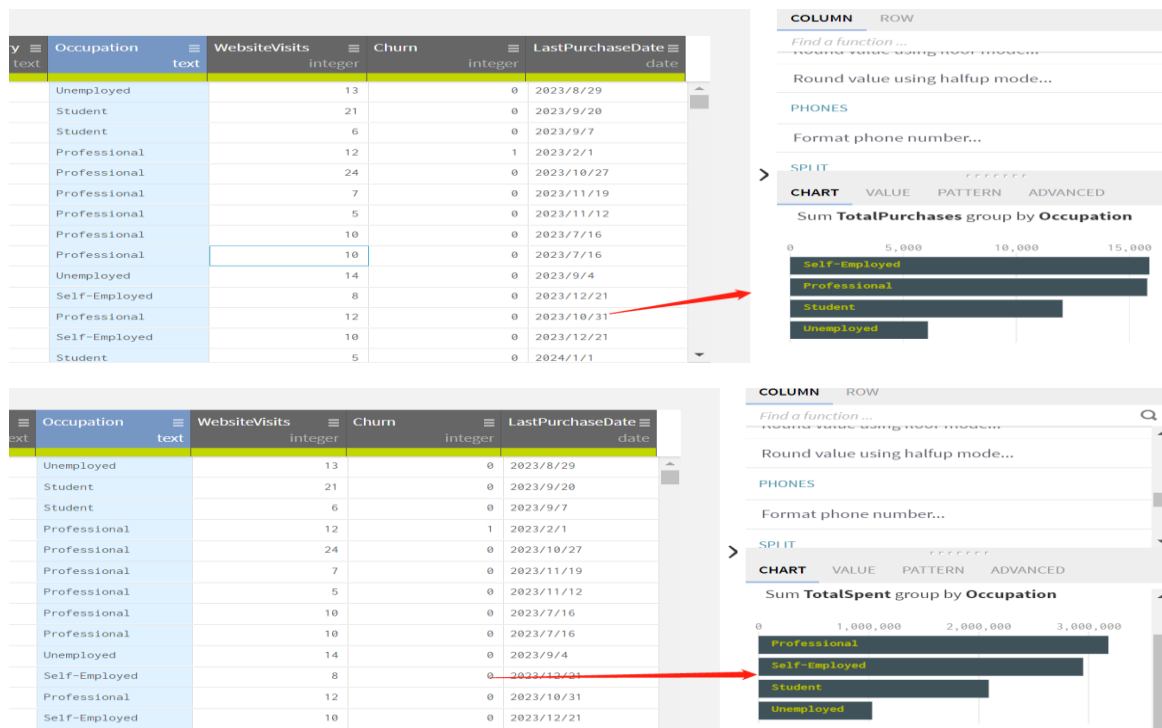


- 2.6 It can be observed that major cities like Beijing, Guangzhou, and Shanghai exhibit significantly higher purchase frequencies and total expenditures compared to other cities.



- 2.7 It is noticeable that individuals with employment or student status tend to have significantly higher purchase frequencies and total expenditures compared

to those who are unemployed.



EXPORT TO CSV

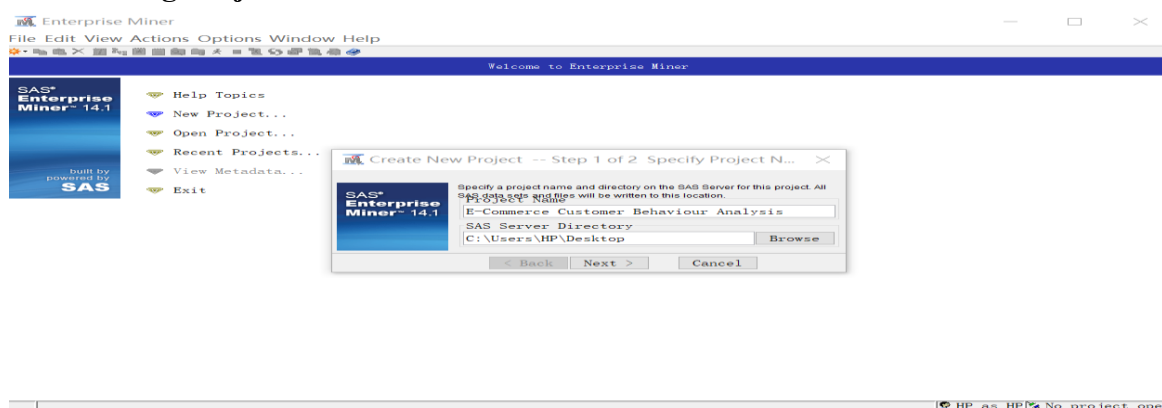
Delimiter: Semicolon

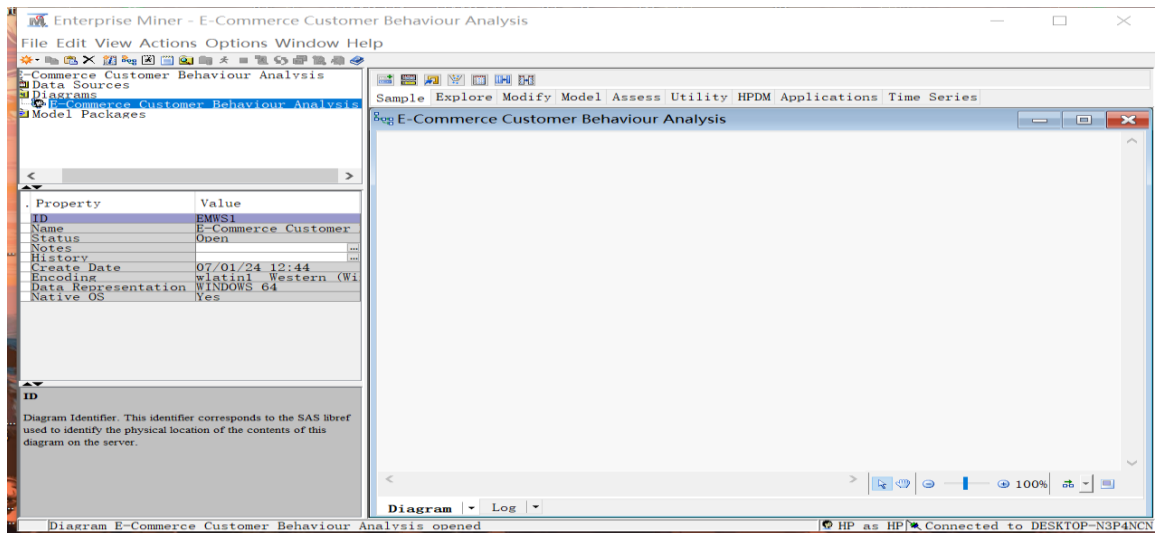
Filename: Aggregate customer information PREPARATION

