



UNIVERSITI MALAYA

MASTER OF DATA SCIENCE

FACULTY OF COMPUTER SCIENCE & INFORMATION
TECHNOLOGY

WQD7005 Data Mining

ALTERNATIVE ASSESSMENT

Name	Sze Hui Ling
Matric Number	22075528
Lecturer	Prof. Dr. Teh Ying Wah

Project Overview

The case study focus on data mining techniques using a Kaggle E-commerce Customer Behaviour Dataset [E-commerce Customer Behavior Dataset \(kaggle.com\)](https://www.kaggle.com/datasets/aliakbar1343/e-commerce-customer-behavior-dataset). Talend Data Preparation Tool used to handle missing values and adding attributes like "Churn" and "Age Group." The project use SAS Enterprise Miner for data import and setting variable roles. In-depth data exploration and analysis are conducted to identify key variables influencing 'Total Spend.' Various models, including regression trees, Random Forest, and Gradient Boosting, are implemented to understand and predict customer spending patterns. The goal is to enhance customer segmentation and targeting in marketing by leveraging statistical and machine learning techniques to analyze customer behaviors and spending patterns.

Tasks

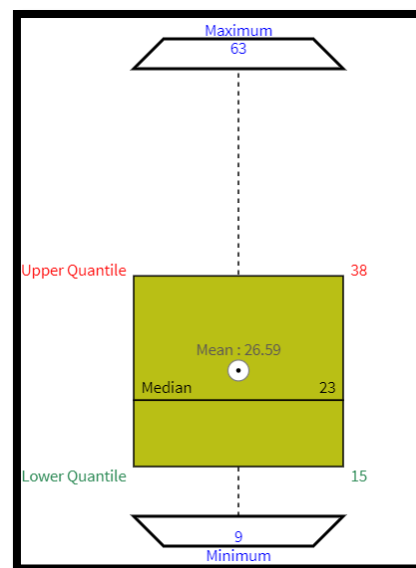
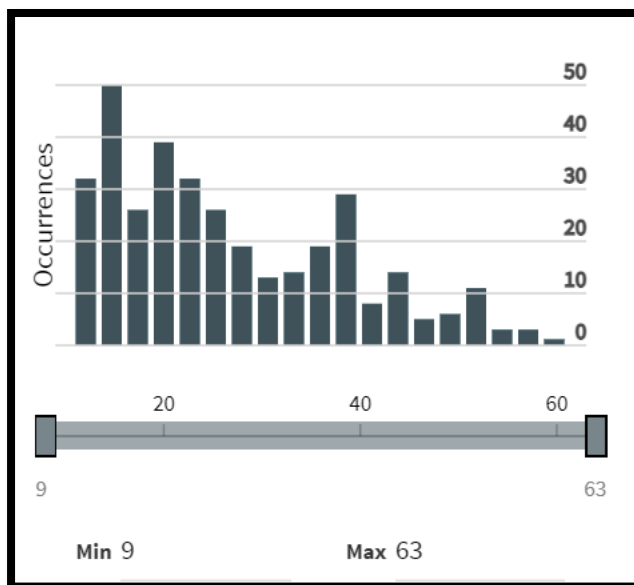
Data Import and Preprocessing: Import your dataset into SAS Enterprise Miner, handle missing values, and specify variable roles.

[15 marks]

1. Handle missing value in Talend Data Preparation Tool

Remove the row contain missing value in attribute "Satisfied level" in using function "Delete the row with empty cell"

2. Add attribute "Churn"



The figures above show the histogram chart and box plot of attribute "Days Since Last Purchase". Define the churn based on attribute "Days Since Last Purchase", label the upper quartile as churned. The 75 percentiles of "Days Since Last Purchase" is 38 days. This targets the 25% customers with the highest inactivity and flags them for churn. Use the "Compare numbers" function to evaluate whether the 'Days Since Last Purchase' for each customer is greater than 38 days. Then, convert the resulting Boolean values: 'True' is replaced with '1' to indicate churn, and 'False' is replaced with '0' to indicate no churn.

E-commerce Customer Behavior Preparation

- Delete the rows with empty cell on column Satisfaction Level
- Compare numbers on column Days Since Last Purchase
- Replace the cells that match on column Days Since Last Purchase_gt_38?
- Replace the cells that match on column Days Since Last Purchase_gt_38?
- Rename column on column Days Since Last Purchase_gt_38?

Filters

Add a filter ...

	Total Spend	Items Purchased	Average Rating	Discount Applied	Days Since Last Purchase	Churn	Satisfaction Level
	decimal	integer	decimal	boolean	integer	integer	text
1	1120.2	14	4.6	TRUE	25	0	Satisfied
2	780.5	11	4.1	FALSE	18	0	Neutral
3	510.75	9	3.4	TRUE	42	1	Unsatisfied
4	1480.3	19	4.7	FALSE	12	0	Satisfied
5	720.4	13	4	TRUE	55	1	Unsatisfied
6	440.8	8	3.1	FALSE	22	0	Neutral
7	1150.6	15	4.5	TRUE	28	0	Satisfied
8	800.9	12	4.2	FALSE	14	0	Neutral

3. Add attribute “Age Group”

Set the age group to 3 categories’: below 30, 30-40 and above 40.

Use function “compare” for cell that value less than 30. Then replace true to “below 30”.

Use function “compare” for cell that value greater than 40. Then replace true to “above 40”.

Then only view above 40, and replace in column churn, after that replace all the false to “30-40”

E-commerce Customer Behavior Preparation

- Replace the cells that match on column Days Since Last Purchase_gt_38?
- Rename column on column Days Since Last Purchase_gt_38?
- Compare numbers on column Age
- Replace the cells that match on column Age_lt_30?
- Compare numbers on column Age
- Replace the cells that match on column Age_gt_40?
- Replace the cells that match on column Age_lt_30?
- Replace the cells that match on column Age_lt_30?
- Delete column on column Age_gt_40?
- Rename column on column Age_lt_30?

Filters

Add a filter ...

	Customer ID	Gender	Age	Age Group	City	Membership Type	Total Spend
	integer	gender	integer	text	city	last_name	decimal
1	181	Female	29	below 30	New York	Gold	1120.2
2	182	Male	34	30-40	Los Angeles	Silver	780.5
3	183	Female	43	above 40	Chicago	Bronze	510.75
4	184	Male	30	30-40	San Francisco	Gold	1480.3
5	185	Male	27	below 30	Miami	Silver	720.4
6	186	Female	37	30-40	Houston	Bronze	440.8
7	187	Female	31	30-40	New York	Gold	1150.6
8	188	Male	35	30-40	Los Angeles	Silver	800.9
9	189	Female	41	above 40	Chicago	Bronze	495.1
10	118	Male	28	below 30	San Francisco	Gold	1520
11	111	Male	32	30-40	Miami	Silver	690
12	112	Female	36	30-40	Houston	Bronze	470
13	113	Female	30	30-40	New York	Gold	1280
14	114	Male	33	30-40	Los Angeles	Silver	820.1
15	115	Female	42	above 40	Chicago	Bronze	530
16	116	Male	29	below 30	San Francisco	Gold	1360
17	117	Male	26	below 30	Miami	Silver	780
18	118	Female	36	30-40	Houston	Bronze	450
19	119	Female	32	30-40	New York	Gold	1170
20	120	Male	34	30-40	Los Angeles	Silver	790
21	121	Female	43	above 40	Chicago	Bronze	585.1
22	122	Male	30	30-40	San Francisco	Gold	1470
23	123	Male	27	below 30	Miami	Silver	710
24	124	Female	37	30-40	Houston	Bronze	430
25	125	Female	31	30-40	New York	Gold	1140
26	126	Male	35	30-40	Los Angeles	Silver	810
27	127	Female	41	above 40	Chicago	Bronze	485.1
28	128	Male	28	below 30	San Francisco	Gold	1500

Age Group

COLUMN **ROW**

Find a function ...

BOOLEAN

Negate value

COLUMNS

Concatenate with...

Delete column

Swap columns...

CONVERSIONS

Convert distance...

Convert duration...

CHART **VALUE** **PATTERN** **ADVANCED**

ROW COUNT

0 50 100 150 200

below 30

30-40

above 40

4. Add attribute “Favorite Category”

Table below is the rule that set for favorite category based on membership type, age group, and gender.

