

Part 2: Data Preparation Using Talend Data Integration

Before loading both datasets (D7005AA1_DataA and D7005AA1_DataB) into SAS Enterprise Miner (SAS EM) for analysis, the preprocessing of data is performed on Talend Data Integration (TDI) and Talend Data Preparation (TDP). First of all, the dataset is loaded into Talend Data Integration using the following schema:

The figure shows two screenshots of the Talend Data Integration schema configuration windows. The top window is titled 'Schema of tFileInputDelimited_3' and the bottom window is titled 'Schema of tFileInputDelimited_1'. Both windows display a table of columns with their respective types and other properties.

Schema of tFileInputDelimited_3

Column	K...	Type	<input checked="" type="checkbox"/> N.	Date Pattern (Ctrl+Space ...)	Length	Precision	Default	Comment
CustomerID	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>					
MembershipLevel	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>					
TotalPurchases	<input type="checkbox"/>	Integer	<input checked="" type="checkbox"/>					
TotalSpent	<input type="checkbox"/>	Double	<input checked="" type="checkbox"/>					
FavoriteCategory	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>					
LastPurchaseDate	<input type="checkbox"/>	Date	<input checked="" type="checkbox"/>	"M/dd/yyyy"				
Sastification	<input type="checkbox"/>	Double	<input checked="" type="checkbox"/>					
MinutePerVisit	<input type="checkbox"/>	Integer	<input checked="" type="checkbox"/>					
PromotionPercent	<input type="checkbox"/>	Integer	<input checked="" type="checkbox"/>					
Comments	<input type="checkbox"/>	Boolean	<input checked="" type="checkbox"/>					
Churn	<input type="checkbox"/>	Integer	<input checked="" type="checkbox"/>					

Schema of tFileInputDelimited_1

Column	K...	Type	<input checked="" type="checkbox"/> N.	Date Pattern (Ctrl+Space av...	Length	Precision	Default	Comment
CustomerID	<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"/>					
Location	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>					

Figure 2: Schema when loading the dataset into Talend Data Integration.

Although the CustomerID is made up of 10 to 11-digit integers, however when the data type of CustomerID is set as Integer, some errors exist as the system identifies some of the numbers as strings. As the CustomerID is only used to check for duplicates, which is not involved much in analysis, the string data type is acceptable for Customer ID. The error is shown in Figure 3.

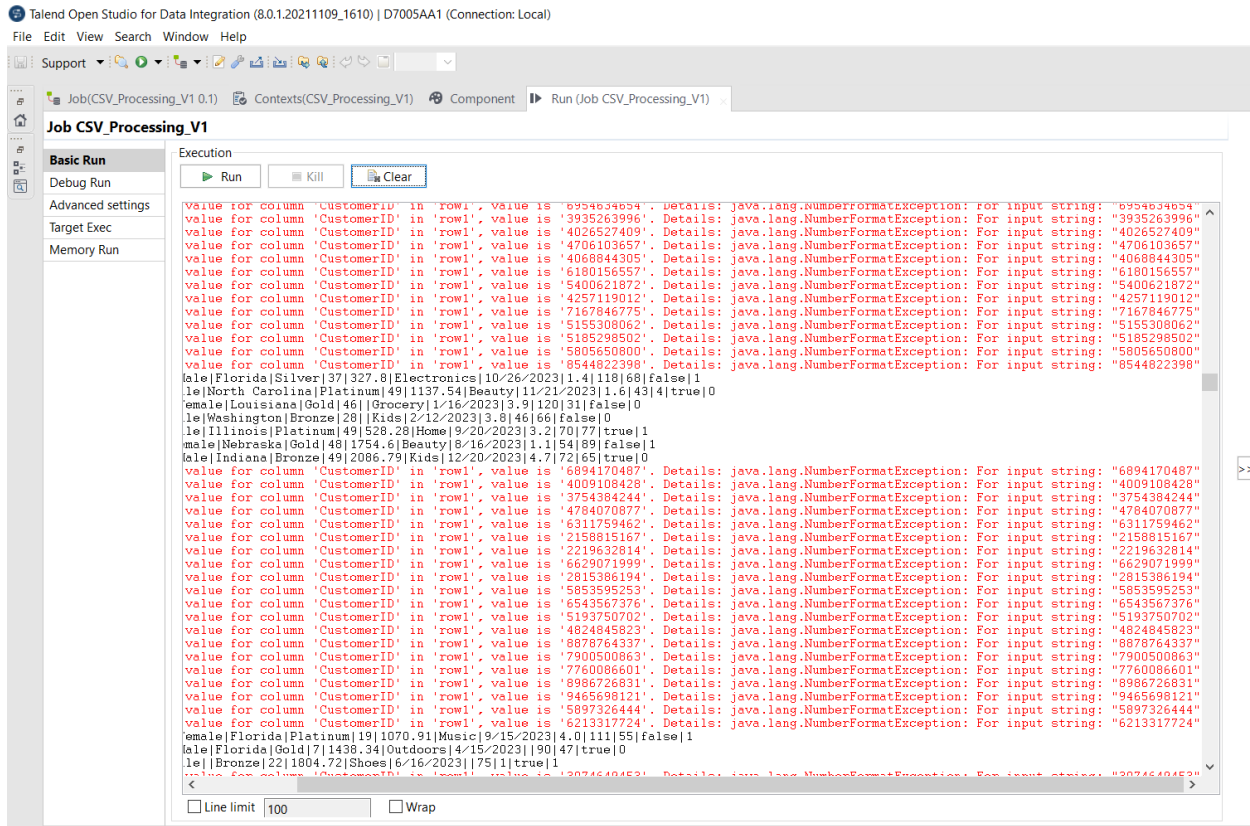


Figure 3: Error when loading the dataset to TDI when CustomerID type is set as integer.

First of all, both datasets are related so they needed to be joined together. The tLogRow component provides an output that enables us to determine whether the data under processing is correct. The tMap is used to join both datasets.

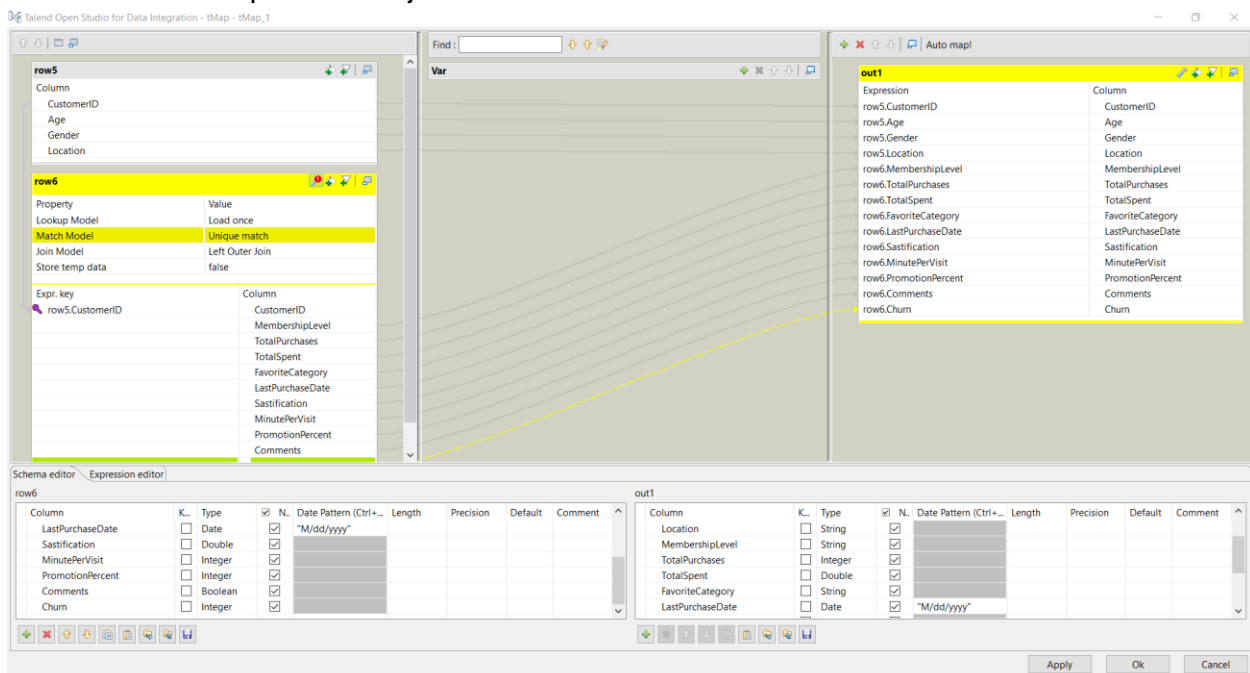


Figure 4: Setting of tMap.

Figure 4 shows both datasets are joined together based on the unique CustomerID match. After the output schema is determined, the input column is linked to the desired output column.

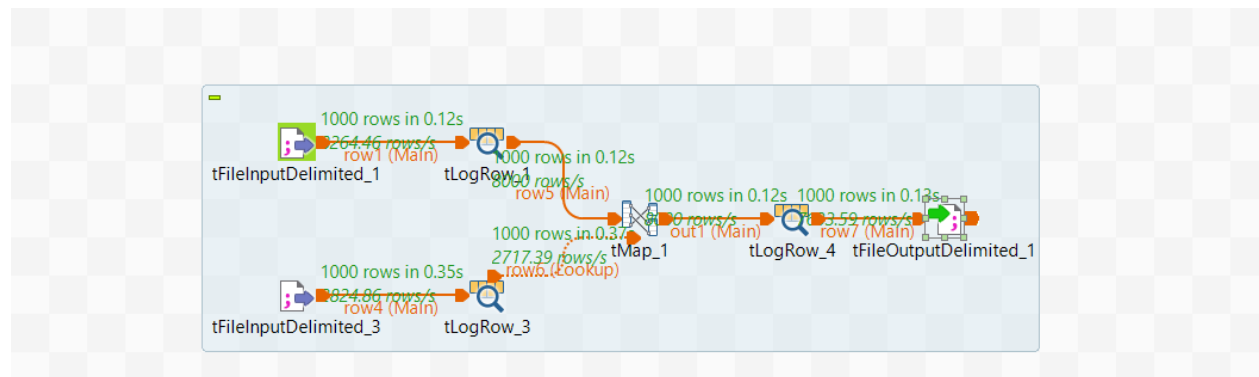


Figure 5: Workflow of producing combined CSV dataset.

As shown in Figure 5, after the datasets are joined, the combined dataset is now ready to be exported for further usage. To ensure smoother loading of data in the further step, the settings as shown in Figure 6 are applied.

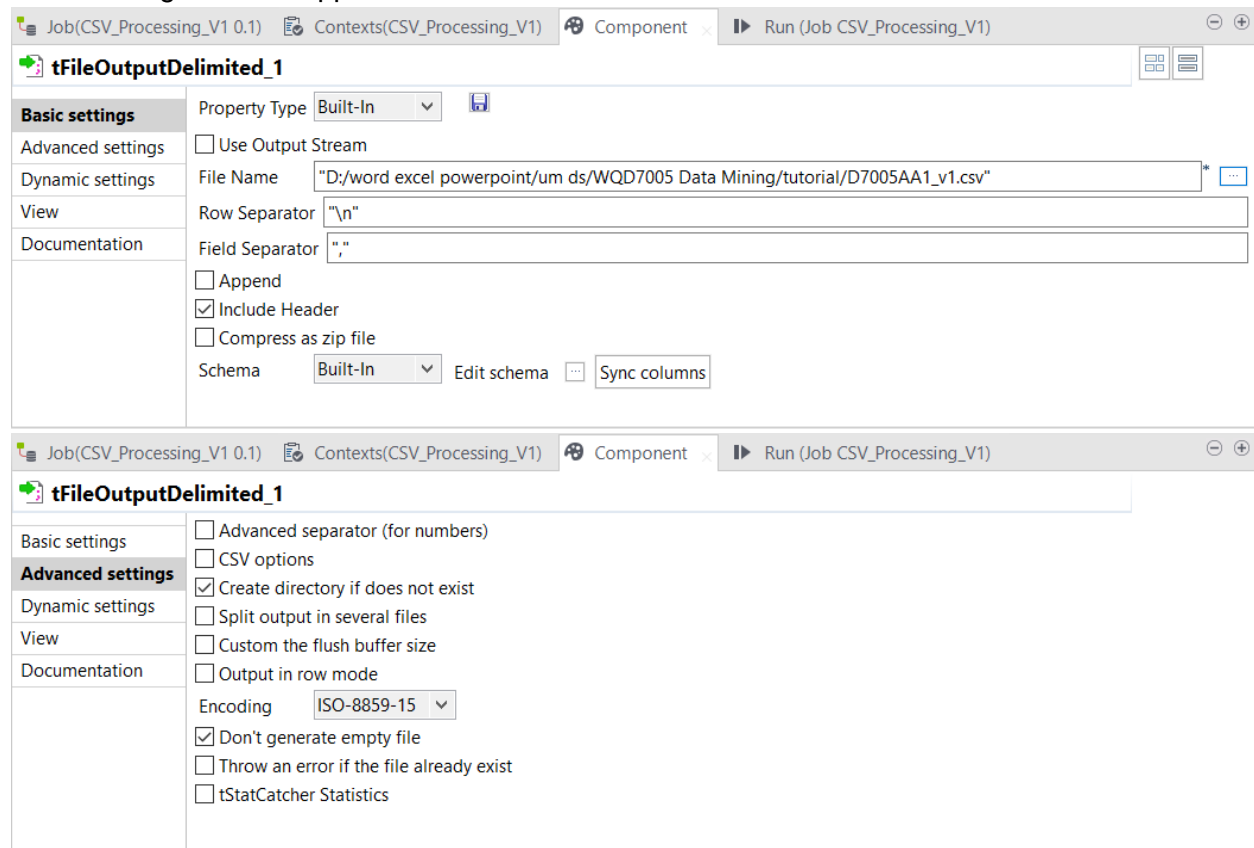


Figure 6: Setting of export of combined CSV dataset.

The newly exported dataset is now able to be used. After loading the combined dataset into TDI, the first attempt of handling missing values is operated, which is dropping all the rows with missing values. By using the tFilterRow component with the setting as shown in Figure 7, the remaining number of rows is shown.

The screenshot shows the Talend Data Integration (TDI) Designer interface. At the top, a workflow diagram is visible with components: tFileInputDelimited_2 (1000 rows in 0.17s), tFilterRow_1 (854 rows in 0.17s), and tLogRow_2. Below the diagram, the 'tFilterRow_1' component is selected, and its configuration is shown in the 'Basic settings' tab. The 'Schema' is set to 'Built-In'. The 'Logical operator used to combine conditions' is set to 'And'. The 'Conditions' table lists 15 columns, all with the function 'Empty' and operator 'Not equal to', and the value 'null'. The 'Basic settings' tab is selected. The 'tFilterRow_1' component is connected to 'tFileInputDelimited_2' and 'tLogRow_2'.

InputColumn	Function	Operator	Value
CustomerID	Empty	Not equal to	null
Age	Empty	Not equal to	null
Gender	Empty	Not equal to	null
Location	Empty	Not equal to	null
MembershipLevel	Empty	Not equal to	null
TotalPurchases	Empty	Not equal to	null
TotalSpent	Empty	Not equal to	null
FavoriteCategory	Empty	Not equal to	null
LastPurchaseDate	Empty	Not equal to	null
Sastification	Empty	Not equal to	null
MinutePerVisit	Empty	Not equal to	null
PromotionPercent	Empty	Not equal to	null
Comments	Empty	Not equal to	null
Churn	Empty	Not equal to	null

Figure 7: Filtering out all rows with missing values.

As shown in the figure above, only 854 rows are remaining in the dataset, which means 146 records (14.6%) are dropped. As the proportion of rows with missing values is not small enough to be removed, data imputation seems to be the better option to handle missing values.

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