CANCEL EXPORT

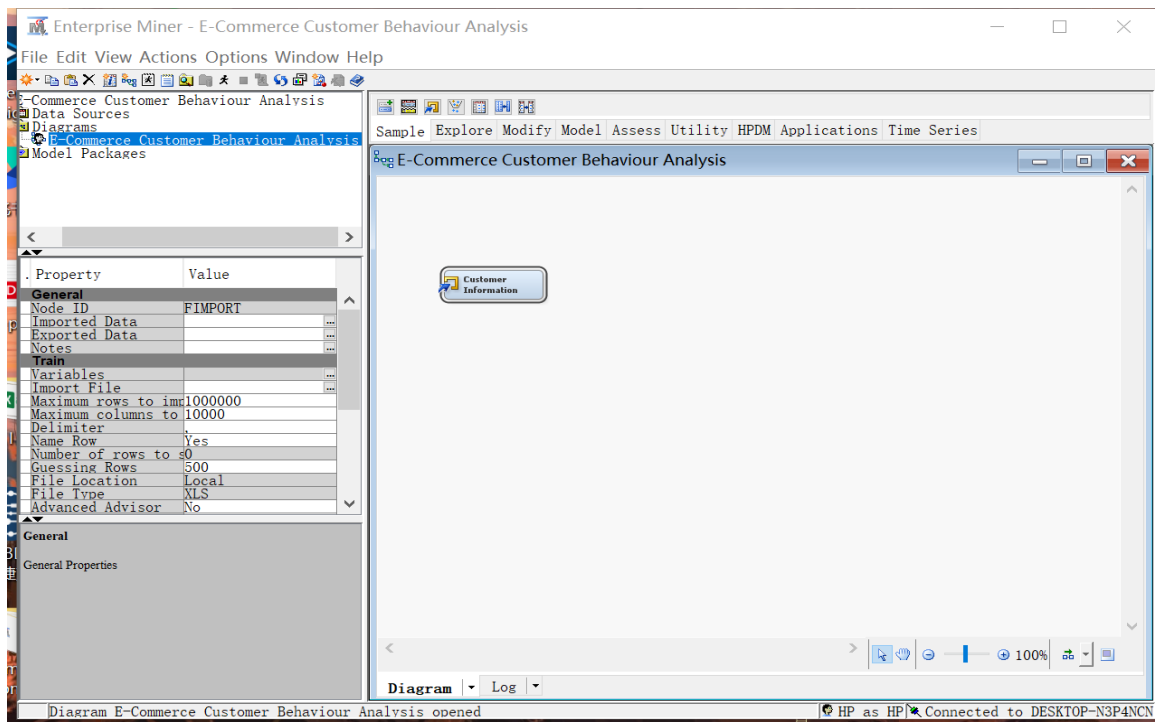
### 3. SAS EM

#### 3.1 Creating Project Files





### 3.2 Import of data and initial checking



**Output**

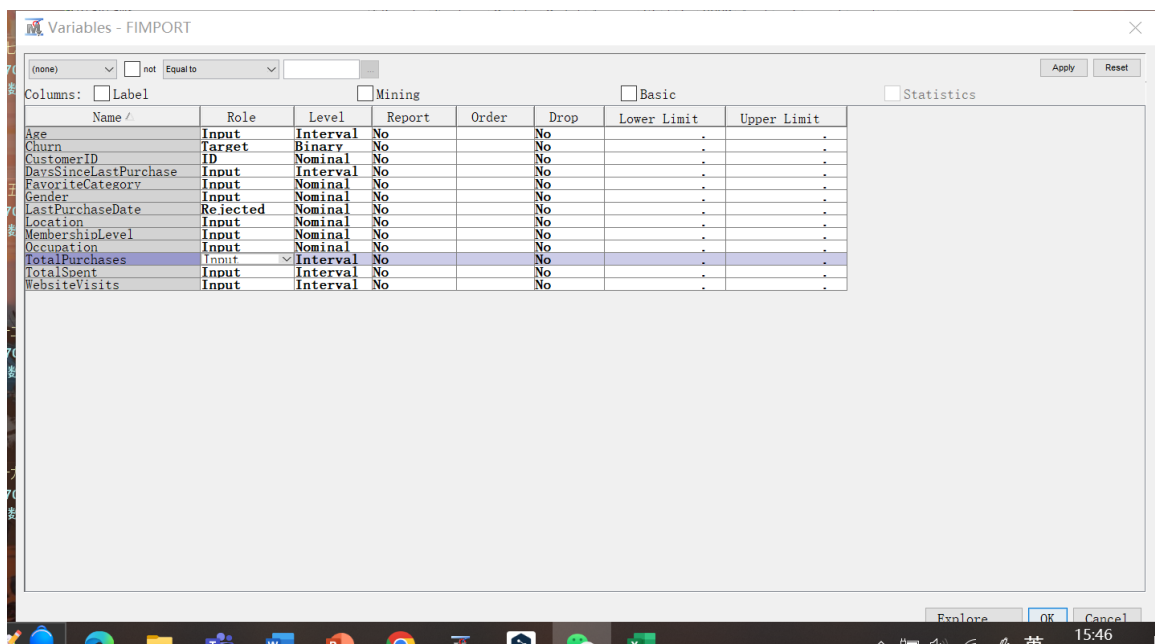
46	Number of Data Set Repairs	0
47	ExtendObsCounter	YES
48	Filename	C:\Users\HP\Desktop\E-Commerce Customer Behaviour Analysis\Workspaces\EMWS1\import_data.sas7bdat
49	Release Created	9.0401M3
50	Host Created	X64_SHOME
51		
52		
53	Alphabetic List of Variables and Attributes	
54		
55	#	Variable
56		Type
57		Len
58		Format
59		Informat
60	2	Age
61	11	Churn
62	1	CustomerID
63	13	DaysSinceLastPurchase
64	8	FavoriteCategory
65	3	Gender
66	12	LastPurchaseDate
67	4	Location
68	5	MembershipLevel
69	9	Occupation
70	6	TotalPurchases
71	7	TotalSpent
72	10	WebsiteVisits
73		
74		
75		



### 3.3 specify variable roles

#### reason:

- CustomerID: This is a unique identifier for each customer and holds no predictive value for the model. It should be set as ID and not used as a feature in decision tree or random forest models.
- Age: This is a continuous numerical variable and should be set as Interval, serving as an input feature.
- Gender, Location, MembershipLevel, FavoriteCategory, Occupation: These are categorical variables and should be set as Nominal, serving as input features.
- TotalPurchases, TotalSpent, and WebsiteVisits: These are continuous numerical variables and should be set as Interval, serving as input features.
- Churn: This is the target variable indicating whether a customer has churned or not. Considering its binary nature (0 or 1), it should be set as Target and Binary.
- LastPurchaseDate: This is a date variable. For constructing decision trees and random forest models, date variables are typically not directly used. Since it has already been transformed into days since the last purchase, it can be set as Rejected.