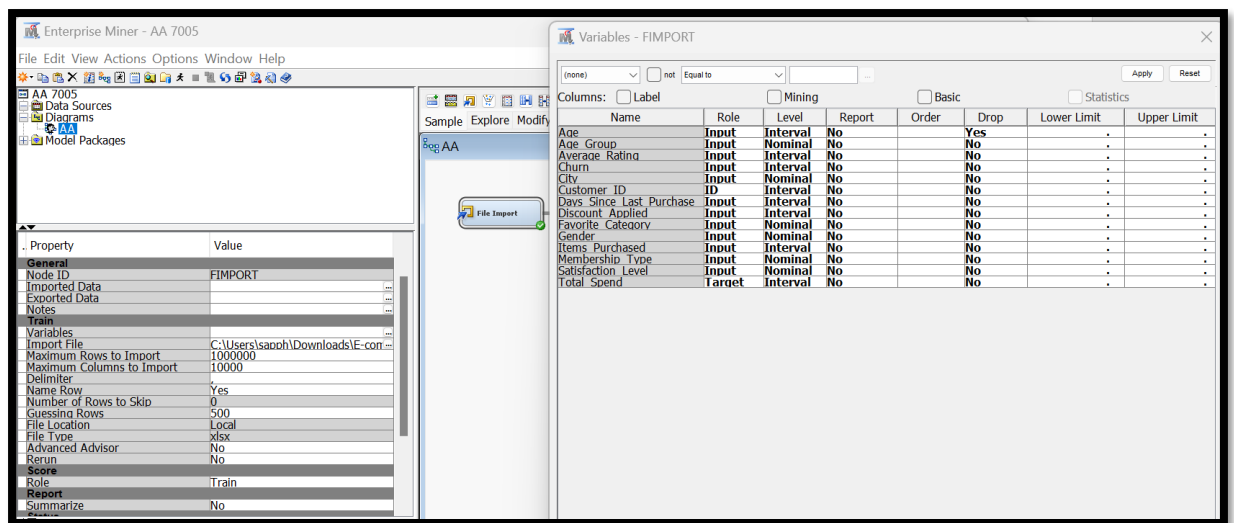
Membership Type	Age Group	Gender	Favorite Category
Gold	Below 30	Male	Electronics
Gold	Below 30	Female	Clothing
Gold	30 - 40	Male	Sports and Outdoors
Gold	30 - 40	Female	Beauty Products
Gold	Above 40	Male	Gourmet Foods
Gold	Above 40	Female	Home Goods
Silver	Below 30	Male	Toys and Games
Silver	Below 30	Female	Books and Media
Silver	30 - 40	Male	Kitchen Appliances
Silver	30 - 40	Female	Gardening and DIY
Silver	Above 40	Male	Home Goods
Silver	Above 40	Female	Gourmet Foods
Bronze	Below 30	Male	Electronics
Bronze	Below 30	Female	Beauty Products
Bronze	30 - 40	Male	Sports and Outdoors
Bronze	30 - 40	Female	Clothing
Bronze	Above 40	Male	Books and Media
Bronze	Above 40	Female	Kitchen Appliances

Then create a new column called “Favorite Category”, then replace each category one by one based on membership type, age group, and gender criteria. Finally, the favorite category got 7 category, Kitchen Application, Clothing, Beauty Products, Toys and Games, Electronics, Sport and Outdoor, Book and Media.

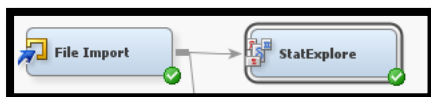
The screenshot displays the Alteryx Designer interface for the 'E-commerce Customer Behavior Preparation' workflow. The central data table shows the following columns: Customer ID, Gender, Age, Membership Type, City, and Favorite Category. The data is filtered by 'Membership Type = Bronze', 'Age Group = above 40', and 'Gender = Female'. The right-hand panel shows the 'Favorite Category' column configuration, with a dropdown menu for selecting the category. The 'Current' value is set to 'Kitchen Appliances'. The 'Replacement' field is empty. The 'Overwrite entire cell' checkbox is unchecked. The 'Apply changes to' dropdown is set to 'All rows'. The 'ROW COUNT' is displayed as 58/348. The bottom right panel shows a bar chart with the following categories: Kitchen Appliances, Clothing, Beauty Products, Toys and Games, Electronics, Sports and Outdoors, and Books and Media.

- Import exported dataset from Talend Data Preparation Tool to SAS Enterprise Miner and set specify variable roles.

The node 'File Import' is drag to diagram and click import file to import exported dataset from Talend Data Preparation Tool to SAS Enterprise Miner. The 'Customer ID' is set as an 'ID' role to uniquely identify each record. The 'Total Spend' has been designated as the 'Target' variable, which allows for the identification of characteristics among high spenders. Understanding these characteristics enables more effective targeting of marketing efforts towards segments with the potential for higher revenue generation or increased sales. All other variables have been assigned the default role of 'Input', serving as predictors in the model. Besides, the age attribute is drop since the dataset already exist age group.



- Data Explore



Explore the data using 'StatExplore', there was 5 interval and 6 nominal input variable and 1 interval target variable 'Total Spend'. The correlation statistics shows that 'Item Purchased' and 'Average Rating' most positive correlation to target variable.

Variable Summary			
Role	Measurement Level	Frequency Count	Correlation Statistics
ID	INTERVAL	1	(maximum 500 observations printed)
INPUT	INTERVAL	5	
INPUT	NOMINAL	6	
TARGET	INTERVAL	1	Data Role=TRAIN Type=PEARSON Target=Total_Spend
			Input
			Correlation
Variable Levels Summary			
(maximum 500 observations printed)			
			Items_Purchased
			0.97228
			Average_Rating
			0.94119
			Discount_Applied
			-0.16853
			Churn
			-0.38441
			Days_Since_Last_Purchase
			-0.54468

7. Remove outlier

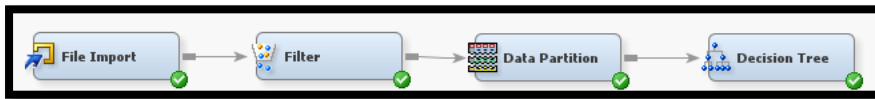
In SAS Enterprise Miner, “Filter” function is used to remove outlier and missing value. “Filter” setting is to 'Rare Values (Percentage)' probably means that the filter will remove values that occur below a certain threshold percentage. Choose to drop the missing values and exclude single-occurrence values as they may be noise. For interval variables, filter out values far from the mean, targeting outliers. As result, there has been 1 observation in Favourite Category been filter out and there is no any outlier presented.

The screenshot shows the SAS Enterprise Miner interface. On the left, the 'Property' pane displays the configuration for the 'Filter' node. The 'General' tab is selected, showing the 'Node ID' as 'Filter'. The 'Train' tab is also visible, showing the 'Default Filtering Method' set to 'Rare Values (Percentage)'. The 'Score' tab is also visible, showing the 'Default Filtering Method' set to 'Standard Deviations from the Mean'. On the right, the 'Diagram' pane shows a workflow diagram with three nodes: 'File Import', 'StatExplore', and 'Filter'. Arrows indicate the flow from 'File Import' to 'StatExplore' and from 'File Import' to 'Filter'.

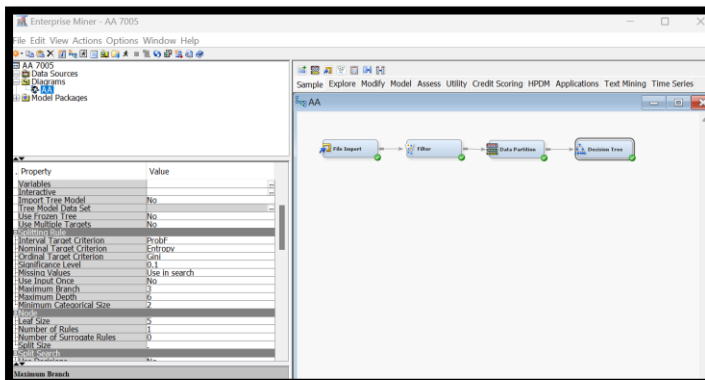
Excluded Class Values						
(maximum 500 observations printed)						
Variable	Role	Level	Train Count	Train Percent	Label	Filter Method
Favorite_Category	INPUT	BOOKS AND MEDIA	1	0.28736	Favorite Category	MINPCT
Number Of Observations						
Data Role	Filtered	Excluded	DATA			
TRAIN	347	1	348			

