The data set of this project was established through AI and personal modification. The data set can truly represent actual e-commerce customer behavior. The data set contains 3,000 items and 12 attributes.

The dataset for an e-commerce website contains several key attributes related to customer transactions. Here is an overview of the dataset structure:

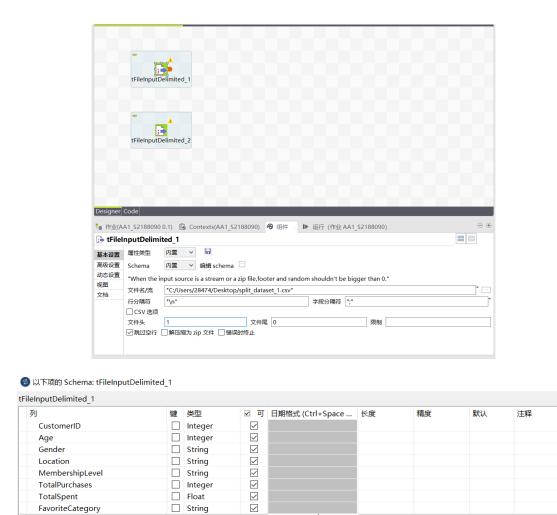
CustomerID	A unique identifier for each customer.			
Age	The age of the customer.			
Gender	The gender of the customer.			
Location	The customer's geographical location.			
MembershipLevel	The customer's membership level (e.g.,			
	Silver, Gold, Bronze).			
TotalPurchases	The total number of purchases made by a			
	customer on the website.			
TotalSpent	The total amount spent by the customer.			
FavoriteCategory	The customer's favorite shopping category.			
LastPurchaseDate	The date of the customer's last purchase.			
Churn	Indicates whether the customer has stopped			
	purchasing (1 means churn, 0 means active).			
Occupation:	The client's occupation.			
WebsiteVisitFrequency	How often customers visit the website.			

# **Data Integration**

Talend Data Integration (Talend DI) is a tool that helps organizations manage, move and transform data. It lets you easily connect different types of data sources, clean and transform data, and move them to where you need them, all through a visual interface without writing code. This tool supports big data processing and has strong community support.

The following is how I use Talend Data Integration to integrate the AI data set with the data set I collected.

By searching for the component of tFileInputDelimited\_1, this is typically used for reading CSV files. Then confirm that the line separator is correct, set the field separator, and enter 1 in the file header. After finishing, edit the schema of the tFileInputDelimited\_1 file by adding the name and data type of each data field in the data set (such as integer, string, date etc.), as well as attributes such as whether to allow null values and whether to serve as a primary key. After defining the schema, Talend can process the data based on this information to ensure data accuracy and consistency. For example, the LastPurchaseDate field is set to a date type, with the date format 'yyyy-MM-dd', so Talend will expect this field in each record to represent the date in this format when processing the data.



yyyy-MM-dd

×

OK Cancel

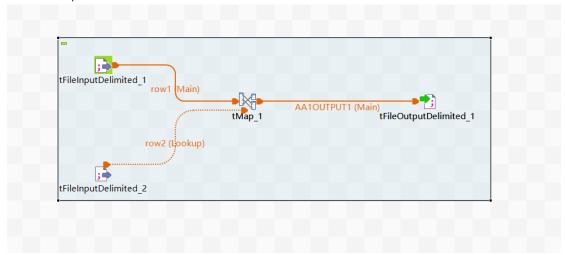
In the same way, edit the schema for the data set tFileInputDelimited\_2.

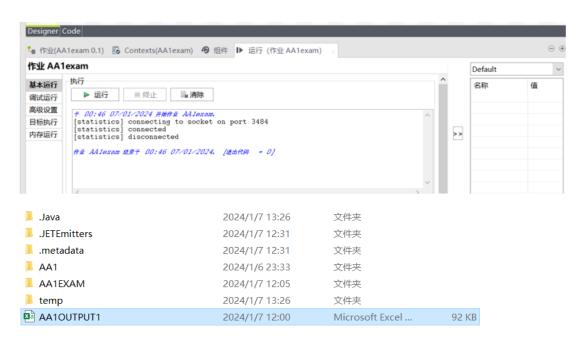
☐ Date

Last Purchase Date



This figure shows fileInputDelimited\_1 and fileInputDelimited\_2, connected through the tMap component. The tMap\_1 component may be used to integrate the data in fileInputDelimited\_1 and fileInputDelimited\_2, and perform operations such as merging fields, filtering records, converting data types, etc. row1(Main) represents the main data flow of the first input file, while row2(lookup) represents the data flow of the second input file as a lookup. In tMap, the search flow is usually used to query or supplement data associated with the main data flow. After processing, the output is connected to tFileOutputDelimited\_1, which indicates that the processed data will be written to a file in a new delimited format. The result of this process is the integrated output of the two datasets. The execution log of the job shows that the export was successful.



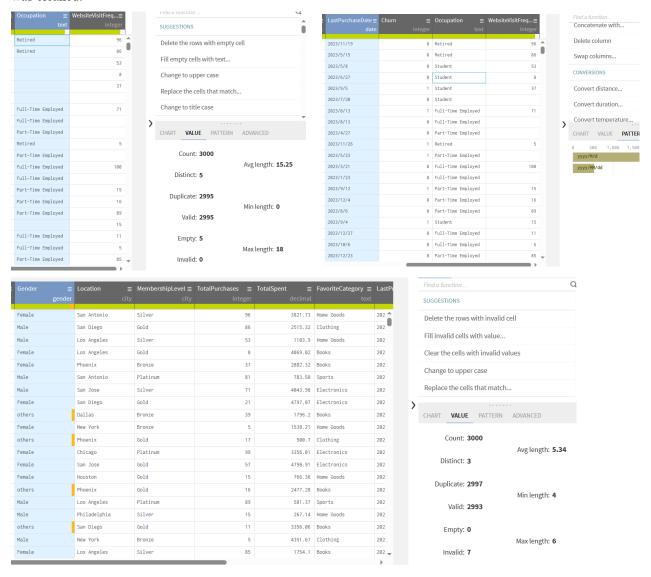


# **Talend Data Preparation**

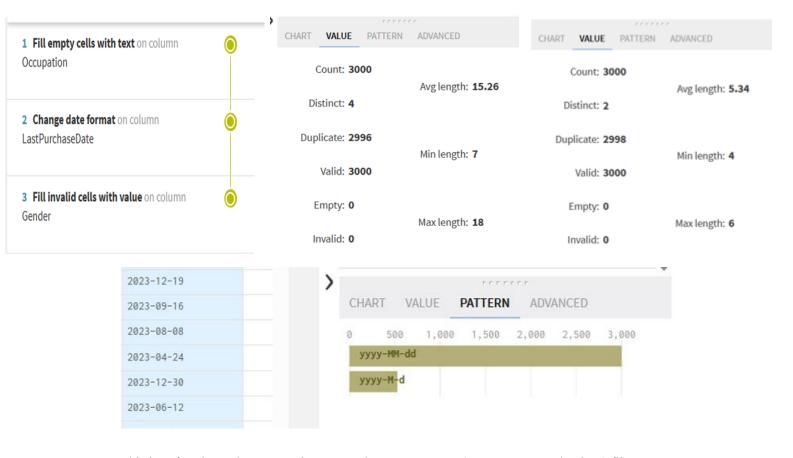
Talend DP is a data preparation tool from Talend Inc. that cleans, transforms, and prepares data for analysis and reporting. It features a visual interface, automation capabilities, multi-data source integration and advanced data analytics to help improve data quality and preparation efficiency.

By uploading the data set to the talend system, we can see that the career attribute of the data set also has a "VALUE" panel on the right, which contains some statistical information about the data: "Empty: 5" Indicates that there are 5 rows of empty data.

In addition, we can see in the date pattern that the format of the date is not yyyy-MM-dd. At the same time, we can see that the gender attribute contains 7 invalid values. We fill in the null value as stduent, change the date to the standard format, and fill in the invalid value. for female. The data set was cleaned.



As you can see from the picture below, the data set has been successfully cleaned, there are no null values, no invalid values, and the date format has been corrected successfully.

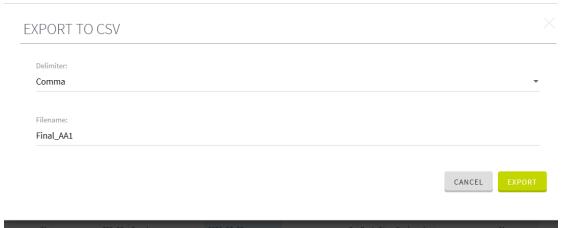


This interface is used to export the current data set to a CSV (comma separated values) file.

Delimiter: The delimiter is set to comma to separate the data in the CSV file.

Filename: The file name has been set to "Final AA1".

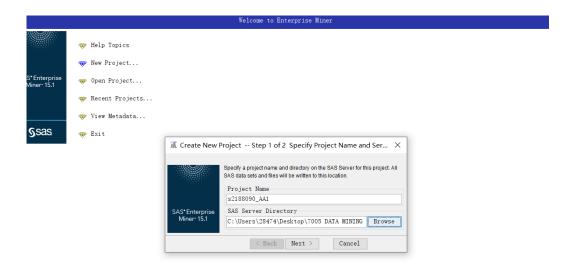
Creates a CSV file based on the selected delimiter and specified file name.



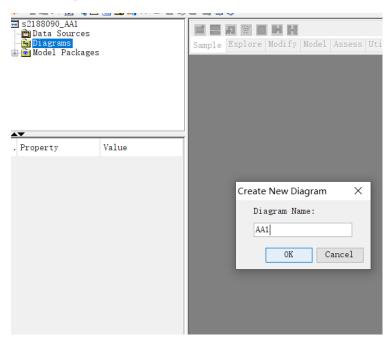
# **Tasks**

1. Data Import and Preprocessing: Import your dataset into SAS Enterprise Miner, handle missing values and specify variable roles.

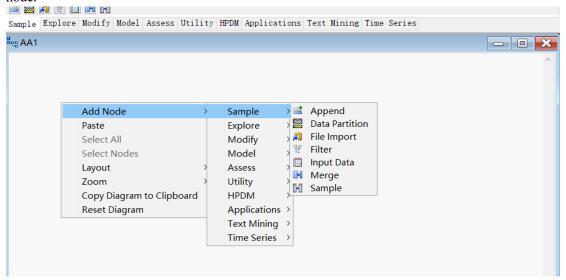
Data import part: Start SAS Enterprise Miner, create project: enter the name S2188090\_AA1, select the server.



