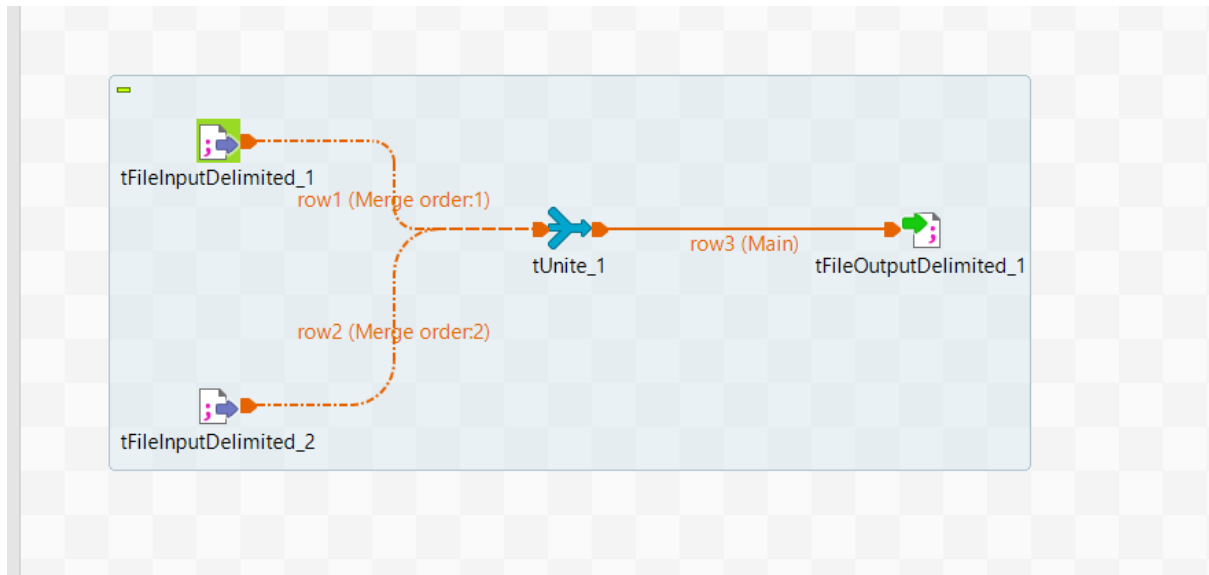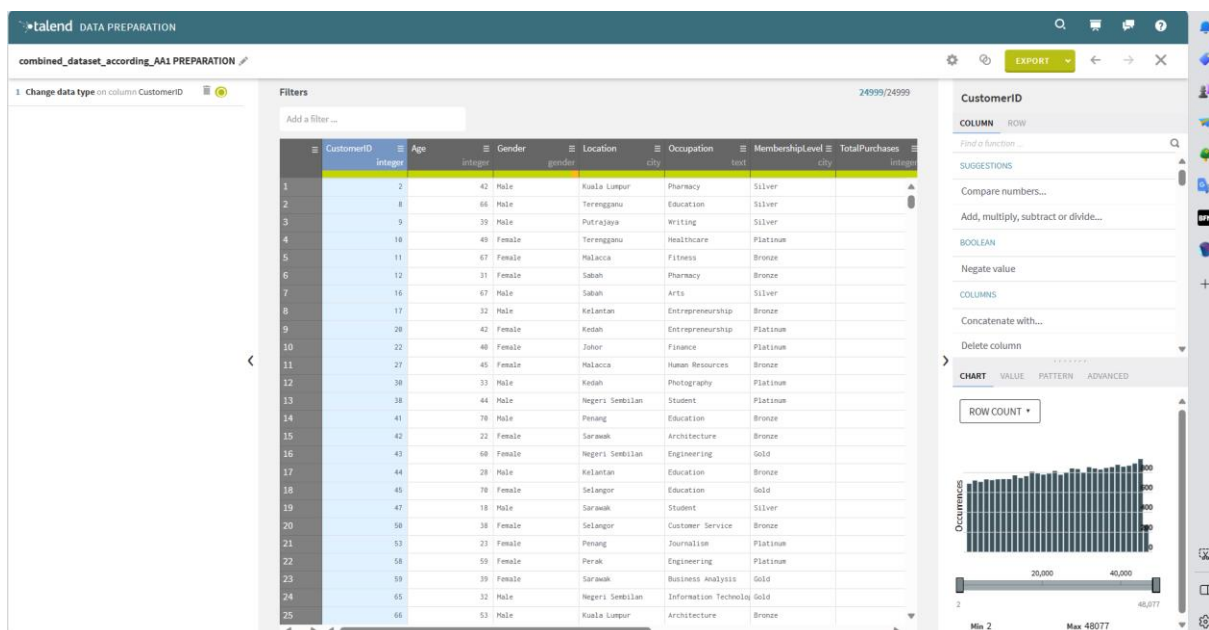# Documentation for Talend Data Integration

- Open Talend Data Integration
- Create new project named wqd7005-AA1
- Create a new job
- Add in the tFileInputDelimited node, tUnite_1 node and tFileOutputDelimited



- Set the fil path for input .csv file and file path for output .csv file
- Set the schema and synchronize it
- Run the job

# Documentation for Talend Data Preparation

- Open Talend Data Preparation
- Add datasets from local device to the Talend Data Preparation



- Check for invalid values and data formatting

- Correcting the Gender using "find and group similar text…"



- Correcting the wrongly labelled Kuala Lumpur at Location attribute using "replace the cells that match…"
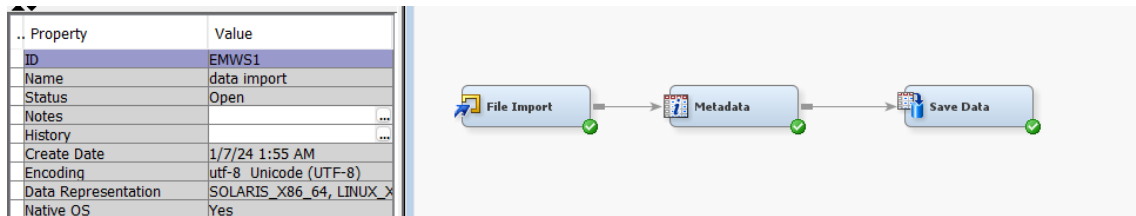


- Check for date format whether it is standardized.



- Missing data is found in both TotalSpent and Churn Attribute. However, it is ignored as the missing value will be handled by SAS Enterprise Miner.

| TotalSpent integer | FavouriteCategory text | LastPurchaseDate date | PaymentMethod text | Churn integer |
|---|---|---|---|---|
| 2162 | Home | 4/25/2023 | Cash | 0 |
| 9181 | Electronics | 5/6/2023 | Credit Card | 0 |

- Export the dataset as basic_clean_dataset_2023.csv

# Documentation for SAS Enteprise Miner

- Open SAS Enterprise Miner 15.2
- Create a new project
- Create a new diagram, add in the "File Import", "Metadata" and "Save Data" node



- Create a new library from the save data
- Create a new data source from the created table



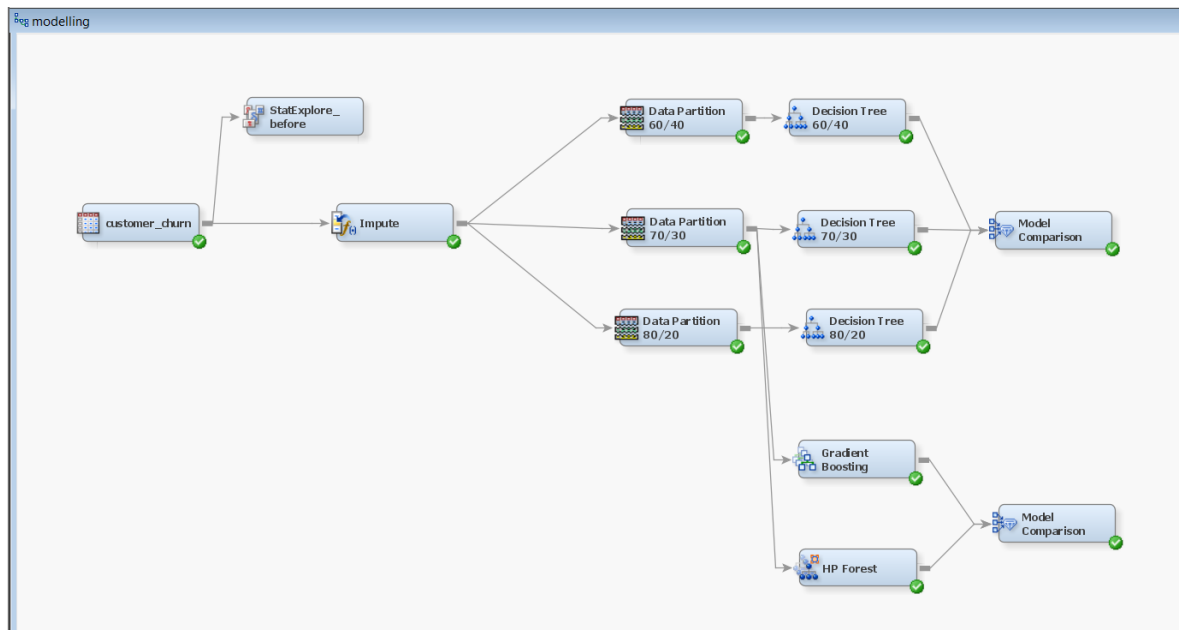- Specifying the variable role and level



- Use StatExplore to check for missing values and "Impute" to fill up the missing values

- Perform variable selection via correlation analysis





- Identify the data splitting ratio
- Perform Decision Tree modelling using Decision Tree node, Bagging modelling using HP Forest node and Boosting modelling using Gradient Boosting node



- Evaluate the results