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

[20 marks]



The data partition is set to a 70-30 split to separate the dataset into training and testing subsets. However, such a split is typically more suited for larger datasets. Therefore, to prevent overfitting and ensure the model's robustness and generalizability, I have chosen to use cross-validation instead of a simple data split. For the decision tree settings, the usual significance level for a splitting rule is 0.05, but I have selected a significance level of 0.1. This more lenient threshold permits splits that may not meet conventional standards of statistical significance but are appropriate for an exploratory model like a decision tree, which seeks to identify potential patterns. I have limited the maximum number of branches to three, aligning with the majority of categorical variables in the dataset that have 2 or 3 categories, with the exception of 'city', which has 6.



Cross Validation	
Perform Cross Validation	Yes
Number of Subsets	10
Number of Repeats	1
Seed	12345

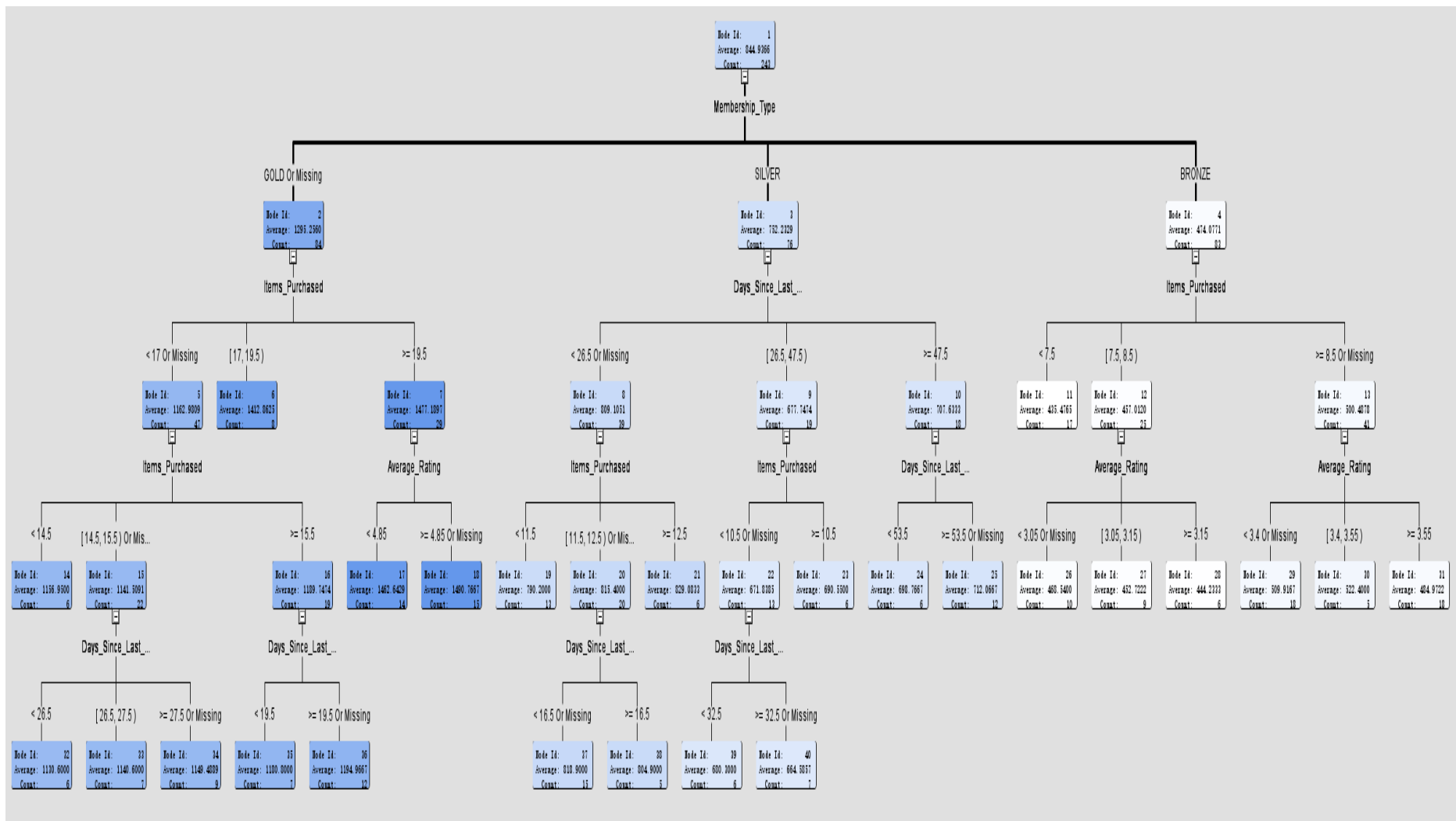
The regression tree analysis, with 'Total Spend' as the target variable, demonstrates how customer behaviors and satisfaction influence spending amounts. The model developed using key customer related features such as 'Items Purchased', 'Days Since Last Purchase' and 'Average Rating'. The model has facilitated the identification of patterns and average spending figures within distinct customer segments, which can inform targeted marketing strategies and customer relationship management.

For example, the 'Membership Type' was initial split, indicating its significant influence on target variable 'Total Spend'. In the 'GOLD' category, the regression tree reveals three spending tiers: members with fewer than 17 purchases have a lower average spend, indicating a base spending level. Those with 17 to 19.5 purchases occupy a mid-tier spending bracket, while members with over 19.5 purchases emerge as the top spenders, with the highest average spend within the group, illustrating a clear positive correlation between the frequency of purchases and spending levels among 'GOLD' members.

In the 'SILVER' category, the regression tree identifies spending patterns based on how recently they've made purchases: those buying within 26.5 days, those between 26.5 to 47.5 days, and those with more than 47.5 days since their last purchase, each with distinct average spending levels, implying more recent shoppers spend differently than less frequent ones.

In the 'BRONZE' category, fewer purchases (less than 7.5) correlate with lower spending, while higher customer ratings suggest distinct spending behaviors, indicating that both purchase frequency and satisfaction levels are key indicators of spending within this group.

Overall, the 'GOLD' category demonstrate that increased purchase frequency aligns with higher average spending within 'Gold' members. The 'SILVER' category demonstrates that the timing of purchases is a predictive factor for spending, with longer intervals between purchases potentially indicating a different type of spending behavior. Besides, the 'BRONZE' category demonstrate that both the quantity of items purchased and the satisfaction level are closely linked to their spending habits. These insights about 'Total Spend' suggest that spending isn't just about how much members buy or how satisfied they are, but also about how these factors interact across different membership levels. This nuanced understanding of spending behavior is vital for tailoring strategies to enhance customer engagement and increase total spending across each segment.

