

WQD7005 DATA MINING 2023/2024 S1

ALTERNATIVE ASSESSMENT 1

Name: Xin Dong

Matric Num: 22060696

Submitted Answer:

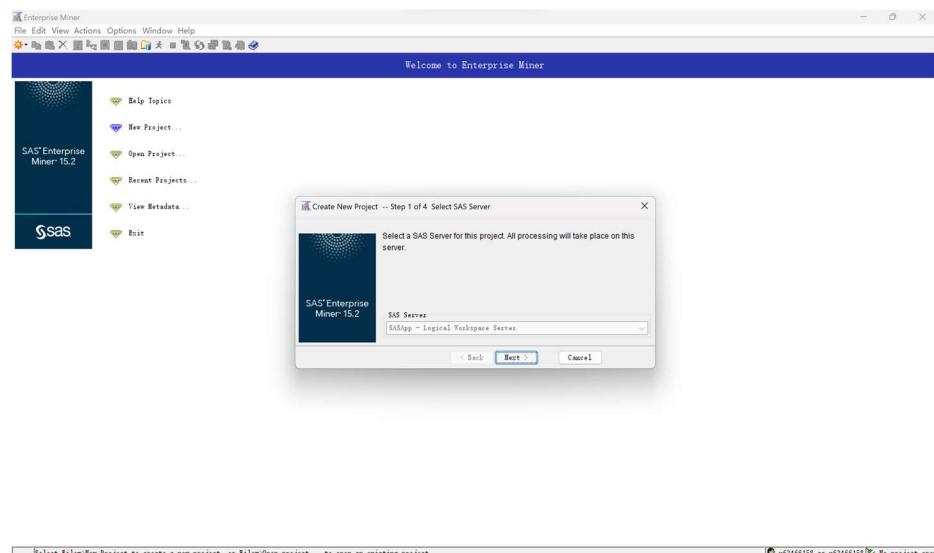
Github Link Repository: https://github.com/xinburg/Final-Assessment-1_Xin-Dong_22060696

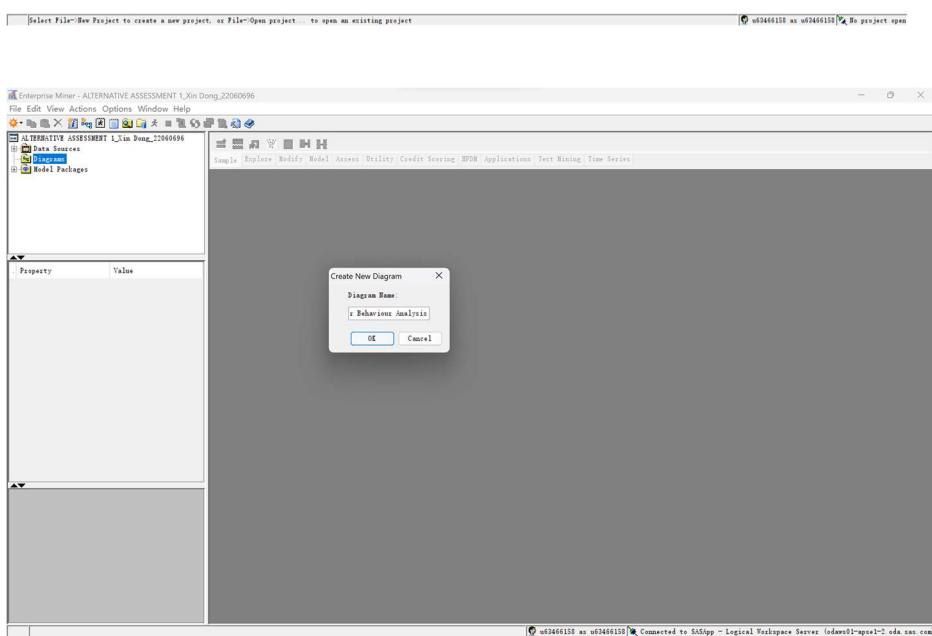
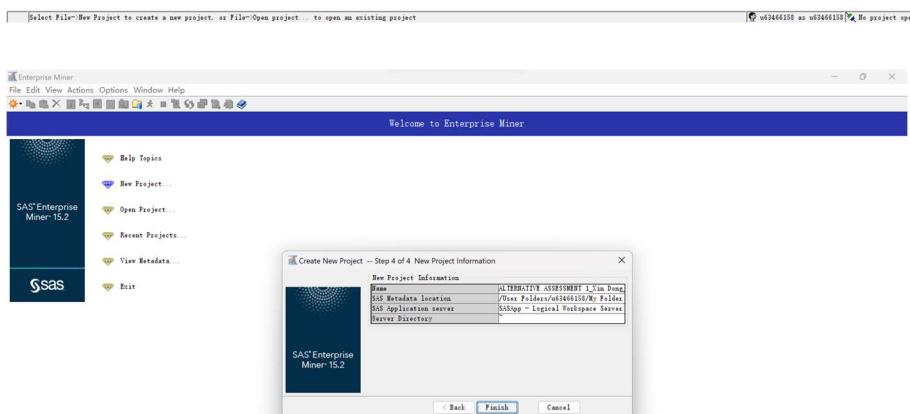
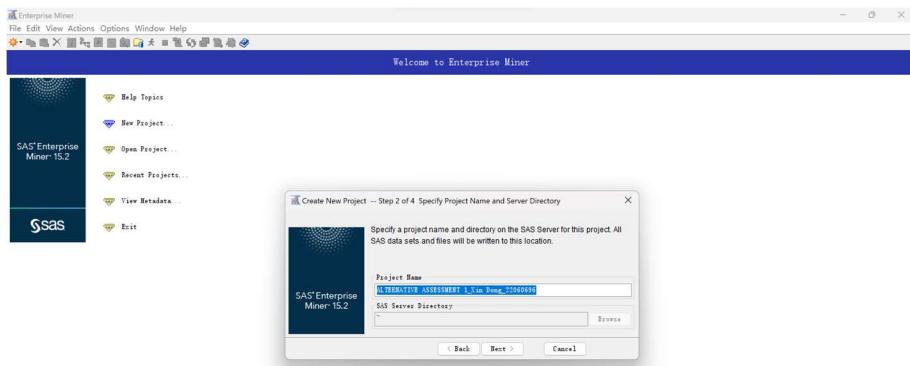
Dataset reference: <https://www.kaggle.com/datasets/shriyashjagtap/e-commerce-customer-for-behavior-analysis/data>

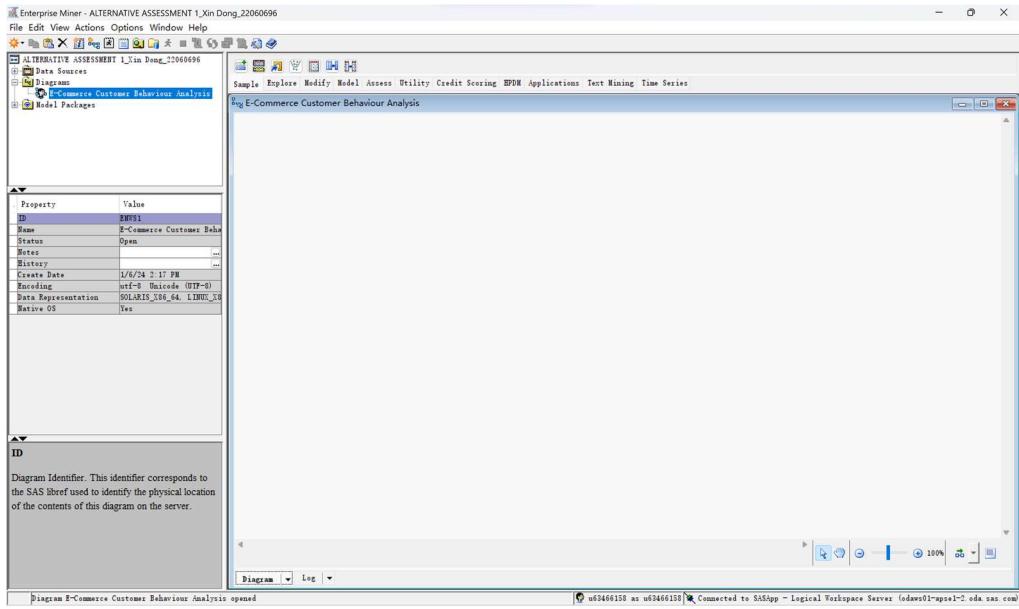
In my study, a new practical dataset was generated based on the real-world reference dataset so that it can closely match the AA1's criteria and allows for a controlled environment to test various analytical methods. The generated data is realistic and representative of actual e-commerce customer behavior.

1. Create Project and Data Import in SAS Enterprise Miner

Start SAS Enterprise Miner: Open the software and create a new project named ALTERNATIVE ASSESSMENT 1_Xin Dong_22060696. And then create a new Diagram and set the name as: E-Commerce Customer Behaviour Analysis.

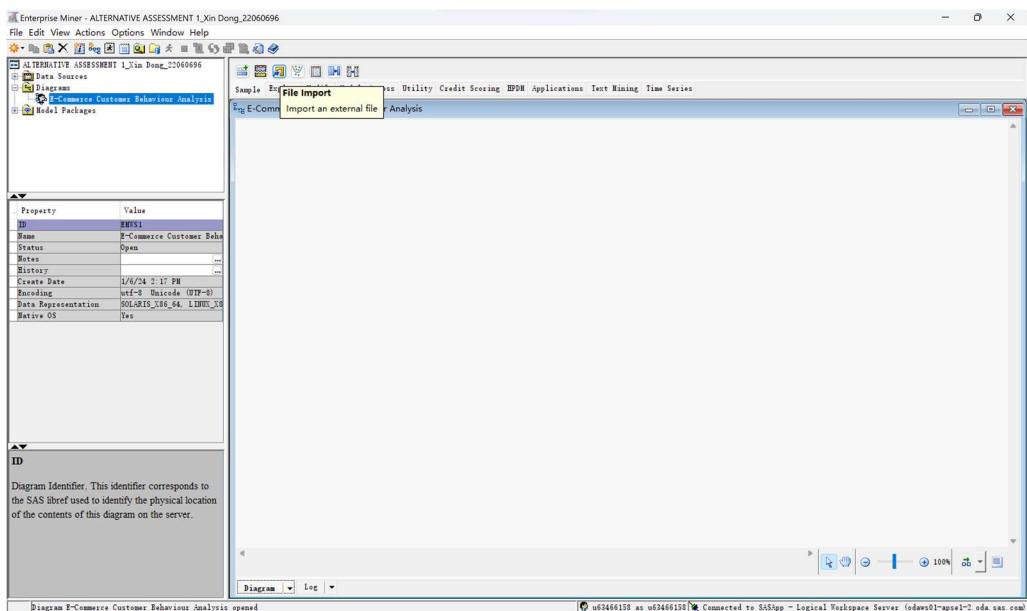


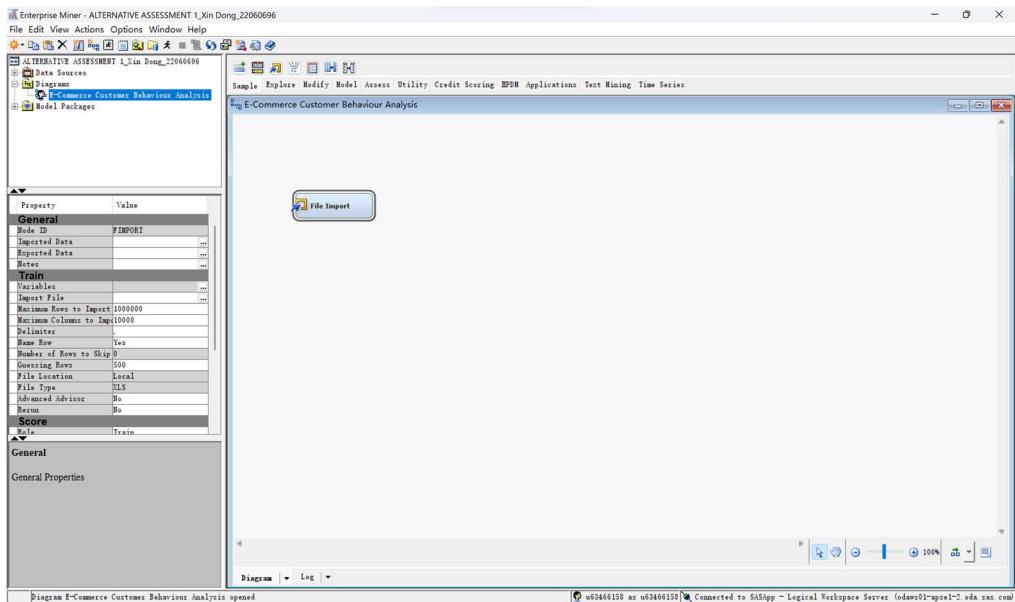




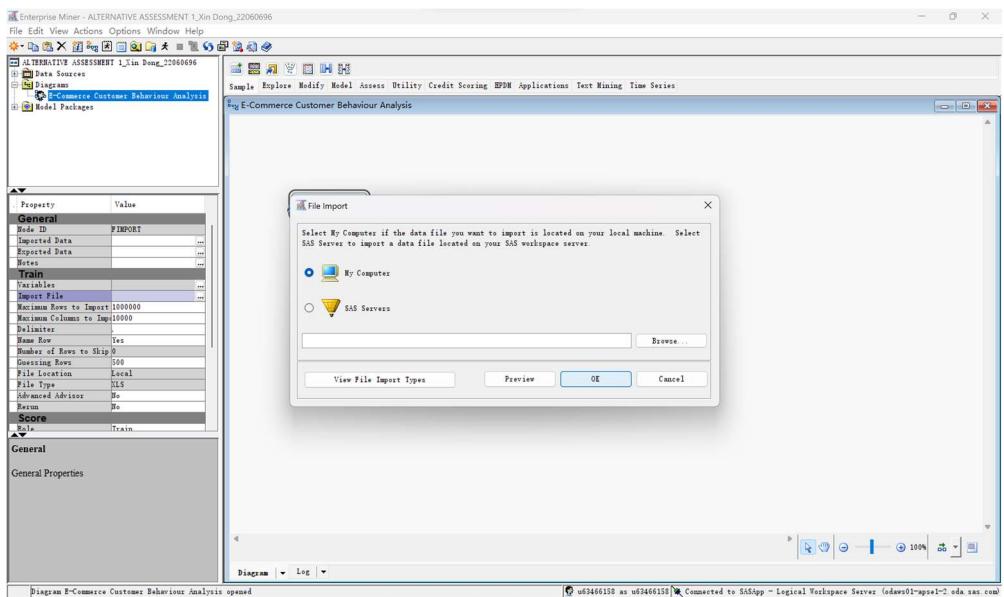
Import Dataset:

Navigate to the File menu and select Import Data.

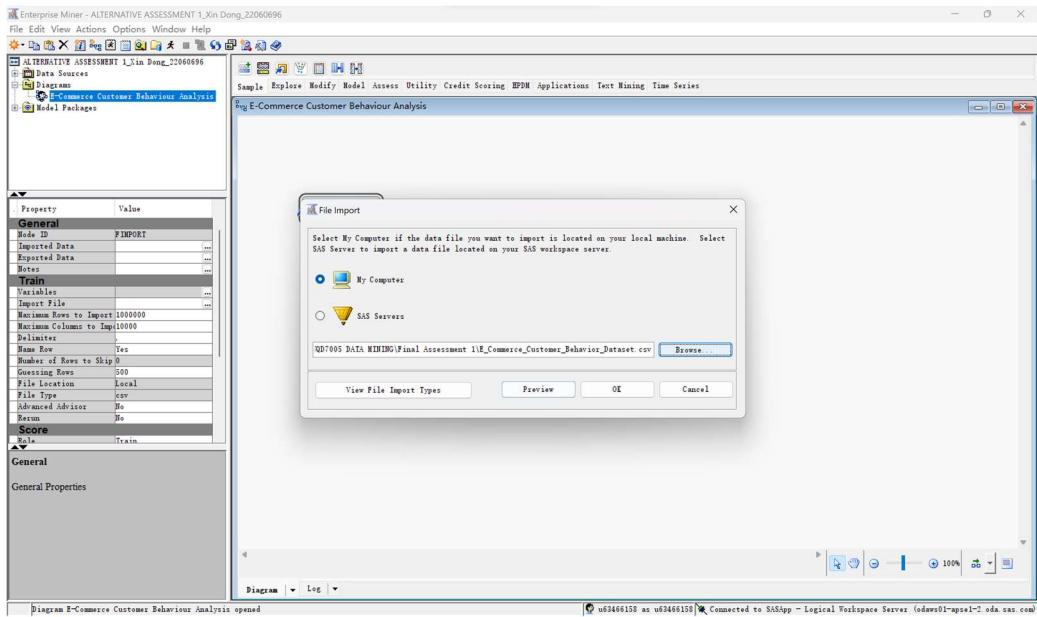




Click on File Import icon and select ‘Import the file’ from the left side configuration bar.



Browse and select the CSV file as the dataset we are going to analyze.



Follow the prompts in the Import Wizard to complete the import process.

2. Data Exploration

View Data: Once imported, view the dataset to ensure it's correctly loaded.

Summary Statistics: Generate summary statistics to understand distributions, averages, and other statistical measures for each variable. This step is crucial for identifying anomalies or patterns that might require addressing during preprocessing.

The details of the dataset as shown in below pictures:

| Sample Properties | |
|-------------------|--------------|
| Property | Value |
| Rows | 1000 |
| Columns | 14 |
| Library | EMWS1 |
| Member | FIMPORT_DATA |
| Type | DATA |
| Sample Method | Top |
| Fetch Size | Default |
| Fetched Rows | 1000 |
| Random Seed | 12345 |