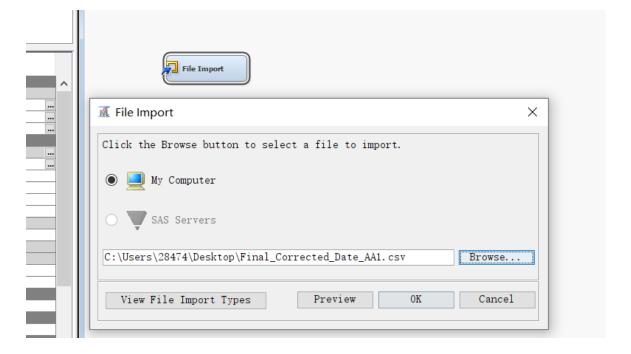
Create New Diagram. Here the user entered "AA1" as the name of the new chart and clicked "OK" to create the chart process.



Right-click the chart interface, Add node appears, then select sample, and then select the file import node.

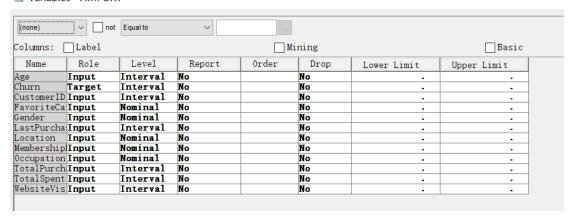


Click file import, and then there is import file in the left column. By selecting this button, you can import data from the local computer or the connected SAS server. File path: Displays the currently selected file path, such as C:\Users\28474\Desktop\Final\_Corrected\_Date\_AA1.csv, indicating that the CSV file named Final\_Corrected\_Date\_AA1.csv is being imported from the desktop path.



Variable settings section: Set churn rate as target and other roles as inputs, which means all these variables will be used in the analysis. Click OK when finished. Since the data set has already been preprocessed in talend DI and talend DP, there is no need to repeat this part.

#### M Variables - FIMPORT



# 2. Decision Tree Analysis: Create a decision tree model in SAS Enterprise Miner to analyse customer behaviour.

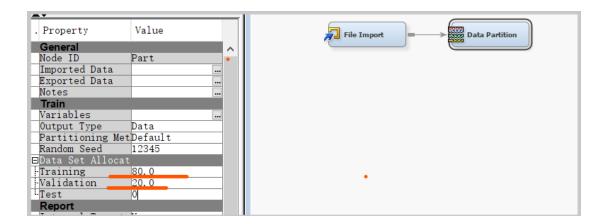
We capture the "Data Partition" node to split the data into different sets, which is a common practice in statistical modeling and machine learning.

Data Partition node configuration: Train: Here defines the proportion of the data set allocated to the training set, validation set and test set. Training: Set to 80.0%, which means 80% of the data will be used for model training.

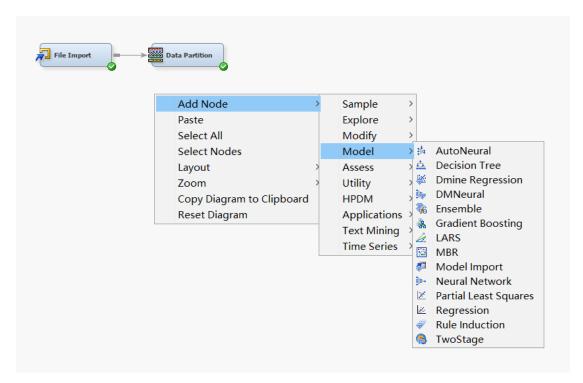
Validation: The display has been set to 20.0%, which means that the remaining 20% of the data will be used for the model validation process.

Random Seed: The random seed is set to 12345 to ensure the reproducibility of the data partition. Using the same random seed ensures the same results for each partition.

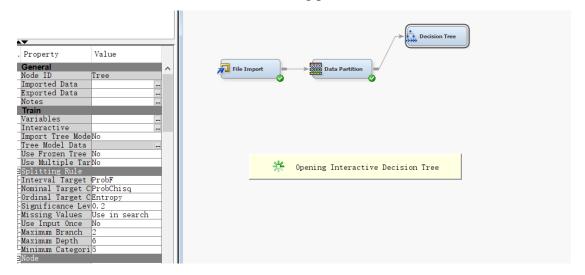
Partitioning Method: Displayed as "Default", which indicates that the data partition will use the SAS Enterprise Miner default partitioning method. Then click Run.



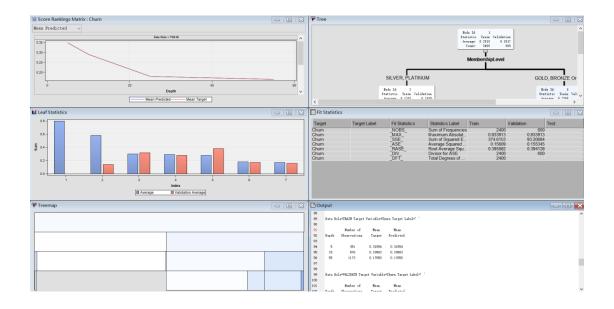
After running, select the decision tree model, right-click on the interface, select add node, then click model, and finally select the decision tree model.

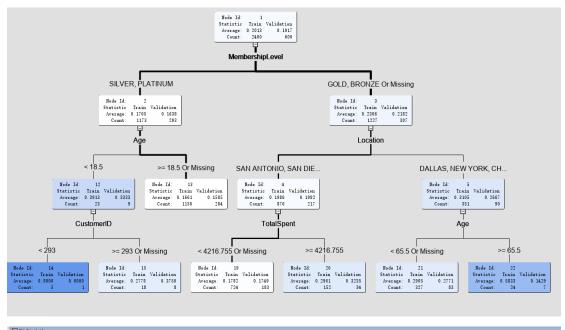


The "Interactive" option in the "Train" configuration of a decision tree model allows the user to manually participate in the decision tree building process. Users can make decisions in real time, such as selecting split points, adjusting tree sizes, applying pruning strategies, and adjusting other model parameters. This interactive approach increases user control over the model building process.

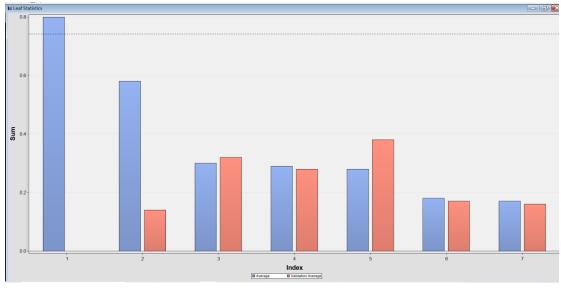


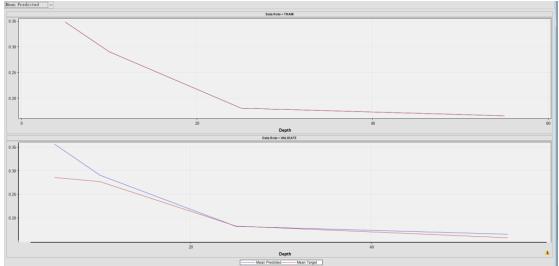
After running, view the results. Below is the overall report of the decision model.





Target	Target Label	Fit Statistics	Statistics Label	Train	Validation
Churn		_NOBS_ _MAX_ _SSE_	Sum of Frequencies	2400	60
Chum		MAX _	Maximum Absolute Error	0.833913	
Chum		SSE_	Sum of Squared Errors	374.6153	
Churn		ASE	Average Squared Error	0.15609	
Churn		RASE_	Root Average Squared Error	0.395082	
Churn		DIV _	Divisor for ASE	2400	6
Churn		DFT	Total Degrees of Freedom	2400	





We can see that it performs data splitting by attributes such as MembershipLevel, Age, Location, etc., where some nodes have high average values on the training data but perform poorly on the validation set, which may mean overfitting.

Signs and Strategies of Overfitting:

Inconsistent node performance: For example, if a certain leaf node has a high mean on the training data, but this value drops significantly on the validation data (such as Node Id 14), this may be a sign of overfitting. This means that the decision tree may have overlearned specific patterns in the training data that do not necessarily hold true for unseen data.

Complex tree structure: The depth of the tree and the number of branches are large, especially when some leaf nodes contain only few data points (such as Node Id 14), which may cause the model to be overly sensitive to noise in the training data.

Solution strategy:

Pruning: Simplify the model by reducing the depth of the tree or setting a higher minimum number of leaf node samples.

Ensemble methods: Use ensemble methods such as random forests or gradient boosted

trees to reduce the risk of overfitting by averaging the results of multiple decision trees. Regularization: Add regularization terms during model training to penalize complex tree structures.

Pruning the tree before the validation error increases can improve the model's generalization ability.

Cross-validation can be used to determine the optimal tree depth that balances model complexity with predictive accuracy.

Incorporating additional data or using ensemble techniques such as random forests can also help improve model performance on unseen data.

# Analyze customer behavior through decision trees.

By analyzing the above report, we can derive some insights about customer behavior: Customer churn (Churn) prediction: The target variable is Churn, which means the model aims to predict whether a customer will churn. Based on the P\_Churn (predicted churn) and R\_Churn (residual of churn) variables, the model attempts to score and predict the likelihood of customer churn.

Variable importance:

Location and MembershipLevel are variables of the model splitting rule, which means that these two variables are important in predicting customer churn.

The importance of Location on the training set is 1.0000, but drops to 0.6270 on the validation set, indicating that the prediction contribution of location to the model is not as significant on the validation set as on the training set.

The validation importance of MembershipLevel is 1.0000, indicating that membership level is a very important predictor variable on the validation data set. Its significance ratio on the validation set relative to the training set exceeds 1 (1.2029), possibly indicating that membership levels have a higher predictive value for predicting churn in real-world data.

Tree depth and number of observations:

As the tree depth increases, the model's predictions tend to be consistent. This may indicate that the model is able to differentiate between customers who are at high risk of churn and those who are at low risk of churn.

In the training data, nodes with depth 5 observed a churn rate of 0.31054, while in the validation set, this number was 0.26667. This suggests that within a specific group of customers, the accuracy of predicting churn decreased on the validation set.

Leaf reports and fit statistics:

The Tree Leaf Report shows the average churn rate for different nodes of the tree (splits based on specific rules).

For nodes with a depth of 2, the churn rates of training data and validation data are similar, indicating that the model's prediction of churn at this depth is consistent on the training and validation sets.

\_RASE\_ (Root Average Squared Error) in Fit Statistics shows the model's error on the training and validation data. A lower error indicates that the model's predictions are more accurate.

At the same time we can also conclude:

The impact of location on churn: A customer's geographic location may be related to their likelihood of churn, with customers in certain locations being more likely to churn. The impact of membership levels: Different membership levels may affect customer loyalty and churn rates. Higher-level members may have lower churn rates.

Possible influence of age: Although age is not mentioned directly in the report, the decision tree may have used age as one of the splitting rules since it is often one of the factors that influence customer behavior.

Relevance of consumption behavior: Consumption behavior is not explicitly mentioned in the report, but since the TotalSpent variable appears in the decision tree model, we can speculate that the customer's consumption level may be related to its churn risk.

Together, these analytics can help business teams better understand and predict which customers are more likely to churn, so they can adopt retention strategies accordingly.

### **Business strategy:**

Based on the above analysis, here are some possible business strategies:

Target focus areas: Since Location appears as an important predictor, companies can develop specific marketing strategies based on location. For example, target areas with high churn rates, increase customer engagement campaigns or offer customized offers. Membership level retention policy:

The importance of MembershipLevel indicates that customers with different membership levels may have different churn risks. Companies can offer loyalty rewards or upgrade offers to membership tiers with a high risk of churn to increase their stickiness.

For those membership levels that exhibit lower churn rates, loyalty can be further enhanced through exclusive offers or customized services.

Personalized communications and promotions:

Use data analysis to identify customer groups at high risk of churn and target them with personalized communications and promotions to increase satisfaction and retention.

Analyzing the nodes of a decision tree can reveal which customer characteristics are associated with churn, helping to design more targeted marketing messages.

Improve customer experience:

If certain nodes are showing unusually high churn rates, this could be a sign of poor customer experience. Investigate these problem areas and take steps to improve the service or product experience.

Optimize resource allocation:

Allocate more resources and attention to those customer groups most likely to churn, for example, by providing additional support or services to reduce their churn rate.

Customer lifecycle management:

Predictive models identify different stages of the customer lifecycle and then provide corresponding interventions at key moments, such as new customer welcome programs, loyal customer rewards or churn prevention programs.

Product and service improvements:

Use churn prediction data to understand why customers leave and improve your product

or service accordingly.

#### Risk Management:

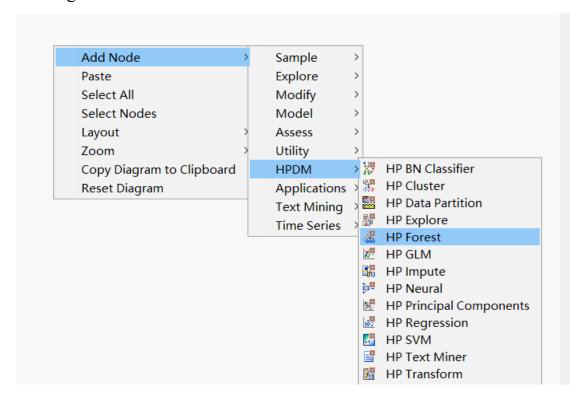
Implement a risk-based pricing strategy to provide different pricing options for customer groups at high risk of churn.

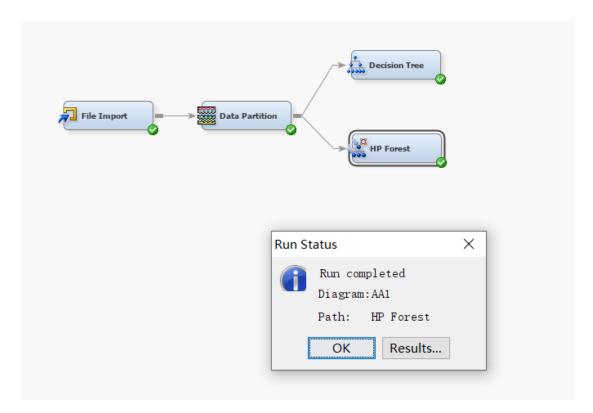
In summary, by analyzing customer behavior using predictive models, companies can more precisely target marketing and service strategies to reduce customer churn and increase customer lifetime value.

# 3. Ensemble Methods: Apply Bagging and Boosting, using the Random Forest algorithm as a Bagging example.

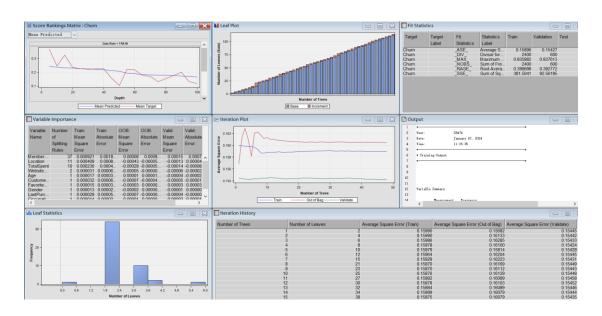
#### Random Forest

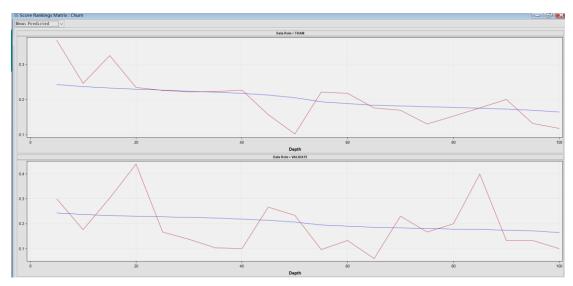
HPforest can be used instead of random forest by selecting the node and clicking HPDM.



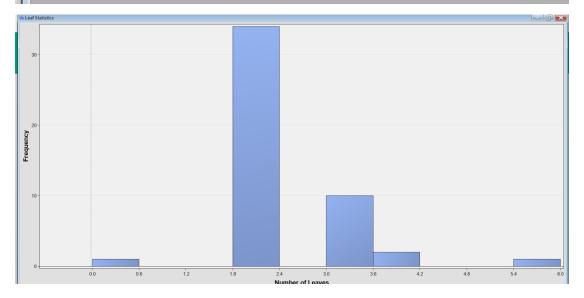


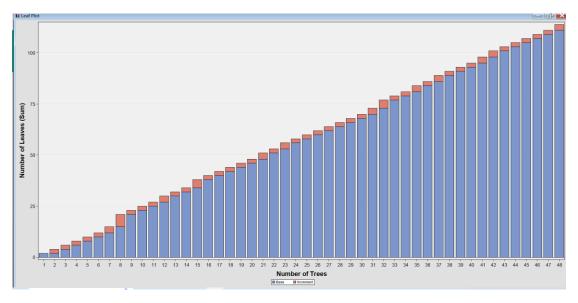
Below is the result report of HPforest.

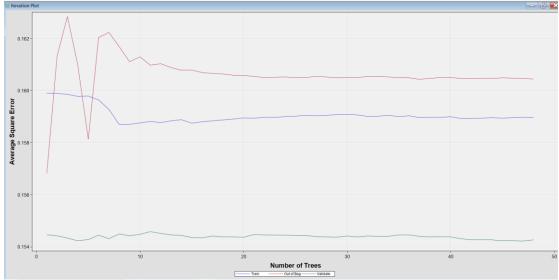




Variable	Importanc	e						
Variable	Number	Train:	Train:	OOB:	OOB:	Valid:	Valid:	Label
Name	of	Maan	Absolute	Mean	Absolute	Mean	Absolute	
	Variable N	Square	Error	Square	Error	Square	Error	
	Rules	Error		Error		Error		
Member	37	0.000921	0.0018	0.00008	0.0009	0.00015	0.0007	
Location	11	0.000409	0.0008	-0.00043	-0.00005	-0.00013	0.00004	
TotalSpent	10	0.000236	0.0004	-0.00028	-0.00005	-0.00014	-0.00008	
Website	2	0.000031	0.00006	-0.00005	-0.00000	-0.00006	-0.00002	
Age	1	0.000017	0.00003	0.00001	0.00001	-0.00004	-0.00002	
Custome	1	0.000032	0.00006	-0.00007	-0.00004	-0.00005	-0.00001	
Favorite	1	0.000015	0.00003	-0.00003	0.00000	0.00003	0.00003	
Gender	1	0.000013	0.00002	-0.00002	0.00000	-0.00001	0.00000	
LastPurc	1	0.000028	0.00005	-0.00007	-0.00000	-0.00004	-0.00000	
Occupati	1	0.000014	0.00002	-0.00001	-0.00000	-0.00003	-0.00001	
TotalPur	0	0.000000	0	0.00000	0	0.00000	0	







Target	Target Label	Fit Statistics	Statistics Label	Train		Validation		Test	
Churn Churn		ASE	Average Squared Error		0.15896		0.15427		
Churn		_DIV_	Divisor for ASE		2400		600		
Churn Churn		MAX	Maximum Absolute Error		0.835982		0.837613		
hum		NOBS	Sum of Frequencies		2400		600		
Churn Churn		RASE	Root Average Squared Error		0.398698		0.392772		
Chum		_SSE_	Sum of Squared Errors		381.5041		92.56195		

Iteration History						
lumber of Trees	Number of Leaves	Average Square Error (Train)		Average Square Error (Out of Bag)	Average Square Error (Validate)	
	1	2	0.15990	0.15682		0.1
	2	4	0.15990	0.16133	3	0.1
	3	6	0.15986	0.16285		0.1
	4	8	0.15978	0.16100	)	0.1
	5	10	0.15979	0.15814		0.1
	6	12	0.15964	0.16204		0.1
	7	15	0.15929	0.16223		0.
	8	21	0.15870	0.16169		0.
	9	23	0.15870	0.16112		0.
	10	25	0.15876	0.16129		0.
	11	27	0.15882	0.16099		0.
	12	30	0.15878	0.16103	3	0.
	13	32	0.15884	0.16089		0.
	14	34	0.15889	0.16079		0
	15	38	0.15875	0.16079		0.
	16	40	0.15881	0.16070 0.16086	)	0.
	17	42	0.15885 0.15887	0.1606 0.1606		Ö
	18	44	0.15897	0.1606		0.
	19 20	46	0.15891	0.1606L 0.1605E	}	0
		48	0.15895	0.16058		0.
	21	51	0.15894 0.15897	0.16056 0.16050		0.
	22	53 56	0.15897	0.1605t 0.1605t	1	0
	23	56	0.15897	0.1605		0.
	24 25	58 60	0.15900	0.1605		0.
	25 26	60	0.15902	0.1605		0
		62	0.15905	0.16050 0.16055		0
	27	64	0.15903 0.15905	0.16050 0.16053		0.
	28 29	66 68	0.15905	0.1605		0
	30	68 70	0.15908	0.1605t 0.1605t		0
	31	73	0.15909	0.1605		0
	31	73	0.15906	0.1605		0
	33	79	0.15903	0.1605		0
	34	81	0.15905	0.1605		0
	35	84	0.15901	0.1605		0
	36	86	0.15903	0.16049		0.
	37	89	0.15897	0.1604		0
	38	91	0.15898	0.1604		0
	39	91	0.15898	0.1605		0
	40	93 95	0.15900	0.16050		ő
	41	98	0.15893	0.16047		Ö
	42	101	0.15893	0.16049		0
	43	103	0.15895	0.16047		ŏ
	44	105	0.15896	0.1604		0
	45	107	0.15895	0.16050		ő
	46	109	0.15896	0.16048		ŏ
	47	111	0.15897	0.16047		0
	48	114	0.15896	0.16045		ő

#### Analysis of results of HPforest model.

Through HPforest's model report, we learned:

The average squared errors (ASE) of the training and validation datasets are very close (0.161 for the training dataset and 0.155 for the validation dataset), indicating that the model has good generalization ability and is not overfitting. Overfitting usually manifests itself as lower training error and higher validation error.

The average squared error for training, out-of-bag (OOB), and validation changes slightly as more trees are added, indicating that the model's performance does not improve significantly after the first few trees. This suggests that a relatively small number of trees are needed to capture patterns in the data.

#### variable importance

MembershipLevel appears to be the most important variable in predicting churn, contributing the most to reducing mean squared error (MSE). Its high importance in the training and validation datasets indicates that it is a strong predictor of churn.

Other variables, such as Age, TotalPurchases, Occupation, and Gender, also contribute to the model, but to a lesser extent.

#### fit statistics

Detailed fit statistics show that as the number of trees increases, so does the number of leaves (decision points in each tree). However, the error rate improved only slightly, suggesting that adding trees beyond a certain point does not necessarily improve the model's predictive power.

The errors during training, OOB, and validation are consistent across the number of trees, again indicating that the model is stable and not overfitting.

Assessment score ranking

This report provides the tree's in-depth performance on the training and validation sets. The average predicted values at different depths indicate that the model is consistent across depths and has no signs of overfitting or underfitting.

#### **Customer behavior:**

- 1. Membership Level Importance: The report indicates that 'MembershipLevel' is highly predictive of churn. This suggests that customers' engagement and satisfaction may vary significantly across different membership tiers, influencing their likelihood to churn.
- 2. Location Influence: 'Location' also plays a role in predicting churn but seems to have a negative impact on the out-of-bag and validation MSE in the variable importance analysis. This could indicate that customers in certain locations are more prone to churn, or there may be regional factors at play affecting customer retention.
- 4. Other Variables: While 'Age', 'TotalPurchases', 'Occupation', 'Gender', 'FavoriteCategory', 'WebsiteVisitFrequency', 'LastPurchaseDate', and 'TotalSpent' have lesser importance compared to 'MembershipLevel', they still contribute to churn prediction. This suggests that demographics, purchasing behavior, and

engagement with the company's online presence are relevant factors in customer churn.

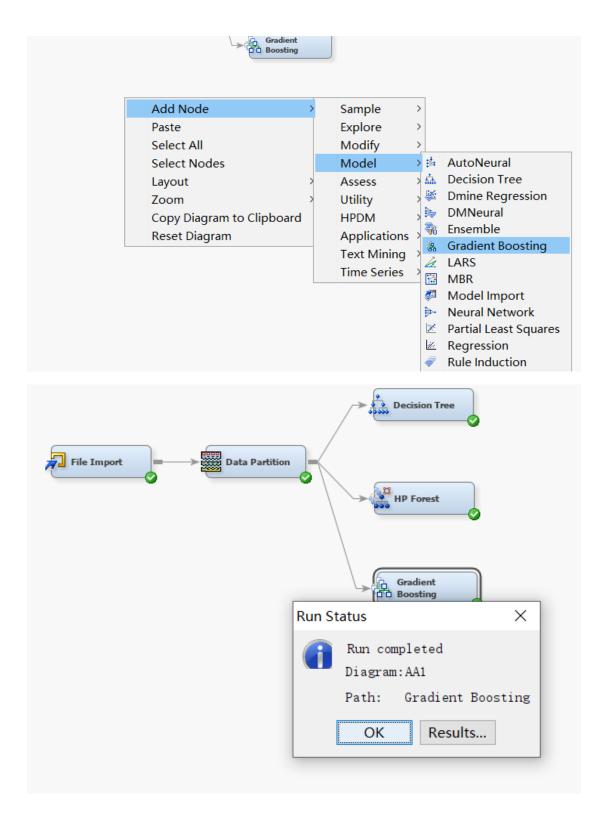
### **Business Strategy Recommendations:**

- 1. Tier-based Engagement Programs: Develop retention programs tailored to different membership levels, especially focusing on the tiers that contribute most to churn. This could include personalized communication, special offers, or loyalty rewards.
- 2. Location-specific Initiatives: Since location seems to affect churn, consider regional marketing initiatives that address the specific needs or pain points of customers in high-churn areas.
- 3. Customer Lifecycle Value Enhancement: For variables with less importance but still contributing to churn prediction, design initiatives to enhance the overall customer experience, such as personalized product recommendations, improved customer service, or loyalty programs.
- 4. Targeted Communication: Use the insights from the model to craft targeted messages for customer segments identified as at risk. This could involve outreach campaigns that address specific behaviors leading to churn.
- 5. Improve Online Experience: Given the inclusion of 'WebsiteVisitFrequency' and 'LastPurchaseDate', improving the online experience and simplifying the purchase process could be beneficial.
- 6. Data-Driven Product and Service Development: Utilize the insights from 'TotalPurchases' and 'FavoriteCategory' to refine product offerings and services that are better aligned with customer preferences.
- 7. Cross-Selling and Up-Selling Opportunities: Leverage information on 'TotalSpent' to identify customers who could be interested in premium offerings or additional services, using up-selling or cross-selling tactics.
- 8. Retention Analytics: Continuously monitor and analyze customer behavior using the model's insights to quickly identify any shifts in behavior that could signal an increased risk of churn.
- 9. Model Monitoring: Since the random forest model used seems to generalize well, it's important to continuously monitor its performance over time with new data to ensure its predictions remain accurate.

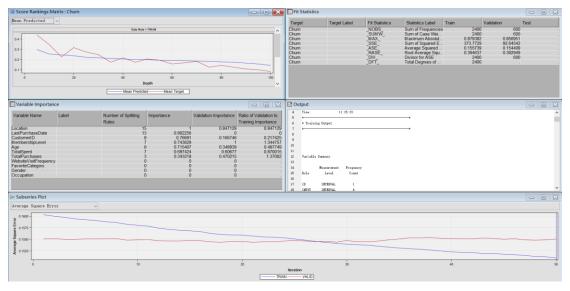
Implementing these strategies can help to reduce churn rates, increase customer satisfaction, and ultimately drive higher customer lifetime value.

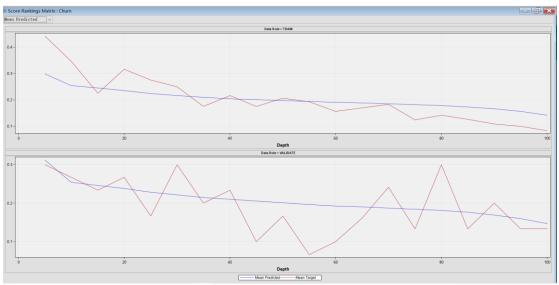
# **Gradient boosting**

For Gradient boosting, we continue the previous method, grab the interface from the model, and then run it to test the performance of the model.

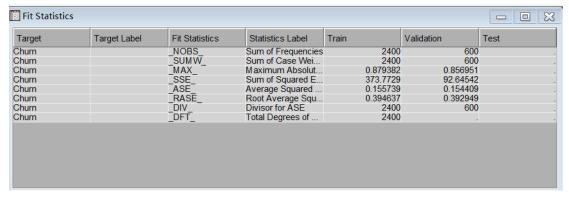


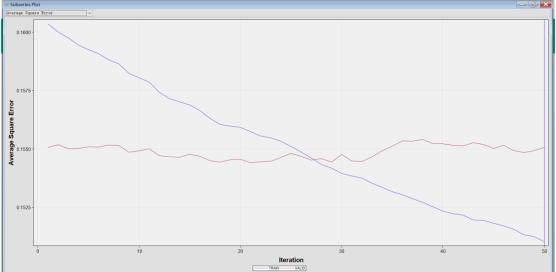
The following is the result report of running the test:





Variable Name	Label	Number of Splitting	Importance	Validation Importance	Ratio of Validation to
variable rearrie	Label		Importance	validation importance	
		Rules			Training Importance
ocation		15	1	0.847129	0.84712
.astPurchaseDate		13	0.902256	0	
CustomerID		8	0.76691	0.166746	0.21742
MembershipLevel		7	0.743629	1	1.34475
Age		8	0.715407	0.348939	0.48774
FotalSpent		7	0.697424	0.60677	0.87001
TotalPurchases		3	0.343218	0.470215	1.3700
NebsiteVisitFrequency		0	0	0	
avoriteCategory		0	0	0	
Gender		0	0	0	
Occupation		0	0	0	





The provided report details the outcomes of a gradient boosting model focused on predicting customer churn. Here is an analysis of the report's results:

#### Variable Importance:

- **Location** is the most critical predictor of churn, with the highest importance score and ratio. This suggests that customer churn is significantly influenced by geographic factors.
- LastPurchaseDate has a high importance score, but a ratio of 0.00000 in the validation importance, indicating it may not be as predictive on unseen data.
- **CustomerID** has a moderate level of importance, which is unusual as IDs are generally not predictive. This could indicate overfitting or data leakage.
- **MembershipLevel** also shows a high importance, suggesting that the type of membership a customer has is predictive of churn likelihood.
- **Age** and **TotalSpent** are also relevant, indicating that demographic factors and spending behavior are associated with churn risk.
- **TotalPurchases** has a lower importance score but a high validation ratio, suggesting it may become a more relevant predictor in the validation dataset.

#### **Fit Statistics:**

- The model has been trained on 2400 observations and validated on 600 observations.
- The maximum absolute error and average squared error are very close between

- the training and validation datasets, which suggests that the model generalizes well and is not overfitting.
- The root average squared error (RASE) is consistent between the training and validation datasets, further supporting the model's generalizability.

#### **Assessment Score Rankings:**

- The mean predicted churn decreases as the number of observations increases, suggesting the model becomes more certain about customers not churning as it has more data to learn from.
- There is a notable decrease in the mean predicted churn from the 5th to the 100th observation in both the training and validation datasets, which could indicate that the model is effectively capturing the underlying patterns associated with churn.

#### **Assessment Score Distribution:**

- The score distribution shows a wide range of mean predicted values, especially in the training dataset, which suggests that the model is identifying different levels of churn risk across the customer base.
- The validation score distribution follows a similar pattern, though with fewer observations in each predicted range, which is expected given the smaller size of the validation dataset.

In summary, the gradient boosting model seems to be performing well, with key variables identified that influence the likelihood of churn. The close fit between training and validation datasets suggests that the model is robust and not overfitting, and the variable importance rankings provide clear indicators of which features are most predictive of churn.

#### **Customer Behavior Insights:**

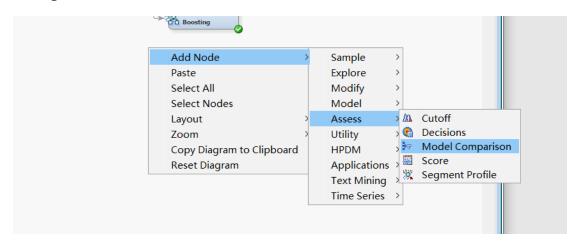
- 1. **Geographic Trends**: Location is highly indicative of churn likelihood, which suggests geographic or regional factors significantly impact customer retention. This could be due to market saturation, local competition, or economic conditions.
- 2. **Engagement Levels**: Membership Level's importance indicates that customers' engagement, as defined by their membership status, plays a crucial role in their decision to stay with or leave the company. Higher or more premium membership levels might correlate with lower churn.
- 3. **Recency of Interaction**: The LastPurchaseDate's initial high importance score may reflect that the more recently a customer has engaged with the company's services or products, the less likely they are to churn.
- 4. **Customer Lifetime Value**: TotalSpent's importance signifies that customers who have spent more over time are key to focus on for retention strategies as their churn could represent a higher loss of revenue.
- 5. **Activity Level**: TotalPurchases, despite its lower overall importance score, is more predictive in the validation set, which might indicate that the frequency of transactions is a relevant factor in understanding churn risk.

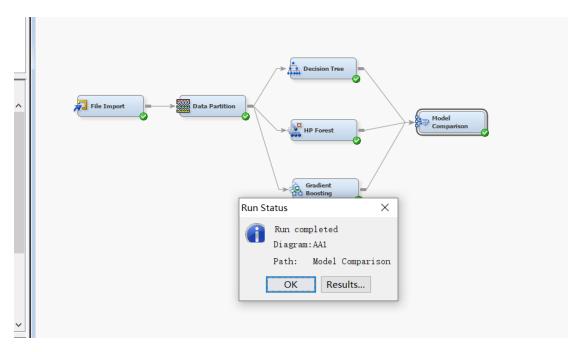
### **Business Strategy Recommendations:**

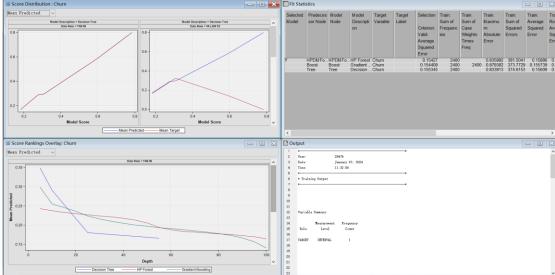
- 1. **Regional Marketing Initiatives**: Develop targeted marketing strategies that address the unique needs and competition in high-churn locations.
- 2. **Loyalty Programs**: Enhance loyalty programs to encourage higher engagement and transition customers to higher membership levels, which are less likely to churn.
- 3. **Customer Engagement**: Implement re-engagement campaigns targeting customers who have not made recent purchases to reduce churn.
- 4. **Customer Value Maximization**: Identify high-value customers based on their total spend and devise retention strategies tailored to them, such as exclusive offers or personalized services.
- 5. **Transaction-Based Engagement**: For customers with fewer purchases, consider strategies to increase their purchase frequency, such as offering discounts on subsequent purchases or rewards for frequent shopping.
- 6. **Data-Driven Product Development**: Use insights from the model to inform product development and service offerings that cater to the needs of different customer segments, particularly focusing on those in high-risk churn groups.
- 7. **Continual Monitoring**: Regularly update and monitor the model to ensure retention strategies remain effective and adjust to new patterns as customer behavior evolves.
- 8. **Customer Segmentation**: Use the model's predictions to segment customers based on churn risk and tailor communication, offers, and interventions accordingly.

By acting on these insights, a business can proactively manage customer churn, tailor its marketing efforts, and ultimately improve customer loyalty and retention.

# Compare these three models







The provided report compares three different models: HP Forest (random forest), Gradient Boosting, and Decision Tree, based on their performance statistics.

# Model Comparison:

Average Squared Error (ASE):

- HP Forest: Shows slightly higher ASE on the training set (0.15896) compared to its validation set (0.15427), suggesting the model may be slightly overfitting.
- Gradient Boosting: Has a lower ASE on the training set (0.15574) than HP Forest and a very similar ASE on the validation set (0.15441). This indicates that the model is generalizing well.
- Decision Tree: Presents the highest ASE on both the training (0.15609) and validation sets (0.15534) among the three models, which may indicate it's the least complex model and potentially underfitting compared to the other models.

## Maximum Absolute Error:

• HP Forest: Has the lowest maximum absolute error on the training set (0.84),

- but a slightly higher error on the validation set (0.838) than the Decision Tree.
- Gradient Boosting: Exhibits the highest maximum absolute error on the training set (0.88) but is still competitive on the validation set (0.857).
- Decision Tree: Shows a consistent maximum absolute error between the training (0.83) and validation sets (0.834), which is desirable.

### Root Average Squared Error (RASE):

- All three models have similar RASE on the training set, around 0.39-0.40, indicating that their overall performance is quite close.
- For the validation set, RASE is also similar across all models, hovering around 0.393-0.394.

### Sum of Squared Errors (SSE):

- HP Forest: Has SSE of 381.50 on the training set and 92.562 on the validation set
- Gradient Boosting: Is slightly better with an SSE of 373.77 on the training set and 92.645 on the validation set.
- Decision Tree: Has an SSE of 374.62 on the training set, which is close to Gradient Boosting, but the highest SSE on the validation set (93.207).

#### Analysis:

- HP Forest and Gradient Boosting are performing similarly, with Gradient Boosting having a slight edge in terms of generalization as indicated by its lower training ASE. The random forest model is slightly overfitting, which is evident from the higher training ASE compared to validation ASE.
- The Decision Tree is the simplest model with the least ability to generalize, as indicated by its highest ASE on both the training and validation sets.

#### Model Selection:

- The model selection based on Valid ASE is indicating Gradient Boosting as the selected model (indicated by the 'Y'), probably due to its balance between training and validation error, showing good generalization without significant overfitting.
- While the Decision Tree model is simpler and may be faster to train and score, its slightly worse performance metrics suggest it may not capture the complexities of the data as well as the other two models.
- The HP Forest model, while slightly overfitting, may still be desirable in scenarios where the ensemble approach of multiple decision trees could provide more robust predictions, especially if the model's complexity can be tuned to reduce overfitting.

In conclusion, Gradient Boosting seems to be the preferred model out of the three due to its performance balance, followed by HP Forest and then the Decision Tree. These insights would be used to choose the best model for deployment in a production environment, considering both performance and computational efficiency.

#### The following are key reflections and learning outcomes:

- 1. Dataset Creation and Integration: I used artificial intelligence and personal modification to build a dataset that truly reflected e-commerce customer behavior. This involves using tools such as Talend Data Integration and Talend Data Preparation to integrate and clean the data.
- 2. Selection of analysis methods: I used decision trees, random forests (HP Forest), and gradient boosting models in SAS Enterprise Miner. These methods help extract meaningful insights from complex data sets.
- 3. Model evaluation and optimization: The article discusses how to identify and deal with overfitting problems, such as simplifying the model through pruning, using ensemble methods, and regularization. This demonstrates an understanding of the data and continued focus on model performance during model building.
- 4. Insights into customer behavior: Through the model, I was able to identify key factors that affect customer churn, such as membership level, geographical location, consumption behavior, etc. These insights help develop targeted business strategies.

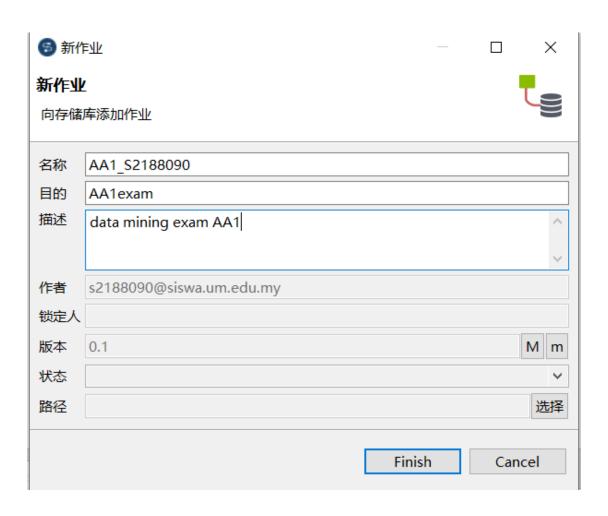
# During this project I faced several major challenges and took some key steps to overcome them:

- 1. Data quality and integration issues: In an e-commerce environment, data is often scattered across different systems and formats. To solve this problem, I use data integration tools (such as Talend Data Integration) to unify data from different sources into a formatted dataset. Additionally, by using Talend Data Preparation, you can clean and transform your data, ensuring its quality and consistency.
- 2. Model selection and optimization: In the data mining process, selecting an appropriate algorithm and optimizing it is a challenge. I address this challenge by comparing the performance of decision tree, random forest, and gradient boosting models. Using SAS Enterprise Miner, I am able to test different algorithms and choose the one that works best for my data.
- 3. Overfitting problem: Overfitting is a common problem when building a prediction model. I combat this problem by using cross-validation and regularization techniques during model building. Additionally, I prune the model appropriately to avoid excessive complexity.
- 4. Understand and interpret model results: A key aspect of data mining is being able to interpret model results. I successfully explained the factors that influence customer behavior through in-depth analysis of the model output and the importance of key variables.

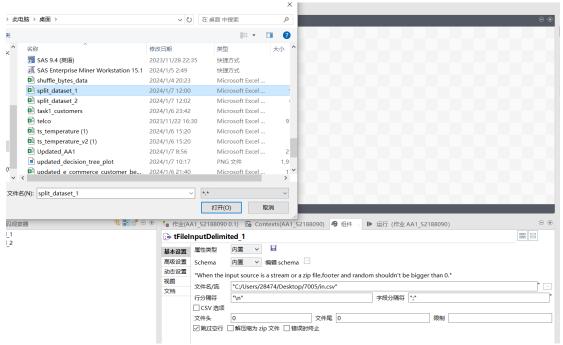
By working through these challenges, I not only improved my data processing and analysis skills, but also enhanced my understanding of complex data patterns. These experiences provide you with a solid foundation for future work in data science.

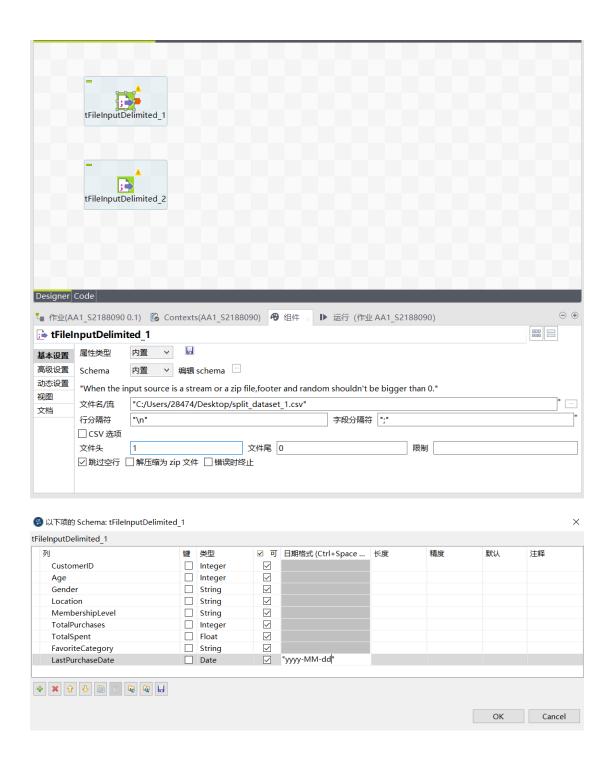
# **Talend DI**

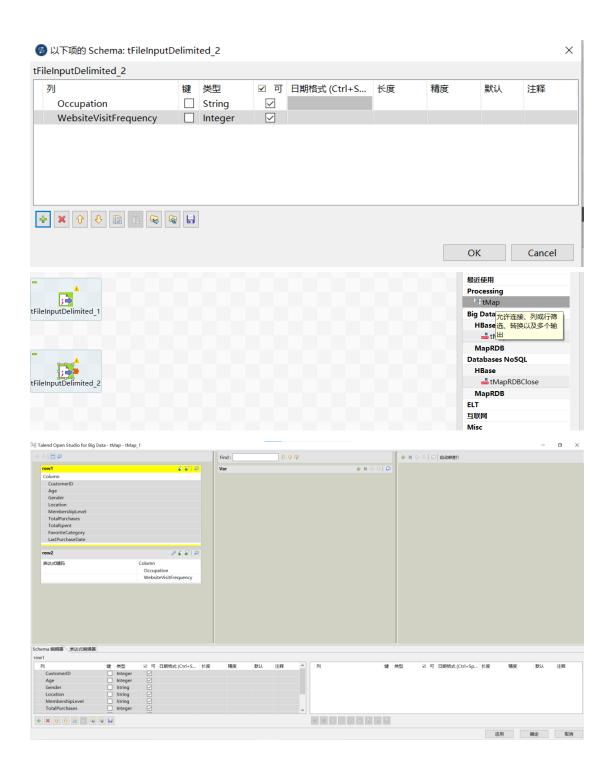


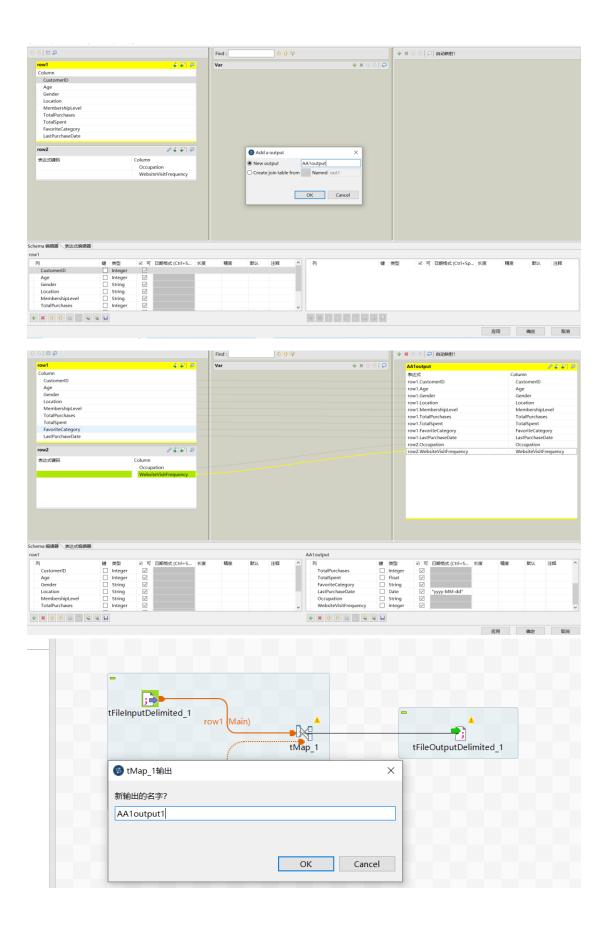


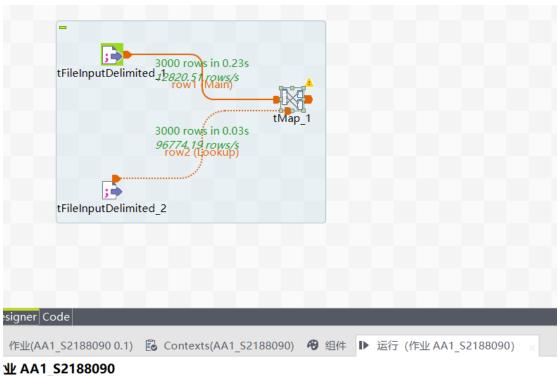




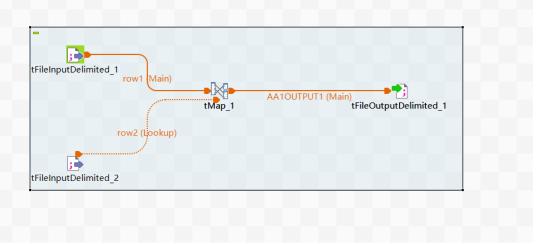


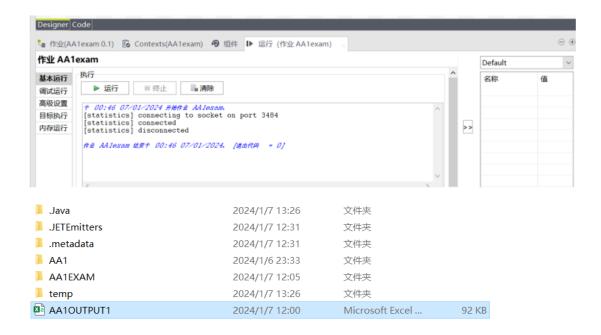




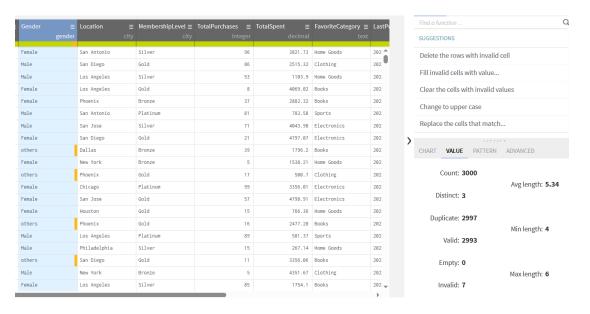






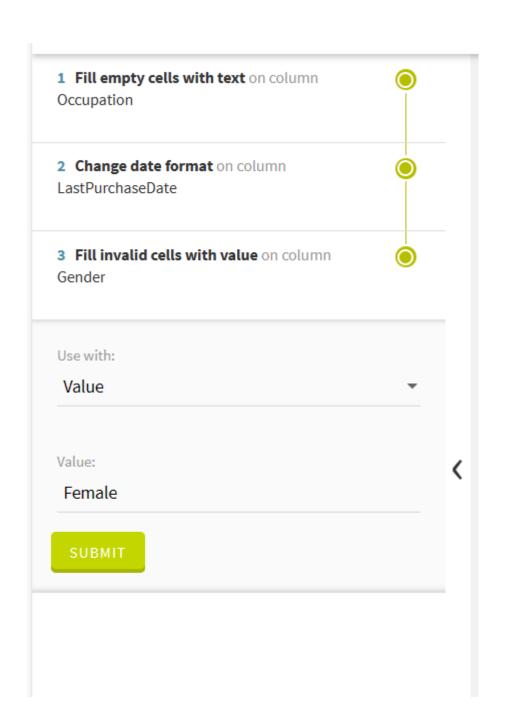


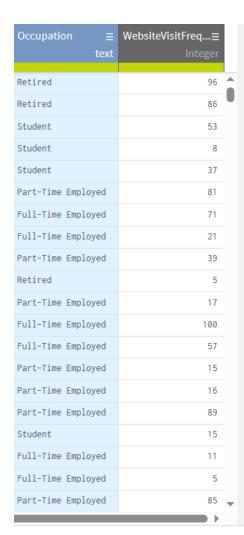
### **Talend DP**

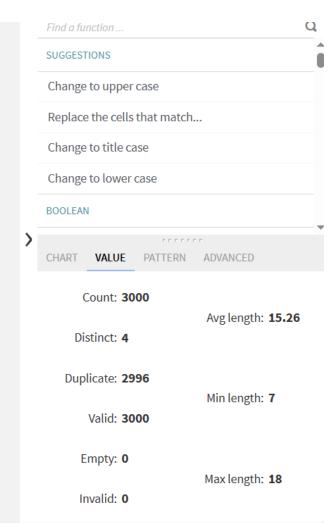


LastPurchaseDate <u></u> date	Churn ≡ integer	Occupation ≡	WebsiteVisitFreq≡ integer
33.0	eger	tent	
2023/11/19	0	Retired	96
2023/5/15	0	Retired	86
2023/5/8	0	Student	53
2023/6/27	0	Student	8
2023/9/5	1	Student	37
2023/7/20	0	Student	
2023/8/13	1	Full-Time Employed	71
2023/8/13	0	Full-Time Employed	
2023/4/27	0	Part-Time Employed	
2023/11/26	1	Retired	5
2023/5/23	1	Part-Time Employed	
2023/3/21	0	Full-Time Employed	100
2023/1/23	0	Full-Time Employed	
2023/9/12	1	Part-Time Employed	15
2023/12/4	0	Part-Time Employed	16
2023/8/6	0	Part-Time Employed	89
2023/9/4	1	Student	15
2023/12/27	0	Full-Time Employed	11
2023/10/6	0	Full-Time Employed	5
2023/12/23	0	Part-Time Employed	85

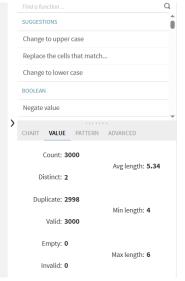


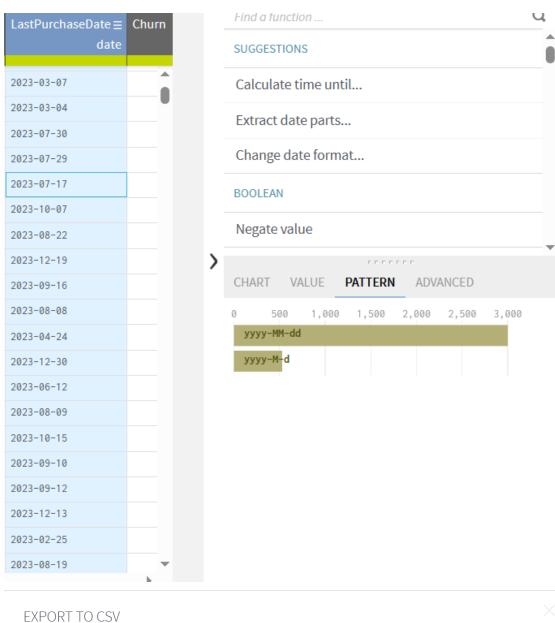




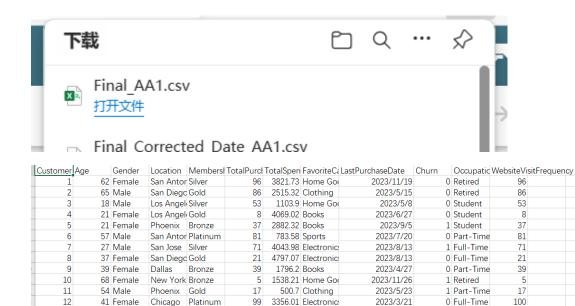








## Delimiter: Comma Filename: Final\_AA1 CANCEL EXPORT



Platinum

Gold

Gold

San Jose

Phoenix

Houston Gold

Philadelph Silver

San Diego Gold

New York Bronze

Los Angel Bronze

New York Bronze

Philadelph Gold

San Antor Silver

Philadelph Gold

Phoenix Gold

San Antor Gold

New York Gold

Los Angel Gold

New York Gold

San Jose Bronze

New York Platinum

Los Angel Platinum

Los Angel Silver

Los Angel Platinum

99

57

15

16

15

11

85

53

76

54

94

68 45

91

92

67

24

29 54

4798.91 Electronic

2477.28 Books

3356.06 Books

4351.67 Clothing

1754.1 Books

4391.25 Sports

673.56 Books

1266.01 Books

4517.87 Home God

286.29 Electronic

729.11 Electronic

373.37 Clothing

4534.5 Clothing

760.5 Sports

1820.07 Electronic

4263.4 Books

4022 93 Clothine

2591.7 Sports

3080.14 Sports

581.37 Sports

766.36 Home God

267.14 Home God

100

57

15

16

89

15

11

85

53

76

54

94

68

45

91

92

67

24

29 54

0 Full-Time

1 Part-Time

0 Part-Time

0 Part-Time

1 Student

0 Full-Time

0 Full-Time

0 Part-Time

1 Part-Time

1 Part-Time

0 Full-Time

0 Part-Time

1 Full-Time

0 Part-Time

0 Part-Time

0 Part-Time

0 Part-Time

0 Part-Time

0 Full-Time

0 Retired

0 Retired

1 Retired

2023/1/23

2023/9/12

2023/12/4

2023/8/6

2023/9/4

2023/12/27

2023/10/6

2023/12/23

2023/7/25

2023/3/13

2023/4/28

2023/3/31

2023/3/7

2023/3/4

2023/7/30

2023/7/29

2023/7/17

2023/10/7

2023/8/22

2023/12/19

2023/9/16

2023/8/8

Sas enterprise miner

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

24 Female

42 Female

42 Female

56 Female

30 Male

19 Male

57 Male

41 Female

64 Female

42 Female

35 Female

55 Female

43 Female

31 Female

26 Female

27 Male

38 Male

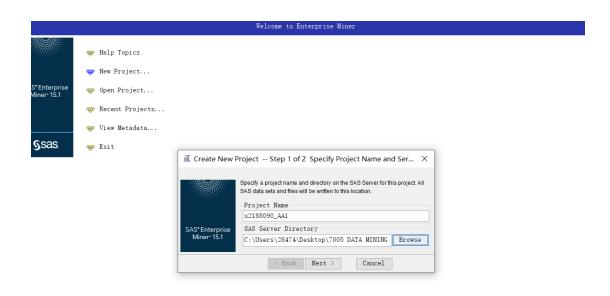
34 Male

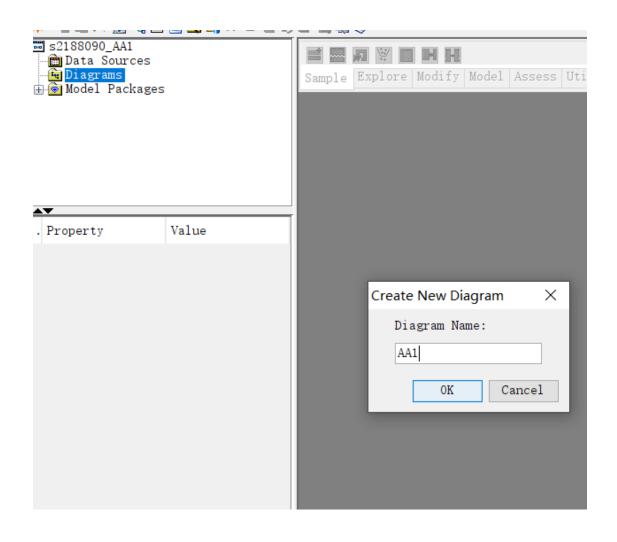
23 Male

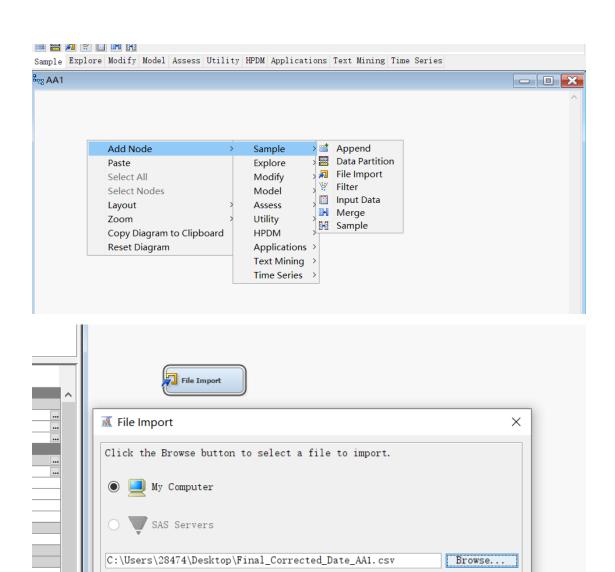
23 Mala

69 Female

69 Female

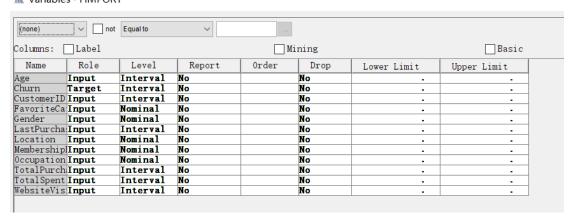






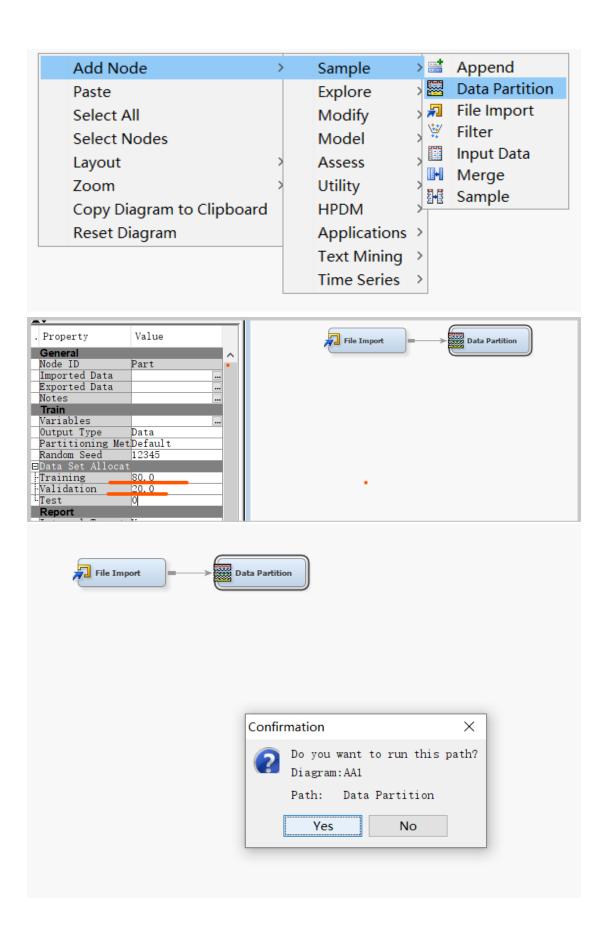
## 

View File Import Types

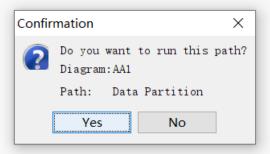


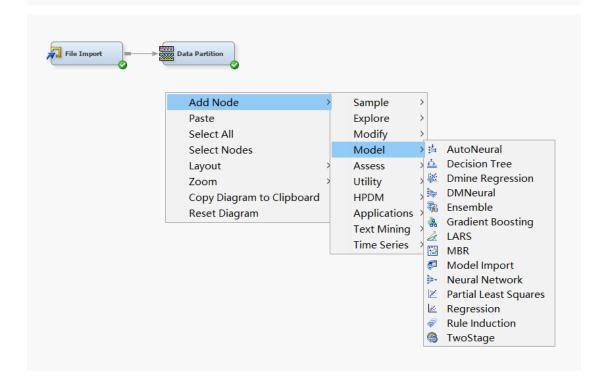
Preview

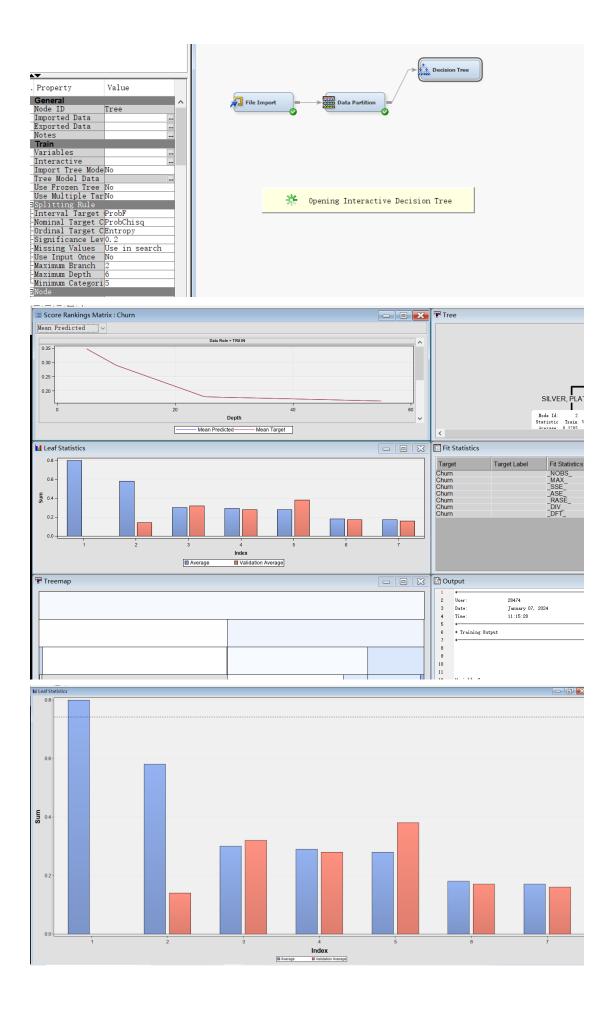
Cancel

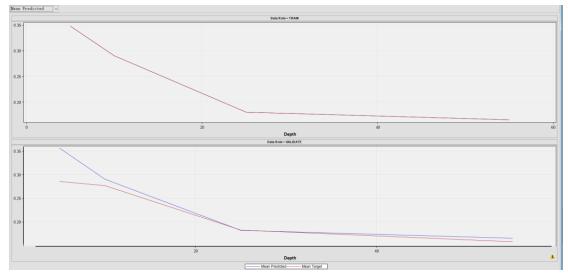




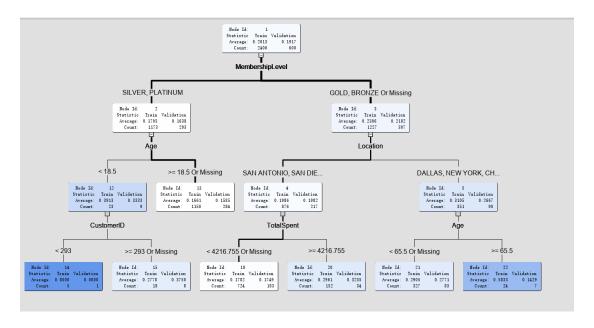


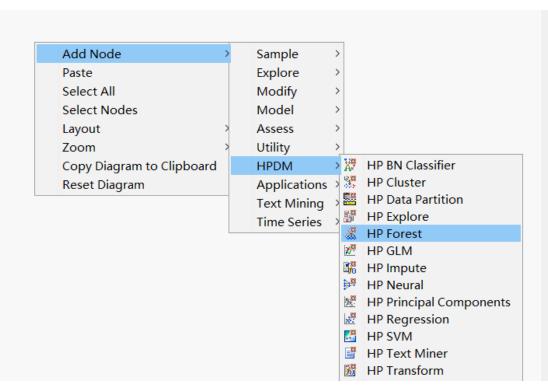


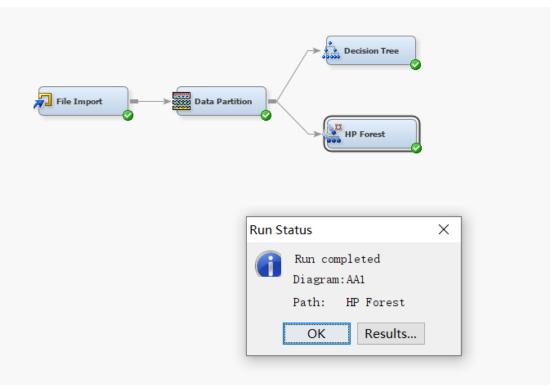




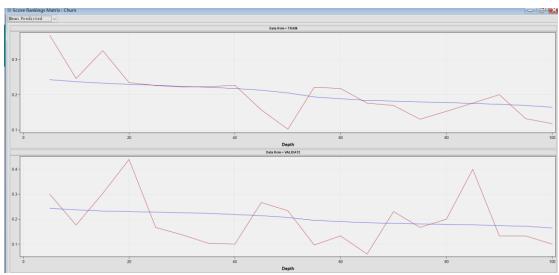




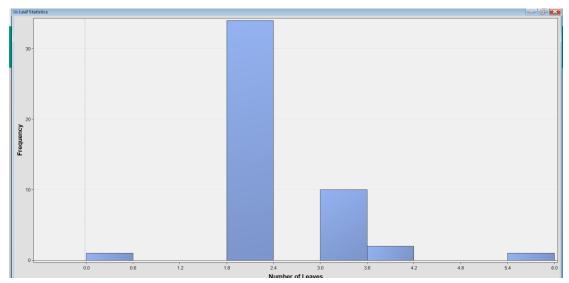


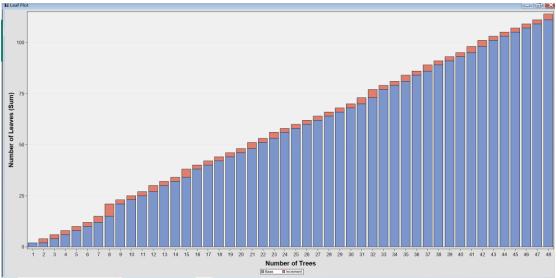


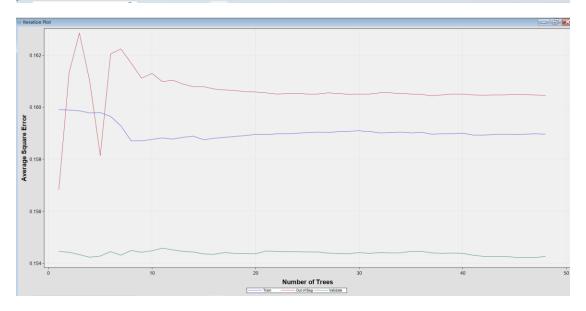


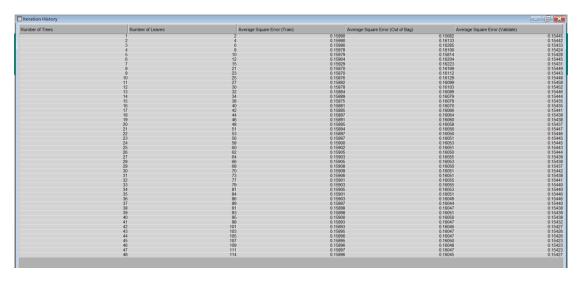


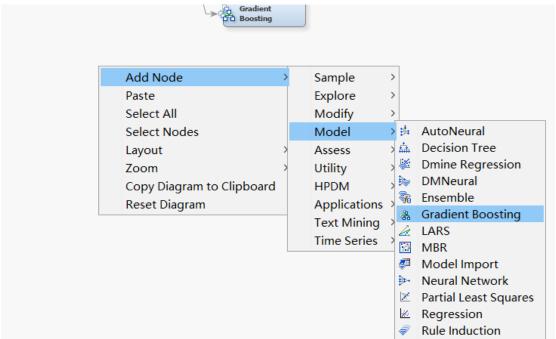
☐ Variable Importance								
Variable	Number	Train:	Train:	OOB:	OOB:	Valid:	Valid:	Label
Name	Variable N	Moon Jamo	Absolute	Mean	Absolute	Mean	Absolute	
	Spilling	Square	Error	Square	Error	Square	Error	
	Rules	Error		Error		Error		
Member	37	0.000921	0.0018	0.00008	0.0009	0.00015	0.0007	
Location	11	0.000409	0.0008	-0.00043	-0.00005	-0.00013	0.00004	
TotalSpent	10	0.000236	0.0004	-0.00028	-0.00005	-0.00014	-0.00008	
Website	2	0.000031	0.00006	-0.00005	-0.00000	-0.00006	-0.00002	
Age	1	0.000017	0.00003	0.00001	0.00001	-0.00004	-0.00002	
Custome	1	0.000032	0.00006	-0.00007	-0.00004	-0.00005	-0.00001	
Favorite	1	0.000015	0.00003	-0.00003	0.00000	0.00003	0.00003	
Gender	1	0.000013	0.00002	-0.00002	0.00000	-0.00001	0.00000	
LastPurc	1	0.000028	0.00005	-0.00007	-0.00000	-0.00004	-0.00000	
Occupati	1	0.000014	0.00002	-0.00001	-0.00000	-0.00003	-0.00001	
TotalPur	0	0.000000	0	0.00000	0	0.00000	0	

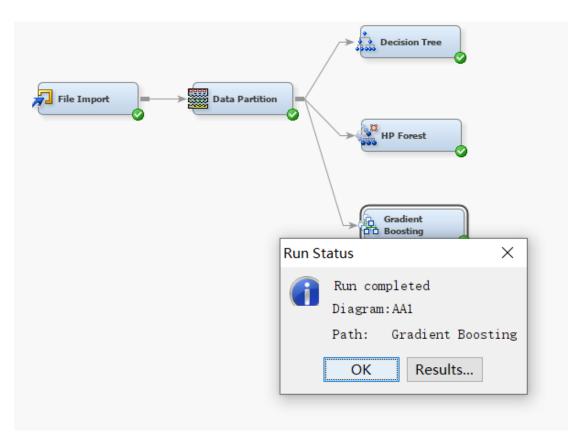


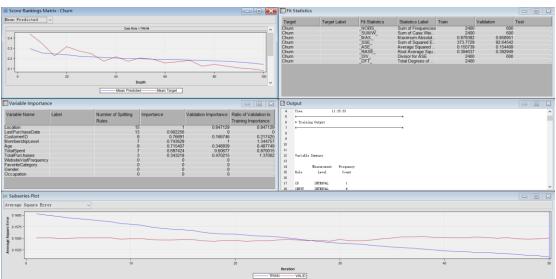


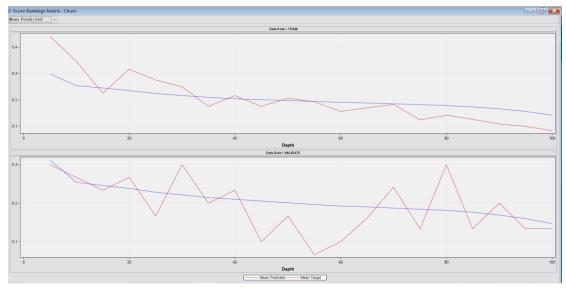












Variable Name	Label	Number of Splitting	Importance	Validation Importance	Ratio of Validation to
		Rules			Training Importance
Location		15	1	0.847129	0.84712
_astPurchaseDate		13	0.902256	0	
CustomerID		8	0.76691	0.166746	0.21742
MembershipLevel		7	0.743629	1	1.3447
Age		8	0.715407	0.348939	0.48774
FotalSpent		7	0.697424	0.60677	0.8700
FotalPurchases		3	0.343218	0.470215	1.370
NebsiteVisitFrequency		0	0	0	
avoriteCategory		0	0	0	
Gender		0	0	0	
Occupation		0	0	0	

Fit Statistics						
Target	Target Label	Fit Statistics	Statistics Label	Train	Validation	Test
Churn		NOBS	Sum of Frequencies	2400	600	
Churn		SUMW	Sum of Case Wei	2400	600	
Churn		_MAX_ _SSE_ _ASE_	Maximum Absolut	0.879382	0.856951	
Churn		SSE	Sum of Squared E	373.7729	92.64542	
Churn		ASE	Average Squared	0.155739	0.154409	
Churn		RASE	Root Average Squ	0.394637	0.392949	
Churn		DIV	Divisor for ASE	2400	600	
Churn		_DIV_ _DFT_	Total Degrees of	2400		