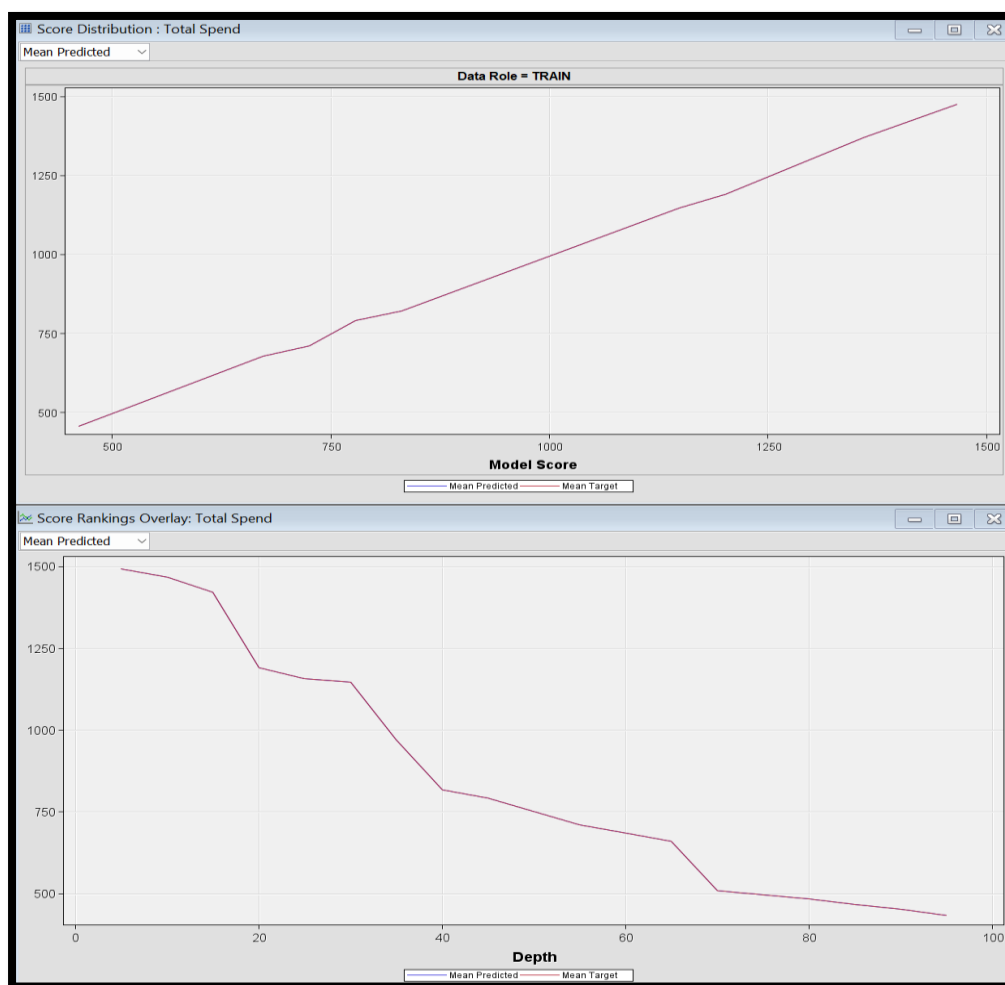


Ensemble Methods: Apply Bagging and Boosting, using the Random Forest algorithm as a Bagging example.

[10 marks]

Random Forest algorithm as a Bagging, the bagging involves creating multiple decision tree, each trained on different sample of the dataset and combined output. The score distribution plot shows the relationship between the predicted total spend (Model Score) and the actual total spend (Mean Target) for the training set. The plot is a 45-degree line plot, which suggests a perfect prediction where every point's predicted value matches the target value. Score ranking overlay plot indicates the model's accuracy at different levels of confidence in predictions. The descending line shows model is less accurate at predicting spending customers. The model is better at predicting higher spenders than lower spenders.

The model achieves accurate predictions with a low ASE and RASE of 7.87, indicating precision in forecasting spending. Besides, the model captures complex spending patterns by utilizing deeper trees, enable to distinguish between different customer behaviors. The model consistently predicts both high and low spending across various depths, demonstrating its reliability in capturing diverse spending behaviors.



			Assessment Score Rankings				Assessment Score Distribution			
			Data Role=TRAIN Target Variable=Total_Spend Target Label=Total_Spend				Data Role=TRAIN Target Variable=Total_Spend Target Label=Total_Spend			
			Depth	Number of Observations	Mean Target	Mean Predicted		Mean Target	Mean Predicted	Number of Observations
Fit Statistics			5	24	1492.18	1492.18				
			10	12	1467.17	1467.17				
			15	22	1420.38	1420.38				
			20	25	1189.99	1189.99				
			25	10	1158.29	1158.29				
Fit Statistics			30	17	1146.48	1146.48				
			35	15	970.01	970.01				
			40	19	818.27	818.27				
			45	32	791.90	791.90				
			55	18	710.96	710.96				
Fit Statistics			60	30	695.40	695.40				
			65	9	661.44	661.44				
			70	34	510.80	510.80				
			80	24	484.42	484.42				
			95	12	468.87	468.87				
Fit Statistics			90	16	453.20	453.20				
			95	28	435.38	435.38				

For Gradient Boosting model, it builds one tree at a time, where each new tree helps to correct errors made by previously trained trees. The Gradient Boosting model is set to run for 50 iterations to build enough trees for accurate predictions while avoiding overfitting and excessive computation time. A learning rate of 0.1 helps the model learn steadily without making drastic changes that could lead to overfitting. Training on 60% of the data allows for a robust learning process while reserving 40% for testing the model's ability to predict new data, ensuring it performs well not just on the training data but also on unseen data. These settings aim to create a well-generalizing model that balances learning complexity and predictive reliability.