81	*-----*		
82			
83			
84			
85			
86	Exported Attributes for TRAIN Port		
87			
88		Measurement	Frequency
89	Role	Level	Count
90			
91	ID	NOMINAL	1
92	INPUT	INTERVAL	5
93	INPUT	NOMINAL	5
94	REJECTED	NOMINAL	1
95	TARGET	BINARY	1
96			

### 3.4 Creating Data Sources

Enterprise Miner - E-Commerce Customer Behaviour Analysis

File Edit View Actions Options Window Help

Commerce Customer Behaviour Analysis

Sample Explore Modify Model Assess Utility HPDM Applications Time Series

E-Commerce Customer Behaviour Analysis

Model Packages

Library Wizard -- Step 3 of 3 Confirm Action

Property	Value
Action	Create New
Name	AAEM2
Engine	BASE
Path	C:\Users\HP\Desktop\E-Commerce Custom
Options	

Status  
Action Succeeded!  
The Library "AAEM2" was created.

< Back Finish

Diagram Log 100%

Data Source Wizard -- Step 1 of 8 Metadata Source

Select a metadata source

Source: SAS Table

< Back Next > Cancel

Select a SAS Table

Name	Type
Em save train	Table

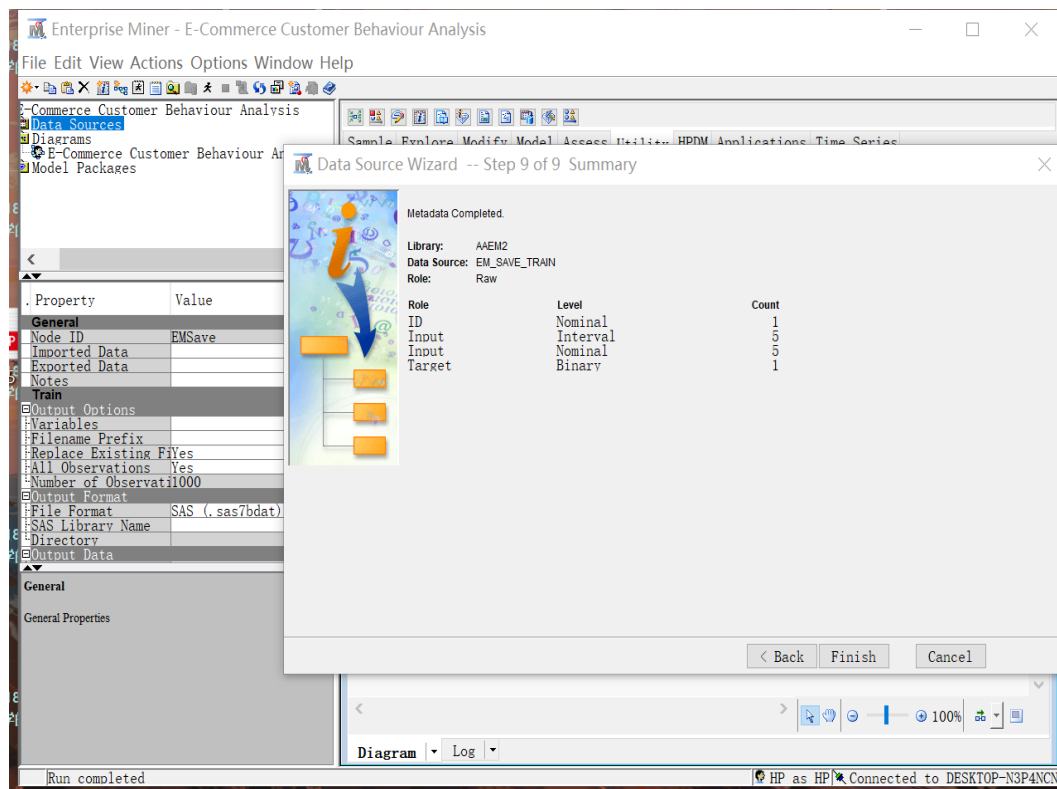
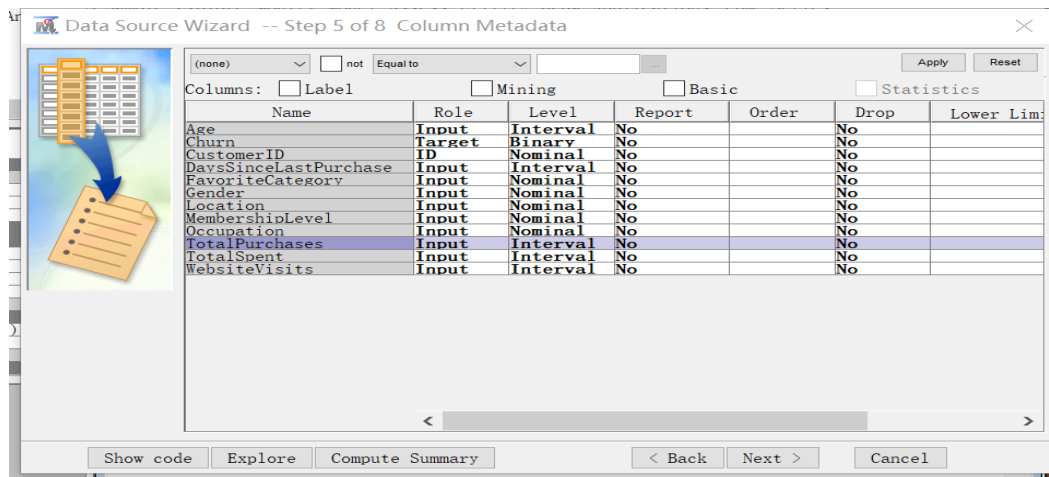
Get Details Properties... Refresh OK Cancel

Data Source Wizard -- Step 3 of 8 Table Information

Table Properties

Property	Value
Table Name	AAEM2. EM SAVE TRAIN
Description	
Member Type	DATA
Data Set Type	DATA
Engine	BASE
Number of Variables	12
Number of Observations	1000
Created Date	07 January 2024 15:57:58 CST
Modified Date	07 January 2024 15:57:58 CST

