



**UNIVERSITI
MALAYA**
K U A L A L U M P U R

Master of Data Science

Faculty of Computer Science & Information Technology

WQD7005 – Data Mining

Instructor: Dr Teh Ying Wah

ASSESSMENT 1

Name	Matric No
Vickneswary Perumal	S2150313

Introduction

Data Source: <https://www.kaggle.com/datasets/shriyashjagtap/e-commerce-customer-for-behavior-analysis>

Dataset Overview:

The "E-commerce Customer Behavior and Purchase Dataset" is a fabricated dataset generated through the use of the Faker Python library. It emulates a comprehensive e-commerce setting, encompassing diverse facets of customer behavior and purchase history within a digital marketplace. This dataset has been specifically crafted for both data analysis and predictive modeling applications in the realm of e-commerce. It proves to be well-suited for tasks such as forecasting customer churn, conducting market basket analysis, implementing recommendation systems, and performing trend analysis.

Details of the Columns:

This dataset encompasses the following columns:

Customer ID: A distinct identifier assigned to each customer.

Customer Name: The customer's name, generated using the Faker library.

Customer Age: The age of the customer, simulated with Faker.

Gender: The gender of the customer, generated using Faker.

Purchase Date: The date of each customer's purchase.

Product Category: The category or type of the purchased product.

Product Price: The cost of the purchased product.

Quantity: The quantity of the product purchased.

Total Purchase Amount: The overall amount spent by the customer in each transaction.

Payment Method: The method of payment employed by the customer (e.g., credit card, PayPal).

Returns: An indicator of whether the customer returned any products from the order (binary: 0 for no return, 1 for return).

Churn: A binary column denoting whether the customer has churned (0 for retained, 1 for churned).

Objective:

Understanding Purchase Patterns:

Analyze the dataset to identify trends in customer purchasing behavior. This includes popular product categories, peak purchasing times, and the average amount spent per transaction. Insights into these patterns can guide inventory management and marketing strategies.

Customer Segmentation:

Group customers based on common characteristics such as age, gender, and purchasing frequency. This segmentation helps tailor marketing campaigns and promotions to specific customer segments, enhancing the relevance of communication.

Churn Prediction:

Leverage the "Churn" column to predict and understand customer churn. Identify factors leading to customer attrition, such as dissatisfaction or changing preferences. Implement targeted retention strategies to reduce churn and foster customer loyalty.

Payment Method Preferences:

Examine the preferred payment methods among customers. This insight can influence payment processing strategies, partnerships with payment providers, and the development of secure and convenient payment options.

Return Analysis:

Investigate the "Returns" column to understand the frequency and reasons for product returns. Insights into return patterns can guide improvements in product quality, customer service, and overall customer satisfaction.

Personalization Opportunities:

Utilize customer names, ages, and genders for personalized marketing initiatives. Personalization enhances the customer experience and fosters a sense of connection with the brand.

Recommendation Systems:

Explore the potential for implementing recommendation systems based on historical purchase data. Recommending relevant products to customers can increase cross-selling and upselling opportunities.

Trend Analysis:

Identify emerging trends in product categories, customer preferences, and market demand. Stay ahead of industry trends to proactively adapt business strategies and offerings.

Optimizing Marketing Channels:

Evaluate the effectiveness of different marketing channels. Focus efforts on channels that generate the highest customer engagement and return on investment.

Customer Engagement Strategies:

Develop strategies to enhance customer engagement, such as loyalty programs, exclusive promotions, or interactive content. Engaged customers are more likely to remain loyal and contribute to the business's success.

By delving into these aspects of customer behavior, businesses can formulate a well-informed business strategy that addresses customer needs, enhances satisfaction, and maximizes profitability. The goal is to create a customer-centric approach that fosters long-term relationships and sustainable business growth.

Talend Data Preparation

There are several data cleaning steps taken in Talend Data Preparation

- Selecting Rows with Missing Cells in Total Purchase Amount:

The first step identifies and isolates rows where there's incomplete data in the "Total Purchase Amount" column. This aims to address missing values before further analysis.

- Finding and Grouping Similar Text on Payment Method:

This step involves analyzing the "Payment Method" column to detect and group together similar variations in text entries. This is to standardize payment method names for consistency or categorization.

- Changing to Title Case on Payment Method:

The third step involves converting the text in the "Payment Method" column to title case (e.g., "CASH" becomes "Cash"). This aims to improve readability and standardization.

- Filling Cells with Value on Gender:

This step addresses missing values in the "Gender" column by filling them with specific values.

- Deleting Rows with Invalid Cell on Total Purchase Amount:

The final step removes rows that contain invalid or unusable data in the "Total Purchase Amount" column.

The screenshot displays the Talend Data Preparation interface for a dataset named "ecommerce_customer_data_large (dirty) PREPARATION". The workflow on the left includes steps for deleting rows with empty cells, finding and grouping similar text, deleting rows with invalid cells, and filling cells with values. The central pane shows a data table with columns: Customer ID, fr_postal_code, Purchase Date, Product Category, and Product. The table is filtered to show "rows with invalid values". The right pane shows the "Customer ID" column selected, with a summary of statistics for the filtered rows.

Customer ID	fr_postal_code	Purchase Date	Product Category	Product
68	7796	16/9/2022 13:46	Clothing	
69	7796	22/10/2022 12:12	Electronics	
71	7796	21/2/2022 3:50	Clothing	
72	7796	21/2/2022 3:50	Clothing	
94	2642	4/5/2021 2:04	Books	
95	2642	18/5/2020 13:08	Clothing	
96	2642	6/8/2020 2:42	Books	
97	2642	18/3/2022 18:13	Clothing	
98	2642	24/8/2021 15:36	Home	
99	1254	15/8/2023 3:19	Home	
100	1254	21/3/2023 20:00	Home	
140	6829	4/2/2020 20:12	Home	
141	6829	24/6/2020 1:56	Home	
142	6829	1/8/2020 18:26	Books	

Summary statistics for the filtered rows:

- Count: 24350
- Min: 7
- Distinct: 5024
- Max: 50000
- Duplicate: 19326
- Mean: 24954.31
- Valid: 19470
- Variance: 208906141
- Empty: 0
- Median: 24864.5
- Lower quantile: 12475

Talend Open Studio

The cleaned CSV file import into Talend Open Studio.

The screenshot shows the Talend Open Studio Designer interface. The 'Designer' tab is active, displaying a job named 'Job(DataCleaningJob 0.1)'. The 'Component' palette is open, showing the 'tFileInputDelimited_1' component. The 'Basic settings' tab is selected, showing the following configuration:

- Property Type: Built-In
- Schema: Built-In
- Dynamic settings: "When the input source is a stream or a zip file, footer and random shouldn't be bigger than 0."
- File name/Stream: "C:/Users/Vickneswary Perumal/Desktop/sem 4/DM/AA1/ecommerce_customer_c..."
- Row Separator: "\n"
- Field Separator: ","
- CSV options: ☐ CSV options
- Header: 0
- Footer: 0
- Limit:

Ensure data types are corresponding to data

The screenshot shows the 'Schema of tFileInputDelimited_1' dialog box. The table below lists the columns and their data types, with checkboxes for 'K...' and 'N.'.

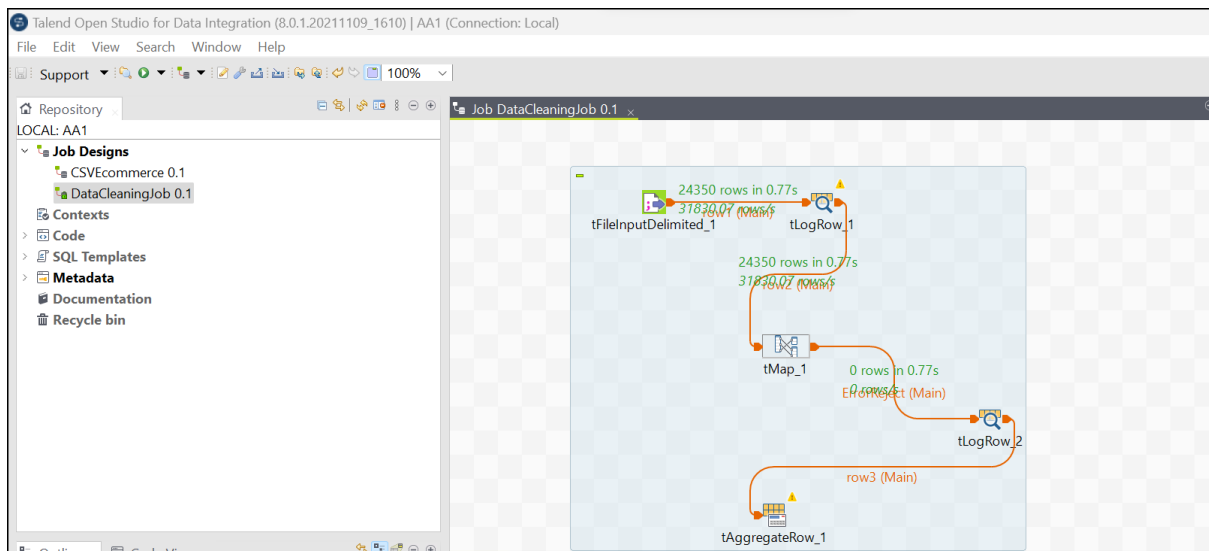
Column	K...	Type	<input checked="" type="checkbox"/> N.	Date Pattern (Ctrl...	Length	Precision	Default	Comment
Customer_ID	<input type="checkbox"/>	Integer	<input checked="" type="checkbox"/>					
Purchase_Date	<input type="checkbox"/>	Date	<input checked="" type="checkbox"/>	"dd/MM/yyyy"				
Product_Category	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>					
Product_Price	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>					
Quantity	<input type="checkbox"/>	Integer	<input checked="" type="checkbox"/>					
Total_Purchase_Amount	<input type="checkbox"/>	Integer	<input checked="" type="checkbox"/>					
Payment_Method	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>					
Customer_Age	<input type="checkbox"/>	Integer	<input checked="" type="checkbox"/>					
Returns	<input type="checkbox"/>	Integer	<input checked="" type="checkbox"/>					
Customer_Name	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>					
Age	<input type="checkbox"/>	Integer	<input checked="" type="checkbox"/>					
Gender	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>					
Churn	<input checked="" type="checkbox"/>	Integer	<input checked="" type="checkbox"/>					

At the bottom of the dialog, there are buttons for 'OK' and 'Cancel'.

tMap is used to remove identical column Customer_Age and Age.

The screenshot displays the Talend Open Studio interface for configuring a tMap component. The top section shows the 'row2' input table with columns: Customer_ID, Purchase_Date, Product_Category, Product_Price, Quantity, Total_Purchase_Amount, Payment_Method, Customer_Age, Returns, Customer_Name, Gender, and Churn. The 'ErrorReject' component is configured to reject rows with errors, showing columns: ErrorMessage and errorStackTrace. The 'out1' output table lists the columns to be written, including Customer_ID, Purchase_Date, Product_Category, Product_Price, Quantity, Total_Purchase_Amount, Payment_Method, Customer_Age, Returns, Customer_Name, Gender, and Churn. The bottom section shows the 'Schema editor' and 'Expression editor' tabs. The 'Schema editor' shows the 'row2' table schema with columns: Customer_ID (Integer), Purchase_Date (Date), Product_Category (String), Product_Price (String), and Quantity (Integer). The 'Expression editor' shows the 'ErrorReject' component schema with columns: ErrorMessage (String) and errorStackTrace (String).

Final Output of data cleaning process in Talend Open Studio



SAS Enterprise Miner

Import Data:

Import your dataset into SAS Enterprise Miner.

Create a New Project:

Start a new project or use an existing one.

Create a Decision Tree Node:

Drag and drop a "Decision Tree" node from the "Common" tab onto the process flow canvas.

Connect Data Source:

Connect the Decision Tree node to the dataset by dragging the arrow from the dataset node to the Decision Tree node.

Specify Training and Validation Percentages:

Data Partition: Determine the percentage of data want to allocate for training and validation. Enter these percentages in the corresponding fields (e.g., 70% for training, 30% for validation).

Configure Decision Tree Node:

Right-click on the Decision Tree node and select "Properties."

Set the target variable to "Churn."

Run the Decision Tree Node:

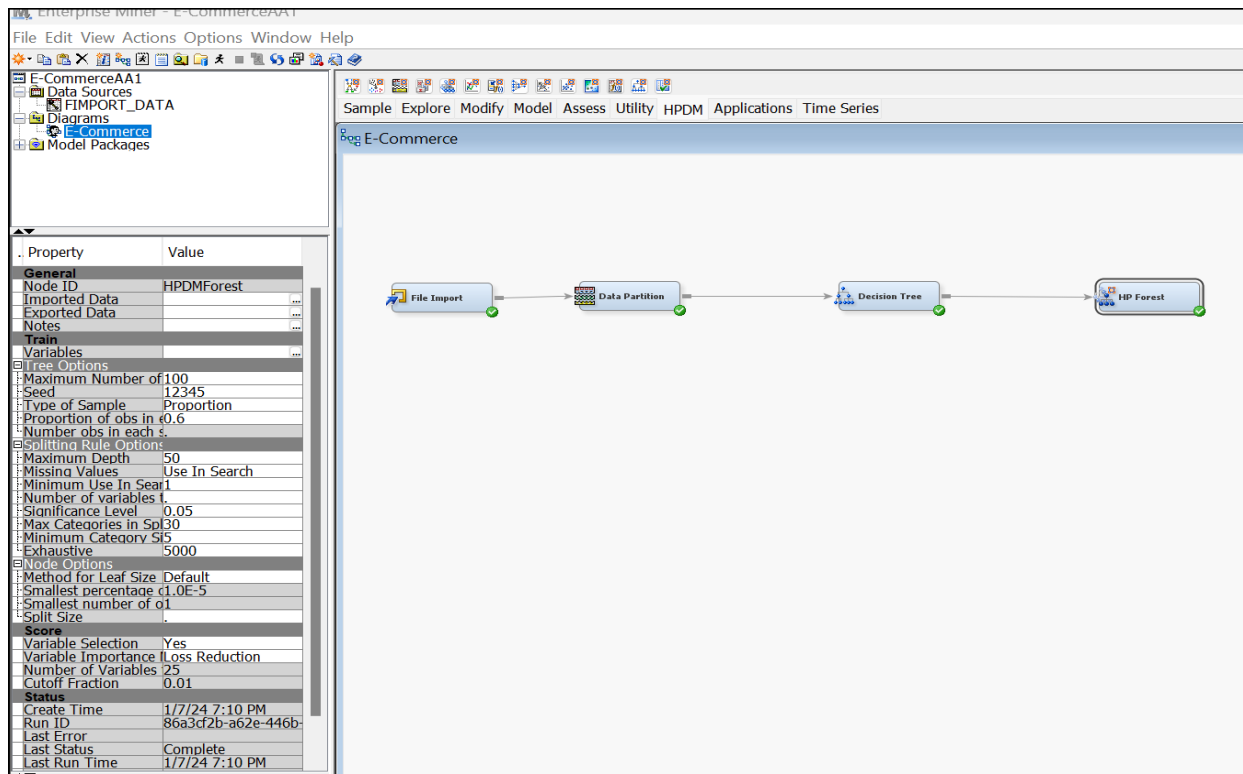
Run the Decision Tree node to build the model. Right-click on the Decision Tree node and select "Run."

