**Talend Data Integration**:Merge two data sets into one.

Step:

Open Talend Data Integration Studio:

Launch the Talend Data Integration tool and open your project.

Create a New Job:

Right-click in the Repository panel on the left under Job Designs.

Choose Create job.

Provide a name for the job and a description, if desired. Click Finish.
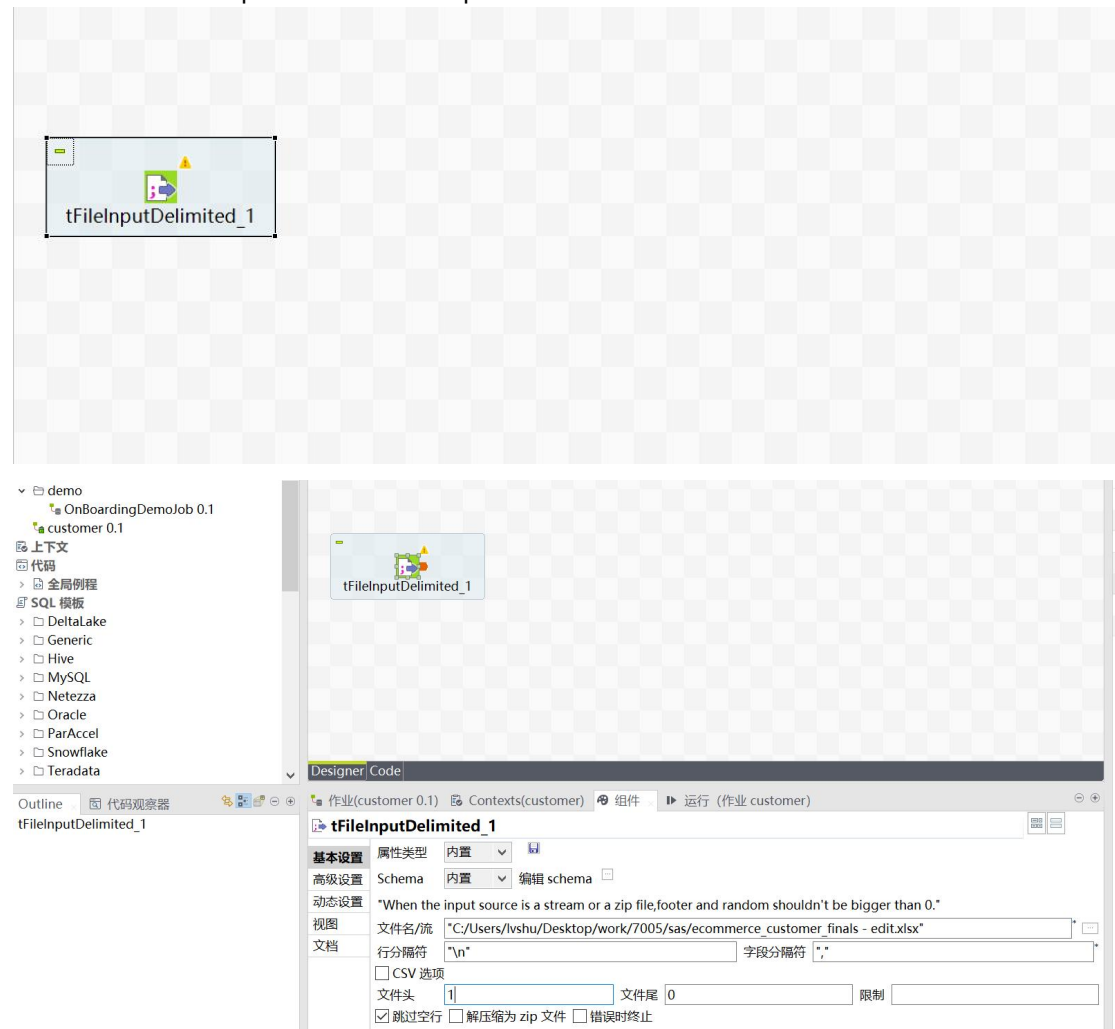
Add a tFileInputDelimited Component:

This component is used to read delimited files like CSV.

From the Palette panel on the right, type "tFileInputDelimited" into the search bar.

Drag the tFileInputDelimited component to the design workspace.

Configure the tFileInputDelimited Component:

Click on the tFileInputDelimited component to select it.

文件设置

服务器  Localhost 127.0.0.1

文件  C:/Users/lvshu/Desktop/work//005/sas/dataSet/ecommerce_customer_edit01.xlsx  浏览...

☑ 读取 excel2007 文件格式 (xlsx)

生成模式  Memory-consuming(User mode)

文件查看器和表单设置

设置表单参数

☑ All sheets/☐Select sheet
　☑ sheet1

填选择工作表 (工作表结构作为 schema 指南)  sheet1 ▼

| A | B | C | D | E | F | G | H | I | J |
|---|---|---|---|---|---|---|---|---|---|
| Cust... | Last... | Fav... | Tot... | Tot... | City | Pay... | Cus... | Cus... | Chu... |
| 2.01... | 202... | Ho... | 177... | 1.0 | Ne... | Pay... | 31.0 | Joh... | 0.0 |
| 2.01... | 202... | Elec... | 174... | 3.0 | Los... | Pay... | 31.0 | Joh... | 0.0 |
| 2.01... | 202... | Boo... | 413... | 1.0 | Chi... | Cre... | 31.0 | Joh... | 0.0 |
| 2.01... | 202... | Elec... | 396... | 3.0 | San... | Cash | 31.0 | Joh... | 0.0 |
| 2.01... | 202... | Boo... | 259... | 4.0 | Mia... | Pay... | 31.0 | Joh... | 0.0 |
| 2.01... | 202... | Ho... | 191... | 3.0 | Ho... | Cre... | 27.0 | Lau... | 0.0 |
| 2.01... | 202... | Elec... | 205... | 1.0 | Ne... | Cre... | 27.0 | Lau... | 0.0 |
| 2.01... | 202... | Boo... | 370... | 5.0 | Los... | Cash | 27.0 | Lau... | 0.0 |
| 2.01... | 202... | Ho... | 13.0 | 2.0 | Chi... | Cash | | Lau... | 0.0 |

< Back　Next >　Finish　Cancel

---

文件锐道

编码  UTF-8

☐ 高级分隔符 (用于数字)

千位分隔符:

小数分隔符:

元数据列设置

第一列:  1

最后一列:

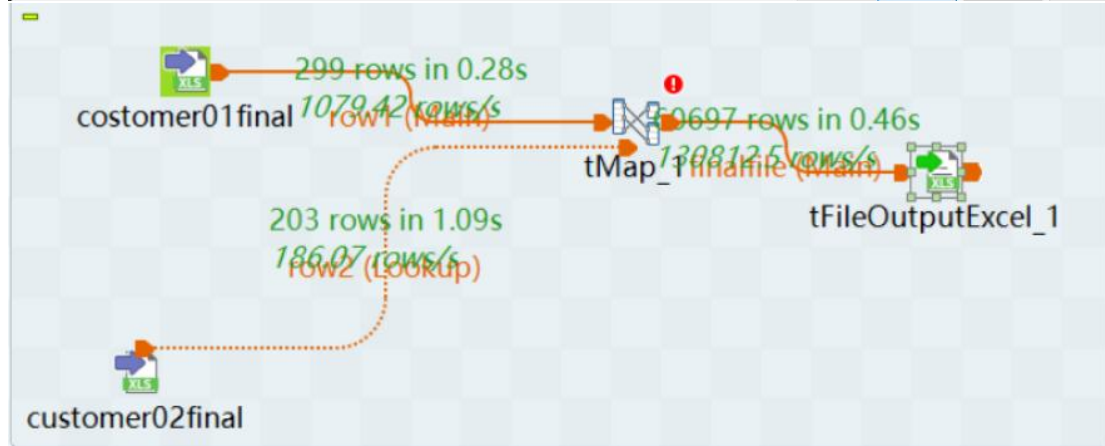要减过的行数
如果必须忽略任何行，请指定以下参数

文件头 ☐

文件足 ☐

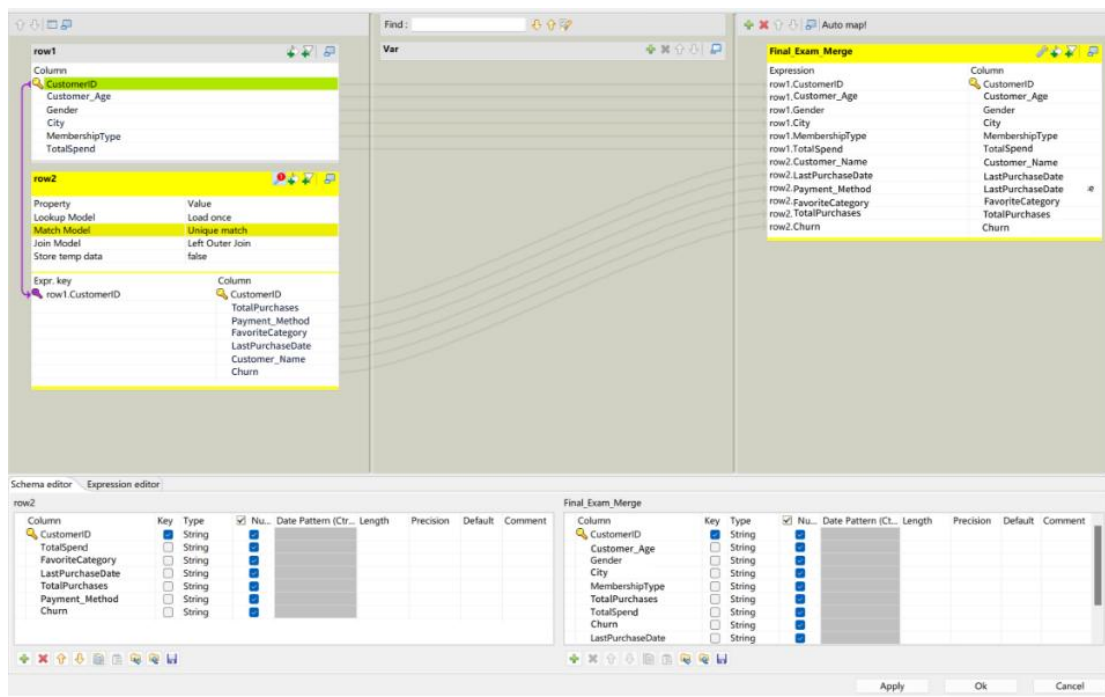行数限制
如果必须限制行数，则指定此数字。

限制 ☐

预览  输出

☐ 将标题行设为列名  刷新预览

| A | B | C | D | E | F | G | H | I | J |
|---|---|---|---|---|---|---|---|---|---|
| Customer ID | LastPurchaseDate | FavoriteCategory | TotalSpent | TotalPurchases | City | Payment Method | Customer Age | Customer Name | Churn |
| 20140001 | 2023/5/3 21:30 | Home | 177 | 1 | New York | PayPal | 31 | John Rivera | 0 |
| 20140002 | 2021/5/16 13:57 | Electronics | 174 | 3 | Los Angeles | PayPal | 31 | John Rivera | 0 |
| 20140003 | 2020/7/13 6:16 | Books | 413 | 1 | Chicago | Credit Card | 31 | John Rivera | 0 |
| 20140004 | 2023/1/17 13:14 | Electronics | 396 | 3 | San Francisco | Cash | 31 | John Rivera | 0 |

导出为上下文　恢复上下文

---

**Talend Data Prep**: Remove missing values

When clicking on each column you can see the number of missing values in that column.

The city names with the highest frequency of occurrence usually represent the major cities in the dataset, and populating such city names helps to maintain the consistency of the overall data distribution. For missing addresses, this experiment uses taking the city name with the highest frequency of occurrence to fill it.



As shown, New York is the most frequent. So fill in the four missing city names as New York.

For critical attributes such as missing churn value, where there are few missing values, the deletion of missing values is used to deal with it.Churn value is usually a critical attribute in user churn prediction as it directly reflects whether a user is churned or not. In this case, it is important to ensure the accuracy of this attribute as it is the target variable for model training. Removing missing values avoids introducing uncertainty about the accuracy of user churn prediction during the modeling process.





**SAS e-Miner**:

decision trees：

1. To create a project: Select the "File" menu and then select "New Project". Name the project and set the properties of the project.

2. Create a flowchart: In the project, select the "Diagram" menu, and then select "Create Diagram".

3. Import data: In the project, import the dataset that contains the data you want to analyze.

4. Set the ratio of test set to validation set

5. Configure the decision tree nodes.

6. Run the entire flowchart. Select the "Run" menu and then select "Run Diagram".

Cluser: