# Talend Data Preparation

The original Online_Sales.csv dataset has 52924 rows while Talend Data Preparation can only access 30000 rows at most, so I split the dataset into two files, one with 30000 rows and the other with 22924 rows, as shown below:

To make the data type more accurate and convenient for subsequent processing, I change the data type of CustomerID, Transaction_ID into integer, and Delivery_Charges into decimal.

In addition, the format of the Transaction_Date is not uniform, some shows 'M/d/yyyy', the others show 'd/M/yyyy', so I change the date format into 'yyyy-MM-dd' for subsequent processing:

For subsequent analysis, I calculate the 'Total_Amount' spent on each transaction using 'Quantity' mutyplies 'Avg_Price', and rename the new column as 'Total_Amount'.
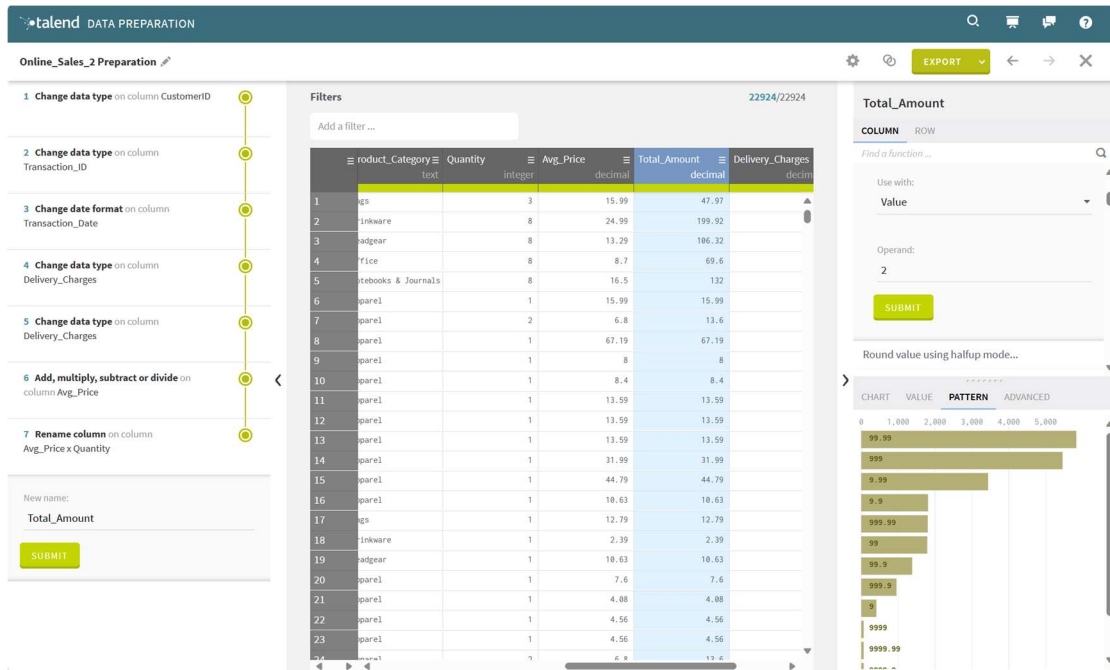
After preprocessing, the two splitted sets of Online_Sales.csv file are cleaned and then combined as Online_Sales_preprocessed.xlsx.

The original CustomerData.xlsx file is as follows, there are no missing value or invalid value at all, and it need no more preprocessing: