**Part 3: Data Preparation Using Talend Data Preparation:**

The data preprocessing is now moved to Talend Data Preparation (TDP) as TDP provide an easier yet robust way to process the data.

After importing the dataset to TDP, the data health bar for every column is inspected. The first step is to change the data type of CustomerID to Integer from Phone Number to eliminate the invalid data issues questioned by the tool. Note that TDP accepts all CustomerID as integers, unlike TDI, indicating the CustomerID values from the original dataset might not have any issues with the data type.



*Figure 8: Data health increased to perfect green after changing the data type of CustomerID.*
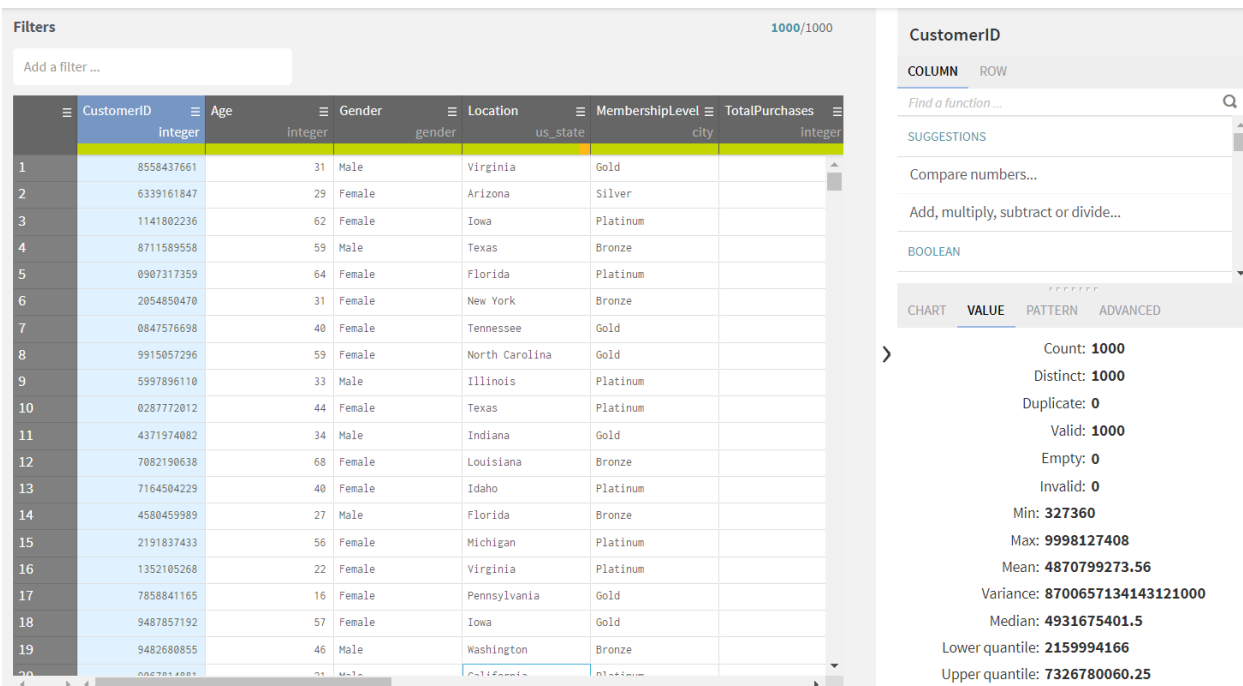


*Figure 9: The right plane shows there is no duplicate CustomerID.*

As shown in the output of Figure 9, as the duplicate check is performed, no duplicates are occurring at the CustomerID column which every row should have unique values. Hence, there is no action taken for the removal of duplicate data.
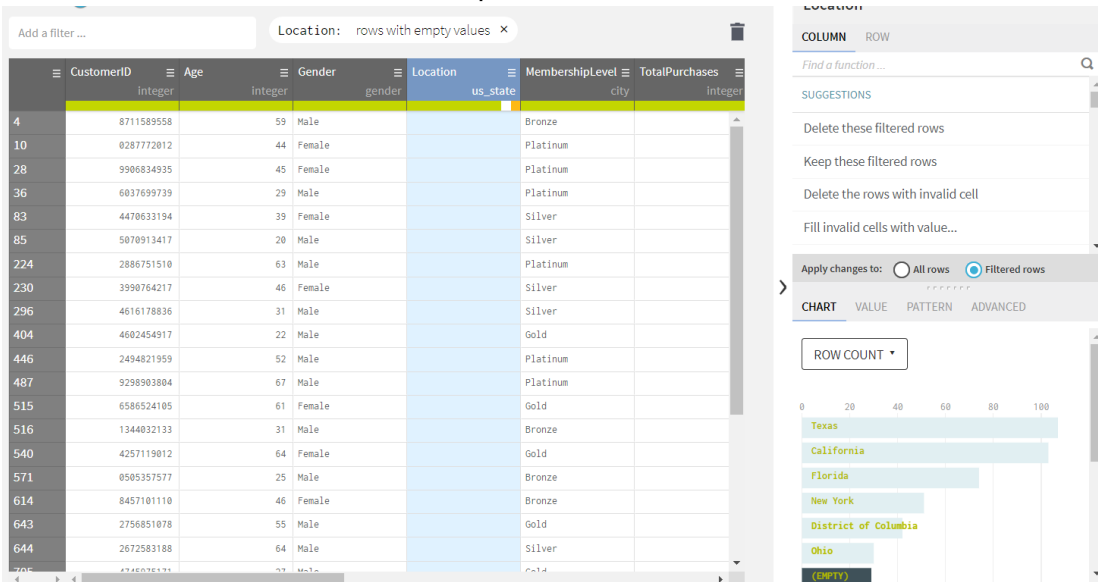


*Figure 10: Handling missing values for Location.*

Referring to Figure 10, there are missing values in location as designated during data generation. For this column, we select the state with the highest occurrence to perform data imputation for missing locations. The way to impute the data is shown in the left pane of Figure 8.
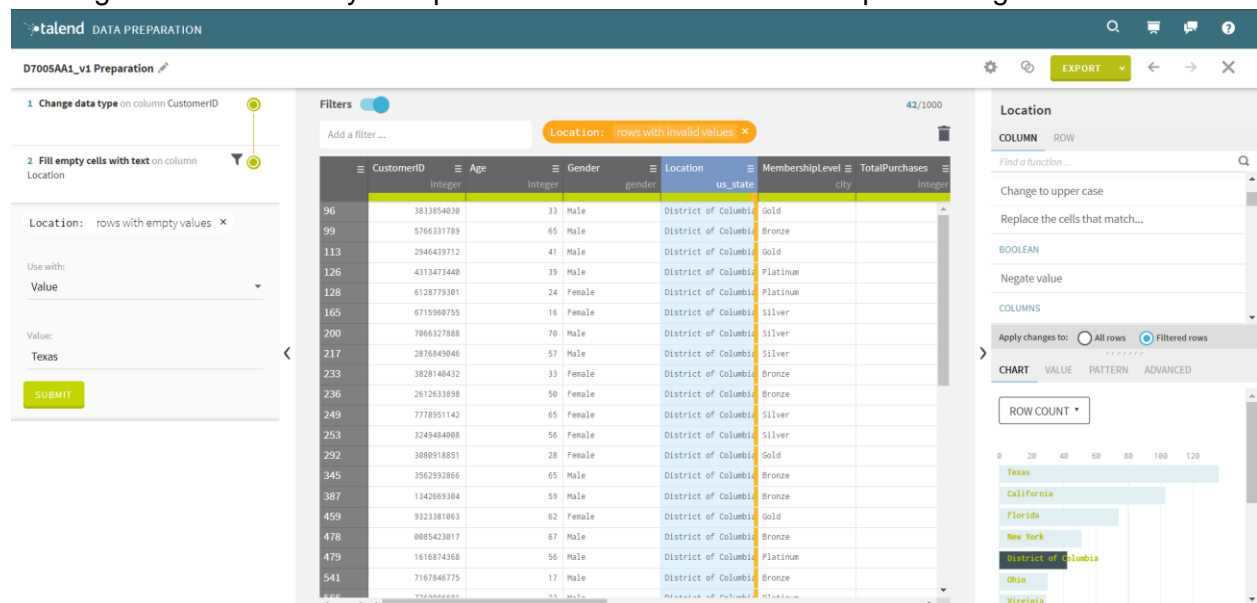


Figure 11: Location imputation and potential invalid location.

Besides missing values, there exist invalid values too. However, these values are the same, which is the District of Columbia, which is not a US state, but a federal district of the United States which not belong to any state. Hence, it is safe to ignore this flag as it is still an area of the United States.
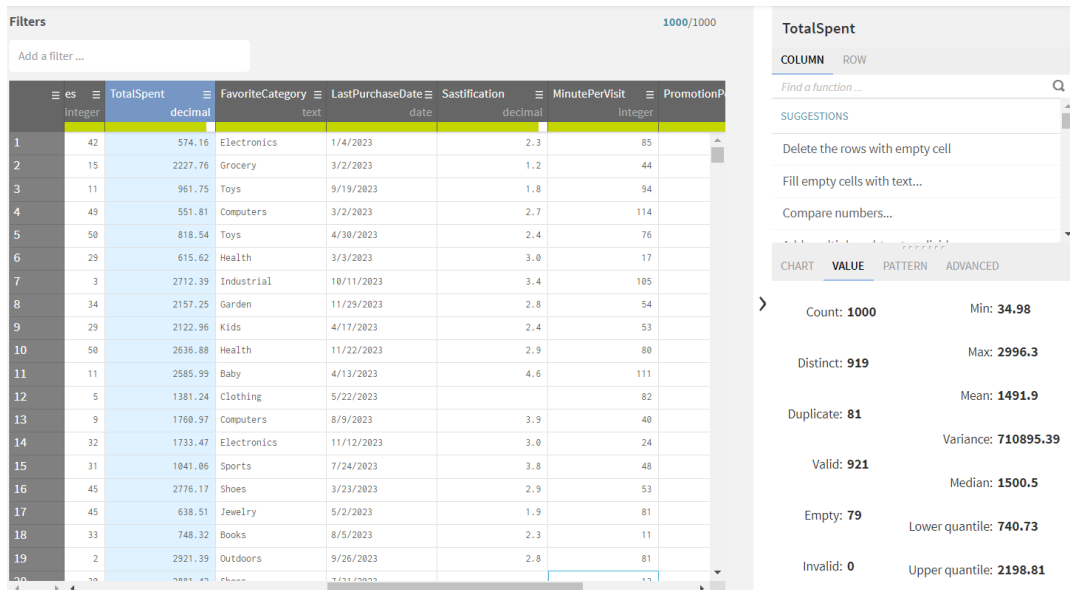
*Figure 12: Deciding the imputation method of missing TotalSpent values.*

For the TotalSpent, since the range of the value is large (34.98 to 2996.3), it is safer to select the median as the imputation of missing values to minimise the effect and bias caused by the potential outliers although the boxplot does not significantly skew to either side.
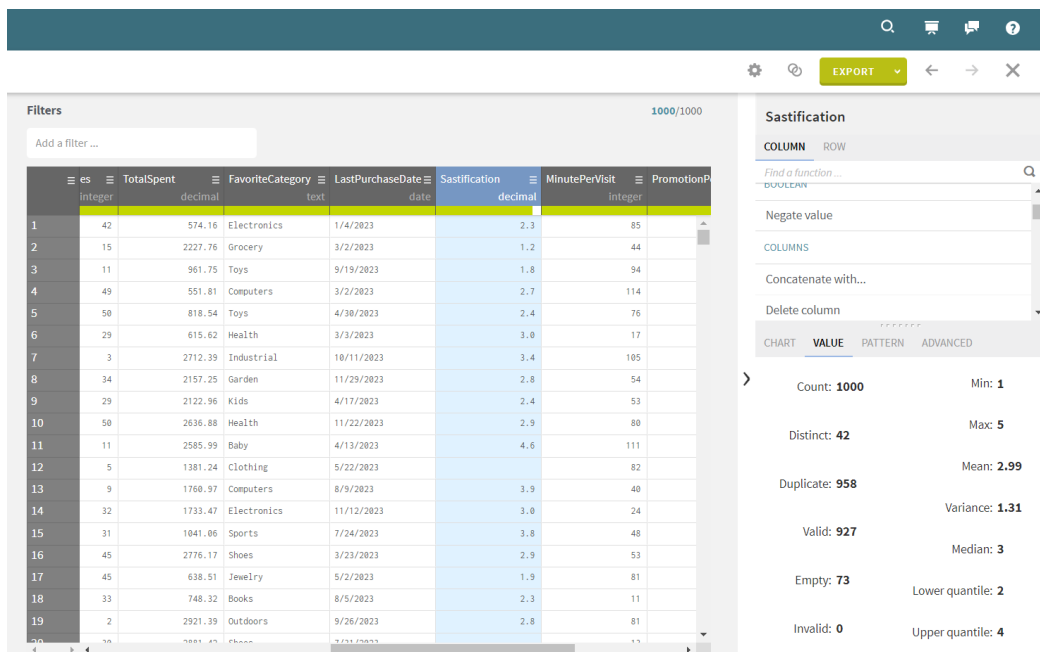


*Figure 13: Deciding the imputation method of missing Satisfaction values.*

Regarding missing values of Satisfaction, the mean and median values are nearly identical which indicates that the distribution of values is possibly evenly distributed and both values are located at the middle of the range. Hence, the mean can be selected to be imputed. However, to ensure

data formalisation which allows only one decimal place for input, the mean value of 2.99 is rounded off to 3.0 for imputation of missing values.
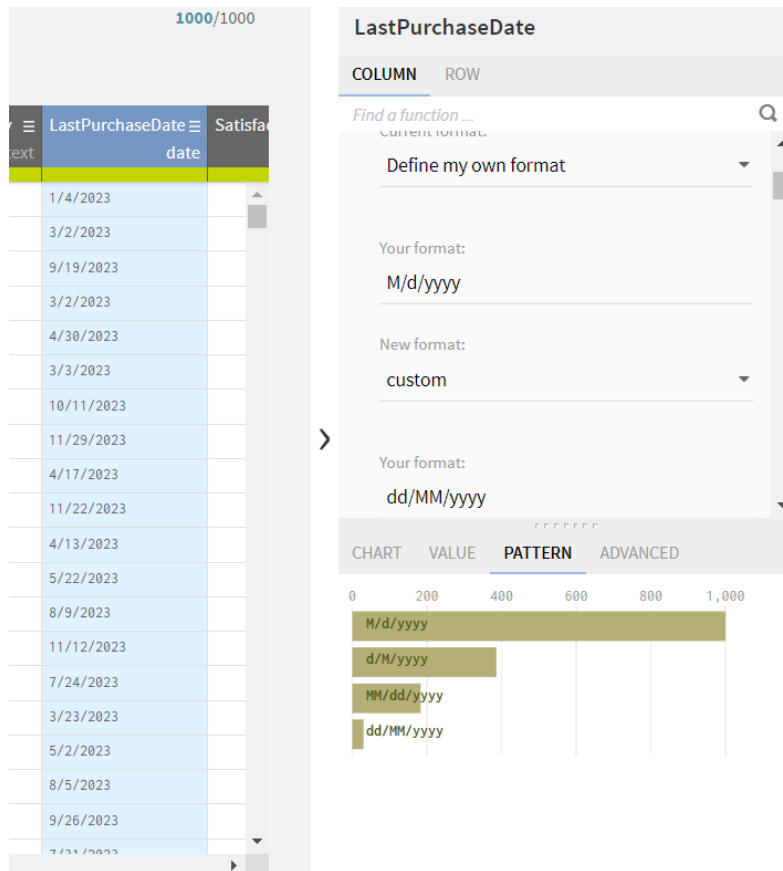


*Figure 14: Transforming data format.*

As this is a US-related dataset, the date format is recorded in American style. To change the date format to a format that is used more frequently in Malaysia, choose the "Change date format" at the right pane of the interface and apply the setting in Figure 14.
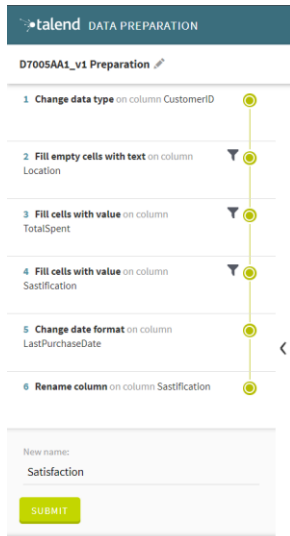
*Figure 15: Data processing workflow on Talend Data Preparation.*

Finally, rename the unsatisfied column name (for example, wrong spelling) during the final check. As the data quality of all columns is satisfied, the processed dataset is exported into CSV format to the local environment as shown in Figure 16.



*Figure 16: Exporting processed dataset.*