Explore - EMWS1.FIMPORT_DATA

File View Actions Window

Detailed View

EMWS1.FIMPORT_DATA

| Obs # | CustomerID | Churn | Tenure | PreferredLoginDevice | Age | Gender | Location | MembershipLevel | TotalPurchases | TotalSpent | FavoriteCategory | LastPurchaseDate | Occupation | FrequencyOfWebsiteVisits |
|-------|------------|-------|------------|----------------------|----------|----------|----------|-----------------|----------------|---------------|------------------|------------------|---------------|--------------------------|
| 1 | 50001 | 1 | 2 Tablet | 58Male | Suburban | Platinum | | | 49 | 10237.70352 | Home Goods | 2023-02-12 | Student | 7 |
| 2 | 50002 | 1 | 1 Computer | 65Other | Rural | Platinum | | | 20 | 6337.57845 | Electronics | 2023-01-11 | Professional | 14 |
| 3 | 50003 | 0 | 1 Mobile | 54Male | Suburban | Silver | | | 87 | 8718.62064 | Home Goods | 2023-03-08 | Employed | 8 |
| 4 | 50004 | 0 | 1 Computer | 52Male | Rural | Platinum | | | 50 | 4759.40720 | Home Goods | 2023-04-19 | Student | 28 |
| 5 | 50005 | 0 | 5 Tablet | 67Male | Rural | Silver | | | 16 | 7425.91514 | Home Goods | 2023-11-09 | Professional | 25 |
| 6 | 50006 | 0 | 4 Computer | 63Other | Urban | Gold | | | 59 | 26163.38785 | Home Goods | 2023-05-09 | Student | 12 |
| 7 | 50007 | 1 | 1 Tablet | 27Other | Urban | Silver | | | 15 | 23055.39606 | Home Goods | 2023-08-22 | Self-Employed | 3 |
| 8 | 50008 | 0 | 1 Computer | 56Female | Rural | Bronze | | | 12 | 4461.20936 | Clothing | 2023-08-13 | Self-Employed | 12 |
| 9 | 50009 | 0 | 1 Computer | 28Female | Urban | Bronze | | | 9 | 4465.93358 | Electronics | 2023-09-06 | Self-Employed | 3 |
| 10 | 50010 | 0 | 5 Computer | 53Male | Rural | Bronze | | | 40 | 18989.02570 | Home Goods | 2023-04-09 | Retired | 25 |
| 11 | 50011 | 1 | 4 Mobile | 41Other | Suburban | Gold | | | 24 | 4123.24401 | Home Goods | 2023-02-06 | Retired | 27 |
| 12 | 50012 | 0 | 2 Mobile | 27Other | Suburban | Silver | | | 87 | 27849.12929 | Clothing | 2023-12-17 | Retired | 30 |
| 13 | 50013 | 0 | 4 Tablet | 39Male | Suburban | Gold | | | 19 | 8552.43176 | 3 Home Goods | 2023-10-12 | Professional | 29 |
| 14 | 50014 | 0 | 4 Computer | 37Female | Rural | Bronze | | | 17 | 18315.74989 | Electronics | 2023-12-05 | Self-Employed | 20 |
| 15 | 50015 | 0 | 5 Tablet | 29Male | Rural | Bronze | | | 96 | 42098.59403 | Electronics | 2023-12-26 | Retired | 10 |
| 16 | 50016 | 0 | 2 Mobile | 33Female | Suburban | Silver | | | 27 | 5169.53803 | Home Goods | 2023-04-02 | Self-Employed | 18 |
| 17 | 50017 | 0 | 5 Mobile | 45Other | Suburban | Bronze | | | 54 | 13883.12982 | Home Goods | 2023-05-06 | Retired | 17 |
| 18 | 50018 | 0 | 5 Mobile | 40Male | Urban | Silver | | | 100 | 2267.51446 | Electronics | 2023-02-03 | Retired | 14 |
| 19 | 50019 | 1 | 4 Mobile | 69Other | Rural | Gold | | | 40 | 16151.14103 | Clothing | 2023-03-22 | Retired | 14 |
| 20 | 50020 | 0 | 1 Mobile | 51Other | Suburban | Platinum | | | 44 | 8007.9316839 | Clothing | 2023-02-07 | Professional | 8 |
| 21 | 50021 | 0 | 1 Tablet | 63Male | Suburban | Platinum | | | 29 | 10193.762346 | Home Goods | 2023-02-27 | Retired | 5 |
| 22 | 50022 | 0 | 2 Mobile | 24Other | Urban | Gold | | | 15 | 19179.77702 | Clothing | 2023-01-10 | Retired | 27 |
| 23 | 50023 | 1 | 5 Mobile | 64Other | Suburban | Silver | | | 49 | 14741.15642 | Clothing | 2023-05-02 | Self-Employed | 10 |
| 24 | 50024 | 0 | 5 Mobile | 44Female | Suburban | Silver | | | 43 | 3662.1917275 | Home Goods | 2023-10-3 | Retired | 25 |
| 25 | 50025 | 0 | 2 Mobile | 27Other | Urban | Silver | | | 68 | 26345.74381 | Electronics | 2023-11-30 | Student | 20 |
| 26 | 50026 | 0 | 1 Tablet | 50Female | Rural | Gold | | | 14 | 3037.89507 | Clothing | 2023-01-18 | Retired | 8 |
| 27 | 50027 | 0 | 2 Mobile | 67Female | Urban | Platinum | | | 89 | 14983.03085 | Home Goods | 2023-08-07 | Retired | 19 |
| 28 | 50028 | 0 | 3 Computer | 64Other | Urban | Platinum | | | 98 | 5264.9702523 | Clothing | 2023-09-29 | Student | 10 |
| 29 | 50029 | 0 | 3 Computer | 31Female | Suburban | Platinum | | | 29 | 2504.9830159 | Clothing | 2023-04-20 | Student | 19 |
| 30 | 50030 | 0 | 5 Mobile | 37Male | Urban | Silver | | | 55 | 1846.33907 | Electronics | 2023-05-02 | Self-Employed | 3 |
| 31 | 50031 | 1 | 2 Mobile | 24Other | Suburban | Platinum | | | 84 | 23930.62225 | Home Goods | 2023-10-22 | Retired | 15 |
| 32 | 50032 | 0 | 2 Mobile | 62Other | Urban | Bronze | | | 84 | 32781.42011 | Electronics | 2023-08-03 | Professional | 3 |
| 33 | 50033 | 0 | 1 Computer | 49Female | Rural | Gold | | | 99 | 18516.76740 | Home Goods | 2023-02-17 | Professional | 8 |
| 34 | 50034 | 0 | 4 Mobile | 58Male | Rural | Silver | | | 99 | 4686.48505 | Clothing | 2023-03-07 | Self-Employed | 1 |
| 35 | 50035 | 0 | 2 Mobile | 59Other | Rural | Platinum | | | 30 | 16711.2900215 | Home Goods | 2023-10-05 | Retired | 15 |
| 36 | 50036 | 0 | 2 Mobile | 42Other | Rural | Platinum | | | 39 | 11761.46368 | Home Goods | 2023-04-05 | Student | 9 |
| 37 | 50037 | 0 | 5 Mobile | 19Female | Rural | Platinum | | | 82 | 23669.35747 | Electronics | 2023-05-14 | Professional | 13 |
| 38 | 50038 | 0 | 4 Tablet | 38Other | Suburban | Gold | | | 51 | 4437.78455 | Clothing | 2023-03-24 | Retired | 20 |
| 39 | 50039 | 0 | 4 Mobile | 30Male | Suburban | Gold | | | 14 | 2776.7559048 | Electronics | 2023-12-17 | Retired | 21 |
| 40 | 50040 | 1 | 1 Mobile | 44Male | Rural | Silver | | | 21 | 2628.7282629 | Home Goods | 2023-10-09 | Retired | 9 |
| 41 | 50041 | 0 | 2 Mobile | 35Other | Rural | Gold | | | 96 | 4454.24441 | Electronics | 2023-07-30 | Student | 16 |
| 42 | 50042 | 1 | 4 Computer | 66Male | Rural | Bronze | | | 10 | 3565.25424 | Electronics | 2023-03-08 | Student | 22 |
| 43 | 50043 | 0 | 2 Tablet | 26Other | Suburban | Gold | | | 83 | 37280.94756 | Electronics | 2023-10-15 | Student | 7 |
| 44 | 50044 | 1 | 3 Computer | 23Female | Suburban | Bronze | | | 32 | 5377.8377204 | Home Goods | 2023-10-19 | Professional | 17 |
| 45 | 50045 | 0 | 3 Computer | 37Other | Rural | Platinum | | | 34 | 16793.82096 | Clothing | 2024-01-09 | Retired | 22 |
| 46 | 50046 | 0 | 1 Tablet | 27Female | Rural | Gold | | | 20 | 7364.47545 | 3 Home Goods | 2023-04-27 | Self-Employed | 11 |

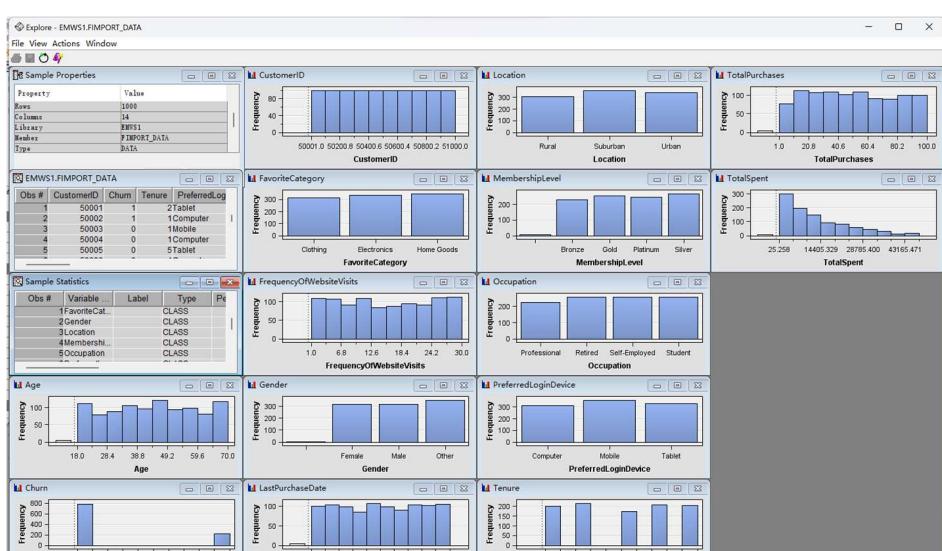
Explore - EMWS1.FIMPORT_DATA

File View Actions Window

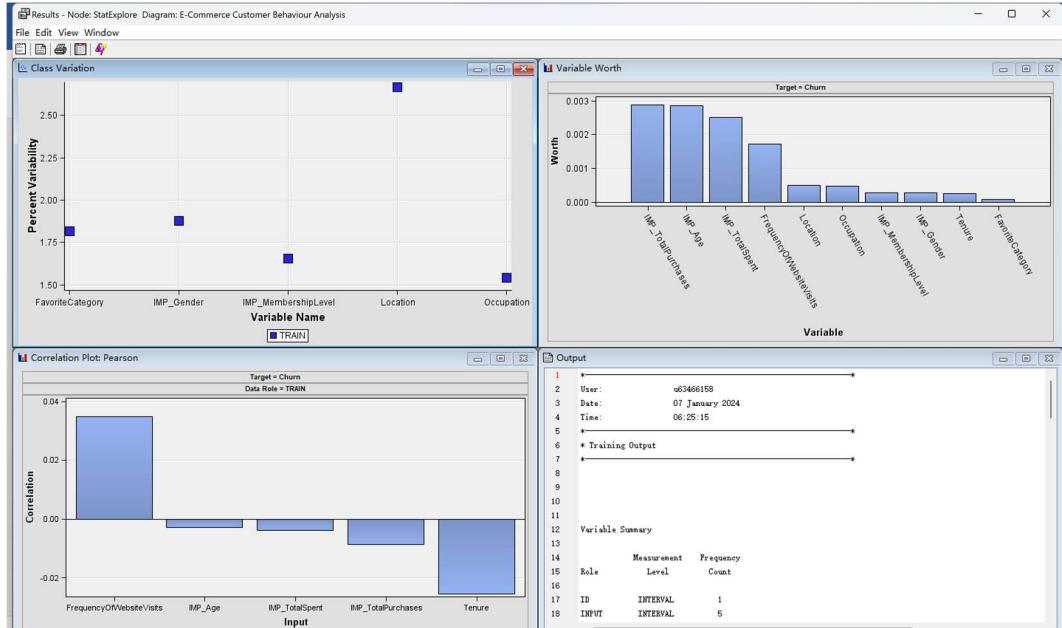
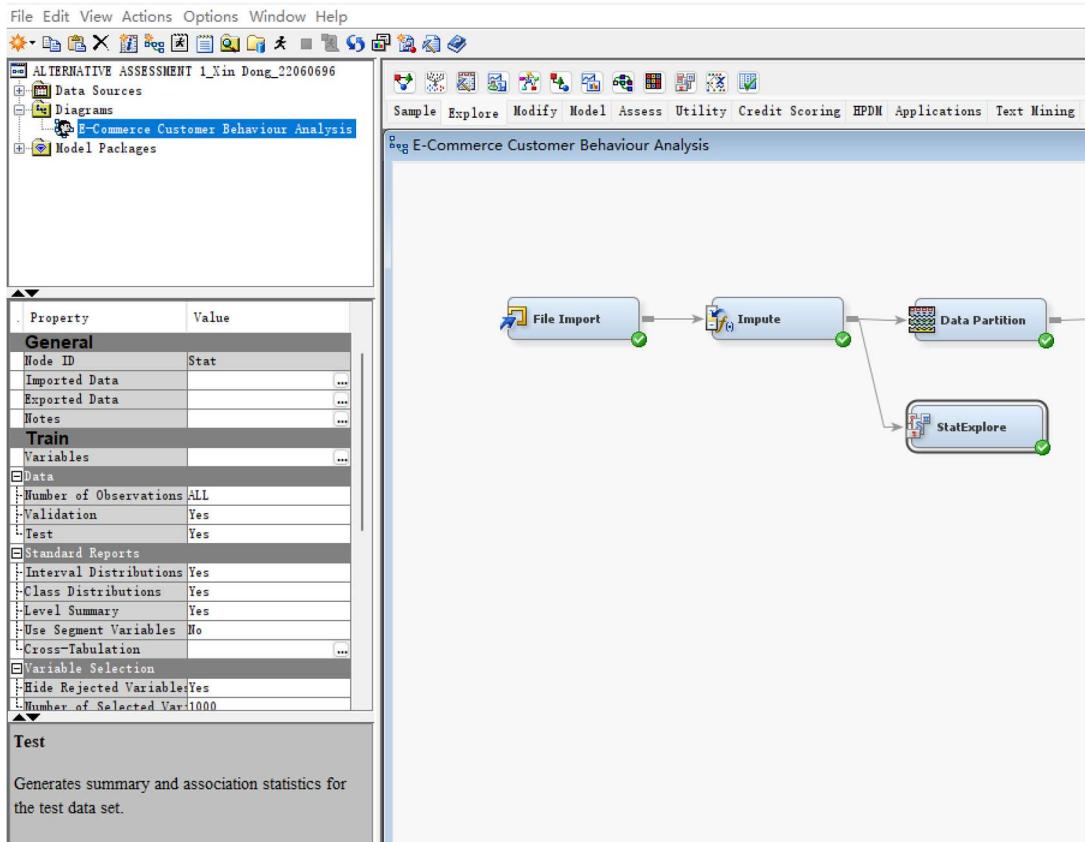
Detailed View

EMWS1.FIMPORT_DATA

| Obs # | Variable ... | Label | Type | Percent ... | Minimum | Maximum | Mean | Number o... | Mode Per... | Mode |
|-------|----------------|-------|-------|-------------|----------|----------|-----------|-------------|----------------|------|
| 1 | FavoriteCat... | CLASS | CLASS | 0 | . | . | . | 3 | 35 HOME GOO... | |
| 2 | Gender | CLASS | CLASS | 0.5 | . | . | . | 4 | 35.5 OTHER | |
| 3 | Location | CLASS | CLASS | 0 | . | . | . | 3 | 35.8 SUBURBAN | |
| 4 | Membershi... | CLASS | CLASS | 0.5 | . | . | . | 5 | 26.8 SILVER | |
| 5 | Occupation | CLASS | CLASS | 0 | . | . | . | 4 | 26 RETIRED | |
| 6 | PreferredLo... | CLASS | CLASS | 0 | . | . | . | 3 | 35.9 MOBILE | |
| 7 | Age | VAR | VAR | 0.5 | 18 | 70 | 44.30452, | . | . | . |
| 8 | Churn | VAR | VAR | 0 | 0 | 1 | 0.214, | . | . | . |
| 9 | CustomerID | VAR | VAR | 0 | 50001 | 51000 | 50500.5, | . | . | . |
| 10 | Frequency... | VAR | VAR | 0 | 1 | 30 | 15.552, | . | . | . |
| 11 | LastPurcha... | VAR | VAR | 0.5 | 23011 | 23376 | 23194.43, | . | . | . |
| 12 | Tenure | VAR | VAR | 0 | 1 | 5 | 2.999, | . | . | . |
| 13 | TotalPurch... | VAR | VAR | 0.5 | 1 | 100 | 50.45226, | . | . | . |
| 14 | TotalSpent | VAR | VAR | 0.5 | 25.25845 | 47958.83 | 12698.71, | . | . | . |



We could also Add one StatExplore node to explore the data:



From the exploration table and plots, we can summarize an attribute table:

| Variable | Role | Level | Rationale |
|----------|------|-------|-----------|
| | | | |

| | | | |
|--------------------------|----------|----------|---|
| Age | Input | Interval | Age is likely to impact customer behavior and could be a predictor of churn. |
| Churn | Target | Binary | This is the variable you are trying to predict. |
| CustomerID | ID | Nominal | Unique identifier for each record, used to track customers but not for prediction. |
| FavoriteCategory | Input | Nominal | Shopping preferences could influence churn decisions. |
| FrequencyOfWebsiteVisits | Input | Interval | Frequency of visits may correlate with customer engagement and churn. |
| Gender | Input | Nominal | Gender may influence purchasing patterns and thus affect churn. |
| LastPurchaseDate | Time ID | Interval | Important for time-series analysis, if applicable, but not used directly in churn prediction. |
| Location | Input | Nominal | Customer location might affect purchasing habits and churn likelihood. |
| MembershipLevel | Input | Nominal | Membership status could impact customer loyalty and churn rate. |
| Occupation | Input | Nominal | Occupation might relate to disposable income and purchase frequency, affecting churn. |
| PreferredLoginDevice | Rejected | Nominal | May have high cardinality or low predictive power for churn. |
| Tenure | Input | Interval | The length of time a customer has been with the company might influence churn. |
| TotalPurchases | Input | Interval | Total number of purchases is directly related to customer engagement and churn. |
| TotalSpent | Input | Interval | Total expenditure could indicate the value of the customer and correlate with churn. |

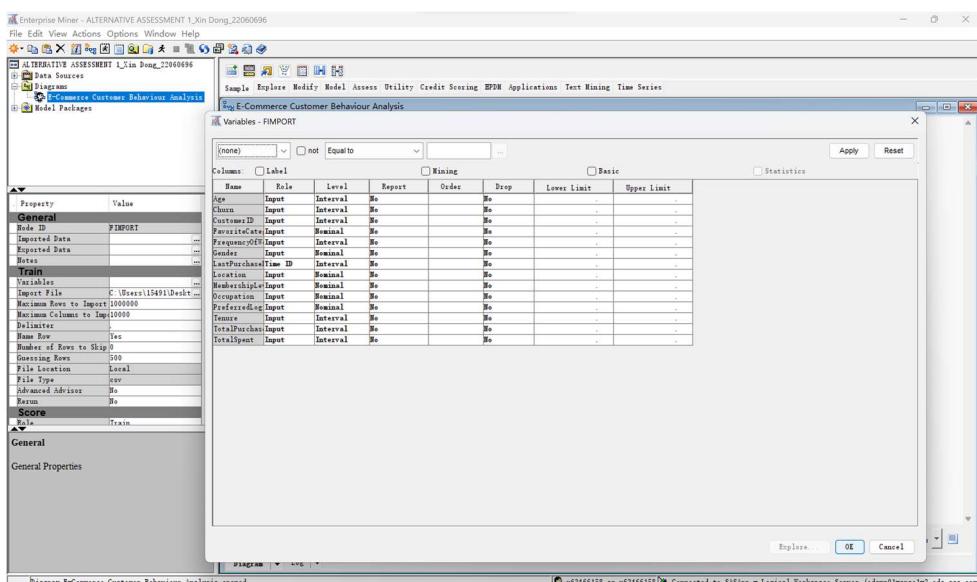
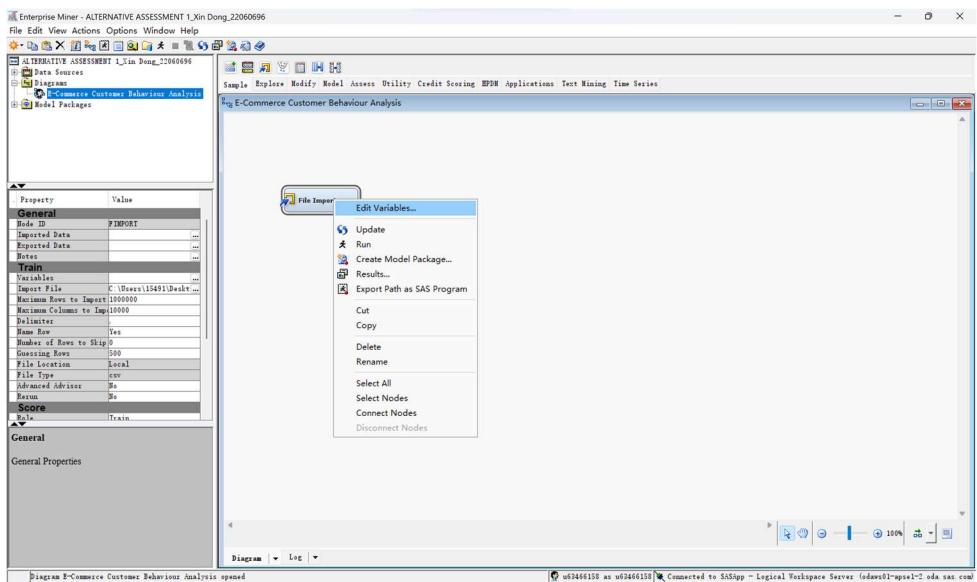
3. Specify Variable Roles and Levels

Set Target Variable: Ensure that the 'Churn' variable is set as the target.

Input Variables: Set other variables like 'Age', 'TotalSpent', etc., as input.

ID Variable: Set 'CustomerID' as the ID variable.

Click on ‘Edit Variables’, as shown in below is the original status of the variables in the dataset:



Specify the variable roles and levels, as shown in the colorful boxes below:

| Variables - FIMPORT | | | | | | | | |
|--------------------------|--------------------------------|----------|----------|----------------------|------------------------------------|---------------------------------|--------------------------------|-------------------------------------|
| Variables - FIMPORT | | | | | | | | |
| Variables - FIMPORT | | | | | | | | |
| (none) ▾ | | | | | | | | |
| (none) | <input type="checkbox"/> | not | Equal to | <input type="text"/> | <input type="button" value="..."/> | <input type="checkbox"/> Mining | <input type="checkbox"/> Basic | <input type="checkbox"/> Statistics |
| Columns: | <input type="checkbox"/> Label | | | | | | | |
| Name | Role | Level | Report | Order | Drop | Lower Limit | Upper Limit | |
| Age | Input | Interval | No | No | - | - | - | |
| Churn | Input | Interval | No | No | - | - | - | |
| CustomerID | Input | Interval | No | No | - | - | - | |
| FavoriteCategory | Input | Nominal | No | No | - | - | - | |
| FrequencyOfWebsiteVisits | Input | Interval | No | No | - | - | - | |
| Gender | Input | Nominal | No | No | - | - | - | |
| LastPurchaseDate | Time ID | Interval | No | No | - | - | - | |
| Location | Input | Nominal | No | No | - | - | - | |
| MembershipLevel | Input | Nominal | No | No | - | - | - | |
| Occupation | Input | Nominal | No | No | - | - | - | |
| PreferredLoginDevice | Rejected | Nominal | No | No | - | - | - | |
| Tenure | Input | Interval | No | No | - | - | - | |
| TotalPurchases | Input | Interval | No | No | - | - | - | |
| TotalSpent | Input | Interval | No | No | - | - | - | |

Churn (Target): This variable has been set from ‘input’ into ‘target’ because it is the outcome we are trying to predict. In customer churn analysis, the goal is to determine which factors influence whether a customer will stop making purchases. Hence, it should be the dependent variable in our predictive model. And the Level should be ‘Binary’.

CustomerID (ID): Customer ID is a unique identifier for each customer and does not contain predictive information that can be generalized to unseen data. Therefore, it's appropriately set as an ‘ID’ variable to index the records without using it as part of the prediction.

PreferredLoginDevice (Rejected): There could be several reasons to reject this variable: If there's a high proportion of missing values that are not easily imputable, it may not be reliable for use in the model. The variable may not show significant variance or impact in exploratory data analysis, suggesting it doesn't contribute much to the prediction of churn. The variable could have too many categories, leading to a sparse representation that is not useful for the model. It might be highly correlated with another variable, and to avoid multicollinearity, choose to exclude it.

LastPurchaseDate (Time ID): This variable can be treated as a time identifier. If we're conducting a time-series analysis or the model needs to account for the time aspect of customer behavior, then it's necessary to indicate when the last purchase was made. It would not be an input variable if the model doesn't use time-series techniques.

The rest of the variables (Input): These variables are likely to contain information that can contribute to predicting customer churn. They are likely to be features that vary between customers and can be used to identify patterns associated with the churn.

4. Data Pre-processing with SAS Enterprise Miner & Talend Data Integration & Talend Data Preparation

4.1 Option one: Utilize SAS Enterprise Miner for data preprocessing

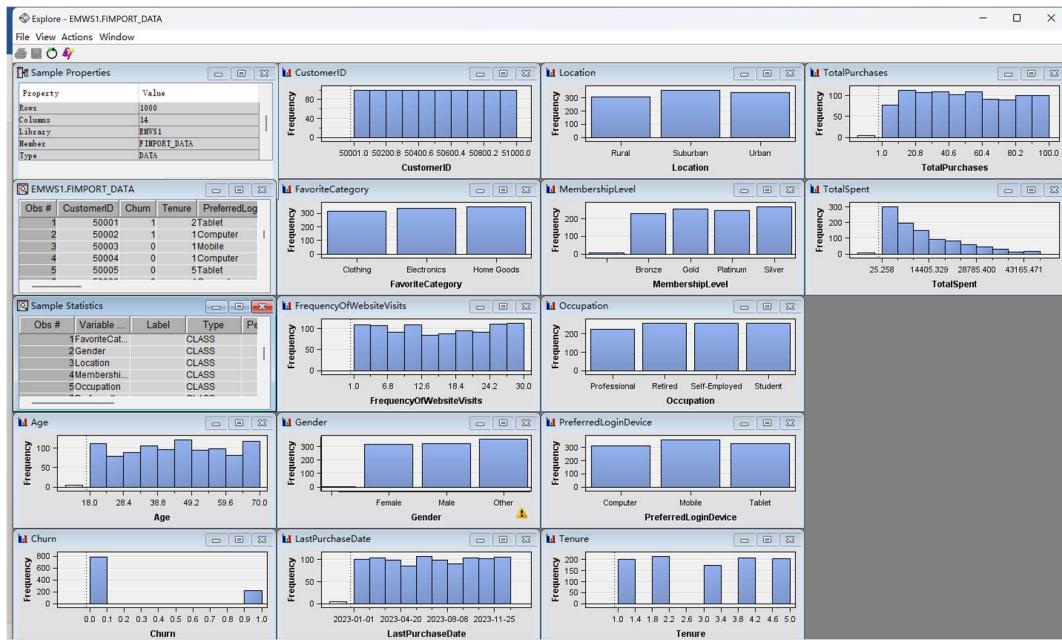
Identify Missing Values: Use the Explore node to identify the missing values in the dataset.

Decide on a Strategy: For handling missing values, decide between different approaches such as deletion, imputation, or using algorithms that can handle missing values.

Deletion: If the missing values are insignificant, we can choose to remove those rows.

Imputation: If the missing values are significant, consider imputation techniques.

For numeric variables, we can use mean/median imputation, and for categorical variables, mode imputation or more sophisticated methods like k-nearest neighbors (KNN) can be used.



As we can observe from the table and plots and summarize the missing values and dealing methods:

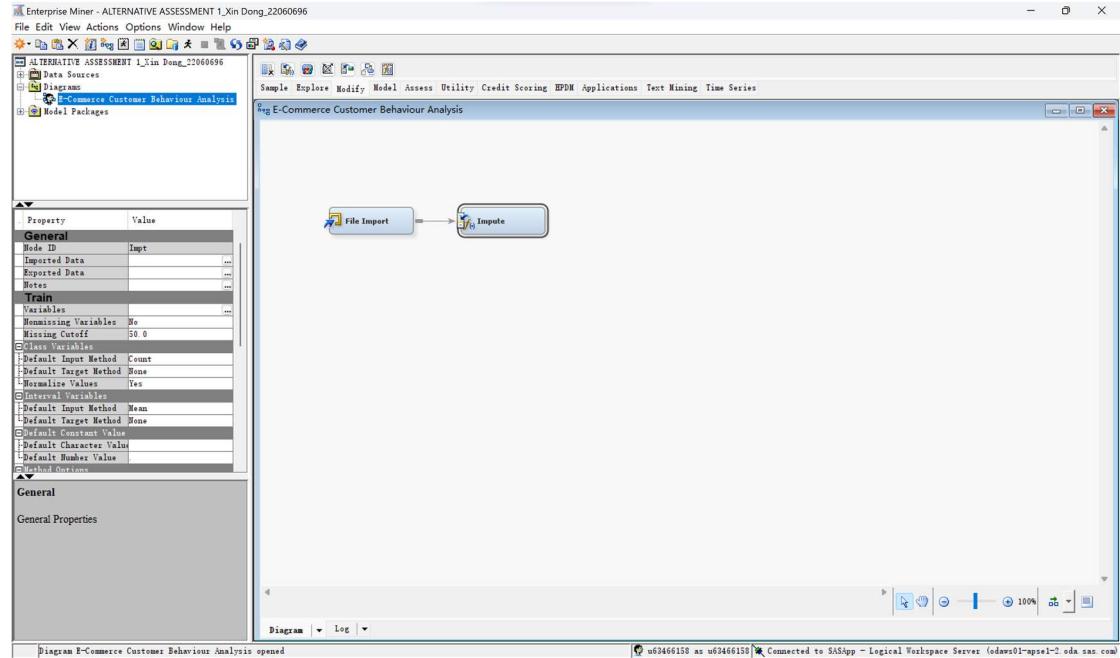
| Variable | Role | Level | Missing Values |
|--------------------------|---------|----------|---------------------|
| Age | Input | Interval | Impute/Median |
| Churn | Target | Binary | None |
| CustomerID | ID | Nominal | None |
| FavoriteCategory | Input | Nominal | Impute/Distribution |
| FrequencyOfWebsiteVisits | Input | Interval | Impute/Mean |
| Gender | Input | Nominal | Impute/Distribution |
| LastPurchaseDate | Time ID | Interval | None |
| Location | Input | Nominal | Impute/Distribution |

| | | | |
|----------------------|----------|----------|---------------------|
| MembershipLevel | Input | Nominal | Impute/Distribution |
| Occupation | Input | Nominal | Impute/Distribution |
| PreferredLoginDevice | Rejected | Nominal | - |
| Tenure | Input | Interval | Impute/Median |
| TotalPurchases | Input | Interval | Impute/Median |
| TotalSpent | Input | Interval | Impute/Mean |

Using the Impute Node:

Drag and drop the Impute node from the Modify tab onto the diagram workspace.

Connect our data source node to the Impute node.



Right-click the Impute node to configure the variables.

For each variable with missing values, you can select an imputation method based on the attribute table:

Mean/Median: For continuous variables like 'Age', 'TotalPurchases', and 'TotalSpent'.

Distribution: For categorical variables like 'FavoriteCategory', 'Gender', and 'MembershipLevel'.

After configuring, run the Impute node to apply the changes.

| Columns: | | <input type="checkbox"/> Label | <input type="checkbox"/> Mining | <input type="checkbox"/> Basic |
|--------------------------|---------|--------------------------------|---------------------------------|--------------------------------|
| Name | Use | Method | Use Tree | Role |
| Age | Yes | Median | Default | Input |
| Churn | Default | None | Default | Target |
| FavoriteCategory | Yes | Distribution | Default | Input |
| FrequencyOfWebsiteVisits | Yes | Mean | Default | Input |
| Gender | Yes | Distribution | Default | Input |
| Location | Yes | Distribution | Default | Input |
| MembershipLevel | Yes | Distribution | Default | Input |
| Occupation | Yes | Distribution | Default | Input |
| PreferredLoginDevice | No | None | Default | Rejected |
| Tenure | Yes | Median | Default | Input |
| TotalPurchases | Yes | Median | Default | Input |
| TotalSpent | Yes | Mean | Default | Input |

And then we can check the results of handling missing values:

| Variable Name | Impute Method | Imputed Variable | Impute Value | Role | Measurement Level | Label | Number of Missing for TRAIN |
|-----------------|---------------|---------------------|---------------|-------|-------------------|-------|-----------------------------|
| Age | MEDIAN | IMP_Age | 44INPUT | INPUT | INTERVAL | | 5 |
| Gender | DISTRIBUTION | IMP_Gender | | INPUT | NOMINAL | | 5 |
| MembershipLevel | DISTRIBUTION | IMP_MembershipLevel | | INPUT | NOMINAL | | 5 |
| TotalPurchases | MEDIAN | IMP_TotalPurchases | 50INPUT | INPUT | INTERVAL | | 5 |
| TotalSpent | MEAN | IMP_TotalSpent | 12698.71INPUT | INPUT | INTERVAL | | 5 |

| Output | | | | | | | |
|--|--|--|--|--|--|--|--|
| 31 | | | | | | | |
| 32 | | | | | | | |
| 33 | | | | | | | |
| 34 Imputation Summary | | | | | | | |
| 35 Number Of Observations | | | | | | | |
| 36 | | | | | | | |
| 37 | | | | | | | |
| 38 Impute Impute Measurement Number of 39 Variable Name Method Imputed Variable Value Role Level Label Missing 40 | | | | | | | |
| 41 Age MEDIAN IMP_Age 44.00 INPUT INTERVAL 5 42 Gender DISTRIBUTION IMP_Gender . INPUT NOMINAL 5 43 MembershipLevel DISTRIBUTION IMP_MembershipLevel . INPUT NOMINAL 5 44 TotalPurchases MEDIAN IMP_TotalPurchases 50.00 INPUT INTERVAL 5 45 TotalSpent MEAN IMP_TotalSpent 12698.71 INPUT INTERVAL 5 46 | | | | | | | |
| 47 | | | | | | | |
| 48 | | | | | | | |
| 49 | | | | | | | |
| 50 Variable Distribution Training Data | | | | | | | |
| 51 | | | | | | | |
| 52 Number of 53 Missing Number of Percent of 54 Obs for TRAIN Variables Variables | | | | | | | |
| 55 | | | | | | | |
| 56 1 5 5 100 | | | | | | | |
| 57 | | | | | | | |

Age:

Method: Median

Imputed Value: 44.00

Interpretation: Missing ages have been replaced with the median age value of 44, which is robust against skewed distributions and outliers.

Gender:

Method: Distribution

Interpretation: The missing values in the 'Gender' variable were filled in following the existing distribution of the 'Gender' categories in the dataset. This preserves the proportionality of each category.

MembershipLevel:

Method: Distribution

Interpretation: Similar to 'Gender', the missing values for 'MembershipLevel' were imputed based on the current distribution of membership levels among customers, maintaining the original category proportions.

TotalPurchases:

Method: Median

Imputed Value: 50.00

Interpretation: The missing values for 'TotalPurchases' were replaced with the median value of 50, which avoids the influence of extreme values and is a typical approach for count data that may be

right-skewed.

TotalSpent:

Method: Mean

Imputed Value: 12,698.71

Interpretation: The mean value of 12,698.71 was used to impute missing values in 'TotalSpent'.

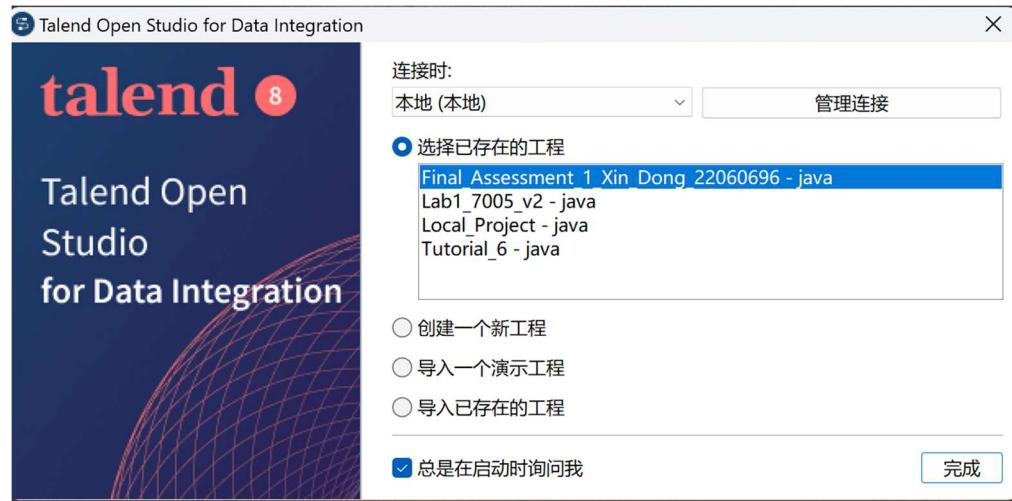
This suggests that on average, customers spent this amount, and it assumes that the spending is normally distributed.

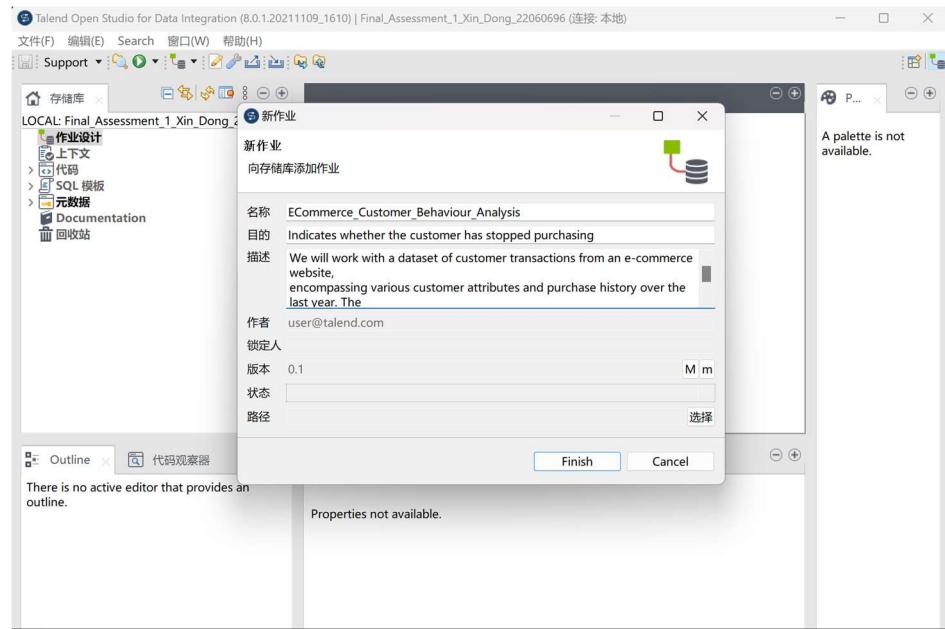
Each of these imputation methods was selected based on the nature of the variable and aimed to minimize the bias that missing data could introduce to the analysis. The number of missing values handled for each variable is consistent 5.

4.2 Another Option: Utilize Talend Data Integration & Talend Data Preparation to do Data Pre-Processing

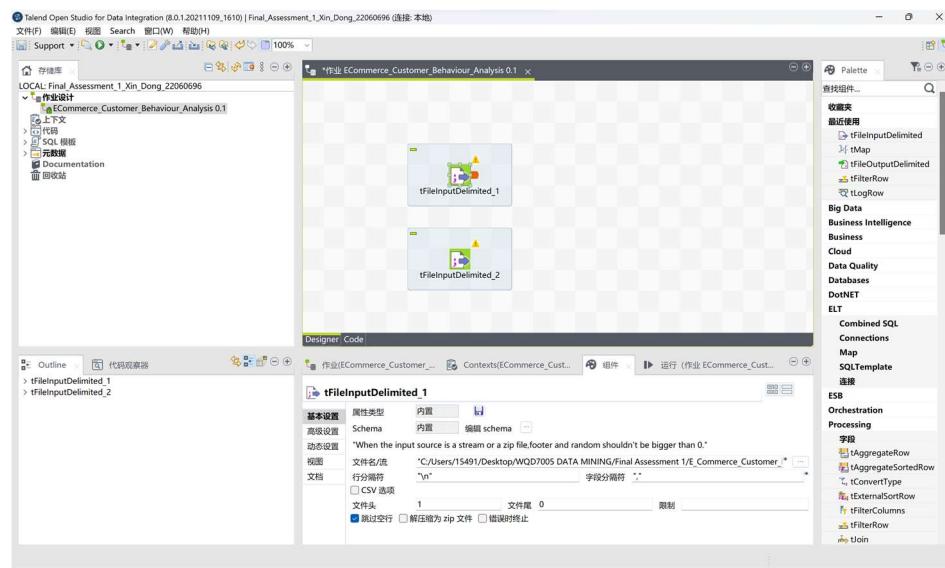
4.2.1 Talend Data Preparation:

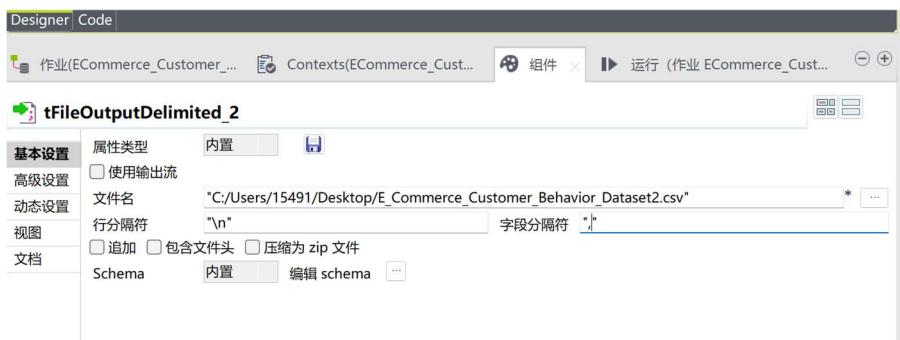
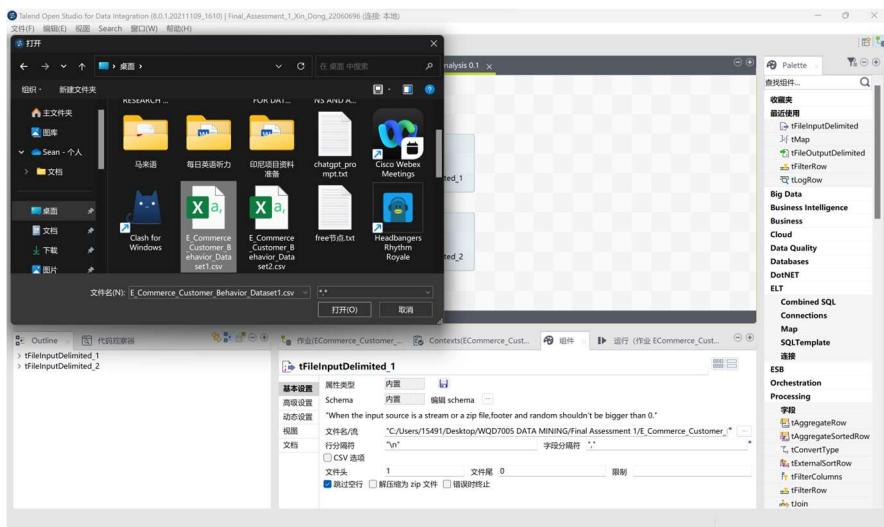
Create a project



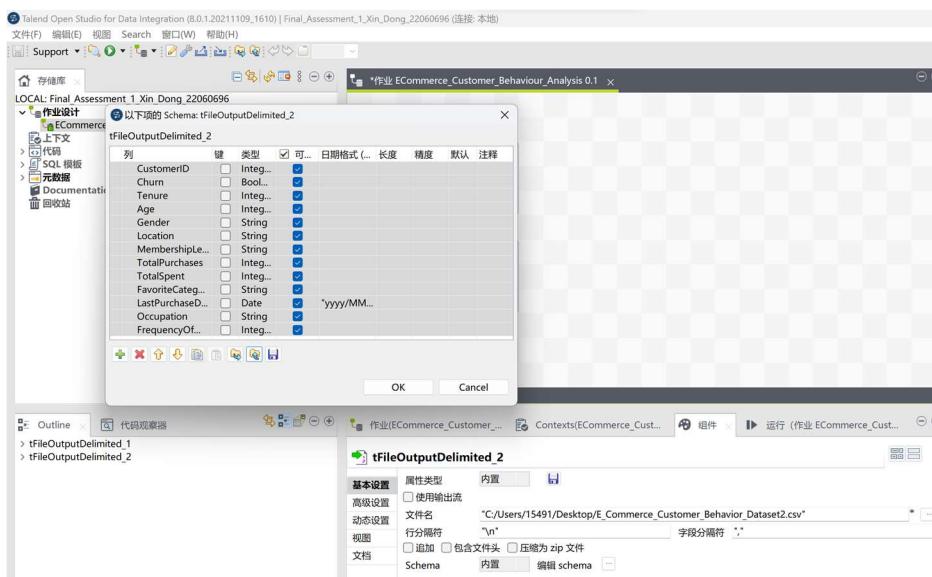
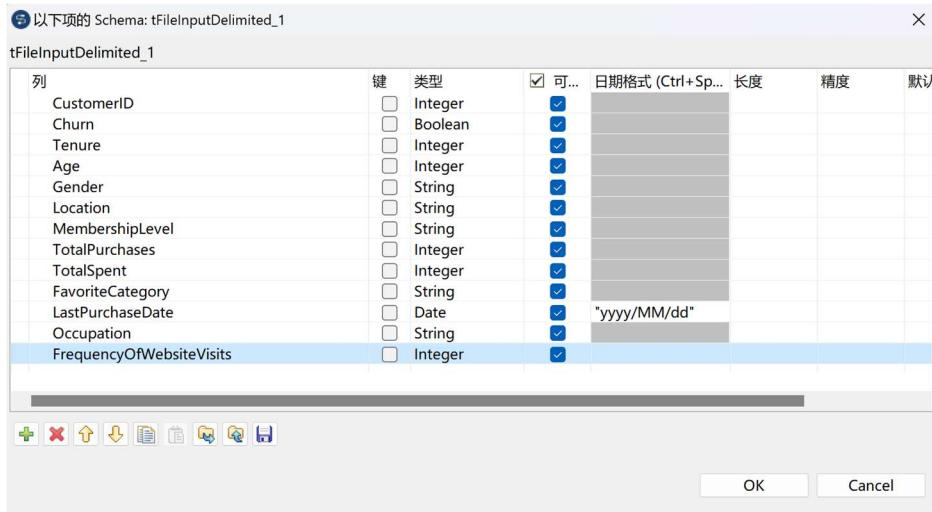


