ID: 17147394 | Fathia Farhana bt Agusalim

<div align="center">

**WQD7005 DATA MINING**
**Alternative Assessment 1**

</div>

**E-commerce Case Study (file in github repository can be retrieved in this link**
**https://github.com/slayerliverave/WQD7005-E-Commerce-Dataset** )

The dataset that has been used for the assessment is obtained from Kaggle titled "Customer Behaviour and Shopping Habits Dataset" (https://www.kaggle.com/datasets/zeesolver/consumer-behavior-and-shopping-habits-dataset) which consists of 3,900 purchases with 18 attributes that similar to the dataset structure proposed for the assessment as follow.

| Required Dataset Structure | Availability in proposed dataset | Remarks |
|---|---|---|
| CustomerID | **Yes** | |
| Age | **Yes** | |
| Gender | **Yes** | |
| Location | **Yes** | |
| MembershipLevel | **Yes** | *Note: initial data structure has no membership level however the dataset is being modified such a way that it meets the required data structure based on Subscription Status and Frequency of Purchase* |
| TotalPurchases | **Yes** | Purchase Amount (USD) |
| TotalSpent | **Yes** | Purchase Amount (USD) |
| FavouriteCategory | **Yes** | Purchase category e.g. Clothing, Footwear, Outerwear, Accessories, |
| LastPurchaseDate | No | However other attributes that can be considered in the dataset is:<br>• **Previous Purchases**: Provides information on the number or frequency of prior purchases made by the customer, contributing to customer segmentation and retention strategies.<br>• **Frequency of Purchases** |
| Additional attributes e.g. customer's occupation, frequency of website visits, etc | **Yes** | Color, Season, Review Rating, Shipping Type, Discount Applied, Promo Code, |
| Churn | **Yes** | Subscription Status (Yes/No) - *Indicates whether the customer has opted for a subscription service, offering insights into their level of loyalty and potential for recurring revenue* |

# WQD7005 DATA MINING | ATERNATIVE ASSESSMENT 1

ID: 17147394 | Fathia Farhana bt Agusalim

## 1. Data Import and Preprocessing

### 1.1. Data Import in Talend Data Preparation

Prior importing the data in SAS E-Miner, the dataset is being imported to Talend Data Preparation in order to inspect the data and do necessary preprocessing to handle the missing values and any either anomalies or non-standardisation observed within the data. It is to be noted that, since the dataset in Kaggle is already cleaned data, hence the data is intentionally dirtied in order to showcase and demonstrate how data cleaning and preprocessing are executed for data mining.

Based on the dataset, there is some non-stanrdardisation of the datatype for attributes Season and Size where both datatype is changed to "text" as shown in Figure 1.1 and 1.2 below



**Figure 1. 1 Original datatype for Season and Size**



**Figure 1. 2: Change datatype for Season and Size to "text" format**

## 1.2. Missing Values and Data standardisation

Based on the data inspection, it can be observed that few attributed have missing values for "Review Rating" where 195 records are missing and similarly for "Payment Method" where 390 records are missing from the dataset as shown in Figure 1.3 and Figure 1.4 respectively. The remaining attributes are acceptable where no missing values have been indicated. For handling missing values, the data will then be imputed using SAS E-Miner.
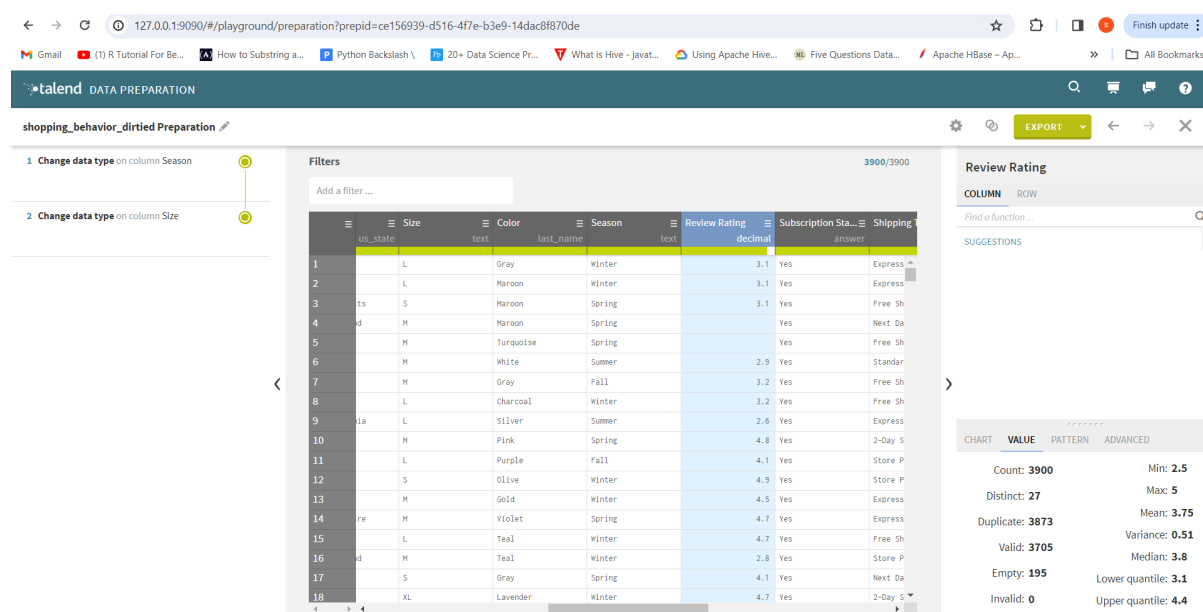


**Figure 1. 3: Total of 195 empty values in Review Rating attributes**
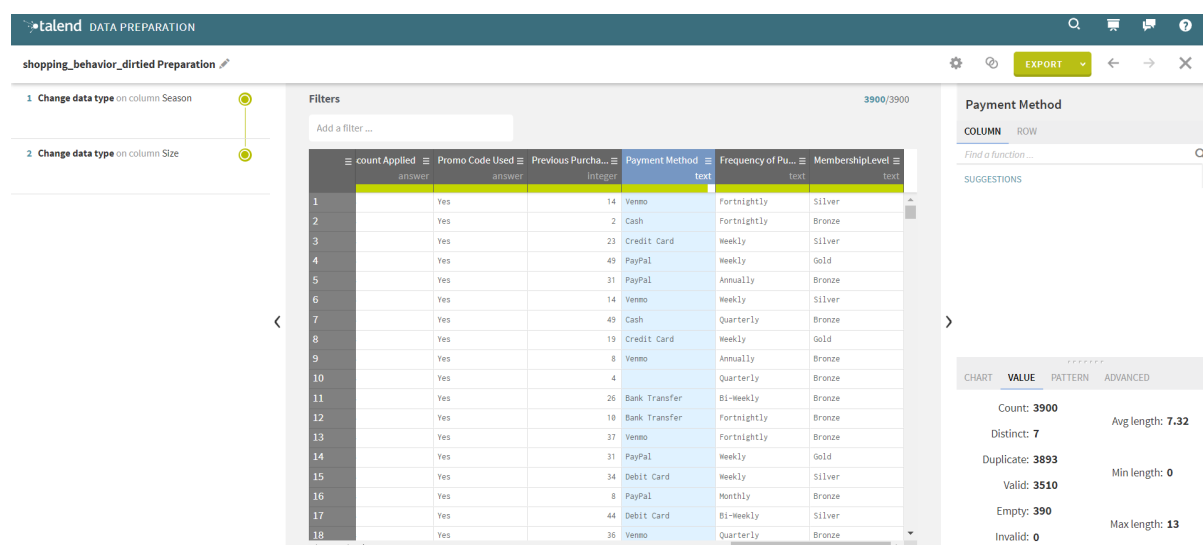


**Figure 1. 4: Total of 390 Empty values in Payment Method**

In order to ensure that Purchase Category is standardized, and more generalisation and customisation can be made to understand the customer behaviour, the Category will be only kept to four (4) main categories such as Clothing, Footwear, Outerwear and Accessories. Hence, the Category type shown in Figure 1.5 for Apparel, Footgear, Electronic and Household will be changed according to the relevant category mentioned above.
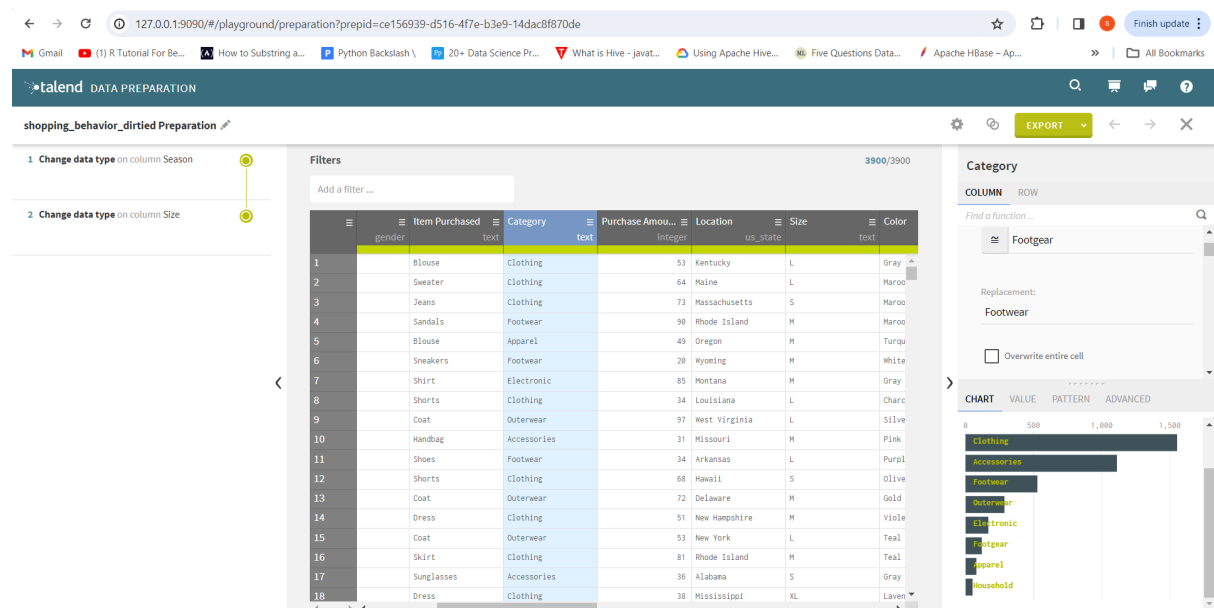
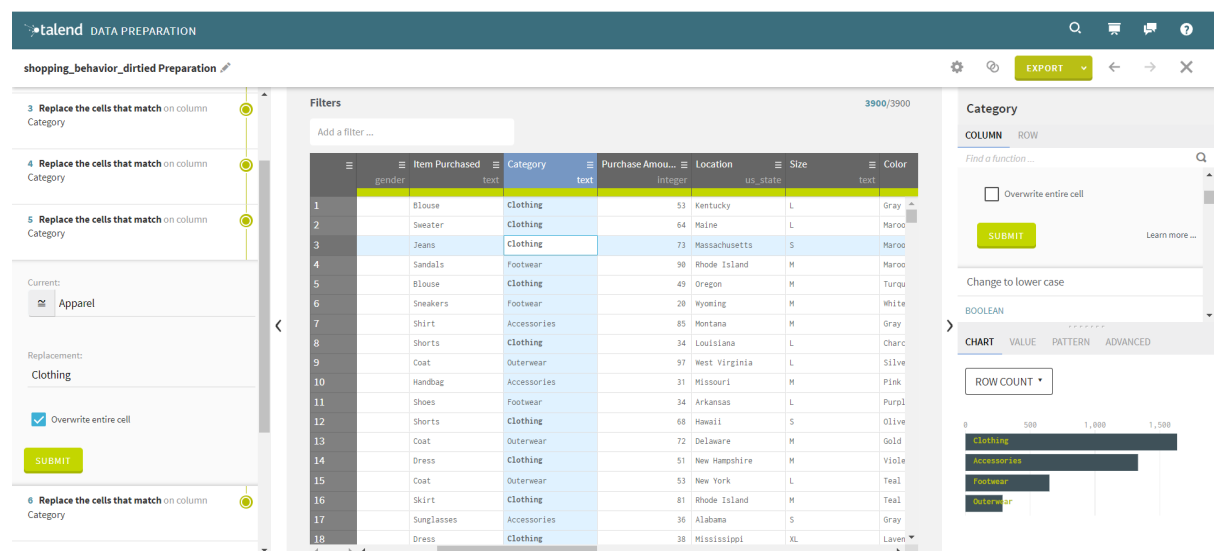**Figure 1. 5: Purchase Category before standardisation**



**Figure 1. 6: Purchase Category after standardisation**

## 1.3. Data Understanding

In Talend Data Preparation, we are also able to understand the data in depth which will assist in getting more insights on customer behaviour. For example, based on Purchase Amount shown in Figure 1.7, it can be summarised that the maximum purchase of the transaction is USD 100 and minimum the purchase is USD 20 whilst average purchase amount for overall records is USD 59.76. With respect to the satisfaction level, maximum Review Rating given by customer is 4.0 over 5.0 with 177 occurrences in which overall satisfaction for the purchase by customer is 3.75 over 5.0.

# WQD7005 DATA MINING | ATERNATIVE ASSESSMENT 1

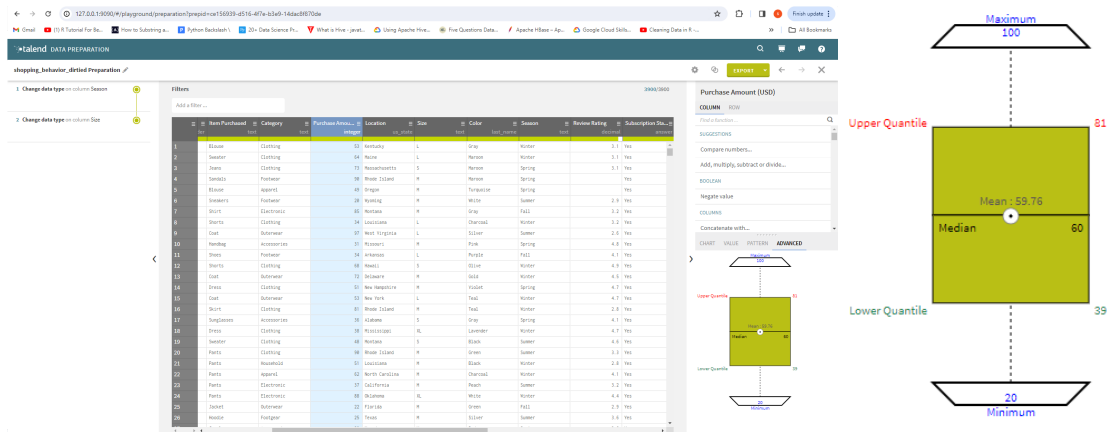ID: 17147394 | Fathia Farhana bt Agusalim



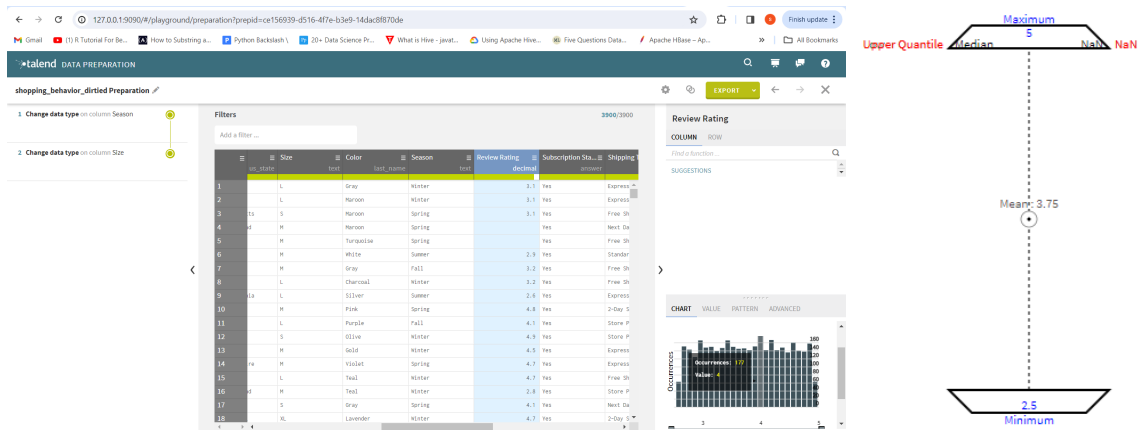**Figure 1. 7: Purchase Amount summary in e-commerce**



**Figure 1. 8: Review Rating summary in e-commerce**

## 1.4. Data Import and Preprocessing with SAS E-Minter

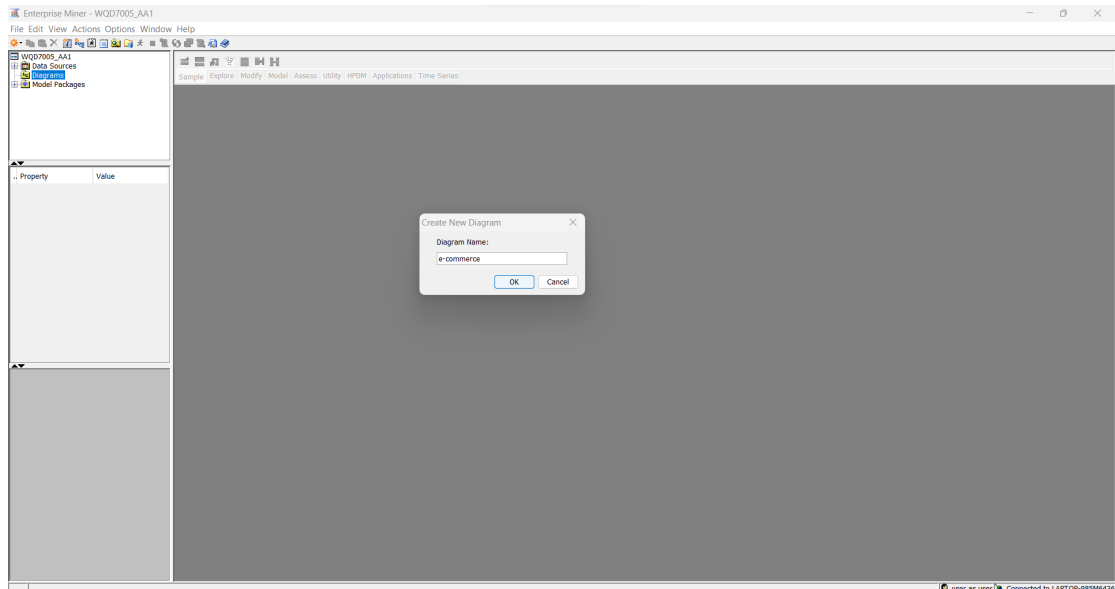To import the dataset in SAS, an 'e-commerce' diagram will be created to New Diagram



**Figure 1. 9: New Diagram "e-commerce" is established**

From the "Sample" tab, select "File Import" and drag the node into the workspace. Right click the File Import and rename the node as "e-commerce"
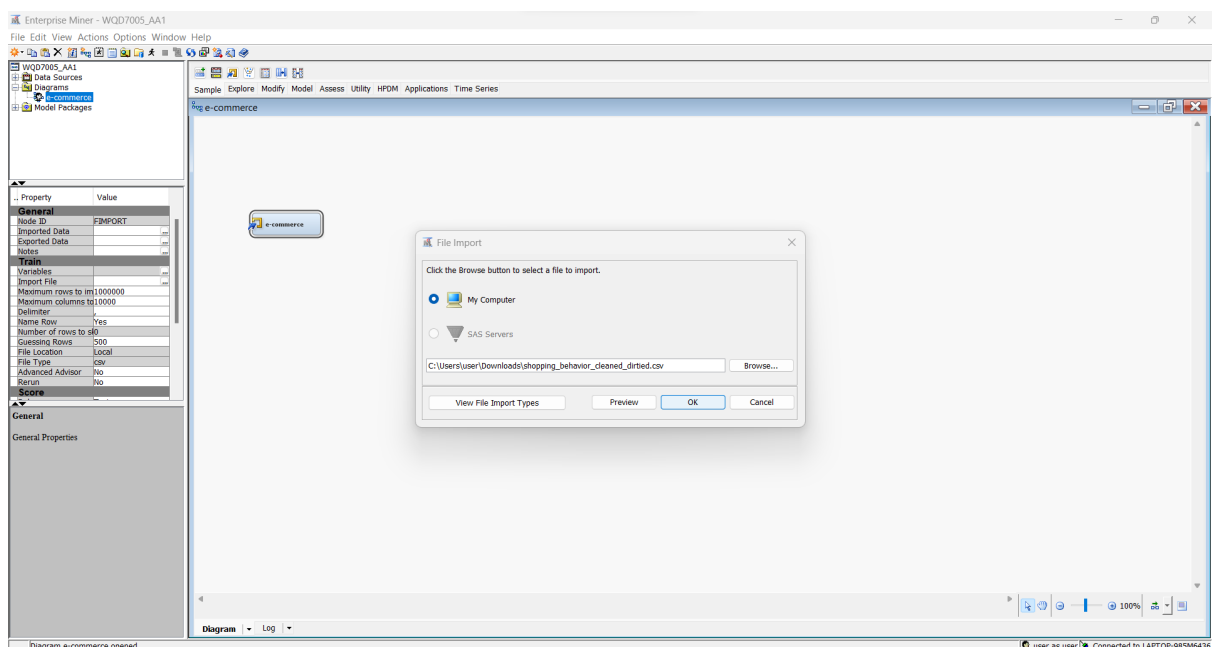


**Figure 1. 10: Importing e-commerce dataset into SAS E-Miner**

It essential to edit the variable for the dataset to understand on customer behaviour. Hence the role of the dataset is amended as follows:

- For Customer_ID, the Role is changed to "Rejected" and select "Yes" in Drop column as this attribute is insignificant for the analysis
- For the Subscription_Status, the Role is changed to "Target" to indicate whether customer will continue the subscription or not. Based on Figure 1.12, we can summarised that majority of the customer opt to "No" subscription with 2897 no. of records (73%)



**Figure 1. 11: Specifying Variable Roles in the dataset**



**Figure 1. 12: Explore Subscription Status data distribution**

Then, Run the "e-commerce" dataset and the summary of Results is demonstrated in Figure 1.13 below:



**Figure 1. 13: Summary of e-commerce dataset**

## 1.5. Data Import and Integration in Talend Data Integration

A new project is created in Talend Data Integration which title as 'e-commerce'

# WQD7005 DATA MINING | ATERNATIVE ASSESSMENT 1

ID: 17147394 | Fathia Farhana bt Agusalim

- To create a New Job:

  - Right-click in the Repository panel on the left under Job Designs.
  - Choose Create job.
  - Provide a name for the job (e.g., "ecommerce ") and a description. Click Finish.



- Add a tFileInputDelimited Component:

  - This component is used to read delimited files like CSV.
  - From the Palette panel on the right, type "tFileInputDelimited" into the search bar.
  - Drag the tFileInputDelimited component to the design workspace.

# WQD7005 DATA MINING | ATERNATIVE ASSESSMENT 1

ID: 17147394 | Fathia Farhana bt Agusalim

- Configure the tFileInputDelimited Component:

  - Click on the tFileInputDelimited component to select it.
  - In the Component panel below, set the properties:
    - File Name/Stream: Browse to and select your CSV file.
    - Row Separator: Set this to "\n" for new lines (most common for CSV files).
    - Field Separator: Set this to  "," depending on file CSV is delimited.
    - Header: Since the dataset obtain Header, hence, set this to 1. .
    - Under the Schema section, click on the Edit schema button. Here, define the columns in the CSV file by adding columns and specifying their data types. This schema should match the structure of CSV file.

- Run the Job to View Data:



Based on the Running the job, the Talend Data Integration encounter an error. This is due to the:

- There was a connection made to a socket on port 3735, which is likely for real-time debugging or statistics gathering.
- Disconnected: The connection to the socket was disconnected normally.
- Job ended: The job ecommerce ended with an exit code of 0 which indicates that the job has completed successfully without errors. If there was an error, we would typically see a non-zero exit code and an error message in the log.

Several attempts on refreshing and creating new job to ensure that job can be successfully Run in the tools such as creating e-commerce2 new job as well as restarting the PC, however similar results achieved.

Hence for execution of Data Preprocessing using Talend Data Integration, the exercise using this tool will be excluded and assessment will be proceeded with Talend Data Preparation and SAS E-Miner.

In summary, three (3) tools have been executed to be used for this e-commerce assessment

## 1.6. Imputing Missing Values

As described in the above section, there are missing values for "Review Rating" and "Payment Method" variables. Hence to handle the missing values, from "Modify" tab, select "Impute" and edit the variables for:

- "Review Rating" to "Mean"
- "Payment Method" to "Distribution". Selecting "Distribution" for imputation means that the missing values will be imputed based on the statistical distribution of the non-missing values of that variable.

# WQD7005 DATA MINING | ATERNATIVE ASSESSMENT 1

ID: 17147394 | Fathia Farhana bt Agusalim



**Figure 1. 14: Select Impute from Modify tab**
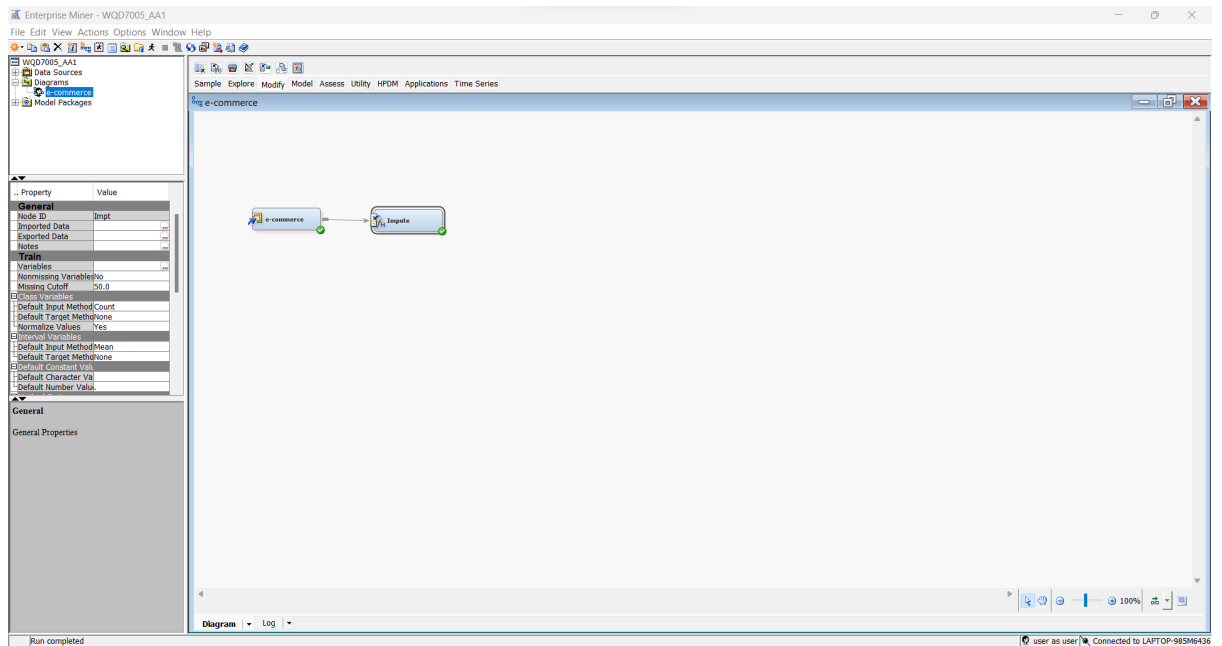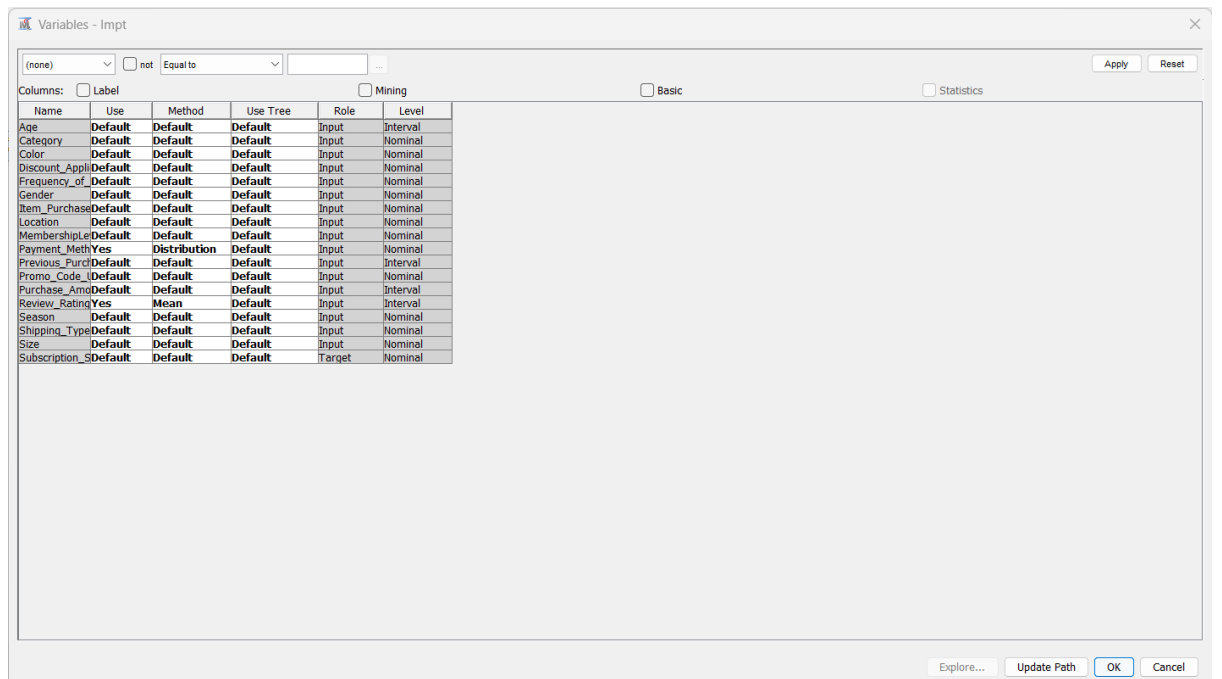


**Figure 1. 15: Impute type for Payment Method and Review Rating**

**Table 1. 1: Representation of Imputation Result**

| Original Dataset | After imputation |
|---|---|
| **Review Rating** – the missing values has been edited with Mean Values | |
|  |  |
| **Payment Method** – the missing values has been edited based on statistical distribution | |
|  |  |

## 2. Decision Tree Analysis

### 2.1. Decision Tree Model

To understand the customer behaviour, decision tree model is established. From "Model" tab, select "Decision Tree".

In order to splitting Role and Node, we need to define the Maximum Branch and Maximum Depth to 2 and 6 and Leaf Size Node to 5. We can maintain all the remaining default value and click "Run" to execute the model

**Figure 2. 1: Decision Tree Model**



**Figure 2. 2: Specifying the Property in Decision Tree Model**

**Figure 2. 3: Overall Result from Decision Tree Model**



**Figure 2. 4: Decision Tree Model based on Prediction Subscription Status**

Based on the above result, the decision tree model appears to predict 'Subscription Status', with different paths in the tree indicating the profiles of customers who are likely to subscribe or not, based on the input variables in "e-commerce" dataset that shown in Figure 2.3 and Figure 2.4 above. It can be seen that each node (representing in the box) represents a decision rule that splits the data based on the best predictor variable at that point.

In the Decision Tree, we can see splits based on 'Discount Applied', 'MembershipLevel', 'Frequency of Purchases', and 'Purchase Amount (USD)'. The splits are binary, dividing records into two groups based on the criterion in each node.
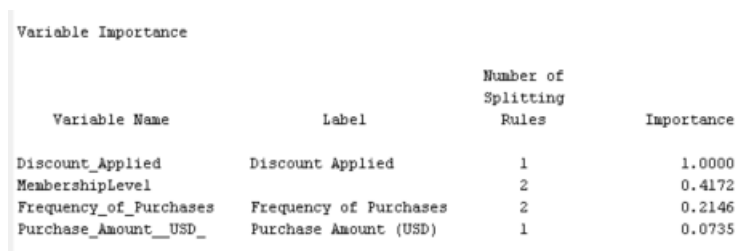
The fact that 'Discount Applied' is the first split suggests it is an important predictor of 'Subscription Status' with important of 1.0000 followed by Membership Level with degree of important 0.4172. The Purchase_Amount_USD_ indicate least important on the customer subscription status with important 0.0735 as shown in Figure 2.5 below.

Variable Importance

| Variable Name | Label | Number of Splitting Rules | Importance |
|---|---|---|---|
| Discount_Applied | Discount Applied | 1 | 1.0000 |
| MembershipLevel | | 2 | 0.4172 |
| Frequency_of_Purchases | Frequency of Purchases | 2 | 0.2146 |
| Purchase_Amount__USD_ | Purchase Amount (USD) | 1 | 0.0735 |

**Figure 2. 5: Variable Importance from Decision Tree Model**

The decision tree model's predictive ability can be evaluated from the lift chart and fit statistics which can be observed from Figure 2.6 and Figure 2.7 respectively. From the chart. it seems that the predictive model to be performing better than random guessing but is not perfect.

The point where the lift line is furthest above the baseline indicates where the model is most effective. The 'Depth' on the x-axis usually represents the percentage of the dataset or population. For example, a depth of 20 could represent the top 20% of predictions from the model. The lift chart typically starts high and then decreases. This is expected because we are capturing the best predictions first (those most likely to be true positives), and as we move through the population, the effectiveness of the model diminishes.

If the lift line is below the baseline, it indicates that the model is performing worse than random chance at that point. Based on the lift chart, the lift line shows that the model has good lift at the beginning, meaning it is effective at identifying positive outcomes, but the effectiveness decreases as we move through the population. This is typical behaviour for a lift chart and indicates that using the model to prioritize the outreach or interventions would be more effective than approaching customers randomly.
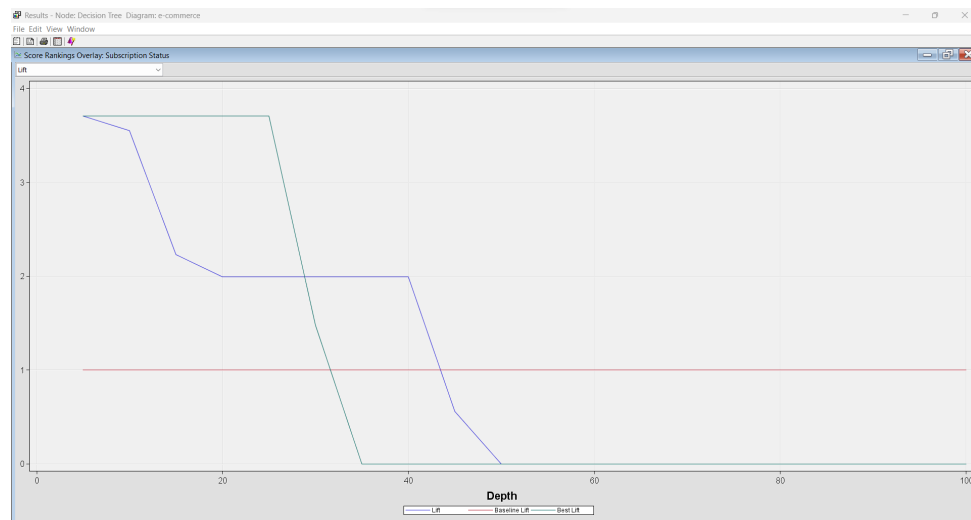
**Figure 2. 6: Lift Chart from Decision Tree Model**

However, it is to be noted that there is misclassification rate in the Fit Statistics with 0.144103 which indicate that the percentage of the predictions were wrong, which can be a straightforward way of evaluating the model's accuracy. The misclassification rate of about 14.41% suggests that the model is relatively accurate, correctly classifying about 85.59% of cases.

The Average Squared Error (ASE) and Root Average Squared Error (RASE) provide measures of the model's error in terms of probability estimations for the classes. Based on the result, it can be seen that the ASE and RASE is lower which indicate the model is making fewer errors on average.

The 'Total Degrees of Freedom' being equal to the 'Number of Observations' suggests that this might be a simple model without many parameters, or it might be an indication of an error in reporting or calculation



| Target | Target Label | Fit Statistics | Statistics Label | Train |
|--------|-------------|----------------|------------------|-------|
| Subscription Status | Subscription Status | NOBS | Sum of Frequencies | 3900 |
| Subscription Status | Subscription Status | MISC | Misclassification Rate | 0.144103 |
| Subscription Status | Subscription Status | MAX | Maximum Absolute Error | 0.854167 |
| Subscription Status | Subscription Status | SSE | Sum of Squared Errors | 613.6688 |
| Subscription Status | Subscription Status | ASE | Average Squared Error | 0.078675 |
| Subscription Status | Subscription Status | RASE | Root Average Squared Error | 0.280492 |
| Subscription Status | Subscription Status | DIV | Divisor for ASE | 7800 |
| Subscription Status | Subscription Status | DFT | Total Degrees of Freedom | 3900 |

**Figure 2. 7: Fit Statistics Result Summary**

In order to improve the prediction model, it would be also essential to look at the performance on a validation set to get a sense of how the model might perform on new or unseen data.

Now, we try to see how does changing the Maximum Branch and Depth and Node for Leaf Size and Number of Rules will affect the result. From the Property, we will specify

the Maximum Branch and Maximum Depth to be 4 and 8 and Node for Lead Size and Number of Rules to 10 (double from the first decision tree model)
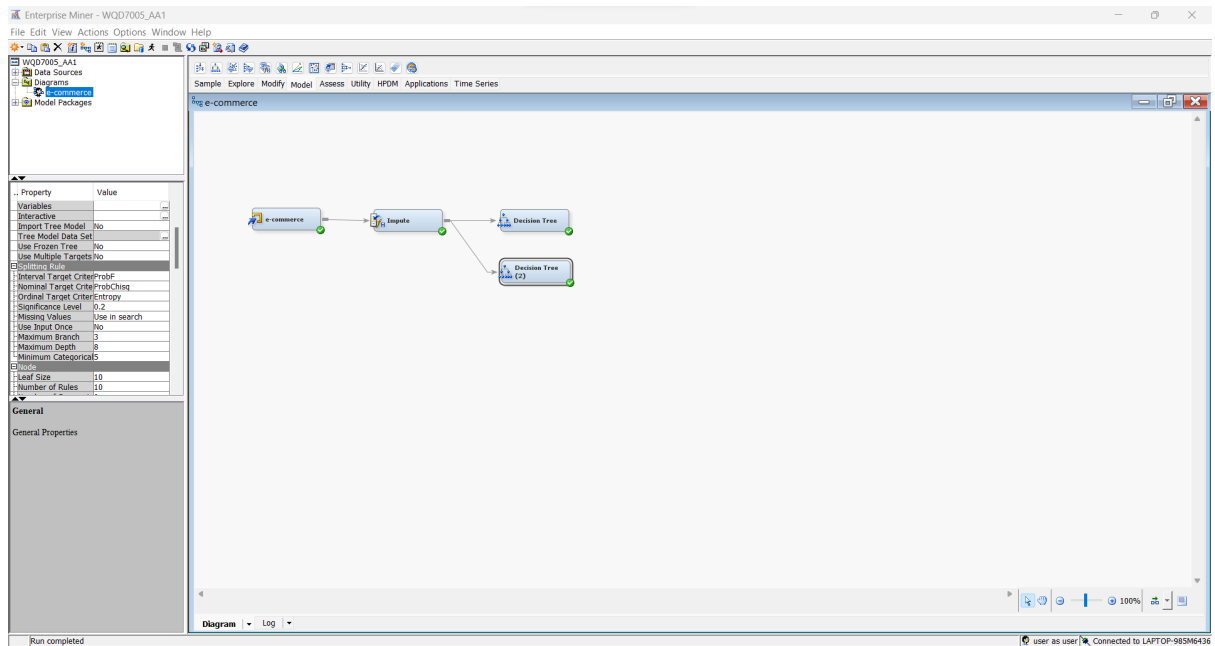


**Figure 2. 8: Decision Tree Model 2**



**Figure 2. 9: Specifying the Property for Decision Tree Model 2**
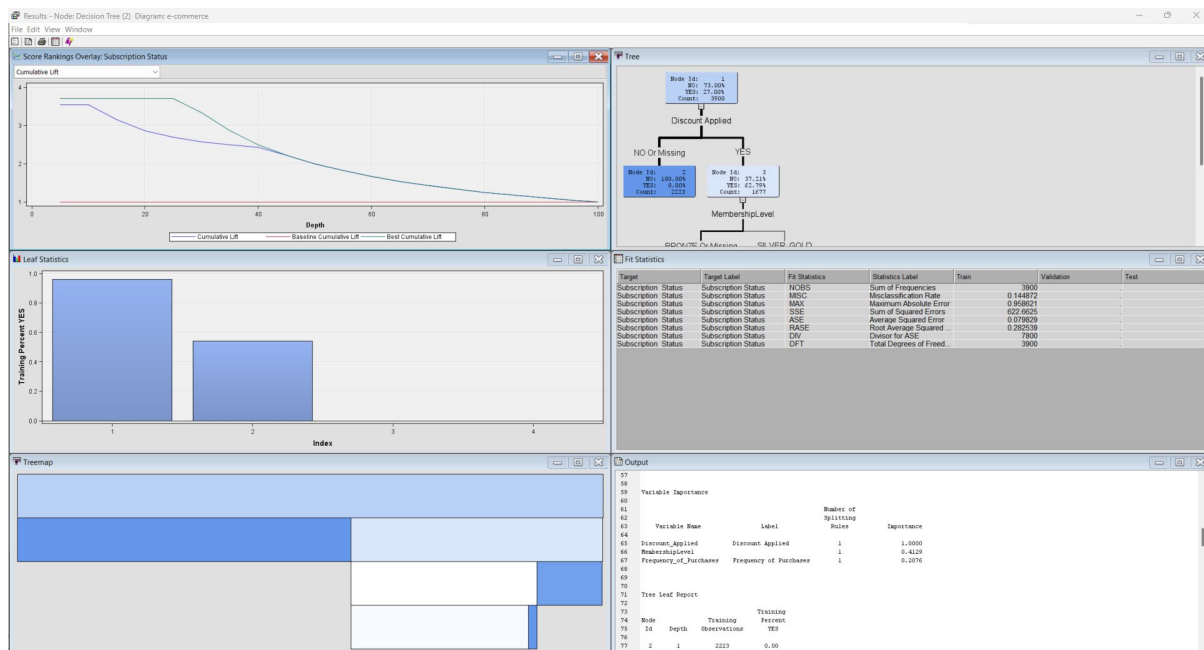
**Figure 2. 10: Overall Result from Decision Tree Model 2**
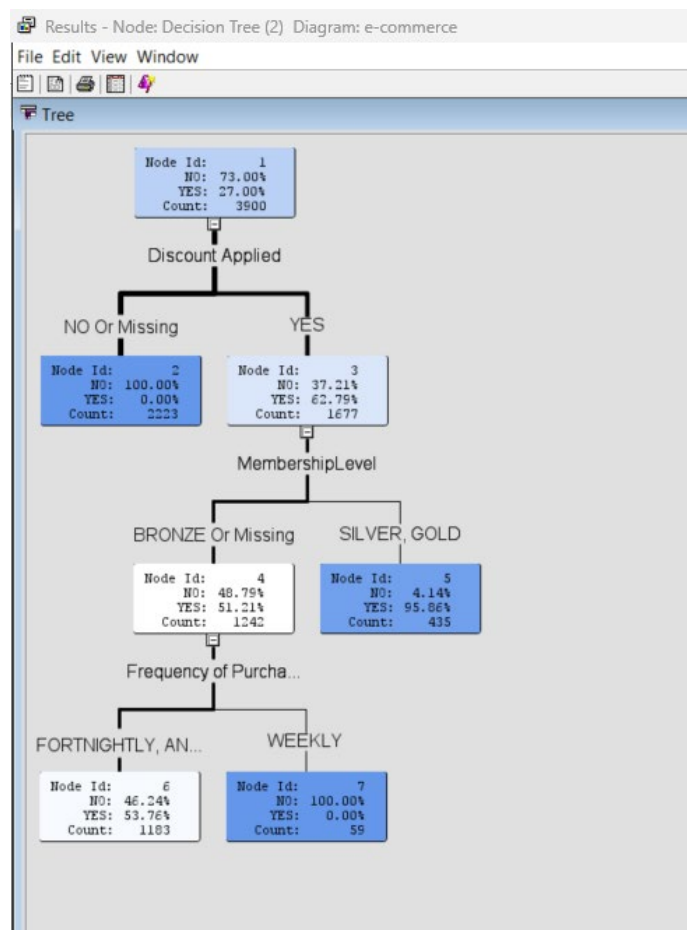


**Figure 2. 11: Decision Tree Model 2 based on Prediction Subscription Status**

Based on the result, the "Discount Applied" remain the first variable importance in the decision tree where it is the most significant predictor of "Subscription Status". For those who received a discount, 'MembershipLevel' is an important next predictor, with 'Frequency of Purchases' also playing a crucial role for customers with higher membership levels. The tree suggests that customers with "SILVER, GOLD" membership and "WEEKLY" purchase frequency have a higher likelihood (95.86% "YES") of subscribing.

In terms of comparison between the first and second Decision Tree Model, the choice between the two depends on the context and the goal of the analysis.

If interpretability and simplicity are priorities (for example, if the rules are to be communicated to a non-technical audience or if the model is to be implemented in a real-world scenario where simplicity is key), the second tree might be preferable.

If the goal is to maximize prediction accuracy and the additional complexity of the first tree offers a significant improvement in performance (as assessed by validation metrics), the first tree might be the better choice. In addition, it is t's also important to consider overfitting whereby more complex trees can fit the training data too closely, capturing noise rather than underlying patterns, which can lead to poor performance on new data. This is why model validation using separate datasets is crucial.

## 3. Ensemble Method Analysis

In performing Ensemble Method, below steps will be executed in the SAS E-Miner. The Ensemble Node creates a new model by taking a function of posterior probabilities (for class targets) or the predicted values (for interval target) from multiple models.

First, we need to drag and drop the "Ensemble" node from the "Model" tab into the workspace. To configure the Ensemble Node, double-click on the "Ensemble" node to configure its properties. We then specify type of ensemble method that need to be executed, in this case for Bagging (Bootstrap Aggregating), we will select "Bagging" as the method. Then, choose "Random Forest" as the algorithm within the bagging method settings. In SAS Enterprise Miner, Random Forest is a separate node, but it can also be part of the ensemble node options.

Connect the data source node to the "Ensemble" node and "Run" the result.

**Figure 3. 1: Ensemble Method Analysis**



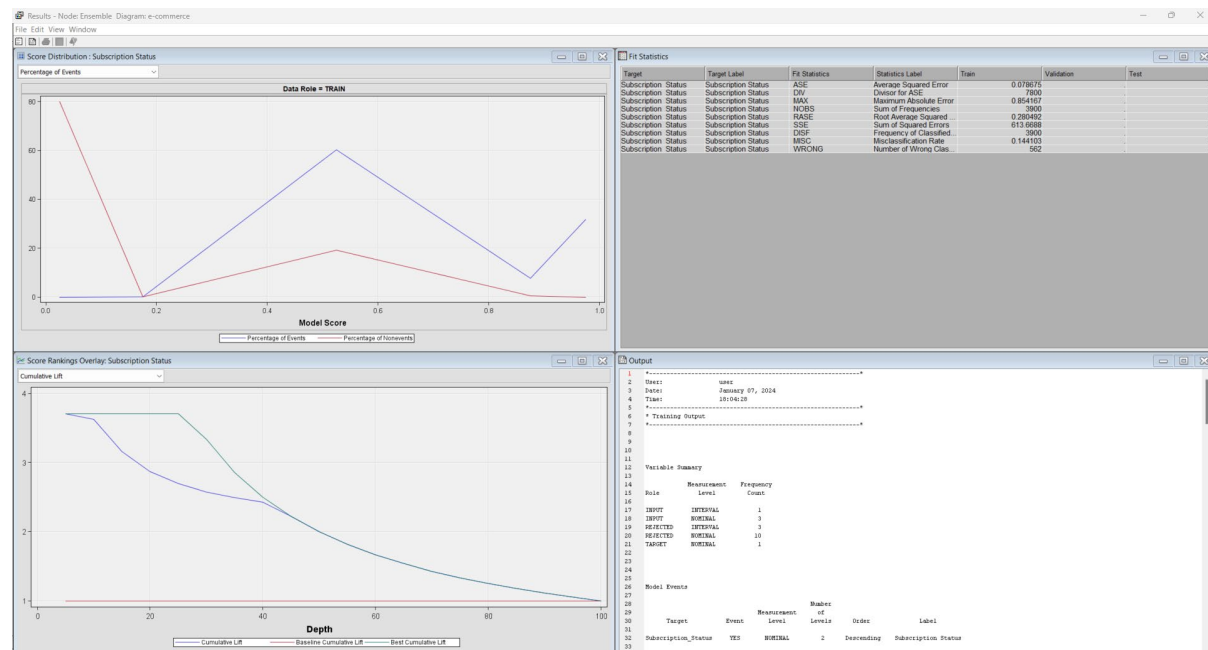**Figure 3. 2: Variables for Ensemble 1 and Ensemble 2**

**Figure 3. 3: Result Summary for Ensemble Method**

Based on the result shown in Figure 3.3, it can be seen be summarised that the model appears to have a moderate level of predictive accuracy, with a misclassification rate of 14% (0.144103). The ensemble method, which seems to be averaging predictions from a set of models, has been able to correctly identify a good proportion of both events (subscriptions – "YES") and non-events (non-subscriptions "NO").

Based on the Score Distribution Plot, it can be seen that the x-axis represents the model score, which is the predicted probability that an observation is an event (YES). Ideally, we want to obtain the blue line chart (events) high on the left (indicating high model scores for actual events) and the red line (non-events) high on the right (indicating low model scores for actual non-events).
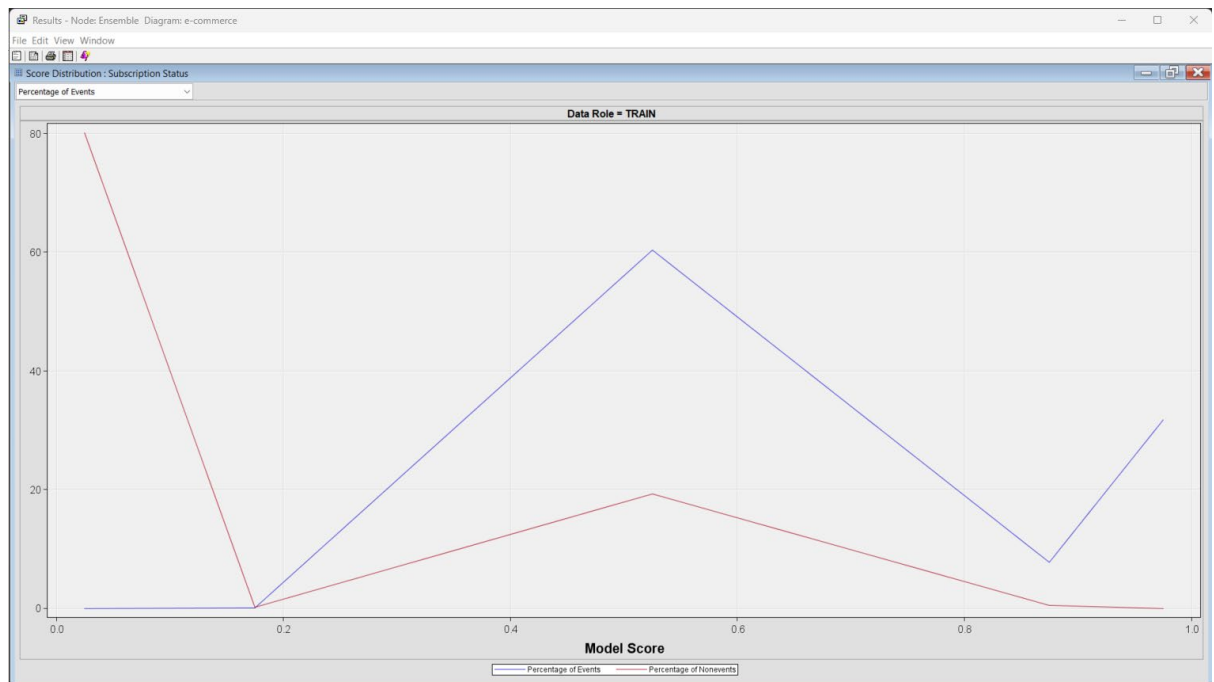
**Figure 3. 4: Score Distribution Plot**

In terms of Cumulative Lift Chart, it can be observed that the chart shows how much better the ensemble model is at predicting the target variable compared to a random guess. The "Cumulative Lift" line should ideally be above the "Baseline Cumulative Lift" line, which indicates the model is better than random chance.
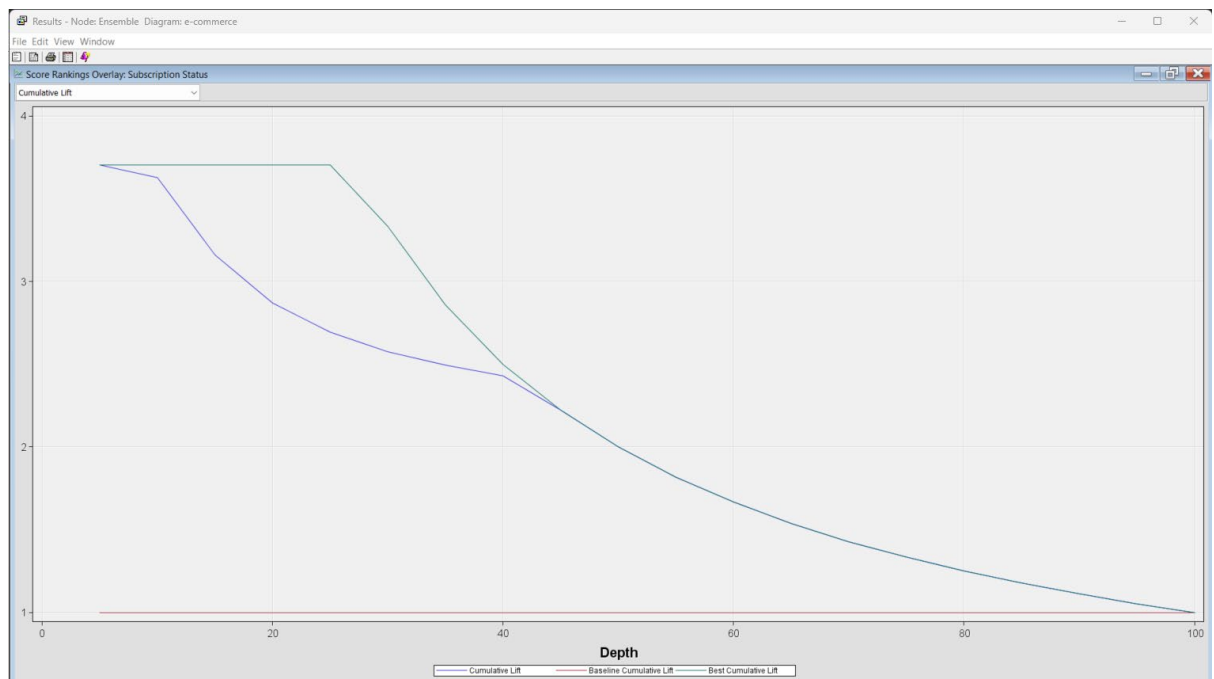


**Figure 3. 5: Cumulative Lift Results**

For the Gain summary result, the model starts off capturing a high number of positive cases (subscriptions) quickly. This is evident in the steep initial slope of the Gain Line.

As more cases are considered, the incremental gains decrease. The Gain Line is above the Baseline, indicating that the model is doing better than random chance at identifying positive outcomes. The Best Gain Line indicates the maximum potential for gain, which is not fully visible in the chart but would typically be at or near the top of the chart, showing the theoretical maximum gain if every positive case was identified before any negative case.

Based on the overall Gain line chart, the ensemble model is effective, especially at the top end, meaning it is good at identifying the most likely positive outcomes. However, as it reach deeper into the population, the effectiveness diminishes, which is typical of such models. The model would be most useful for targeting the top-ranked individuals according to the model's scores, as that's where the most substantial gains are realized compared to random selection.
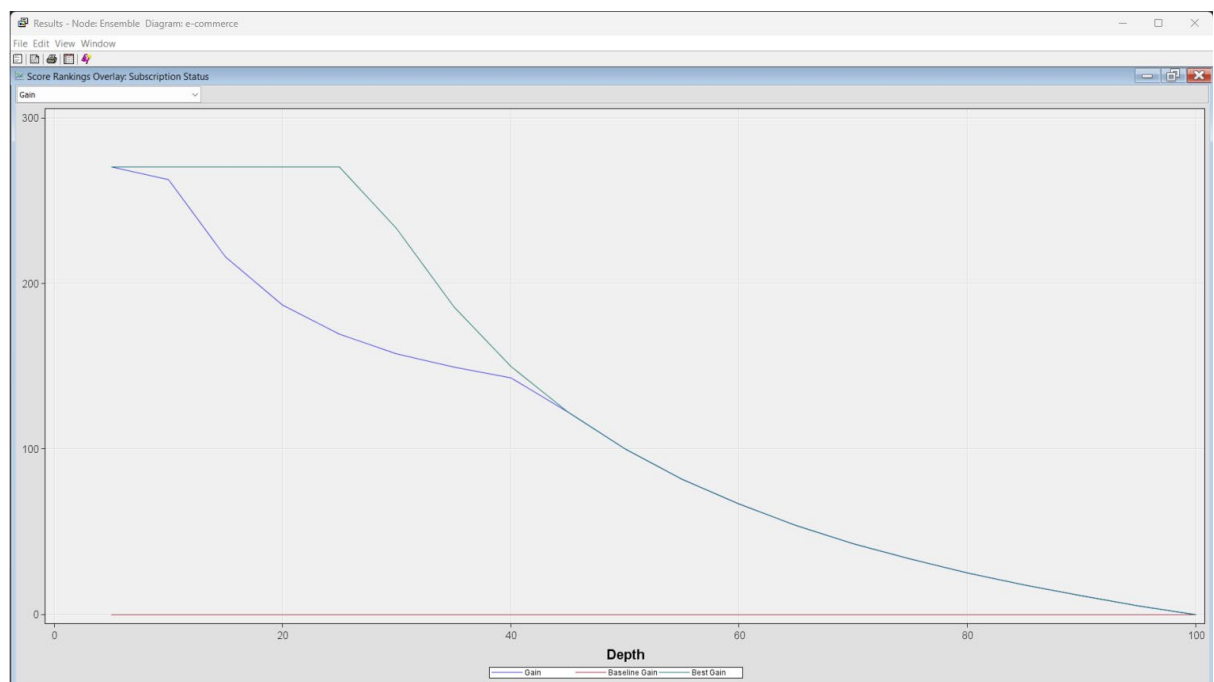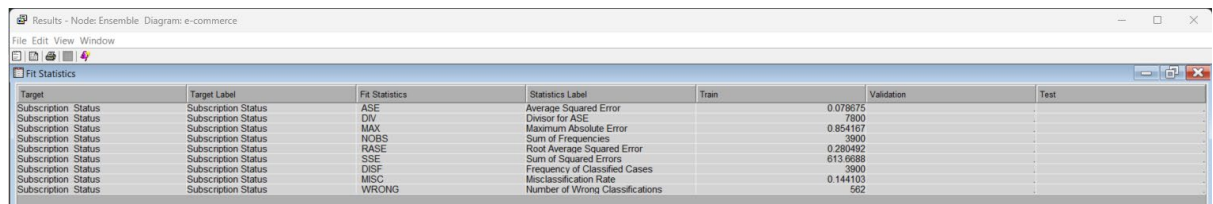


**Figure 3. 6: Gain Results**

# WQD7005 DATA MINING | ATERNATIVE ASSESSMENT 1

ID: 17147394 | Fathia Farhana bt Agusalim



Based on the Fit Statistics, it can be summarised that the ensemble model seems to have a reasonably good fit with an Average Squared Error (ASE) of about 0.08, and a misclassification rate of about 14.41%.

The RASE suggests that the model's predicted probabilities are, on average, about 0.28 away from the actual binary outcomes, which could be considered moderate performance. The number of wrong classifications (562 out of 3900) provides a raw count of errors, which can be useful for assessing the model's performance in more tangible term.