The key findings are that 'Items Purchased' and 'Average Rating' are the most influential factors predicting customer spending. The statistical output from the boosting model demonstrates a strong performance with a low Average Squared Error, indicating that the predictions are generally close to the actual spending amounts. The model performs better at predicting higher spending, as shown by the consistent decrease in predicted values with increasing depth, suggesting higher confidence in these predictions. The model's effectiveness is further suggested by the close match between the mean predicted and actual values across different segments of the data. This consistency is crucial for the model's reliability in practical applications. These insights can guide strategies to enhance customer spending, such as encouraging more purchases or improving product ratings.

File Edit View Actions Options Window Help

AA 7005

Data Sources
Diagrams
AA
AA1
Model Packages

Property Value

General

Node ID Boost

Imported Data ...

Exported Data ...

Notes ...

Train

Variables ...

Series Options

N Iterations 50

Seed 12345

Shrinkage 0.1

Train Proportion 60

Splitting Rule

Huber M-Regression No

Maximum Branch 2

Maximum Depth 2

Minimum Categorical Size 5

Reuse Variable 1

Categorical Bins 30

Interval Bins 100

Missing Values Use in search

Performance Disk

Node

Sample Explore Modify Model Assess Utility Credit Scoring HPDM Applications Text Mining Time Series

AA

File Import StatExplore Filter Decision Tree Ensemble Gradient Boosting

Variable Importance

Variable Name	Label	Number of Splitting Rules	Importance
Items Purchased		67	1
Average Rating		29	0.478943
Days Since Las...		23	0.240761
City	City	17	0.208358
Satisfaction Level		6	0.165609
Discount Applied		3	0.036897
Favorite Category		4	0.023039
Churn	Churn	1	0.00344
Gender	Gender	0	0
Membership Type		0	0
Age Group		0	0

Fit Statistics				Assessment Score Rankings				
Target=Total_Spend Target Label=Total Spend				Data Role=TRAIN Target Variable=Total_Spend Target Label=Total Spend				
Fit				Depth	Number of Observations	Mean Target	Mean Predicted	Assessment Score Distribution
Statistics				Statistics	Statistics Label	Train		
_N_OBS_	Sum of Frequencies	347.00		5	18	1495.10	1484.71	
SUMW	Sum of Case Weights Times Freq	347.00		10	24	1469.15	1466.61	
MAX	Maximum Absolute Error	56.27		15	16	1405.96	1405.50	
SSE	Sum of Squared Errors	47359.35		20	12	1187.45	1189.81	
ASE	Average Squared Error	136.48		25	22	1180.14	1168.60	
BASE	Root Average Squared Error	11.68		30	17	1146.48	1145.62	
DIV	Divisor for ASE	347.00		35	16	979.40	986.64	
DFT	Total Degrees of Freedom	347.00		40	16	817.77	815.89	
				45	35	794.39	793.52	
				55	24	707.90	707.59	
				60	9	689.49	691.50	
				65	19	674.51	672.29	
				70	29	632.36	634.71	
				75	7	520.40	499.64	
				80	27	489.39	498.24	
				85	7	470.50	465.28	
				90	15	460.65	461.71	
				95	17	441.35	452.21	
				100	17	433.15	440.61	

Assessment Score Distribution				
Data Role=TRAIN Target Variable=Total_Spend Target Label=Total Spend				
Range for Predicted	Mean Target	Mean Predicted	Number of Observations	Model Score
1433.032 - 1485.285	1476.02	1470.62	49	1459.16
1328.526 - 1380.779	1371.31	1372.31	9	1354.65
1171.767 - 1224.020	1189.54	1187.66	24	1197.89
1119.514 - 1171.767	1148.23	1145.87	35	1145.64
958.248 - 910.501	920.75	877.02	1	884.37
805.995 - 858.248	821.63	817.68	26	832.12
753.742 - 805.995	791.90	791.58	32	779.87
701.489 - 753.742	707.90	707.59	24	727.62
648.236 - 701.489	676.44	677.27	33	675.36
492.477 - 544.730	499.88	501.78	58	518.60
440.224 - 492.477	447.65	452.87	56	466.35