Import the 2 files to be integrated and set configuration:

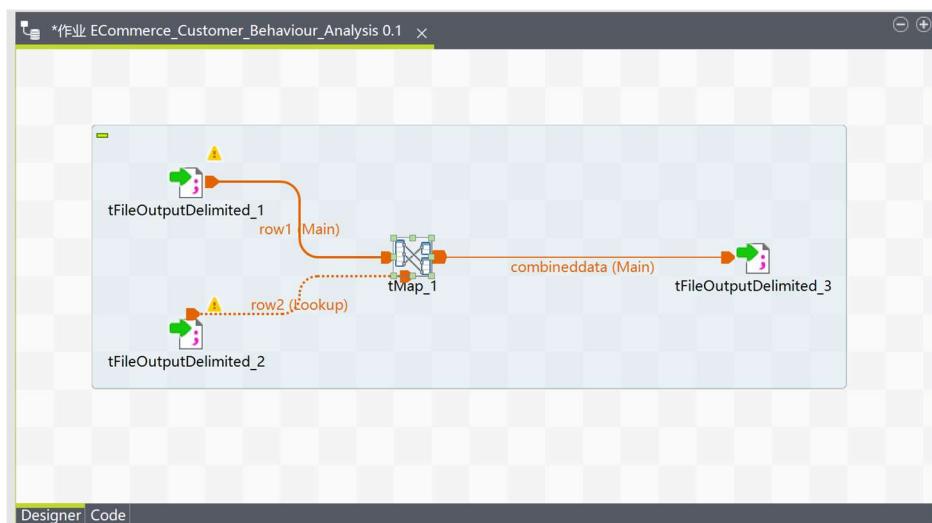




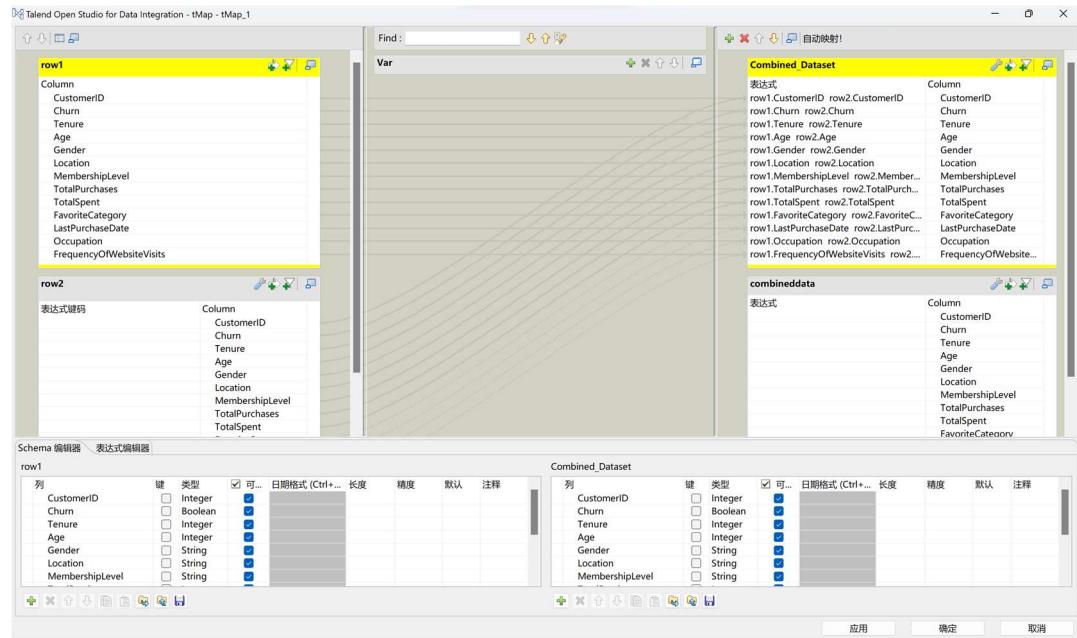
Set schema for the 2 files separately:



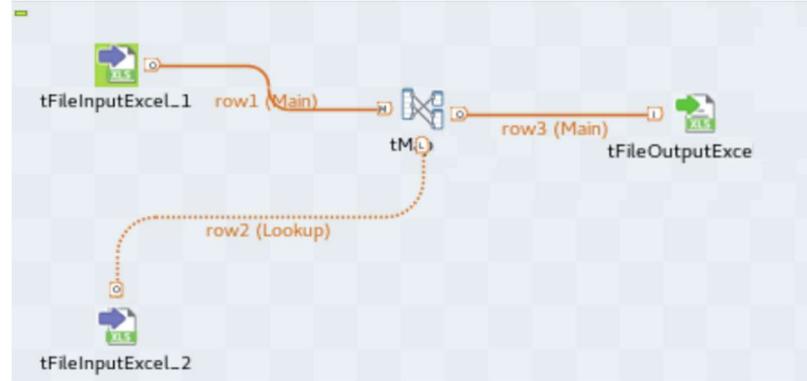
And then we put tMap and set tMap configuration, and finally set the tFileOutput.



tMap configuration:

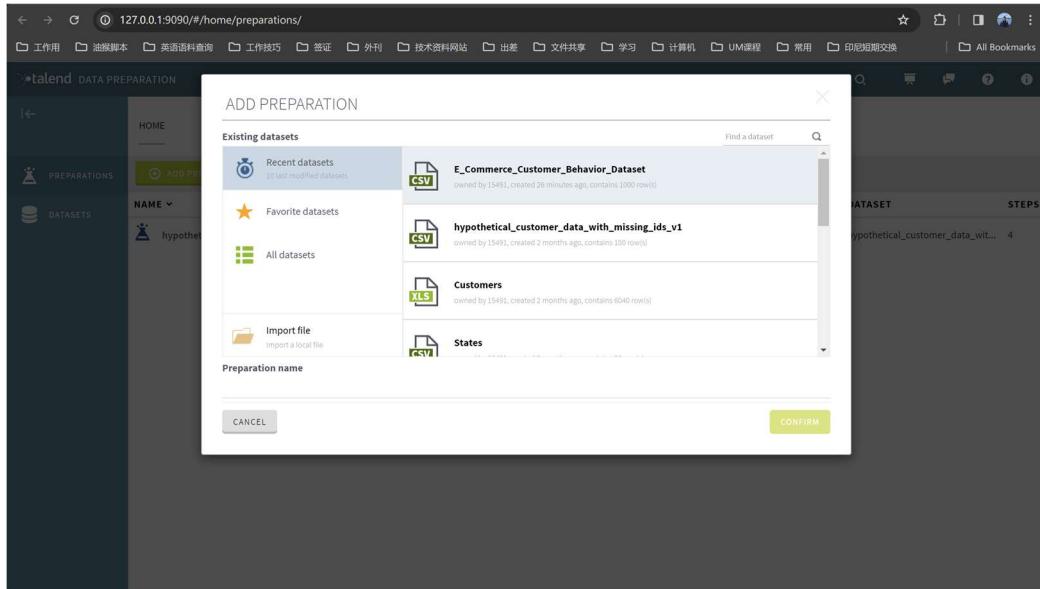


We changed the name and get the final output pic:



So we successfully combined 2 files of dataset into one.

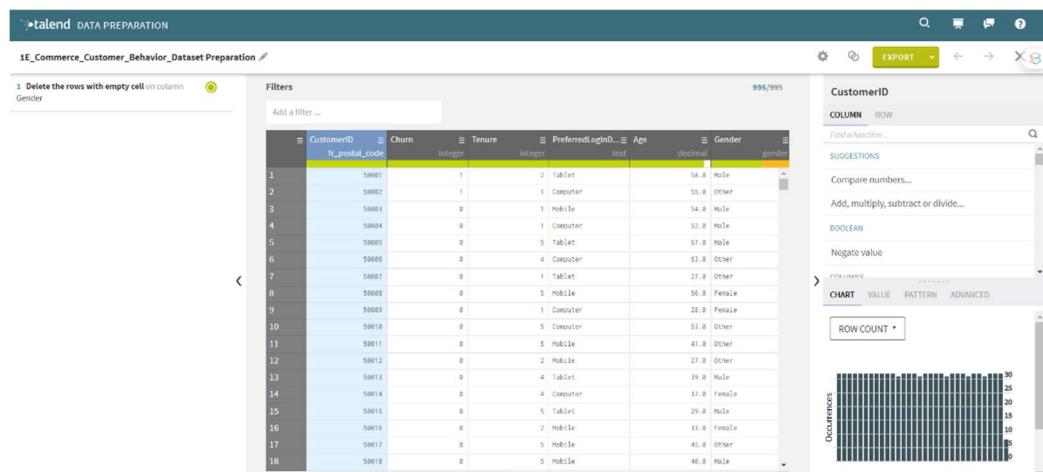
4.2.2 Talend Data Preparation:



| PREPARATIONS | | | | | | |
|--|--------|-------------------|-------------------|-------------------------------|-------|--|
| NAME | AUTHOR | CREATED | MODIFIED | DATASET | STEPS | |
| E_Commerce_Customer_Behavior_Dataset Preparation | 15491 | a few seconds ago | a few seconds ago | E_Commerce_Customer_Beh... | 0 | |
| hypothetical_customer_data_with_missing_ids_v1 Preparation | 15491 | 2 months ago | 2 months ago | hypothetical_customer_data... | 4 | |

Find patterns and visualization for all the columns/variables:

We can observe the white gap for each variables represents the missing values.



Filters

Add a filter ...

| | Age | Gender | Location | MembershipLevel | TotalPurchases | TotalSpent |
|----|------|--------|----------|-----------------|----------------|--------------------|
| 1 | 58.0 | Male | Suburban | Platinum | 49.0 | 10237.707351897 |
| 2 | 55.0 | Other | Rural | Platinum | 20.0 | 633.7578459192 |
| 3 | 54.0 | Male | Suburban | Bronze | 87.0 | 6106.863186448 |
| 4 | 52.0 | Male | Rural | Platinum | 50.0 | 4759.40720763 |
| 5 | 67.0 | Male | Rural | Silver | 16.0 | 7425.91514997 |
| 6 | 63.0 | Other | Urban | Gold | 59.0 | 26165.38785422 |
| 7 | 27.0 | Other | Rural | Silver | 15.0 | 2508.686097404 |
| 8 | 56.0 | Female | Rural | Bronze | 12.0 | 4461.209360231 |
| 9 | 28.0 | Female | Urban | Bronze | 9.0 | 4465.59358993 |
| 10 | 53.0 | Other | Rural | Bronze | 40.0 | 18588.27236745 |
| 11 | 41.0 | Other | Rural | Platinum | 24.0 | 4312.245401668 |
| 12 | 27.0 | Other | Suburban | Silver | 87.0 | 27849.1292795 |
| 13 | 39.0 | Male | Suburban | Gold | 19.0 | 8552.431767280 |
| 14 | 37.0 | Female | Rural | Bronze | 17.0 | 5453.107616315 |
| 15 | 29.0 | Male | Suburban | Bronze | 96.0 | 42098.59480281 |
| 16 | 33.0 | Female | Suburban | Silver | 27.0 | 5169.53980354545 |
| 17 | 45.0 | Other | Suburban | Bronze | 54.0 | 13863.129821054 |
| 18 | 40.0 | Male | Urban | Silver | 100.0 | 22410.543001986385 |

Age

COLUMN ROW

Find a function ...

SUGGESTIONS

Delete the rows with empty cell

Fill empty cells with text...

Compare numbers...

Add, multiply, subtract or divide...

Round value using halfup mode...

CHART VALUE PATTERN ADVANCED

Count: 995 Min: 18

Distinct: 54 Max: 70

Duplicate: 941 Mean: 44.37

Valid: 990 Variance: 228.36

Empty: 5 Median: 44

Invalid: 0 Lower quantile: 32

Lower quantile: 32

Upper quantile: 57

Filters

Add a filter ...

| | Location | MembershipLevel | TotalPurchases | TotalSpent | FavoriteCategory | LastPurchase |
|----|----------|-----------------|----------------|--------------------|------------------|--------------|
| 1 | Suburban | Platinum | 49.0 | 10237.70735189774 | Home Goods | 2023-02- |
| 2 | Rural | Platinum | 20.0 | 633.7578459192816 | Electronics | 2023-07- |
| 3 | Suburban | Bronze | 87.0 | 6106.863186448275 | Home Goods | 2023-10- |
| 4 | Rural | Platinum | 50.0 | 4759.40720763638 | Home Goods | 2023-04- |
| 5 | Rural | Silver | 16.0 | 7425.91514997109 | Home Goods | 2023-11- |
| 6 | Urban | Gold | 59.0 | 26165.3878542678 | Home Goods | 2023-05- |
| 7 | Rural | Silver | 15.0 | 2508.686097404642 | Home Goods | 2023-02- |
| 8 | Rural | Bronze | 12.0 | 4461.209360231154 | Clothing | 2023-08- |
| 9 | Urban | Bronze | 9.0 | 4465.593589931341 | Electronics | 2023-09- |
| 10 | Rural | Bronze | 48.0 | 18588.272306745566 | Home Goods | 2023-04- |
| 11 | Rural | Platinum | 24.0 | 4312.245401668005 | Home Goods | 2023-04- |
| 12 | Suburban | Silver | 87.0 | 27849.12927955158 | Clothing | 2023-12- |
| 13 | Suburban | Gold | 19.0 | 8552.431767280506 | Home Goods | 2023-10- |
| 14 | Rural | Bronze | 17.0 | 5453.107616315044 | Electronics | 2023-02- |
| 15 | Suburban | Bronze | 96.0 | 42098.59480281735 | Electronics | 2023-12- |
| 16 | Suburban | Silver | 27.0 | 5169.539803545292 | Home Goods | 2023-04- |
| 17 | Suburban | Bronze | 54.0 | 13863.129821054747 | Home Goods | 2023-05- |
| 18 | Urban | Silver | 100.0 | 22410.543001986385 | Home Goods | 2023-09- |

MembershipLevel

COLUMN ROW

Find a function ...

SUGGESTIONS

Delete the rows with empty cell

Fill empty cells with text...

Change to upper case

Replace the cells that match...

Change to lower case

CHART VALUE PATTERN ADVANCED

ROW COUNT

Silver

Gold

Platinum

Bronze

(EMPTY)

talend DATA PREPARATION

1_E_Commerce_Customer_Behavior_Dataset Preparation

1 Delete the rows with empty cell on column Gender

Filters

Add a filter ...

| | CustomerID | Churn | Tenure | PreferredLoginDevice | Age | Gender |
|----|------------|-------|--------|----------------------|------|--------|
| 1 | 50001 | 1 | 2 | Tablet | 58.0 | Male |
| 2 | 50002 | 1 | 1 | Computer | 55.0 | Other |
| 3 | 50003 | 0 | 1 | Mobile | 54.0 | Male |
| 4 | 50004 | 0 | 1 | Computer | 52.0 | Male |
| 5 | 50005 | 0 | 5 | Tablet | 47.0 | Male |
| 6 | 50006 | 0 | 4 | Computer | 63.0 | Other |
| 7 | 50007 | 0 | 1 | Tablet | 27.0 | Other |
| 8 | 50008 | 0 | 5 | Mobile | 56.0 | Female |
| 9 | 50009 | 0 | 1 | Computer | 28.0 | Female |
| 10 | 50010 | 0 | 5 | Computer | 51.0 | Other |
| 11 | 50011 | 0 | 5 | Mobile | 41.0 | Other |
| 12 | 50012 | 0 | 2 | Mobile | 27.0 | Other |
| 13 | 50013 | 0 | 4 | Tablet | 33.0 | Male |
| 14 | 50014 | 0 | 4 | Computer | 37.0 | Female |
| 15 | 50015 | 0 | 5 | Tablet | 29.0 | Male |
| 16 | 50016 | 0 | 2 | Mobile | 31.0 | Female |
| 17 | 50017 | 0 | 5 | Mobile | 45.0 | Other |
| 18 | 50018 | 0 | 5 | Mobile | 48.0 | Male |

Churn

COLUMN ROW

Find a function ...

SUGGESTIONS

Compare numbers...

Add, multiply, subtract or divide...

BOOLEAN

Negate value

CHART VALUE PATTERN ADVANCED

ROW COUNT

Occurrences

700

600

500

400

300

200

100

0

Filters

Add a filter ...

| | gender | Location | MembershipLevel | TotalPurchases | TotalSpent | FavoriteCategory | LastPurchase |
|----|--------|----------|-----------------|----------------|--------------------|------------------|--------------|
| | | text | city | decimal | decimal | text | |
| 1 | | Suburban | Platinum | 49.0 | 10237.70735189774 | Home Goods | 2023-02-12 |
| 2 | | Rural | Platinum | 20.0 | 633.7578459192816 | Electronics | 2023-07-11 |
| 3 | | Suburban | Bronze | 87.0 | 6106.8631806448275 | Home Goods | 2023-10-25 |
| 4 | | Rural | Platinum | 50.0 | 4759.40720763638 | Home Goods | 2023-04-19 |
| 5 | | Rural | Silver | 16.0 | 7425.91514097109 | Home Goods | 2023-11-01 |
| 6 | | Urban | Gold | 59.0 | 26165.38785422678 | Home Goods | 2023-05-01 |
| 7 | | Rural | Silver | 15.0 | 2508.606097404642 | Home Goods | 2023-02-12 |
| 8 | | Rural | Bronze | 12.0 | 4461.209360291154 | Clothing | 2023-08-13 |
| 9 | | Urban | Bronze | 9.0 | 4465.59358993141 | Electronics | 2023-09-01 |
| 10 | | Rural | Bronze | 40.0 | 18588.272306745566 | Home Goods | 2023-04-01 |
| 11 | | Rural | Platinum | 24.0 | 4312.245401668005 | Home Goods | 2023-04-19 |
| 12 | | Suburban | Silver | 87.0 | 27849.129297995158 | Clothing | 2023-12-17 |
| 13 | | Suburban | Gold | 19.0 | 8552.431767280506 | Home Goods | 2023-10-01 |
| 14 | | Rural | Bronze | 17.0 | 5453.107616315044 | Electronics | 2023-02-12 |
| 15 | | Suburban | Bronze | 96.0 | 42098.59480281735 | Electronics | 2023-12-01 |
| 16 | | Suburban | Silver | 27.0 | 5169.539803545292 | Home Goods | 2023-04-01 |
| 17 | | Suburban | Bronze | 54.0 | 13863.129821054747 | Home Goods | 2023-05-01 |
| 18 | | Urban | Silver | 100.0 | 22410.543001986385 | Home Goods | 2023-09-01 |

Filters

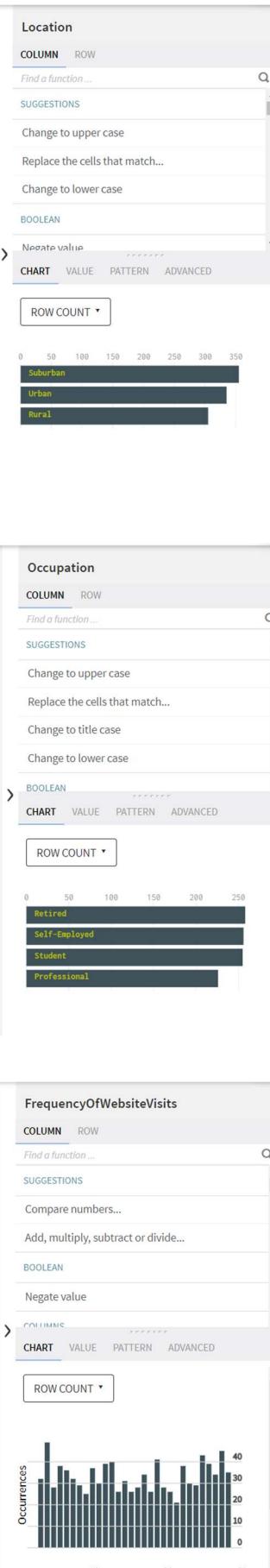
Add a filter ...

| | allPurchases | TotalSpent | FavoriteCategory | LastPurchaseDate | Occupation | FrequencyOfWe... |
|----|--------------|--------------------|------------------|------------------|---------------|------------------|
| | decimal | decimal | text | date | text | Integer |
| 1 | 49.0 | 10237.70735189774 | Home Goods | 2023-02-12 | Student | 7 |
| 2 | 20.0 | 633.7578459192816 | Electronics | 2023-07-11 | Professional | 14 |
| 3 | 87.0 | 6106.8631806448275 | Home Goods | 2023-10-25 | Self-Employed | 8 |
| 4 | 50.0 | 4759.40720763638 | Home Goods | 2023-04-19 | Student | 26 |
| 5 | 16.0 | 7425.91514097109 | Home Goods | 2023-11-09 | Professional | 25 |
| 6 | 59.0 | 26165.38785422678 | Home Goods | 2023-05-08 | Student | 12 |
| 7 | 15.0 | 2508.606097404642 | Home Goods | 2023-02-27 | Self-Employed | 3 |
| 8 | 12.0 | 4461.209360291154 | Clothing | 2023-08-13 | Self-Employed | 12 |
| 9 | 9.0 | 4465.59358993141 | Electronics | 2023-09-08 | Self-Employed | 3 |
| 10 | 40.0 | 18588.272306745566 | Home Goods | 2023-04-06 | Retired | 25 |
| 11 | 24.0 | 4312.245401668005 | Home Goods | 2023-04-23 | Retired | 27 |
| 12 | 87.0 | 27849.129297995158 | Clothing | 2023-12-17 | Retired | 30 |
| 13 | 19.0 | 8552.431767280506 | Home Goods | 2023-10-12 | Professional | 26 |
| 14 | 17.0 | 5453.107616315044 | Electronics | 2023-02-13 | Self-Employed | 20 |
| 15 | 96.0 | 42098.59480281735 | Electronics | 2023-12-30 | Student | 13 |
| 16 | 27.0 | 5169.539803545292 | Home Goods | 2023-04-25 | Self-Employed | 18 |
| 17 | 54.0 | 13863.129821054747 | Home Goods | 2023-05-05 | Retired | 17 |
| 18 | 100.0 | 22410.543001986385 | Home Goods | 2023-09-01 | Retired | 14 |

Filters

Add a filter ...

| | allPurchases | TotalSpent | FavoriteCategory | LastPurchaseDate | Occupation | FrequencyOfWe... |
|----|--------------|--------------------|------------------|------------------|---------------|------------------|
| | decimal | decimal | text | date | text | integer |
| 1 | 49.0 | 10237.70735189774 | Home Goods | 2023-02-12 | Student | 7 |
| 2 | 20.0 | 633.7578459192816 | Electronics | 2023-07-11 | Professional | 14 |
| 3 | 87.0 | 6106.8631806448275 | Home Goods | 2023-10-25 | Self-Employed | 8 |
| 4 | 50.0 | 4759.40720763638 | Home Goods | 2023-04-19 | Student | 26 |
| 5 | 16.0 | 7425.91514097109 | Home Goods | 2023-11-09 | Professional | 25 |
| 6 | 59.0 | 26165.38785422678 | Home Goods | 2023-05-08 | Student | 12 |
| 7 | 15.0 | 2508.606097404642 | Home Goods | 2023-02-27 | Self-Employed | 3 |
| 8 | 12.0 | 4461.209360291154 | Clothing | 2023-08-13 | Self-Employed | 12 |
| 9 | 9.0 | 4465.59358993141 | Electronics | 2023-09-08 | Self-Employed | 3 |
| 10 | 40.0 | 18588.272306745566 | Home Goods | 2023-04-06 | Retired | 25 |
| 11 | 24.0 | 4312.245401668005 | Home Goods | 2023-04-23 | Retired | 27 |
| 12 | 87.0 | 27849.129297995158 | Clothing | 2023-12-17 | Retired | 30 |
| 13 | 19.0 | 8552.431767280506 | Home Goods | 2023-10-12 | Professional | 26 |
| 14 | 17.0 | 5453.107616315044 | Electronics | 2023-02-13 | Self-Employed | 20 |
| 15 | 96.0 | 42098.59480281735 | Electronics | 2023-12-30 | Student | 13 |
| 16 | 27.0 | 5169.539803545292 | Home Goods | 2023-04-25 | Self-Employed | 18 |
| 17 | 54.0 | 13863.129821054747 | Home Goods | 2023-05-05 | Retired | 17 |
| 18 | 100.0 | 22410.543001986385 | Home Goods | 2023-09-01 | Retired | 14 |



We can select the rows with valid values:

The screenshot shows a data analysis interface with a table of 986/995 rows. Three filters are applied:

- FrequencyOfWebsiteVisits:** rows with valid values
- LastPurchaseDate:** rows with valid values
- Age:** rows with valid values

The table columns include: id, Churn, Tenure, PreferredLoginD..., Age, Gender, and Location. The right panel displays a bar chart titled "PreferredLoginDevice" showing the count of rows for different devices: Mobile, Tablet, and Computer.

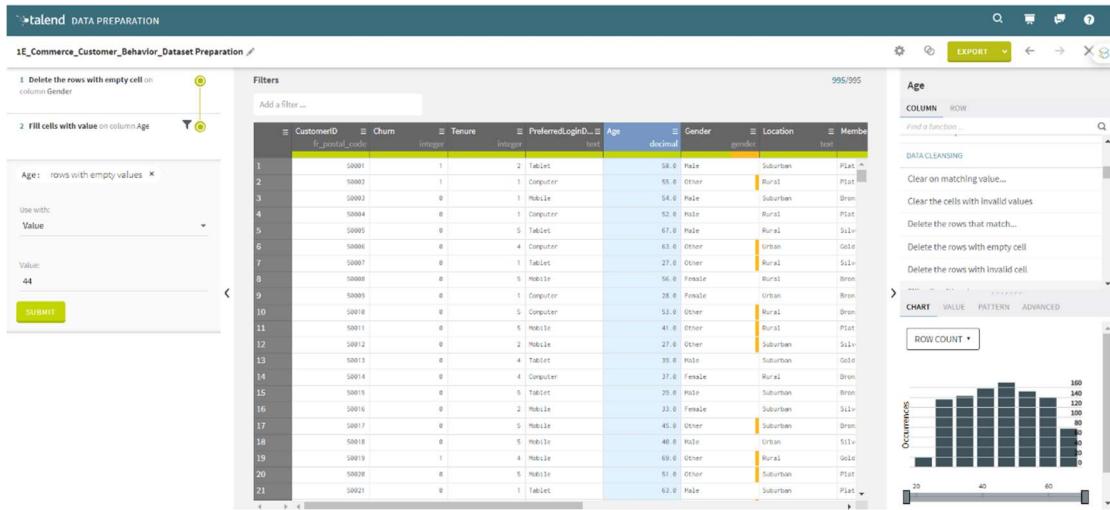
Sample of filling the missing values, including age (numeric), MembershipLevel (string), and other variables (delete unimportant rows):

First select the rows with empty cells.

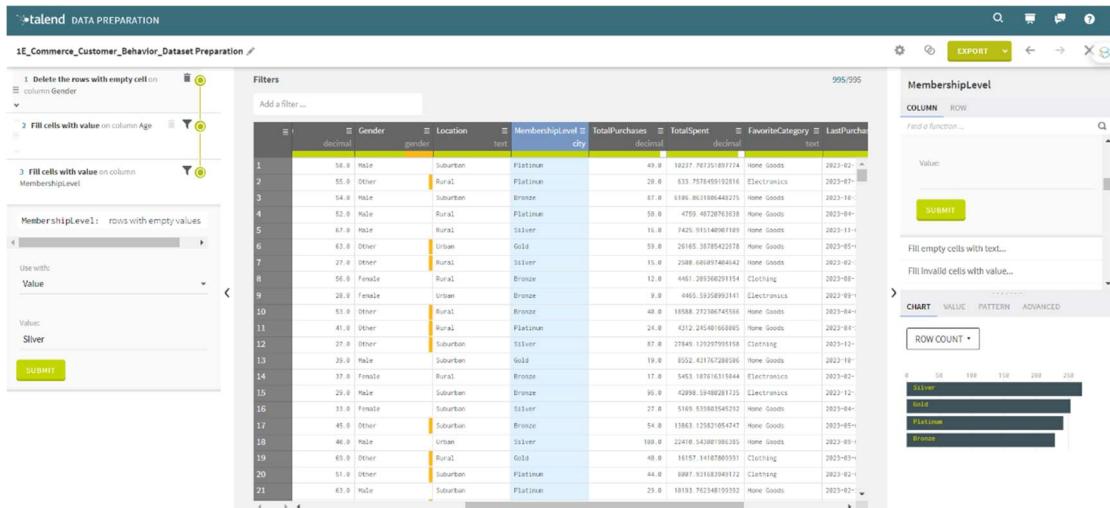
The screenshot shows a data analysis interface with a table of 5/995 rows. A filter is applied for rows with empty values in the Age column.

The table columns include: id, Churn, Tenure, PreferredLoginD..., Age, Gender, Location, and MembershipLevel. The right panel displays a histogram titled "Age" showing the occurrences of different age values.

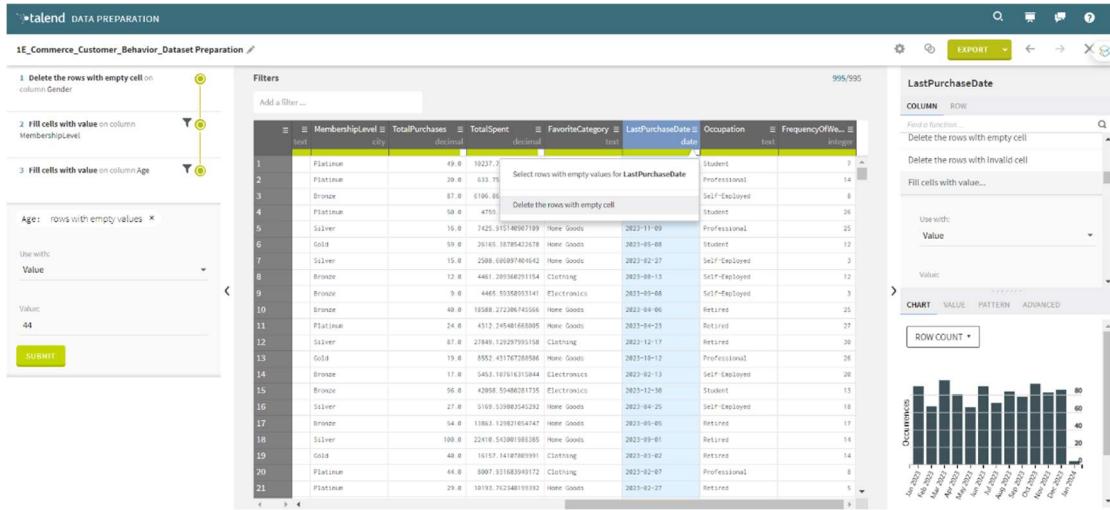
And then fill the missing values with median/mean number.



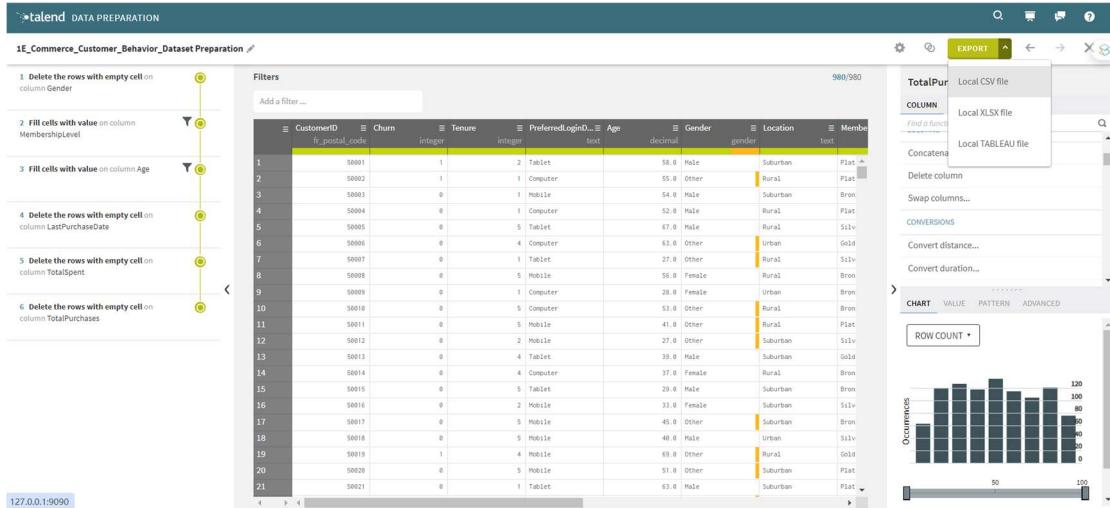
For the variables with missing value type of string/text, use the same method, impute the most frequent value based on the distribution.



Other variables, delete unimportant rows with empty values.

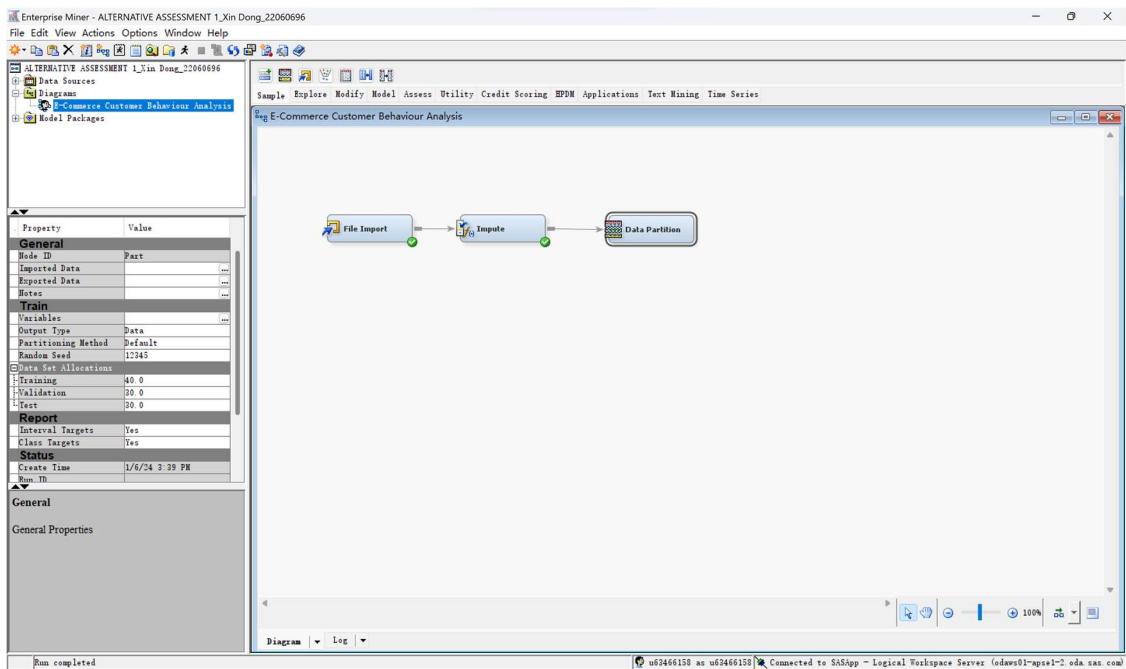


And then we finish the data pre-processing part.



5. Data Partition

Partition the Data: Use the Data Partition node to split your data into training, validation, and testing sets. A common split is 70% training, 15% validation, and 15% testing.



| Property | Value |
|-----------------------------|----------------|
| General | |
| Node ID | Part |
| Imported Data | ... |
| Exported Data | ... |
| Notes | ... |
| Train | |
| Variables | ... |
| Output Type | Data |
| Partitioning Method | Default |
| 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 | 1/6/24 3:39 PM |
| Run ID | ... |

As below shows the results of data partition:

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

File Edit View Window

Output

```

1   *
2   User:          u63466158
3   Date:          06 January 2024
4   Time:          15:41:50
5   *
6   * Training Output
7   *
8
9
10
11
12 Variable Summary
13
14     Measurement  Frequency
15     Role        Level      Count
16
17 ID         INTERVAL    1
18 INPUT      INTERVAL    5
19 INPUT      NOMINAL    5
20 REJECTED   NOMINAL    1
21 TARGET     INTERVAL    1
22 TIMEID    INTERVAL    1
23
24
25
26
27 Partition Summary
28
29
30     Type       Data Set      Number of
31
32 DATA       EMWS1.Impt_TRAIN 1000
33 TRAIN      EMWS1.Part_TRAIN 700
34 VALIDATE   EMWS1.Part_VALIDATE 150
35 TEST       EMWS1.Part_TEST 150
36
37

```

Output

```

46
47
48
49
50 Summary Statistics for Interval Targets
51
52 Data=DATA
53
54     Variable  Maximum  Mean  Minimum  Number of Observations  Standard Deviation  Label
55
56     Churn      1        0.214  0        1000            0           0.410332027
57
58
59
60 Data=TEST
61
62     Variable  Maximum  Mean  Minimum  Number of Observations  Standard Deviation  Label
63
64     Churn      1        0.2        0        150            0           0.4013400373
65
66
67
68 Data=TRAIN
69
70     Variable  Maximum  Mean  Minimum  Number of Observations  Standard Deviation  Label
71
72     Churn      1        0.21       0        700            0           0.4075994836
73
74
75
76 Data=VALIDATE
77
78     Variable  Maximum  Mean  Minimum  Number of Observations  Standard Deviation  Label
79
80     Churn      1        0.2466666667  0        150            0           0.4325151457
81
82

```

6. Decision Tree Analysis

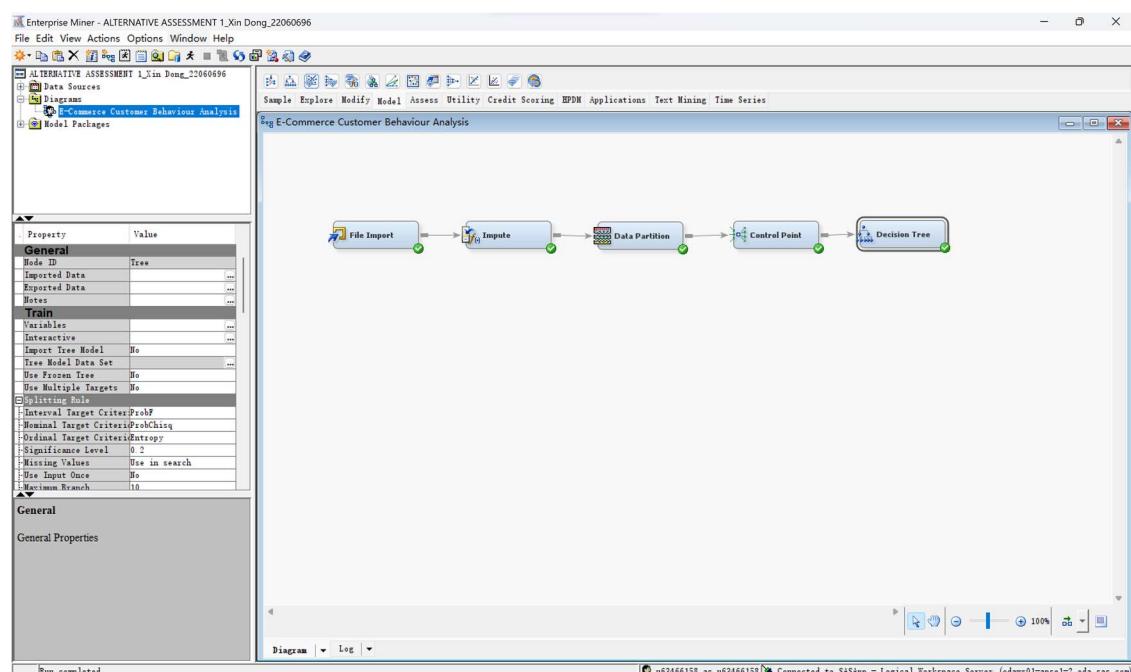
Step 1: Setting Up the Decision Tree Node

Open our project within SAS Enterprise Miner.

Drag and drop the Decision Tree node from the toolbar onto the workspace.

Connect the Data Partition node (which has the preprocessed data) to the Decision Tree node.

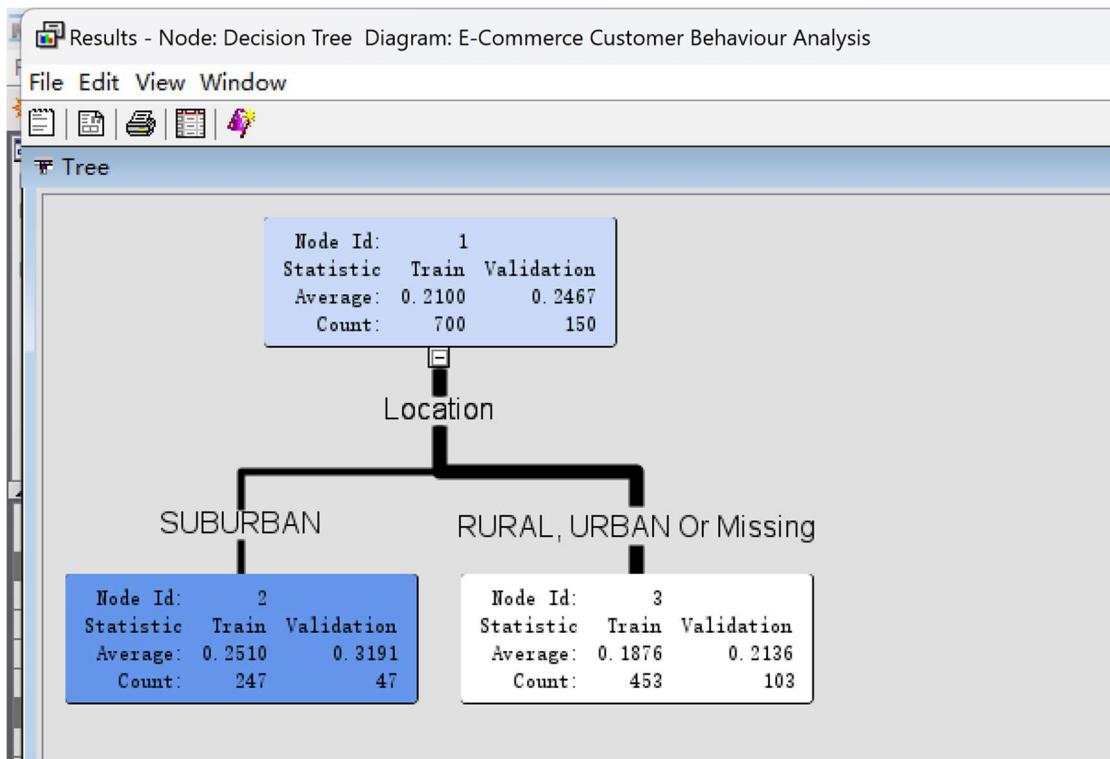
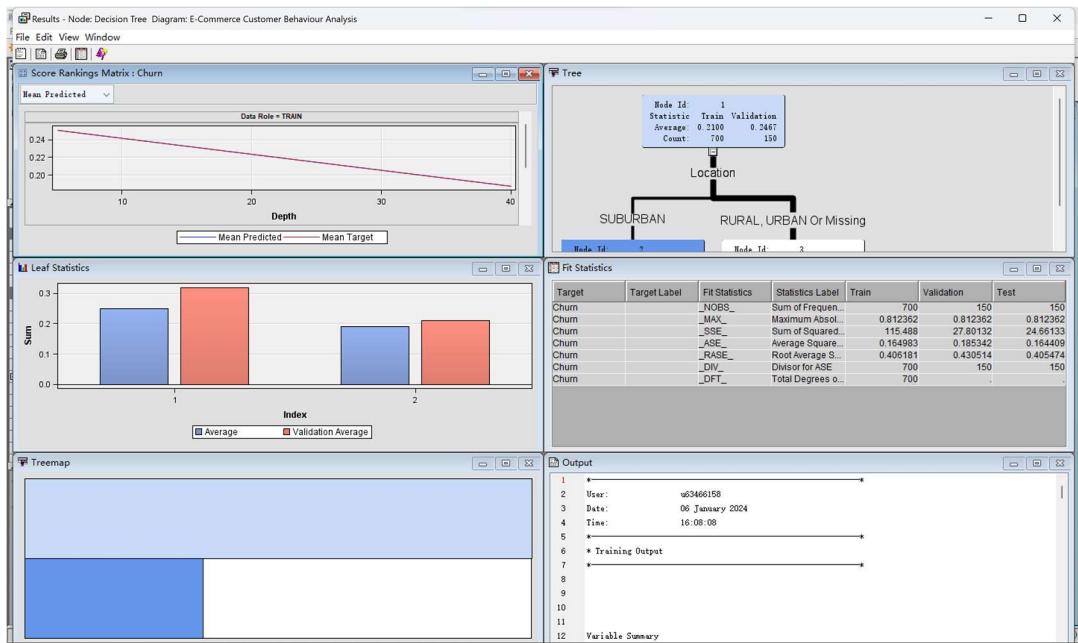
Use a Control point to prepare for more models to implement and simplify the process flow diagram.



Step 2: Configuring the Decision Tree Node

Click on the Decision Tree node to check the left side configuration bar.

Here below is the output results based on the default configure settings:



Results - Node: Decision Tree Diagram: E-Commerce Customer Behaviour Analysis

```

File Edit View Window
Output
1
2 User: u63460158
3 Date: 06 January 2024
4 Time: 16:08:08
5
6 * Training Output
7
8
9
10
11
12 Variable Summary
13
14 Role Measurement Frequency
15 Level Count
16
17 ID INTERVAL 2
18 INPUT INTERVAL 5
19 INPUT NOMINAL 6
20 REJECTED NOMINAL 1
21 TARGET INTERVAL 1
22 TIMEID INTERVAL 1
23
24
25
26 Predicted and decision variables
27
28 Type Variable Label
29
30
31 TARGET Churn
32 PREDICTED P_Churn Predicted: Churn
33 RESIDUAL R_Churn Residual: Churn
34
35
36
37 * Score Output
38
39
40

```

Results - Node: Decision Tree Diagram: E-Commerce Customer Behaviour Analysis

```

File Edit View Window
Output
46
47 Variable Importance
48
49
50 Variable Number of Splitting Ratio of Validation to Training
51 Name Label Rules Importance Importance Importance
52
53 Location 1 1.0000 1.0000 1.0000
54
55
56
57
58 Tree Leaf Report
59
60 Node Training Training Validation Validation Training Validation
61 Id Depth Observations Average Observations Average Root ASE Root ASE
62
63 3 1 453 0.19 103 0.21 0.39042 0.41066
64 2 1 247 0.25 47 0.32 0.43360 0.47110
65
66
67
68
69 Fit Statistics
70
71 Target='Churn Target Label'
72
73 Fit
74 Statistics Statistics Label Train Validation Test
75
76 _NBS_ Sum of Frequencies 700.000 150.000 150.000
77 _MAX_ Maximum Absolute Error 0.812 0.812 0.812
78 _SSE_ Sum of Squared Errors 115.488 27.801 24.661
79 _ASE_ Average Squared Error 0.165 0.185 0.164
80 _RASE_ Root Average Squared Error 0.406 0.431 0.405
81 _DIV_ Divisor for ASE 700.000 150.000 150.000
82 _DFT_ Total Degrees of Freedom 700.000 , ,

```

Results - Node: Decision Tree Diagram: E-Commerce Customer Behaviour Analysis

```

File Edit View Window
Output
87 Assessment Score Rankings
88
89 Data Role='TRAIN Target Variable='Churn Target Label'
90
91 Number of Mean Mean
92 Depth Observations Target Predicted
93
94 5 247 0.25101 0.25101
95 40 453 0.18764 0.18764
96
97
98 Data Role='VALIDATE Target Variable='Churn Target Label'
99
100 Number of Mean Mean
101 Depth Observations Target Predicted
102
103 5 47 0.31915 0.25101
104 35 103 0.21359 0.18764
105
106
107
108
109 Assessment Score Distribution
110
111 Data Role='TRAIN Target Variable='Churn Target Label'
112
113 Range for Mean Mean Number of Model
114 Predicted Target Predicted Observations Score
115
116 0.248 - 0.251 0.25101 0.25101 247 0.24943
117 0.188 - 0.191 0.18764 0.18764 453 0.18922
118
119
120 Data Role='VALIDATE Target Variable='Churn Target Label'
121
122 Range for Mean Mean Number of Model
123 Predicted Target Predicted Observations Score
124
125 0.248 - 0.251 0.31915 0.25101 47 0.24943
126 0.188 - 0.191 0.21359 0.18764 103 0.18922

```

Summary of Results with Default Configuration:

Significance Level:

Set a lower significance level for splits to reduce complexity and overfitting, such as 0.05 or 0.01.

Leaf Size:

Increase the minimum leaf size to 10 or 20 to prevent the model from being too sensitive to noise in the training data.

Number of Rules:

Adjust the number of rules based on the complexity of the dataset. If overfitting is observed, reduce the number of rules to simplify the tree.

Maximum Depth:

Set a maximum depth based on the point where validation performance plateaus or decreases to prevent overfitting. A depth of 5 or 6 might be a good starting point.

Pruning:

Implement post-pruning techniques using a validation set to remove branches that do not contribute significantly to the model's predictive power.

Variable Importance:

Review the variable importance and ensure that the tree uses a diverse set of variables for decision-making, rather than relying predominantly on one.

Max depth 6 to 10, leaf size 5 to 8, number of surrogate rules 0 to 4, make sure the assessment measure is decision, subtree method from Assessment to Largest, Observation Based Importance No to Yes, maximum branch 2 to 10.

Improvements for Better Configuration:

As updated the configuration in the screenshot below:

| Property | Value |
|-------------------------------------|------------------------------|
| General | |
| Node ID | Tree |
| Imported Data | ... |
| Exported Data | ... |
| Notes | ... |
| Train | |
| Variables | ... |
| Interactive | ... |
| Import Tree Model | No |
| Tree Model Data Set | ... |
| Use Frozen Tree | No |
| Use Multiple Targets | No |
| Splitting Rule | |
| Interval Target Criterion | ProbF |
| Nominal Target Criterion | ProbChisq |
| Ordinal Target Criterion | Entropy |
| Significance Level | 0.2 |
| Missing Values | Use in search |
| Use Input Once | No |
| Maximum Branch | 10 |
| Maximum Depth | 10 |
| Minimum Categorical Size | 5 |
| Node | |
| Leaf Size | 8 |
| Number of Rules | 5 |
| Number of Surrogate Rules | 4 |
| Split Size | 1 |
| Split Search | |
| Use Decisions | No |
| Use Priors | No |
| Exhaustive | 5000 |
| Node Sample | 20000 |
| Subtree | |
| Method | Largest |
| Number of Leaves | 1 |
| Assessment Measure | Decision |
| Assessment Fraction | 0.25 |
| Cross Validation | |
| Perform Cross Validation | No |
| Number of Subsets | 10 |
| Number of Repeats | 1 |
| Seed | 12345 |
| Observation Based Imputation | |
| Observation Based Impute | Yes |
| Number Single Var Impute | 5 |
| P-Value Adjustment | |
| Bonferroni Adjustment | Yes |
| Time of Bonferroni Adjustment | Before |
| Inputs | No |
| Number of Inputs | 1 |
| Depth Adjustment | Yes |
| Output Variables | |
| Leaf Variable | Yes |
| Interactive Sample | |
| Create Sample | Default |
| Sample Method | Random |
| Sample Size | 10000 |
| Sample Seed | 12345 |
| Performance | Disk |
| Score | |
| Variable Selection | Yes |
| Leaf Role | Segment |
| Report | |
| Precision | 4 |
| Tree Precision | 4 |
| Class Target Node Col | Percent Correctly Classified |
| Interval Target Node Col | Average |
| Node Text | ... |
| Status | |
| Create Time | 1/6/24 4:56 PM |
| Run ID | 3046bece-6bc3-e84c-ad04 |
| Last Error | |
| Last Status | Complete |
| Last Run Time | 1/6/24 5:23 PM |
| Run Duration | 0 Hr. 0 Min. 4.34 Sec. |
| Grid Host | |
| User-Added Node | No |

Max Depth (6 to 10):

Increasing the maximum depth allows the tree to create more levels of decision nodes. This will

enable the model to capture more complex patterns in the data. It's important to monitor performance on validation data to find the optimal depth. I have tried that start at the lower end of this range and incrementally increase, validating performance at each step.

Maximum Branch (2 to 10):

Adjusting the maximum number of branches from a single decision node from 2 to 10 allows the tree to consider more than binary splits. This is useful for multi-class categorical variables and can provide a more nuanced understanding of the data.

Leaf Size (5 to 8):

Increasing the minimum leaf size makes the tree less sensitive to noise in the data, reducing overfitting. A leaf size of 8 means that any decision node must represent at least 8 instances in the data. This helps ensure that patterns found by the tree are statistically reliable and not due to chance variations in smaller sample sizes.

Number of Surrogate Rules (0 to 4):

Surrogate rules are backup criteria used to split a node when the primary splitting criterion is missing for some cases. This helps handle missing data and maintains the robustness of the tree structure. Introducing up to 4 surrogate rules provides a fallback for handling missing values without relying on imputation alone.

Assessment Measure (ensure it is 'decision'):

The assessment measure should be 'decision' to focus on classification accuracy. This measure will evaluate how well the tree distinguishes between the classes, which is crucial for a binary outcome like churn. This setting aligns the tree's assessment with the goal of making accurate predictions rather than simply fitting the training data.

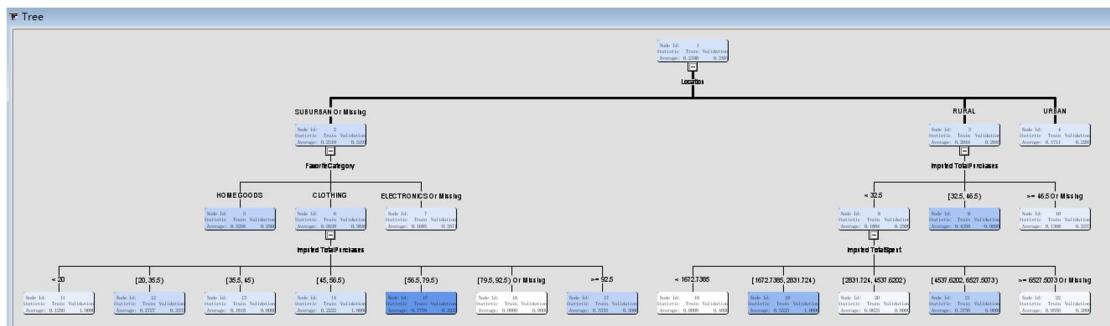
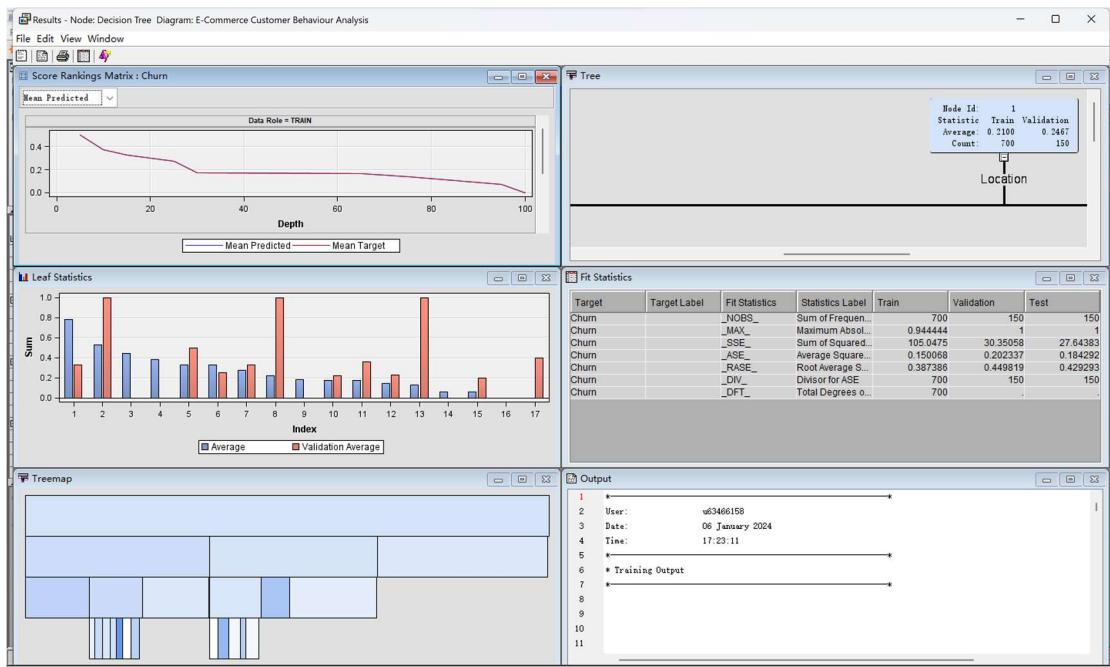
Subtree Method (from 'Assessment' to 'Largest'):

Changing the subtree method from 'Assessment', which might focus on validation performance, to 'Largest' alters the criteria for subtree selection. 'Largest' may retain more of the tree structure, which could capture additional patterns, because the 'Assessment' method is too conservative and we want the model to retain more detail from the training data.

Observation Based Importance (No to Yes):

Enabling observation-based importance allows the tree to weight splits based on how well they improve the prediction for each observation, which can lead to a more accurate representation of variable importance. This setting is particularly useful if we have imbalanced classes or if some observations are more important than others.

And then we got the improved results as shown in the pictures below:



Fit Statistics

| Target | Target Label | Fit Statistics | Statistics Label | Train | Validation | Test |
|--------|--------------|----------------------------|------------------|----------|------------|------|
| Churn | _NOBS_ | Sum of Frequencies | 700 | 150 | 150 | . |
| Churn | _MAX_ | Maximum Absolut. Error | 0.944444 | 1 | 1 | . |
| Churn | _SSE_ | Sum of Squared Errors | 105.0475 | 30.35058 | 27.64383 | . |
| Churn | _ASE_ | Average Squared Error | 0.150068 | 0.202337 | 0.184292 | . |
| Churn | _RASE_ | Root Average Squared Error | 0.387386 | 0.449819 | 0.429293 | . |
| Churn | _DIV_ | Divisor for ASE | 700 | 150 | 150 | . |
| Churn | _DFT_ | Total Degrees of Freedom | 700 | 150 | 150 | . |

```

Output
1   *-----*
2   User:          u63466158
3   Date:          06 January 2024
4   Time:          17:23:11
5   *-----*
6   * Training Output
7   *-----*
8
9
10
11
12 Variable Summary
13
14      Measurement  Frequency
15  Role       Level     Count
16
17  ID        INTERVAL    2
18  INPUT     INTERVAL    5
19  INPUT     NOMINAL     5
20  REJECTED  NOMINAL    1
21  TARGET    INTERVAL    1
22  TIMEID   INTERVAL    1
23
24
25
26
27 Predicted and decision variables
28
29 Type      Variable  Label
30
31 TARGET    Churn
32 PREDICTED P_Churn  Predicted: Churn
33 RESIDUAL   R_Churn  Residual: Churn
34
35
36 *-----*
37 * Score Output
38 *-----*
39

```

```

Output
47 Variable Importance
48
49
50
51
52      Number of      Number of      Ratio of
53      Splitting     Surrogate    Validation
54      Rules         Rules       Importance  to Training
55      Validation   Importance  Importance
56
57 Variable Name   Label
58
59
60
61
62
63

```

| Variable Name | Label | Number of Splitting Rules | Number of Surrogate Rules | Validation Importance | Ratio of Validation to Training Importance |
|--------------------------|-------------------------|---------------------------|---------------------------|-----------------------|--|
| IMP_TotalPurchases | Imputed TotalPurchases | 2 | 3 | 1.0000 | 0.6291 |
| IMP_TotalSpent | Imputed TotalSpent | 1 | 3 | 0.8804 | 0.6164 |
| IMP_Age | Imputed Age | 0 | 4 | 0.6068 | 0.6153 |
| FrequencyOfWebsiteVisits | | 0 | 3 | 0.5515 | 0.6188 |
| IMP_MembershipLevel | Imputed MembershipLevel | 0 | 3 | 0.5384 | 0.0000 |
| Tenure | | 0 | 1 | 0.3996 | 0.0000 |
| FavoriteCategory | | 1 | 0 | 0.3537 | 0.0000 |
| Location | | 1 | 0 | 0.2962 | 1.0000 |
| Occupation | | 0 | 1 | 0.2262 | 0.0000 |

```

Output
66 Tree Leaf Report
67
68 Node      Training   Training   Validation  Validation  Training   Validation
69 Id        Observations Average    Observations Average    Root ASE   Root ASE
70
71 4         228       0.17      59        0.22      0.37655  0.41740
72 10        117       0.14      22        0.23      0.34359  0.42874
73 7         89        0.17      14        0.36      0.37434  0.51494
74 5         86        0.33      20        0.25      0.46859  0.43956
75 9         39        0.44      6         0.00      0.49587  0.43590
76 22        18        0.06      5         0.20      0.22906  0.42528
77 20        16        0.06      3         0.00      0.24206  0.06250
78 19        15        0.53      1         1.00      0.49889  0.46667
79 16        12        0.00      2         0.00      0.00000  0.00000
80 17        12        0.33      2         0.50      0.47140  0.52705
81 18        3         0.00      5         0.40      0.00000  0.63246
82 12        11        0.27      3         0.33      0.44536  0.47528
83 13        11        0.18      1         0.00      0.38569  0.18182
84 14        9         0.22      1         1.00      0.41574  0.77778
85 15        3         0.78      3         0.33      0.41574  0.64788
86 11        8         0.13      1         1.00      0.33072  0.87500
87 21        8         0.38      2         0.00      0.48412  0.37500
88
89
90
91
92 Fit Statistics
93
94 Target=Churn Target Label=' '
95
96 Fit
97 Statistics Statistics Label      Train   Validation   Test
98
99 _NOBS_ Sum of Frequencies    700.000  150.000  150.000
100 _MAX_ Maximum Absolute Error 0.944    1.000    1.000
101 _SSE_ Sum of Squared Errors 105.048  30.351   27.644
102 _ASE_ Average Squared Error 0.150    0.202    0.184
103 _BASE_ Root Average Squared Error 0.387  0.450   0.429
104 _DIV_ Divisor for ASE      700.000  150.000  150.000
105 _DFT_ Total Degrees of Freedom 700.000

```

```

Output
110 Assessment Score Rankings
111
112 Data Role=TRAIN Target Variable=Churn Target Label=' '
113
114 Number of      Mean      Mean
115 Depth   Observations  Target  Predicted
116
117 5          63        0.50794  0.50794
118 10         8         0.37500  0.37500
119 15         98       0.32653  0.32653
120 25         11       0.27273  0.27273
121 30         248      0.17339  0.17339
122 65         89       0.16854  0.16854
123 75         117      0.13675  0.13675
124 95         42        0.07143  0.07143
125 100        24       0.00000  0.00000
126
127
128 Data Role=VALIDATE Target Variable=Churn Target Label=' '
129
130 Number of      Mean      Mean
131 Depth   Observations  Target  Predicted
132
133 5          10        0.20000  0.54821
134 10         24        0.25000  0.33035
135 25         4         0.50000  0.26010
136 30         60       0.21667  0.17123
137 70         14       0.35714  0.16854
138 75         22       0.22727  0.13675
139 90         1        1.00000  0.12500
140 95         8        0.12500  0.05816
141 100        7        0.28571  0.00000

```

| Output | | | | | |
|--|---|----------------|-------------------|---------------------------|----------------|
| 146 Assessment Score Distribution | | | | | |
| 147 | | | | | |
| 148 Data Role=TRAIN Target Variable=Churn Target Label=' ' | | | | | |
| 149 | | | | | |
| 150 | Range for Predicted | Mean Target | Mean Predicted | Number of Observations | Model Score |
| 151 | 0.739 - 0.778 | 0.77778 | 0.77778 | 9 | 0.75833 |
| 152 | 0.506 - 0.544 | 0.53333 | 0.53333 | 15 | 0.52500 |
| 153 | 0.428 - 0.467 | 0.43590 | 0.43590 | 39 | 0.44722 |
| 154 | 0.350 - 0.389 | 0.37500 | 0.37500 | 8 | 0.36944 |
| 155 | 0.311 - 0.350 | 0.32653 | 0.32653 | 98 | 0.33056 |
| 156 | 0.272 - 0.311 | 0.27273 | 0.27273 | 11 | 0.29167 |
| 157 | 0.194 - 0.233 | 0.22222 | 0.22222 | 9 | 0.21389 |
| 158 | 0.156 - 0.194 | 0.17073 | 0.17073 | 328 | 0.17500 |
| 159 | 0.117 - 0.156 | 0.13600 | 0.13600 | 125 | 0.13611 |
| 160 | 0.039 - 0.078 | 0.05882 | 0.05882 | 34 | 0.05833 |
| 161 | 0.000 - 0.039 | 0.00000 | 0.00000 | 24 | 0.01944 |
| 162 | | | | | |
| 163 | | | | | |
| 164 | | | | | |
| 165 | | | | | |
| 166 | Data Role=VALIDATE Target Variable=Churn Target Label=' ' | | | | |
| 167 | | | | | |
| 168 | Range for Predicted | Mean Target | Mean Predicted | Number of Observations | Model Score |
| 169 | 0.739 - 0.778 | 0.33333 | 0.77778 | 3 | 0.75833 |
| 170 | 0.506 - 0.544 | 1.00000 | 0.53333 | 1 | 0.52500 |
| 171 | 0.428 - 0.467 | 0.00000 | 0.43590 | 6 | 0.44722 |
| 172 | 0.350 - 0.389 | 0.00000 | 0.37500 | 2 | 0.36944 |
| 173 | 0.311 - 0.350 | 0.27273 | 0.32629 | 22 | 0.33056 |
| 174 | 0.272 - 0.311 | 0.33333 | 0.27273 | 3 | 0.29167 |
| 175 | 0.194 - 0.233 | 1.00000 | 0.22222 | 1 | 0.21389 |
| 176 | 0.156 - 0.194 | 0.24324 | 0.17072 | 74 | 0.17500 |
| 177 | 0.117 - 0.156 | 0.26087 | 0.13624 | 23 | 0.13611 |
| 178 | 0.039 - 0.078 | 0.12500 | 0.05816 | 8 | 0.05833 |
| 179 | 0.000 - 0.039 | 0.28571 | 0.00000 | 7 | 0.01944 |
| 180 | | | | | |
| 181 | | | | | |

Improved Configuration Summary:

The tree now branches on variables other than 'Location', such as 'FavoriteCategory' and 'TotalPurchases', which suggests a more nuanced analysis of factors contributing to churn.

The tree has a greater depth, indicating that it can capture more complex relationships between variables and churn.

The leaf nodes show a wider spread in the distribution of error rates, suggesting differentiation in the predictive capability across different segments of the tree.

The use of surrogate rules has likely helped to deal with missing data and maintain the integrity of splits, improving the model's robustness.

The tree model has expanded to consider other variables and their interactions, as evidenced by the diverse branches and leaves.

Comparison with Default Configuration:

Model Scoring Improvement highly: The improvement in model score implies that the additional branches and depth have allowed the decision tree to capture more complex patterns that are more predictive of customer churn.

Depth and Complexity: The improved model has a greater depth, allowing for more detailed splits and potentially capturing more complex patterns than the default configuration.

Variable Interaction: Unlike the default model, which primarily used 'Location', the improved model considers other variables, such as 'FavoriteCategory' and 'TotalPurchases', indicating a better-utilized feature space.

Error Rates: The default configuration might have shown a better alignment of error rates between training and validation, but potentially at the cost of oversimplification. The improved model, while having varied error rates, may be capturing more genuine patterns at the risk of some overfitting.

Surrogate Rules: The addition of surrogate rules in the improved model suggests a strategy to handle missing values more effectively than the default setting, which could improve the model's ability to make predictions in the presence of incomplete data.

Assessment Measure and Subtree Method: Switching to 'decision' as the assessment measure and 'Largest' as the subtree method likely resulted in a more complex model that can better capture the true distribution of the target variable.

Based on the results from the updated decision tree configuration, specific inferences can be drawn regarding the factors affecting churn in the e-commerce customer dataset.

Location as a Primary Factor: The initial split on 'Location' in both the default and updated configurations implies that geographic factors are significant in predicting churn. Customers from different locations may have varying experiences or expectations that influence their likelihood to churn.

Favorite Category Significance: The appearance of 'FavoriteCategory' in subsequent splits suggests that the type of products customers frequently purchase is closely linked to their churn behavior. This could indicate satisfaction levels with product offerings or engagement with the platform.

Total Purchases Relevance: The branching on 'TotalPurchases' indicates that the frequency of transactions is a strong indicator of churn, with different purchasing behaviors leading to distinct churn probabilities. Regular purchasers might show more loyalty, while infrequent shoppers could be at higher risk of churn.

Customer Segmentation: Churn appears to be affected by customer segments defined by location and shopping preferences. Tailored strategies can be developed for these segments to address churn.

Purchase Frequency: Customers with lower purchase frequencies may need targeted engagement strategies, such as special offers or reminders, to prevent churn.

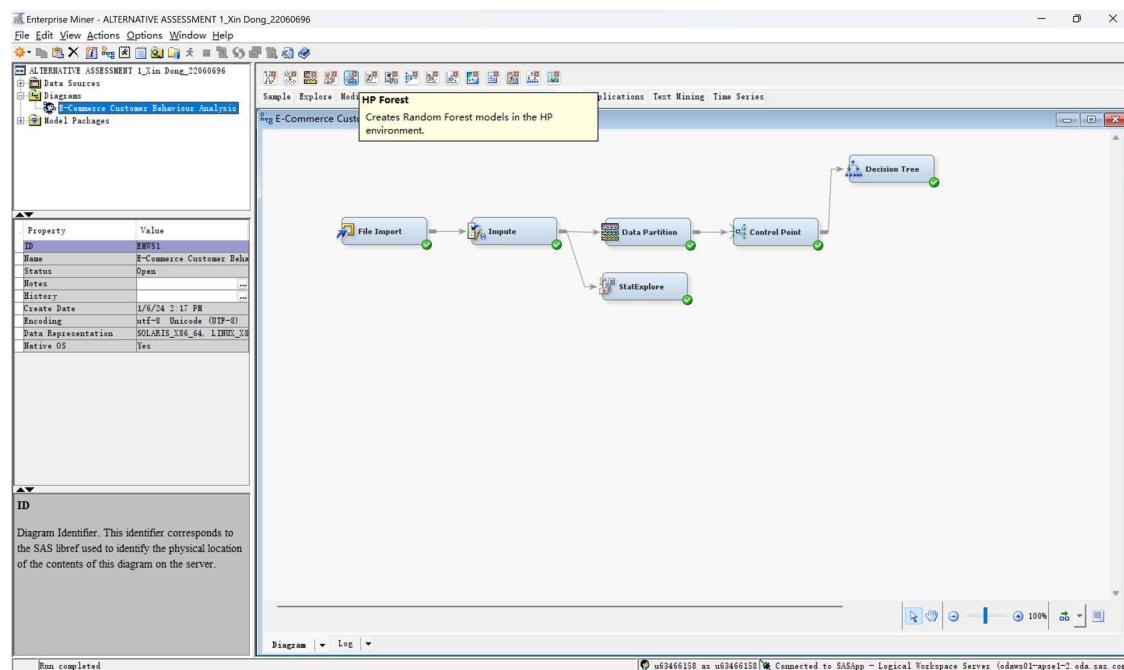
Multi-Dimensional Strategy: Churn prevention requires a multi-dimensional strategy that addresses various factors simultaneously rather than focusing on a single aspect of customer behavior.

Data Completeness: Ensuring data completeness and robust handling of missing data is crucial for accurate churn prediction and subsequent strategy formulation.

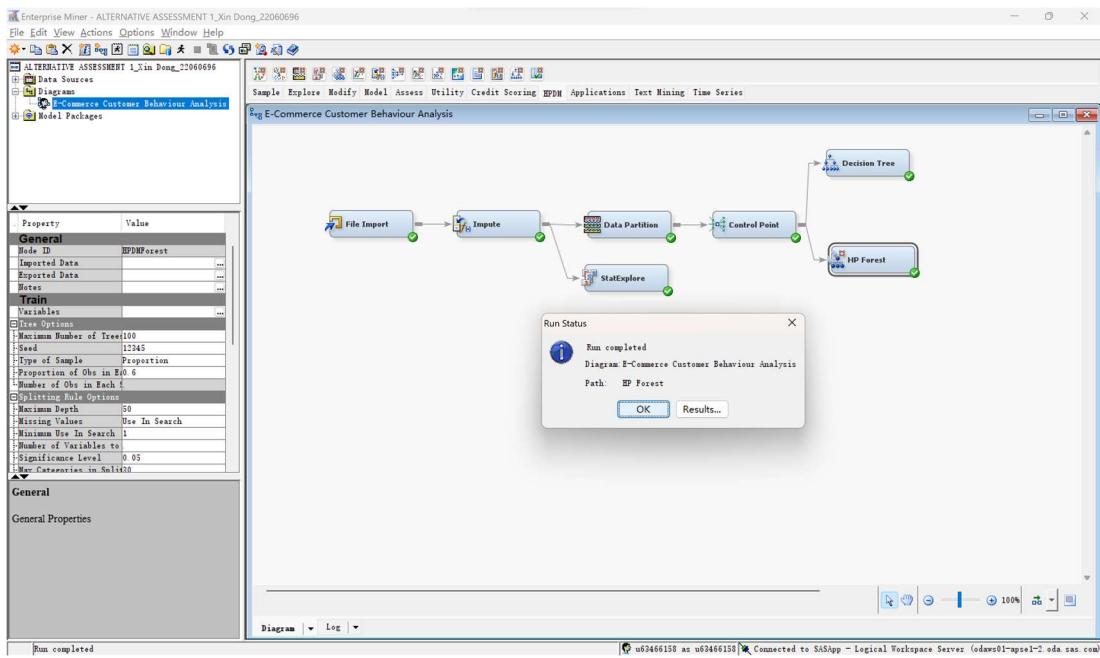
7. Ensemble Methods

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

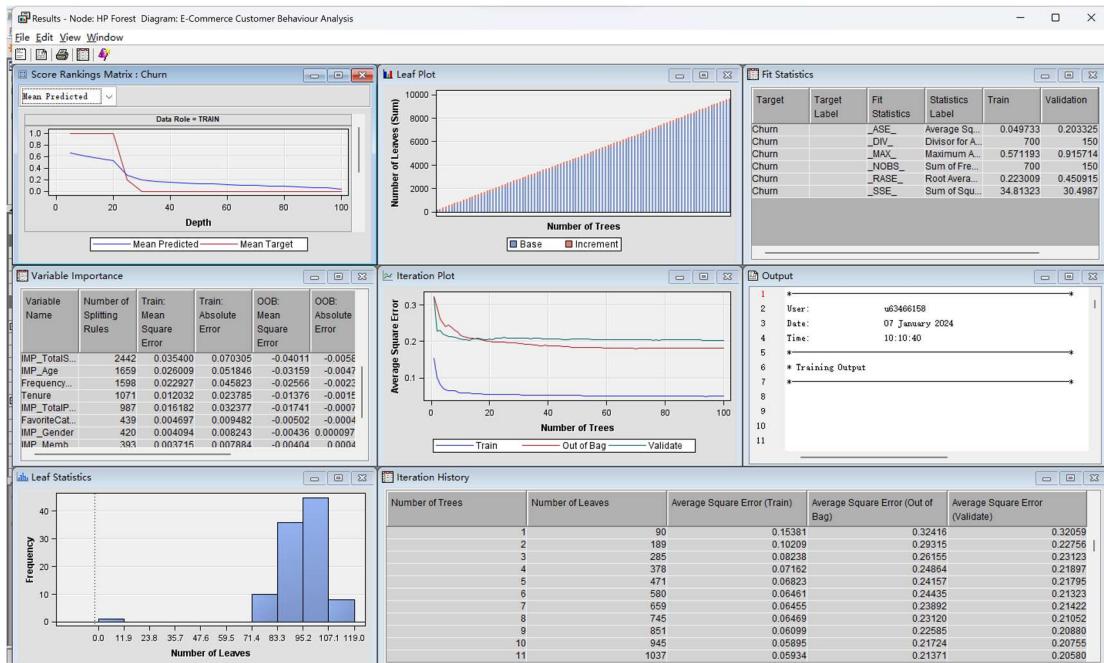
7.1 Apply Random Forest Model:



Connect the processed dataset with HP Random Forest Model by connecting Control Point node with HP Forest, and run:



Results of the Random Forest Model application:



Results - Node: HP Forest Diagram: E-Commerce Customer Behaviour Analysis

File Edit View Window

Output

```
1 *-----*
2 User: u63466158
3 Date: 07 January 2024
4 Time: 10:10:40
5 *-----*
6 * Training Output
7 *-----*
```

8

9

10

11

12 Variable Summary

13

| Role | Measurement Level | Frequency Count |
|----------|-------------------|-----------------|
| ID | INTERVAL | 2 |
| INPUT | INTERVAL | 5 |
| INPUT | NOMINAL | 5 |
| REJECTED | NOMINAL | 1 |
| TARGET | INTERVAL | 1 |
| TIMEID | INTERVAL | 1 |

16

17

18

19

20

21

22

23

24

25

26

27 Predicted and decision variables

28

| Type | Variable | Label |
|-----------|----------|------------------|
| TARGET | Churn | |
| PREDICTED | P_Churn | Predicted: Churn |
| RESIDUAL | R_Churn | Residual: Churn |

30

31

32

33

34

35

36

37

38

Output

```
408 *
409 *-----*
410 * Report Output
411 *-----*
```

412

413

414

415

416 Fit Statistics

417

418 Target=Churn Target Label=' '

419

420 Fit

| Statistics | Statistics Label | Train | Validation | Test |
|------------|----------------------------|---------|------------|---------|
| _ASE_ | Average Squared Error | 0.050 | 0.203 | 0.165 |
| _DIV_ | Divisor for ASE | 700.000 | 150.000 | 150.000 |
| _MAX_ | Maximum Absolute Error | 0.571 | 0.916 | 0.912 |
| _NOBS_ | Sum of Frequencies | 700.000 | 150.000 | 150.000 |
| _RASE_ | Root Average Squared Error | 0.223 | 0.451 | 0.407 |
| _SSE_ | Sum of Squared Errors | 34.813 | 30.499 | 24.796 |

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

771

772

773

774

775

776

777

778

779

780

781

782

783

784

785

786

787

788

789

790

791

792

793

794

795

796

797

798

799

800

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

982

983

984

985

986

987

988

989

990

991

992

993

994

995

996

997

998

999

1000

Output

```

432
433 Assessment Score Rankings
434
435 Data Role=TRAIN Target Variable=Churn Target Label=' '
436
437          Number of      Mean      Mean
438 Depth    Observations   Target   Predicted
439
440      5        35       1.0     0.65856
441     10        35       1.0     0.60556
442     15        35       1.0     0.56894
443     20        35       1.0     0.52559
444     25        35       0.2     0.27955
445     30        35       0.0     0.19396
446     35        35       0.0     0.17221
447     40        35       0.0     0.15505
448     45        35       0.0     0.14328
449     50        35       0.0     0.13190
450     55        35       0.0     0.12299
451     60        35       0.0     0.11443
452     65        35       0.0     0.10657
453     70        35       0.0     0.09774
454     75        35       0.0     0.08856
455     80        35       0.0     0.08201
456     85        35       0.0     0.07384
457     90        35       0.0     0.06455
458     95        35       0.0     0.05443
459    100        35       0.0     0.03883
460

```

Output

```

462 Data Role=VALIDATE Target Variable=Churn Target Label=' '
463
464          Number of      Mean      Mean
465 Depth    Observations   Target   Predicted
466
467      5         8       0.12500  0.41979
468     10        7       0.00000  0.33506
469     15        8       0.12500  0.30641
470     20        7       0.57143  0.28669
471     25        8       0.25000  0.27065
472     30        7       0.28571  0.25874
473     35        8       0.25000  0.24628
474     40        7       0.00000  0.23225
475     45        8       0.12500  0.22084
476     50        7       0.14286  0.20788
477     55        8       0.00000  0.19683
478     60        7       0.28571  0.18919
479     65        8       0.50000  0.17962
480     70        7       0.00000  0.17111
481     75        8       0.37500  0.16408
482     80        7       0.28571  0.15472
483     85        8       0.25000  0.14685
484     90        7       0.85714  0.13873
485     95        8       0.25000  0.11855
486    100        7       0.28571  0.08927
487

```

| Assessment Score Distribution | | | | | |
|--|-------------|----------------|------------------------|-------------|--|
| Data Role=TRAIN Target Variable=Churn Target Label=' ' | | | | | |
| Range for Predicted | Mean Target | Mean Predicted | Number of Observations | Model Score | |
| 0.697 - 0.732 | 1 | 0.71200 | 6 | 0.71447 | |
| 0.661 - 0.697 | 1 | 0.67114 | 8 | 0.67886 | |
| 0.625 - 0.661 | 1 | 0.64293 | 16 | 0.64324 | |
| 0.590 - 0.625 | 1 | 0.61015 | 36 | 0.60762 | |
| 0.554 - 0.590 | 1 | 0.57183 | 37 | 0.57201 | |
| 0.519 - 0.554 | 1 | 0.53704 | 24 | 0.53639 | |
| 0.483 - 0.519 | 1 | 0.50866 | 13 | 0.50078 | |
| 0.447 - 0.483 | 1 | 0.46908 | 5 | 0.46516 | |
| 0.412 - 0.447 | 1 | 0.43429 | 2 | 0.42954 | |
| 0.269 - 0.305 | 0 | 0.28728 | 2 | 0.28708 | |
| 0.234 - 0.269 | 0 | 0.25042 | 8 | 0.25146 | |
| 0.198 - 0.234 | 0 | 0.21502 | 28 | 0.21886 | |
| 0.162 - 0.198 | 0 | 0.17980 | 60 | 0.18023 | |
| 0.127 - 0.162 | 0 | 0.14341 | 105 | 0.14461 | |
| 0.091 - 0.127 | 0 | 0.10993 | 144 | 0.10900 | |
| 0.056 - 0.091 | 0 | 0.07486 | 151 | 0.07338 | |
| 0.020 - 0.056 | 0 | 0.04370 | 55 | 0.03776 | |

| Data Role=VALIDATE Target Variable=Churn Target Label=' ' | | | | | |
|---|-------------|----------------|------------------------|-------------|--|
| Data Role=VALIDATE Target Variable=Churn Target Label=' ' | | | | | |
| Range for Predicted | Mean Target | Mean Predicted | Number of Observations | Model Score | |
| 0.503 - 0.527 | 0.00000 | 0.52651 | 1 | 0.51489 | |
| 0.480 - 0.503 | 0.00000 | 0.48353 | 1 | 0.49164 | |
| 0.457 - 0.480 | 0.00000 | 0.47643 | 1 | 0.46840 | |
| 0.387 - 0.410 | 0.50000 | 0.39491 | 2 | 0.39867 | |
| 0.364 - 0.387 | 0.00000 | 0.37528 | 1 | 0.37542 | |
| 0.341 - 0.364 | 0.00000 | 0.35053 | 4 | 0.35218 | |
| 0.317 - 0.341 | 0.00000 | 0.33002 | 5 | 0.32893 | |
| 0.294 - 0.317 | 0.14286 | 0.30849 | 7 | 0.30569 | |
| 0.271 - 0.294 | 0.41667 | 0.28280 | 12 | 0.28244 | |
| 0.248 - 0.271 | 0.33333 | 0.25908 | 15 | 0.25920 | |
| 0.224 - 0.248 | 0.00000 | 0.23484 | 12 | 0.23596 | |
| 0.201 - 0.224 | 0.16667 | 0.21627 | 12 | 0.21271 | |
| 0.178 - 0.201 | 0.18182 | 0.19127 | 22 | 0.18947 | |
| 0.155 - 0.178 | 0.22727 | 0.16659 | 22 | 0.16622 | |
| 0.131 - 0.155 | 0.52941 | 0.14545 | 17 | 0.14298 | |
| 0.108 - 0.131 | 0.37500 | 0.12132 | 8 | 0.11973 | |
| 0.085 - 0.108 | 0.20000 | 0.10030 | 5 | 0.09649 | |
| 0.062 - 0.085 | 0.33333 | 0.07686 | 3 | 0.07325 | |

Key Summary and Findings from Random Forest Model results:

Fit Statistics Overview:

The model has a low Average Squared Error (ASE) for the training data, indicating good fit.

However, this error is significantly higher on the validation and test sets, suggesting the model may be overfitting and not generalizing as well to unseen data.

The Maximum Absolute Error (MAX) is relatively high on the validation and test sets, which means there are certain predictions that deviate significantly from the actual values.

The Root Average Squared Error (RASE) also increases from training to validation and test sets, which is consistent with the ASE findings and further indicates overfitting.

Assessment Score Rankings (Train):

The model predicts a high likelihood of churn (mean of 1.0) across several depths when trained, but this prediction decreases as the tree depth increases. This could mean the model is very confident in its predictions for shallow trees, but as more complexity is added (i.e., deeper trees), the confidence in the churn prediction decreases.

Assessment Score Rankings (Validate):

There's a significant variance in the mean predictions at different depths. For instance, at a depth of 20, the model predicts a high likelihood of churn (mean target of 0.57143), but this is not consistent across other depths. This suggests that the model may be sensitive to the depth parameter and that certain depths capture the churn behavior better than others.

Assessment Score Distribution (Train):

The model score is highest for the top range of predicted values, indicating good performance for cases with a high likelihood of churn.

As the predicted likelihood of churn decreases, so does the model score, suggesting the model is better at identifying clear cases of churn rather than borderline cases.

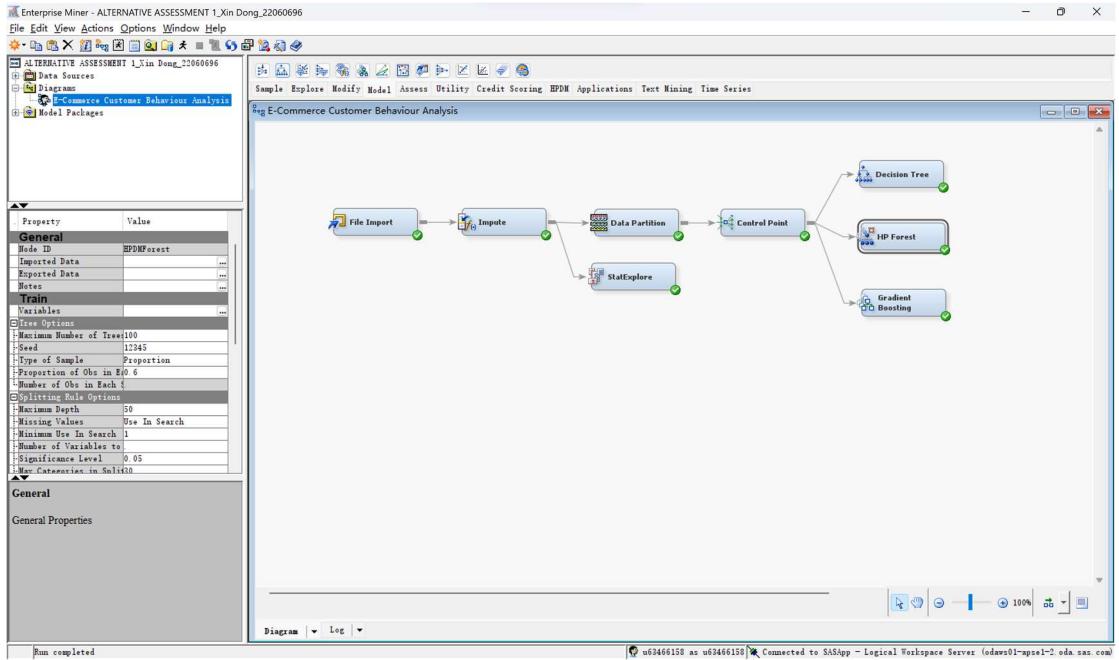
Assessment Score Distribution (Validate):

The model scores are generally lower on the validation set compared to the training set, especially for predictions with a high likelihood of churn, which indicates the model is not performing as well on unseen data.

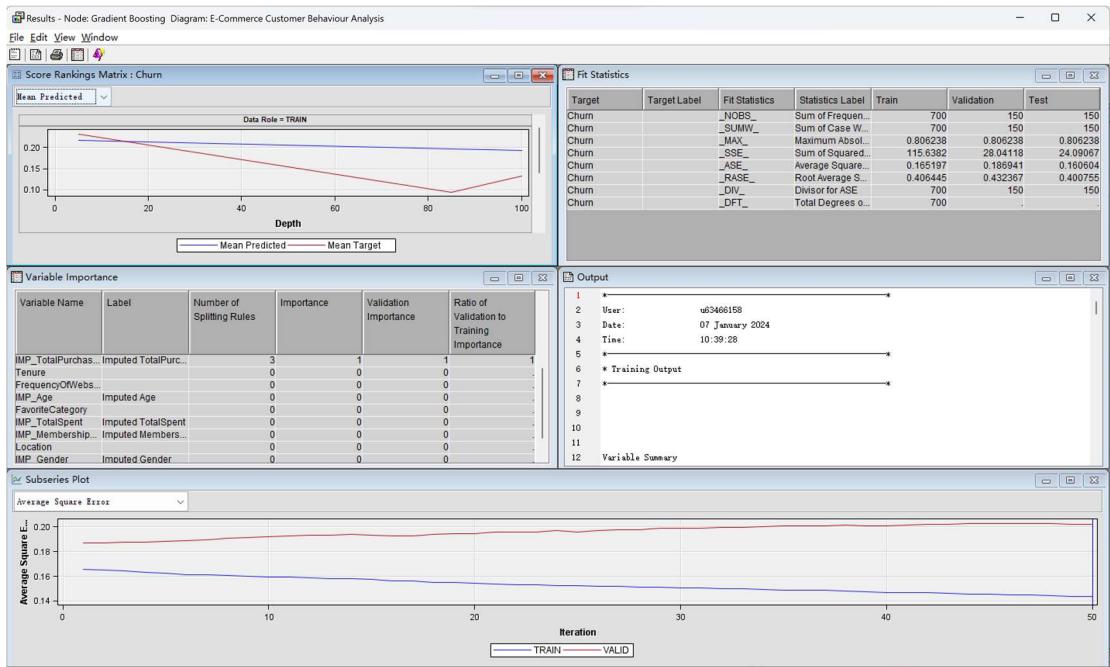
There are instances (e.g., depth of 80 with a mean target of 0.28571) where the model's predicted likelihood of churn does not align well with the actual mean target, which may indicate issues with model bias or variance on unseen data.

7.2 Apply Gradient Boosting Model:

Connect the processed dataset with HP Random Forest Model by connecting Control Point node with Gradient Boosting, and run:



Results of the Gradient Boosting Model application:



```

Output
49   *-----*
50   * Report Output
51   *-----*
52
53
54
55
56   Fit Statistics
57
58   Target=Churn Target Label=' '
59
60   Fit
61   Statistics      Statistics Label      Train    Validation    Test
62
63   _NOBS_      Sum of Frequencies      700.000    150.000    150.000
64   _SUMW_      Sum of Case Weights Times Freq  700.000    150.000    150.000
65   _MAX_       Maximum Absolute Error      0.806      0.806      0.806
66   _SSE_       Sum of Squared Errors      115.638    28.041     24.091
67   _ASE_       Average Squared Error      0.165      0.187      0.161
68   _RASE_      Root Average Squared Error  0.406      0.432      0.401
69   _DIV_       Divisor for ASE          700.000    150.000    150.000
70   _DFT_       Total Degrees of Freedom   700.000    .
71
72
73
74
75   Assessment Score Rankings
76
77   Data Role=TRAIN Target Variable=Churn Target Label=' '
78
79   Number of      Mean      Mean
80   Depth   Observations   Target   Predicted
81
82   5        575        0.23304  0.21660
83   85       95         0.09474  0.19684
84   100      30         0.13333  0.19376
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99   Assessment Score Distribution
100
101  Data Role=TRAIN Target Variable=Churn Target Label=' '
102
103  Range for      Mean      Mean      Number of      Model
104  Predicted     Target   Predicted   Observations   Score
105
106  0.230 - 0.232  0.25000  0.23186    12        0.23090
107  0.215 - 0.217  0.23268  0.21627    563       0.21567
108  0.196 - 0.198  0.09474  0.19684    95        0.19662
109  0.194 - 0.196  0.13333  0.19376    30        0.19471
110
111
112  Data Role=VALIDATE Target Variable=Churn Target Label=' '
113
114  Range for      Mean      Mean      Number of      Model
115  Predicted     Target   Predicted   Observations   Score
116
117  0.230 - 0.232  1.00000  0.23186    2        0.23090
118  0.215 - 0.217  0.22556  0.21627    133       0.21567
119  0.196 - 0.198  0.30769  0.19684    13        0.19662
120  0.194 - 0.196  0.50000  0.19376    2        0.19471
121

```

Key Summary and Findings from Gradient Boosting Model results:

Fit Statistics Overview:

The model processed 700 observations for training, and 150 each for validation and testing.

The Maximum Absolute Error (MAX) is the same across training, validation, and test datasets, which may indicate a consistent outlier or extreme value that the model is uniformly missing across all datasets.

The Average Squared Error (ASE) and Root Average Squared Error (RASE) are relatively low, indicating that on average, the model's predictions are close to the actual values. However, there is a slight increase in ASE from training to validation, which may suggest a small degree of overfitting.

Assessment Score Rankings:

For the training data, the model has a reasonable performance with a slight discrepancy between the mean target and mean predicted values at various depths. This discrepancy increases slightly at higher depths, which could be indicative of the model's increasing uncertainty or decreased performance as the complexity (depth) increases.

On the validation set, there are fewer observations at various depths, but the model shows a higher discrepancy between the mean target and mean predicted, particularly at depths of 95 and 100.

This could indicate that the model's performance is less stable on validation data for those particular depths.

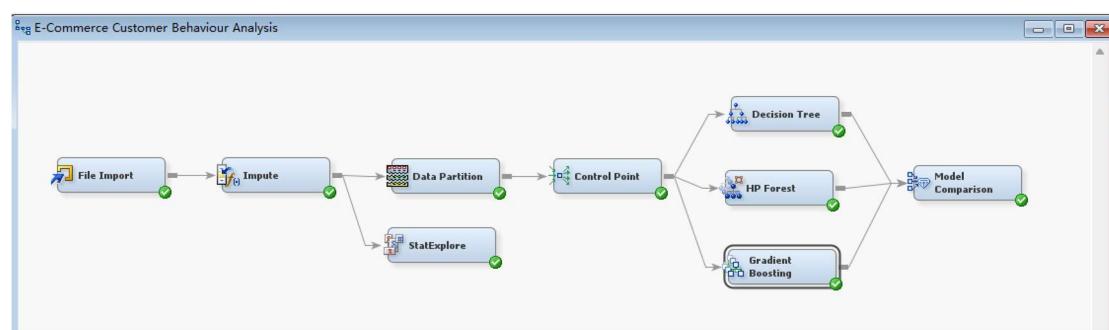
Assessment Score Distribution:

The model scores are closest for predicted ranges around 0.215 - 0.217, showing the smallest gap between the mean target and predicted values for a large number of observations in the training set. This suggests the model is most accurate around this predicted churn probability.

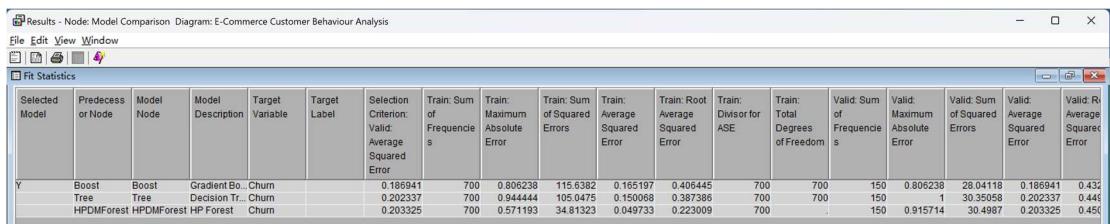
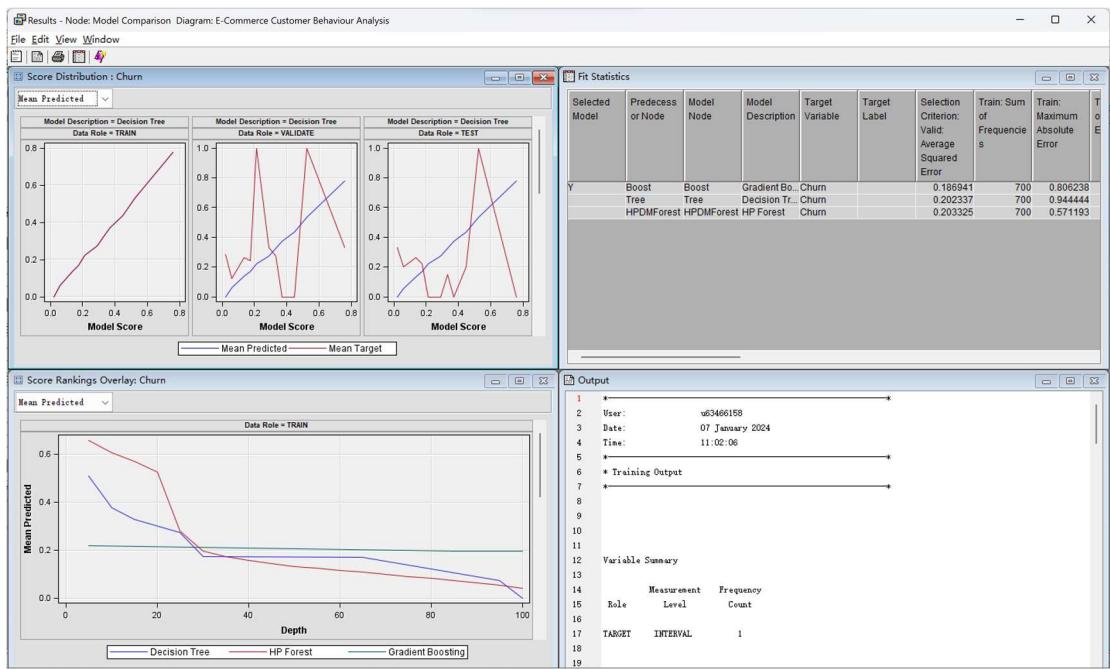
In the validation set, the highest model score occurs at the extreme range of 0.230 - 0.232, where there are very few observations. This could suggest the model is overfitting to specific cases in the validation set.

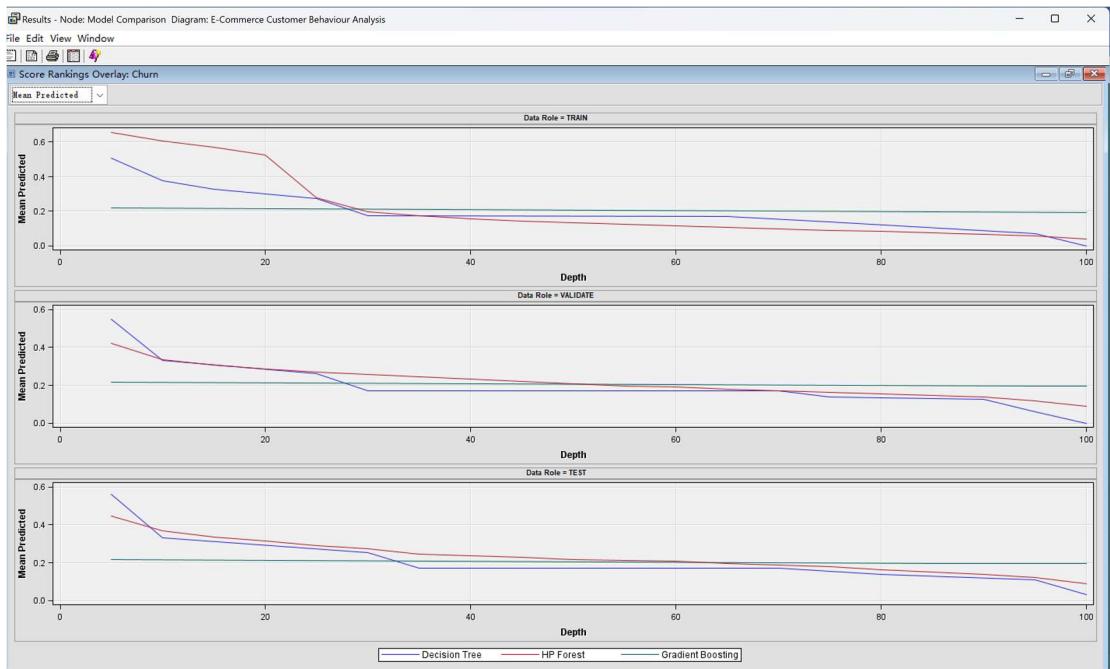
8. Model Comparison and Conclusions

Finally, we use the Model Comparison to compare the three different models and drawn study conclusions



Output Results:





Output

```

28
29 Fit Statistics
30 Model Selection based on Valid: Average Squared Error (_VASE_)
31
32                                         Valid:      Train:
33                                         Average     Average
34 Selected
35 Model   Model Node   Model Description   Error       Error
36
37 Y      Boost        Gradient Boosting    0.18694   0.16520
38 Tree    Decision Tree      0.20234   0.15007
39 HPDFForest  HP Forest      0.20332   0.04973
40
41
42
43
44
45
46
47
48
49
50
51 Fit Statistics Table
52 Target: Churn
53
54 Data Role=Train
55
56 Statistics          Boost      Tree   HPDFForest
57
58 Train: Average Squared Error      0.165     0.150     0.050
59 Selection Criterion: Valid: Average Squared Error  0.187     0.202     0.203
60 Train: Total Degrees of Freedom  700.000   700.000     .
61 Train: Divisor for ASE          700.000   700.000   700.000
62 Train: Maximum Absolute Error   0.806     0.944     0.571
63 Train: Sum of Frequencies      700.000   700.000   700.000
64 Train: Root Average Squared Error  0.406     0.387     0.223
65 Train: Sum of Squared Errors    115.638   105.048   34.813
66 Train: Sum of Case Weights Times Freq  700.000   .         .
67

```

```

Fit Statistics Table
Target: Churn
Data Role=Train
Statistics          Boost      Tree    HPDMForest
Train: Average Squared Error      0.165     0.150     0.050
Selection Criterion: Valid: Average Squared Error      0.187     0.202     0.203
Train: Total Degrees of Freedom   700.000   700.000   .
Train: Divisor for ASE           700.000   700.000   700.000
Train: Maximum Absolute Error    0.806     0.944     0.571
Train: Sum of Frequencies        700.000   700.000   700.000
Train: Root Average Squared Error 0.406     0.387     0.223
Train: Sum of Squared Errors     115.638   105.048   34.813
Train: Sum of Case Weights Times Freq 700.000   .
.
Data Role=Valid
Statistics          Boost      Tree    HPDMForest
Valid: Average Squared Error      0.187     0.202     0.203
Valid: Divisor for VASE           150.000   150.000   150.000
Valid: Maximum Absolute Error    0.806     1.000     0.916
Valid: Sum of Frequencies        150.000   150.000   150.000
Valid: Root Average Squared Error 0.432     0.450     0.451
Valid: Sum of Squared Errors     28.041    30.351    30.499
Valid: Sum of Case Weights Times Freq 150.000   .

```

```

Data Role=Test
Statistics          Boost      Tree    HPDMForest
Test: Average Squared Error      0.161     0.184     0.165
Test: Divisor for TASE           150.000   150.000   150.000
Test: Maximum Absolute Error    0.806     1.000     0.912
Test: Sum of Frequencies        150.000   150.000   150.000
Test: Root Average Squared Error 0.401     0.429     0.407
Test: Sum of Squared Errors     24.091    27.644    24.796
Test: Sum of Weights Times Freqs 150.000   150.000   .

*
* Score Output
*
*
* Report Output
*
```

Based on the model comparison output results as shown in the above pics, here is a summary and key findings:

Model Performance Comparison:

The Gradient Boosting model (Boost) was selected as the best model based on the Valid: Average Squared Error (VASE) with a value of 0.18694, which is the lowest among the three models.

The Decision Tree model (Tree) has a slightly higher VASE of 0.20234, and the HP Random Forest model (HPDMForest) is close behind with a VASE of 0.20332.

Training Data Performance:

The HP Random Forest model shows a remarkably lower Average Squared Error (ASE) of 0.050 on training data compared to the Gradient Boosting (0.165) and Decision Tree (0.150) models. However, the HP Random Forest model seems to have overfitted the training data as indicated by the large discrepancy between its train ASE and validation ASE.

Validation Data Performance:

All models have increased ASE in the validation set compared to the training set, but Gradient Boosting has the smallest increase, indicating better generalization.

The Maximum Absolute Error is notably lower for the HP Forest model (0.916) compared to the Decision Tree model (1.000) on validation data.

Testing Data Performance:

The Gradient Boosting model again performs best on the test data with an ASE of 0.161, followed closely by the HP Forest model with an ASE of 0.165, and lastly the Decision Tree model with an ASE of 0.184.

The Maximum Absolute Error on the test data is lowest for the Gradient Boosting model (0.806), which suggests it may be better at handling outliers or extreme cases.

General Findings:

The Gradient Boosting model not only performed best on the validation set (which was the selection criterion) but also maintained consistent performance on the test set, indicating it has a good balance between bias and variance.

The HP Forest model, despite its exceptional performance on the training data, showed signs of overfitting with a much higher error rate on the validation and test data.

The Decision Tree model, while not outperforming the other models, still holds a respectable performance considering its simplicity compared to the ensemble methods.

In conclusion, the **Gradient Boosting** model is recommended for its strong and consistent performance across all datasets, indicating a robust prediction capability for the churn target variable. It strikes a good balance between fitting the training data and generalizing to unseen data.

Derive Business insights from our analysis:

- Customer Churn Prediction: The Gradient Boosting model has shown to be the most effective at predicting e-commerce customer churn. This insight can be leveraged by the business to identify at-risk customers early and to deploy retention strategies proactively to reduce churn.

- Focus on Model Generalization: The lower validation and test errors of the Gradient Boosting model suggest that it generalizes well to new, unseen data. This implies that the business can expect a reliable performance from this model when applied to real-world e-commerce customer behavior, thus making informed decisions about future marketing and customer service strategies.
- Resource Allocation for Customer Retention: The precision of the Gradient Boosting model allows the e-commerce business to allocate resources more efficiently. Rather than applying broad retention strategies across the entire customer base, the business can target those specific customers identified by the model as being at high risk of churn, optimizing both the impact and cost-effectiveness of these strategies.
- Identifying Key Factors for Churn: The variable importance measures from the models, particularly from the Gradient Boosting model, can highlight the most influential factors contributing to e-commerce customer churn. Understanding these factors can inform the business on which areas (e.g., customer service, product features, pricing strategies) require improvement to enhance customer satisfaction and loyalty.
- Overfitting Awareness: The HP Forest model's overfitting on the training data serves as a reminder to the business that complex models, while powerful, need to be carefully validated to ensure they perform well in practice and not just on historical data.
- Simplicity vs. Complexity in Model Selection: While the Gradient Boosting model outperformed others, the simpler Decision Tree model still provided respectable results. This suggests that for rapid deployment or interpretation, simpler models can still offer valuable insights, especially when time or computational resources are limited.
- Strategic Timing for Interventions: The ability to predict churn at various levels of certainty (as indicated by the range of predictions and mean predicted values) allows the e-commerce business to time their customer interventions more strategically. For example, they can prioritize immediate action for customers predicted with high certainty to churn soon, while planning longer-term strategies for those with a lower immediate risk.
- Tailored Customer Experience: Insights from the model can guide the development of more personalized customer experiences. For example, if frequent purchases are a key variable, loyalty programs or personalized offers can be designed to encourage repeat e-commerce business.
- Feedback Loop for Continuous Improvement: The use of these predictive models creates a feedback loop for the business. By monitoring the outcomes of the strategies implemented based on the model's predictions, the business can further refine its models and strategies, leading to continuous improvement in customer retention efforts.