< Back Next > Cancel

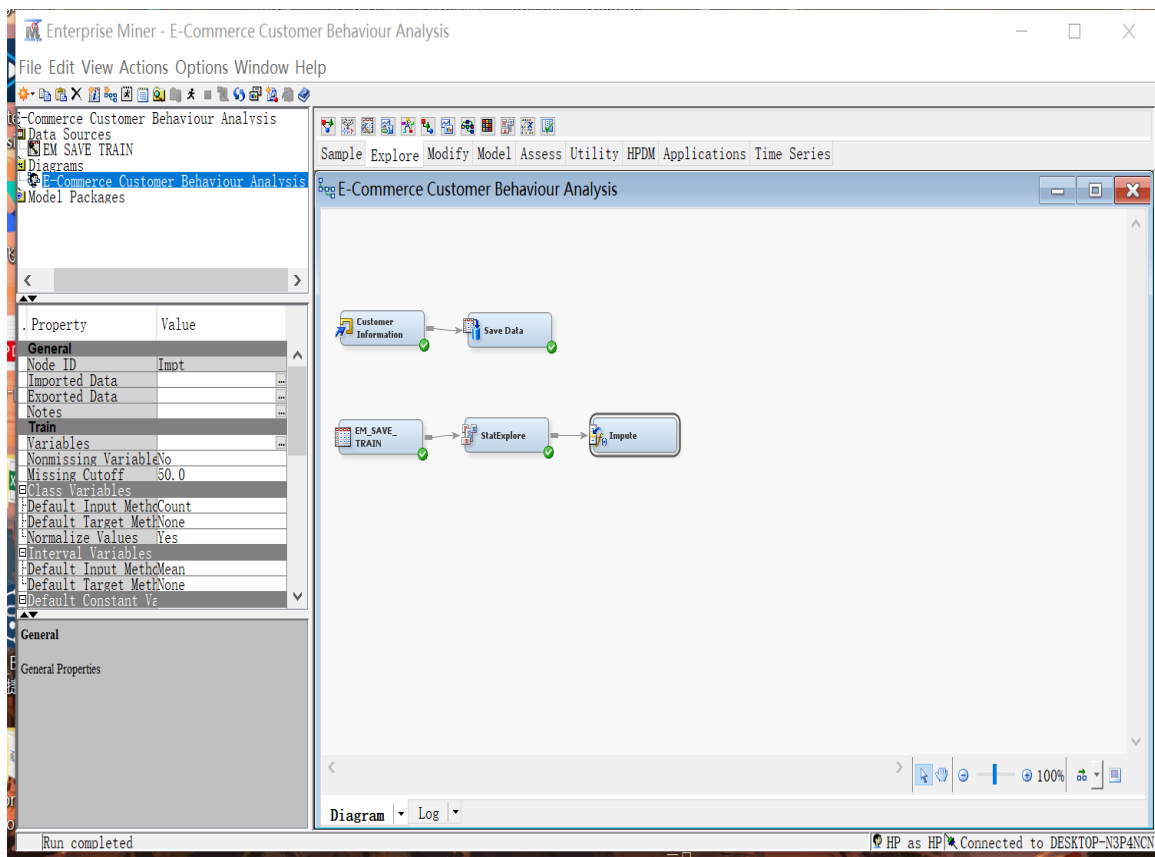
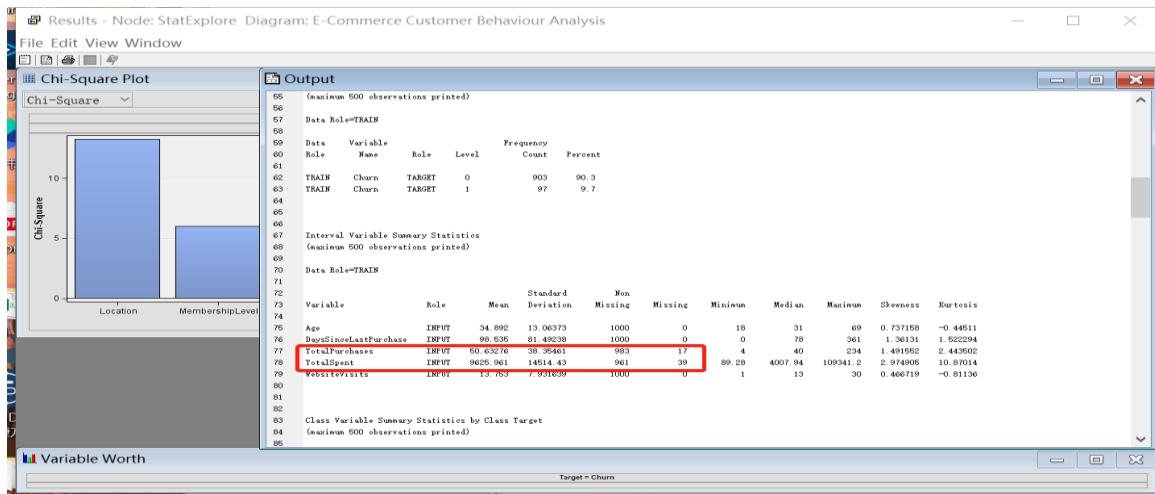


### 3.5 Processing data for missing values

#### Choose median to due missing values

Reason: After checking the data structure found that there are TotalPurchases and

TotalSpent so that there are two columns of missing values, histogram to check the distribution of data is indeed more skewed, so the use of the median to fill in the missing values, which is particularly suitable for dealing with skewed distribution (i.e., asymmetric or long-tailed distribution) of numerical data. The median is less susceptible to extreme values or outliers than the mean.



Enterprise Miner - E-Commerce Customer Behaviour Analysis

File Edit View Actions Options Window Help

E-Commerce Customer Behaviour Analysis

Data Sources

EM\_SAVE\_TRAIN

Diagrams

E-Commerce Customer Behaviour Analysis

Model Packages

Property Value

General

Node ID Impt

Imported Data

Exported Data

Notes

Train

Variables

Nonmissing Variables No

Missing Cutoff 50.0

Class Variables

Default Input Method Count

Default Target Method None

Normalize Values Yes

Interval Variables

Default Input Method Median

Default Target Method None

Default Constant Value

Default Character Value

Default Number Value

Method Options

Random Seed 12345

Tuning Parameters

Tree Imputation

Score

Hide Original Variables Yes

Indicator Variables

Type None

Source Imputed Variables

Role Rejected

Report

Validation and Test Data No

Distribution of missing No

Status

Default Input Method

Sample Explore Modify Model Assess Utility HPDM Applications Time Series

E-Commerce Customer Behaviour Analysis

Customer Information

Save Data

EM\_SAVE\_TRAIN

StatExplore

Impute

Diagram Log

Results - Node: StatExplore (2) Diagram: E-Commerce Customer Behaviour Analysis

File Edit View Window

Chi-Square Plot

Chi-Square

Location

MembershipLevel

Variable Worth

DaysSinceLastPurchase

IMP\_TotalSpent

Age

Location

IMP\_TotalPurchases

WebsiteVisits

MembershipLevel

Occupation

FavoriteCategory

Gender

Output

61	TRAIN	Churn	TARGET	0	903	90.3
62	TRAIN	Churn	TARGET	1	97	9.7
63						
64						
65						
66	Interval Variable Summary Statistics					
67	(maximum 500 observations printed)					
68						
69	Data Role=TRAIN					
70						
71				Standard	Non	
72	Variable	Role	Mean	Deviation	Missing	Missing
73					Minimum	Median
74	Age	INPUT	34.892	13.06373	1000	0
75	DaysSinceLastPurchase	INPUT	98.535	81.49238	1000	0
76	IMP_TotalPurchases	INPUT	50.452	38.05172	1000	0
77	IMP_TotalSpent	INPUT	9408.058	14269.39	1000	0
78	WebsiteVisits	INPUT	13.753	7.931639	1000	0
79						
80						
81						
82	Class Variable Summary Statistics by Class Target					
83	(maximum 500 observations printed)					
84						
85	Data Role=TRAIN Variable Name=FavoriteCategory					
86						
87		Number				
88	Target	of				
89	Level	Levels	Missing	Mode	Percentage	Mode2
90	Percentage					

### 3.6 Data Partitioning

Now, let's proceed with data partitioning before the modeling phase. We'll split the dataset into training, validation, and test sets. These sets are used for training the model, tuning parameters, and evaluating the final model performance. The suggested partitioning is Training: 70.0%, Validation: 15%, Test: 15%, and the Partitioning Method should be set to Stratified.

Reason: Considering the dataset size of only 1000 data points, if you need a larger training set to avoid overfitting, selecting 70% for training is a reasonable choice. It's also wise to use the "Default" partitioning method, especially when dealing with imbalanced class distribution in the target variable (Churn). This helps maintain a closer representation of the overall class proportions within each partition.

The screenshot displays the Alteryx interface. On the left, the 'Properties' pane for the 'Data Partition' node is visible. The 'Partitioning Method' is set to 'Stratified', and the 'Data Set Allocations' are configured as Training: 70.0%, Validation: 15.0%, and Test: 15.0%. The 'Report' section shows 'Interval Targets' and 'Class Targets' both set to 'Yes'. The 'Status' section shows the 'Create Time' as 07/01/24 16:36. On the right, the workflow diagram shows a sequence of nodes: 'Customer Information' (input), 'Save Data', 'EM\_SAVE\_TRAIN', 'StatExplore', 'Impute', 'StatExplore (2)', and 'Data Partition'.

Property	Value
Node ID	Part
Imported Data	
Exported Data	
Notes	
Variables	
Output Type	Data
Partitioning Method	Stratified
Random Seed	12345
Data Set Allocations	
Training	70.0
Validation	15.0
Test	15.0
Report	
Interval Targets	Yes
Class Targets	Yes
Status	
Create Time	07/01/24 16:36
Run ID	
Last Error	
Last Status	
Last Run Time	
Run Duration	
Grid Host	
User-Added Node	No

Results - Node: Data Partition Diagram: E-Commerce Customer Behaviour Analysis

File Edit View Window

Output

Summary Statistics for Class Targets					
Data=DATA					
Variable	Numeric Value	Formatted Value	Frequency Count	Percent	Label
Churn	0	0	903	90.3	
Churn	1	1	97	9.7	

Data=TEST					
Variable	Numeric Value	Formatted Value	Frequency Count	Percent	Label
Churn	0	0	136	90.0662	
Churn	1	1	15	9.9338	

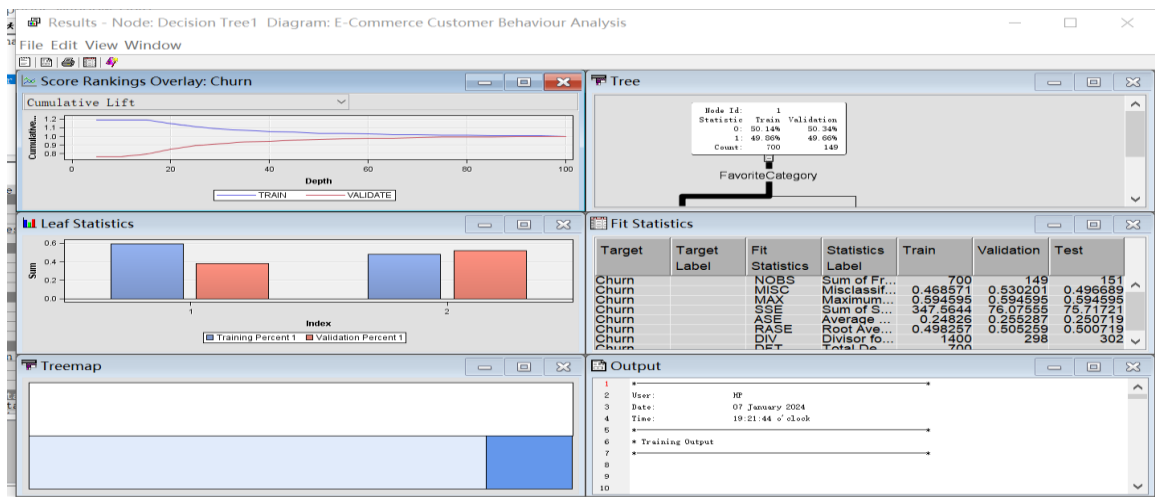
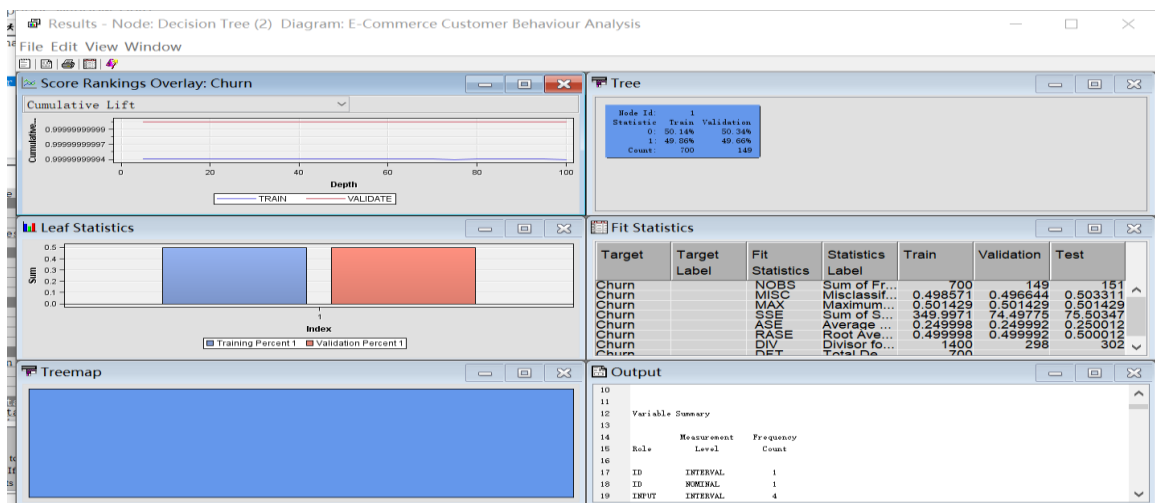
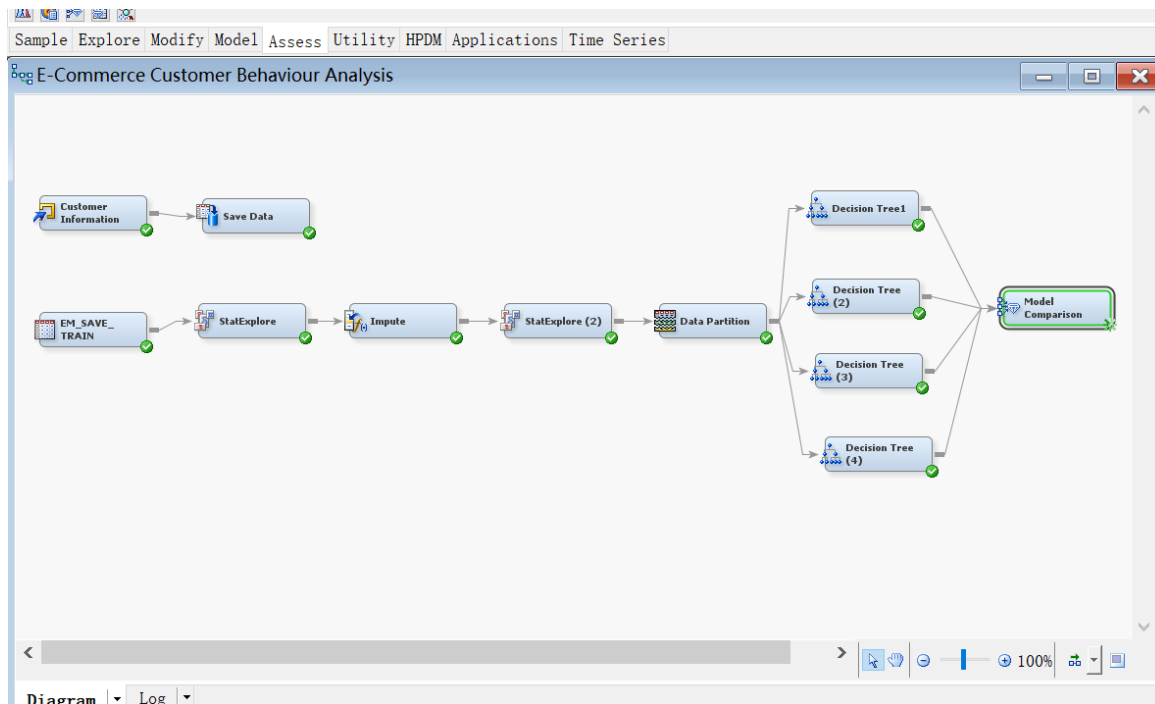
Data=TRAIN					
Variable	Numeric Value	Formatted Value	Frequency Count	Percent	Label
Churn	0	0	632	90.4149	
Churn	1	1	67	9.5851	

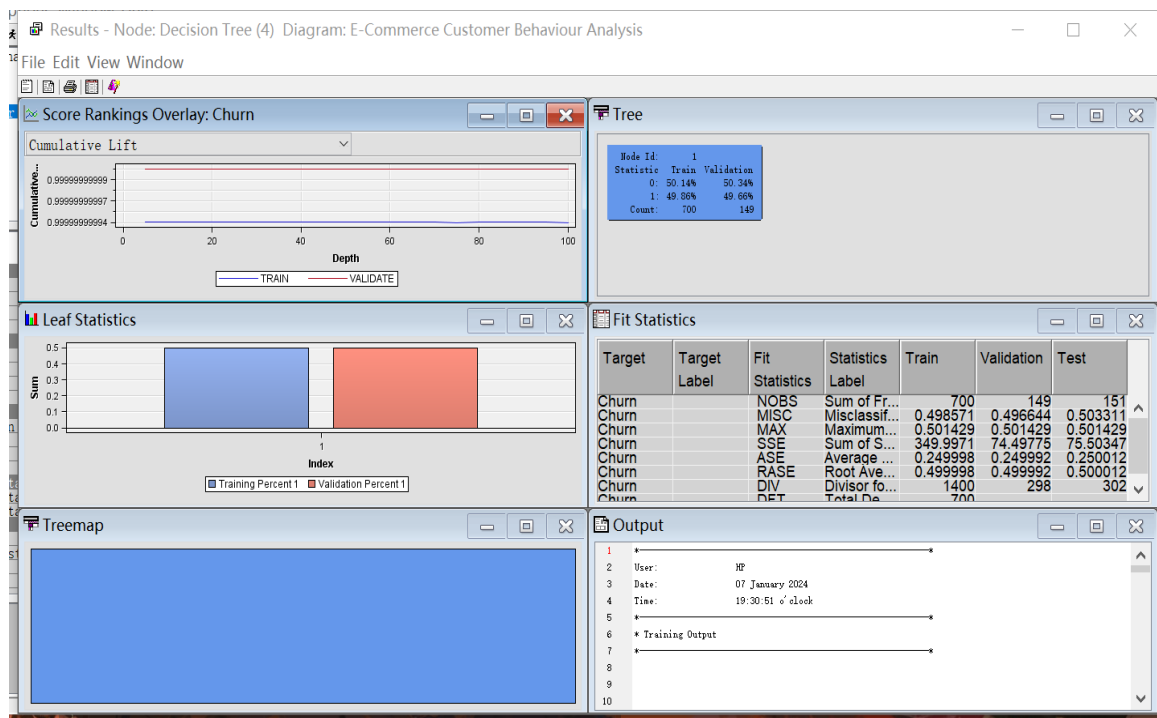
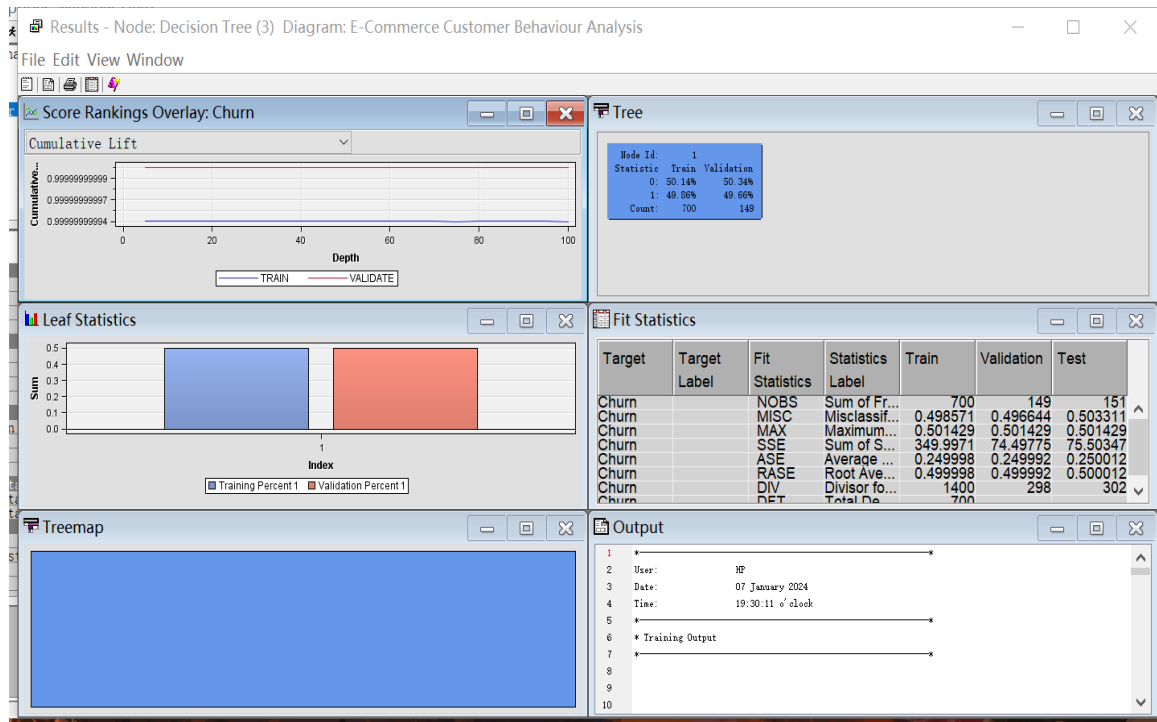
  

Data=VALIDATE					
Variable	Numeric Value	Formatted Value	Frequency Count	Percent	Label

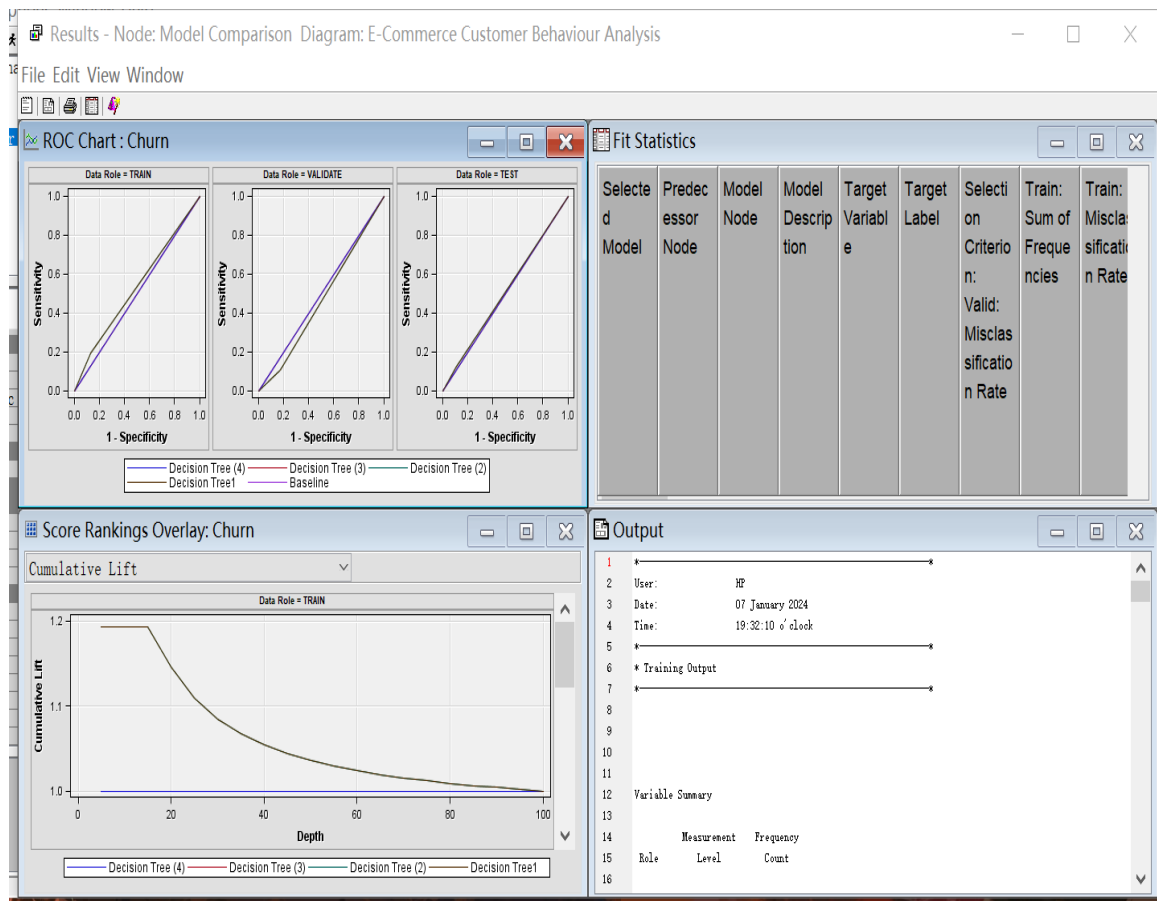
### 3.7 Create a decision tree model

I choose different assessment measure and create four different tree









- **Model Selection:**

Model selection is based on the Validation Misclassification Rate (VMISC), showing the performance of four different decision tree models. Tree3 and Tree4 have the lowest misclassification rates, but according to the event classification table, these models seem to have not captured any positive events (Churn=1). This suggests that the models might simply be predicting all cases as the majority class (Churn=0).

- **Fit Statistics:**

Kolmogorov-Smirnov Statistic (KS Statistic): For Tree3 and Tree4, the KS statistics for the training, validation, and test sets are all 0, indicating that these models lack discriminative power.

Roc Index: The index is approximately 0.5 for all models, implying that the predictive ability of the models is not different from random guessing.

Cumulative Percent Captured Response: In Tree3 and Tree4, this metric is 0, suggesting that the models did not capture any response from positive events.

Event Classification Table: Tree3 and Tree4 did not correctly predict any cases with Churn=1 in both the training and validation sets, while Tree and Tree2, although capturing some Churn=1 cases, still exhibit overall poor performance.

- **Based on these outputs, here is an assessment of how the decision tree models impact customer behavior analysis:**

These models struggled with handling imbalanced datasets, as accuracy is not a good metric in such cases. Even if the models predict no churn for all customers (Churn=0), accuracy may appear high.

Tree3 and Tree4 may be too simplistic, lacking sufficient complexity to learn patterns in the data, and they might be disregarding the minority class (Churn=1).

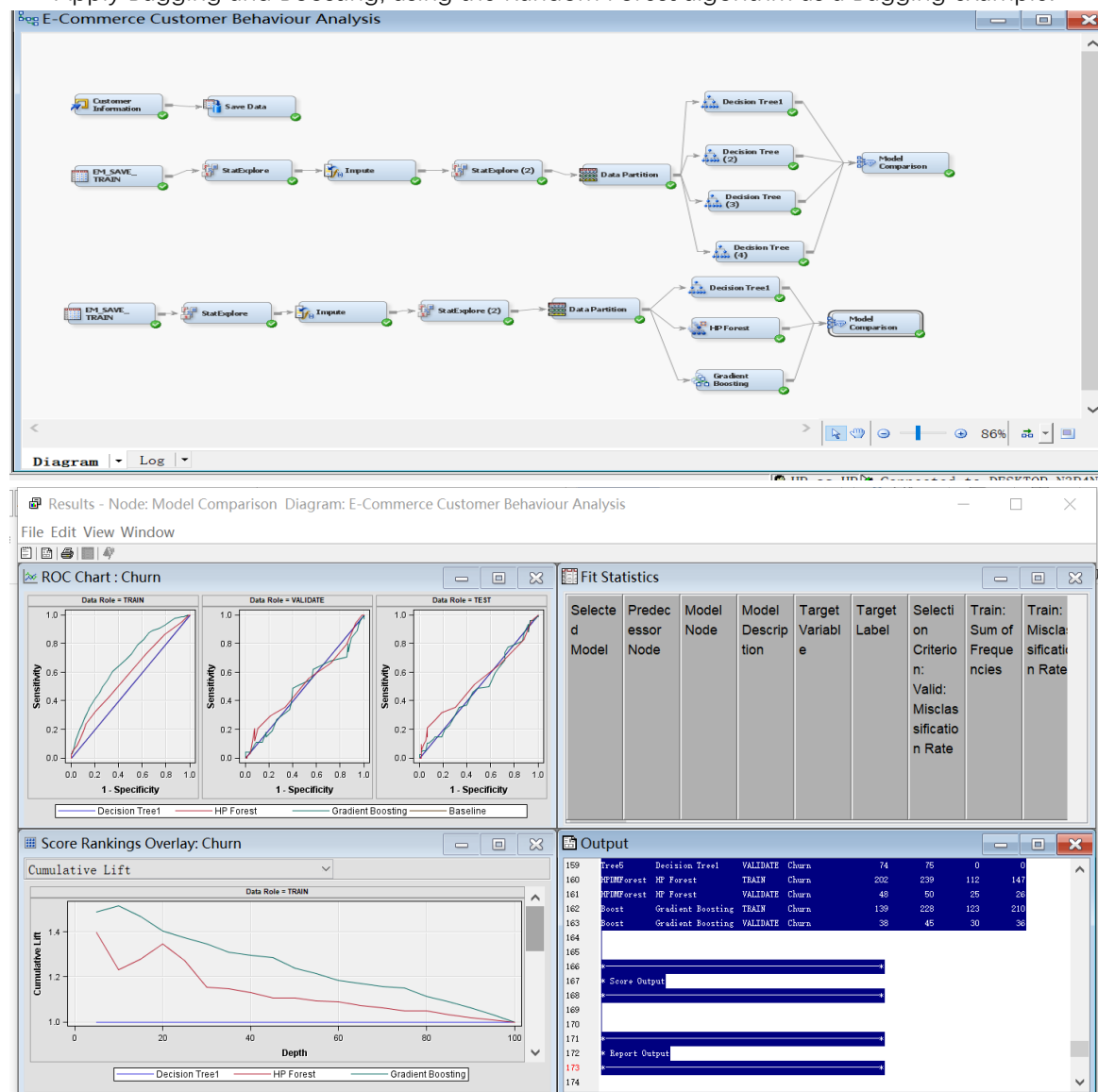
While Tree and Tree2 attempted to distinguish between the two classes, they still need adjustments to improve their ability to identify the minority class.

- **Recommendations:**

To enhance predictive capabilities for customer churn behavior, consider employing resampling techniques to balance classes, adjusting model complexity, or trying different model parameters.

Further fine-tuning and optimization of these decision tree models are necessary. Explore different model structures, parameter settings, or try alternative model types such as random forests or gradient boosting models, which may exhibit better robustness on imbalanced datasets.

### 3.8 Apply Bagging and Boosting, using the Random Forest algorithm as a Bagging example.



We can observe the performance metrics of three models based on different decision tree

algorithms: Gradient Boosting (Boost), HP Forest (HPDMForest), and a single decision tree (Tree5). Here is an analysis of these results:

#### **Model Performance Comparison:**

- **Gradient Boosting (Boost):** It has the lowest misclassification rate on the validation set (0.37429), and its performance in terms of average square error (0.25801) and ROC index (0.479) is also good. This indicates that the Boost model has the overall best performance among the three models.
- **HP Forest (HPDMForest):** This model has a misclassification rate of 0.44857 on the validation set, with a slightly better average square error (0.24902) than Boost, but a lower ROC index (0.515) compared to the Boost model.
- **Decision Tree (Tree5):** It performs the worst among the three models, with a misclassification rate of 0.49857 on the validation set. Both the average square error (0.24999) and ROC index (0.500) are the lowest.

#### **Classification Event Table:**

- The Boost model correctly predicted 210 cases of Churn=1 in the training set and 36 cases in the validation set.
- The HPDMForest model correctly predicted 147 cases of Churn=1 in the training set and 26 cases in the validation set.
- The Tree5 model did not correctly predict any cases of Churn=1 in both the training and validation sets.

#### **Customer Behavior Analysis:**

- The Boost model seems to capture behavioral patterns that may lead to customer churn most effectively. Its predictive ability for Churn=1 is superior to the other two models, suggesting that the Boost model can better identify key factors that may lead to customer churn.
- Although the HPDMForest model has some capability in predicting Churn=1, its performance is not as good as the Boost model. This could be because the model did not effectively capture all relevant behavioral patterns.
- The performance metrics of the Tree5 model suggest that it did not capture any patterns related to customer churn, making it potentially less useful for analyzing customer behavior.

#### **Conclusion:**

In summary, the Boost model performs the best among these three models, indicating that it might be the optimal choice for analyzing and predicting customer behavior. However, it's essential to note that despite the good performance on the validation set, the ROC index is still close to 0.5, suggesting room for improvement in predictive capability. Further exploration may be needed, such as adjusting model parameters, trying other advanced algorithms, and considering strategies to balance the dataset to enhance the model's predictive ability for the minority class.