Explore the Results:

After the model is built, can explore the results by right-clicking on the Decision Tree node and selecting "Results." This will provide information on the model's performance, summary statistics, and the decision tree diagram.


Decision Tree diagram

This process involves importing the dataset, partitioning the data into training and validation sets, building a Decision Tree model, and analyzing the results. Adjustments to the specific settings and configurations based on the characteristics on dataset and the objectives of analysis.



Results - Node: Decision Tree Diagram: E-Commerce

File Edit View Window



Output

55	Gender	1	0.9429	0.5968	0.6329									
56														
57														
58														
59	Tree Leaf Report													
60														
61	Node	Training	Training	Validation	Validation	Training	Validation							
62	Id	Depth	Observations	Average	Observations	Average	Root ASE	Root ASE						
63														
64	2	1	8636	0.22	3675	0.21	0.41276	0.41081						
65	4	2	4185	0.19	1845	0.20	0.39508	0.40254						
66	11	4	2949	0.18	1240	0.18	0.38081	0.38203						
67	6	3	1122	0.12	477	0.16	0.32741	0.36571						
68	10	4	147	0.31	67	0.33	0.46088	0.47014						
69														
70														
71														
72														
73	Fit Statistics													
74	Target=Churn Target Label=' '													
75														
76	Fit													
77														
78	Statistics	Statistics Label	Train		Validation									
79														
80	_NOBS_	Sum of Frequencies	17039.00		7304.00									
81	_MAX_	Maximum Absolute Error	0.88		0.88									
82	_SSE_	Sum of Squared Errors	2703.68		1178.74									
83	_ASE_	Average Squared Error	0.16		0.16									
84	_RASE_	Root Average Squared Error	0.40		0.40									
85	_DIV_	Divisor for ASE	17039.00		7304.00									
86	_DFT_	Total Degrees of Freedom	17039.00		.									
87														
88														
89														
90														
91	Assessment Score Rankings													
92														
93	Data Role=TRAIN Target Variable=Churn Target Label=' '													
94														
95		Number of	Mean	Mean										
96	Depth	Observations	Target	Predicted										
97														
98	5	8783	0.21929	0.21929										
99	55	4185	0.19355	0.19355										
100	80	2949	0.17599	0.17599										
101	95	1122	0.12210	0.12210										
102														
103														
104	Data Role=VALIDATE Target Variable=Churn Target Label=' '													
105														
106		Number of	Mean	Mean										
107	Depth	Observations	Target	Predicted										
108														
109	5	3742	0.21700	0.21939										
110	55	1845	0.20325	0.19355										
111	80	1240	0.17742	0.17588										

The Decision Tree analysis for the provided dataset emphasizes key factors influencing customer churn prediction. According to the variable importance metrics, "Customer_Age" emerges as the most critical determinant, carrying a ratio of 1.0000, indicating its strong impact on predicting churn. Additionally, "Gender" is considered important, albeit to a slightly lesser extent. The tree leaf report showcases nodes with lower root average squared error, such as Node 6 (Depth 3), indicating higher accuracy in predicting churn for specific customer segments. The overall fit statistics, including the Root Average Squared Error, suggest a reasonably accurate fit of the model to the data. Assessment score rankings and distributions offer insights into how well the model differentiates churn likelihood across various depths and predicted value ranges. The model demonstrates proficiency in predicting churn, with the next steps involving a deeper exploration of the decision tree structure and criteria. Visualizing the decision tree can aid in understanding the specific conditions influencing churn predictions, enabling businesses to refine strategies for customer retention.