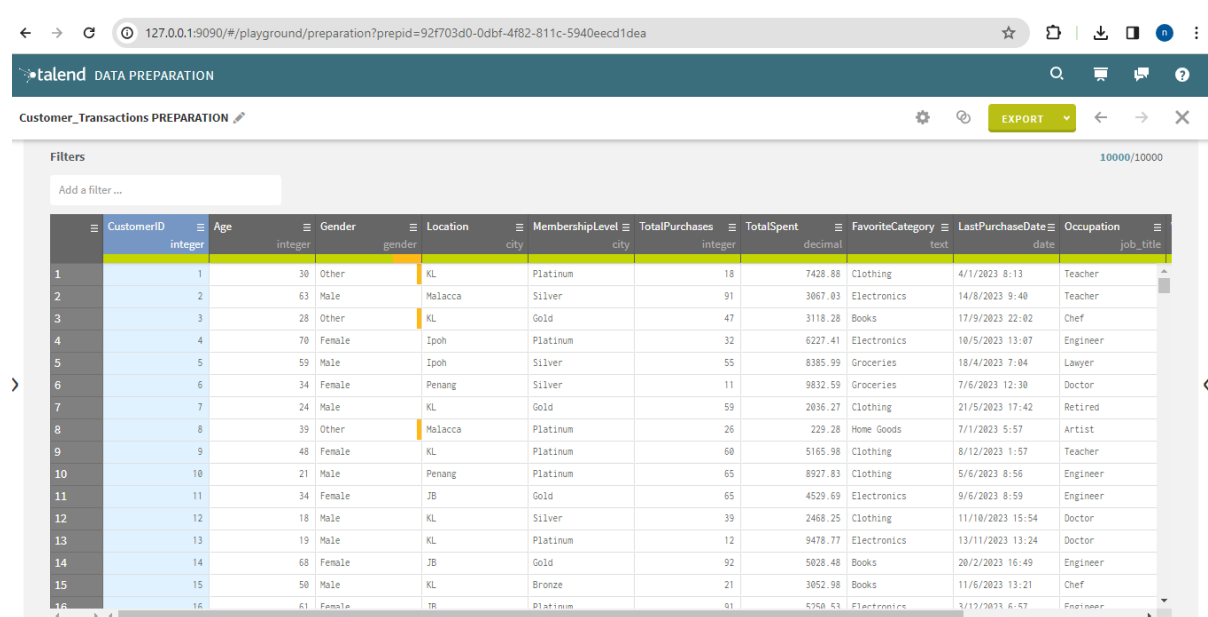**Preprocessing using Talend Data Preparation**

For the next preprocessing, I utilized Talend Data Preparation platform. I uploaded the integrated data from Talend Data Integration, which was produced in CSV delimiter format generated by the tFileOutputDelimited component. The figure below shows the uploaded CSV file in the talend data Preparation.
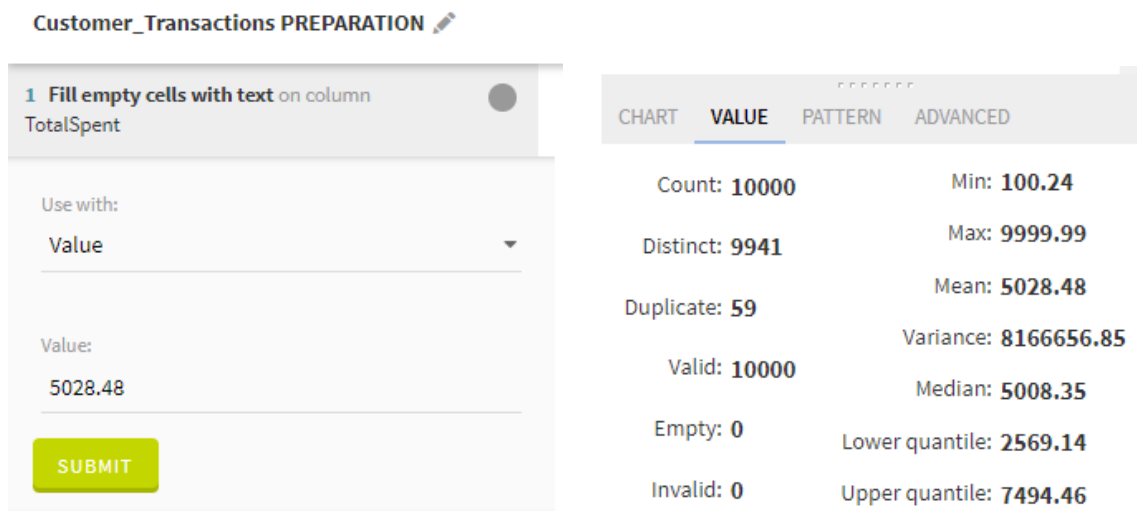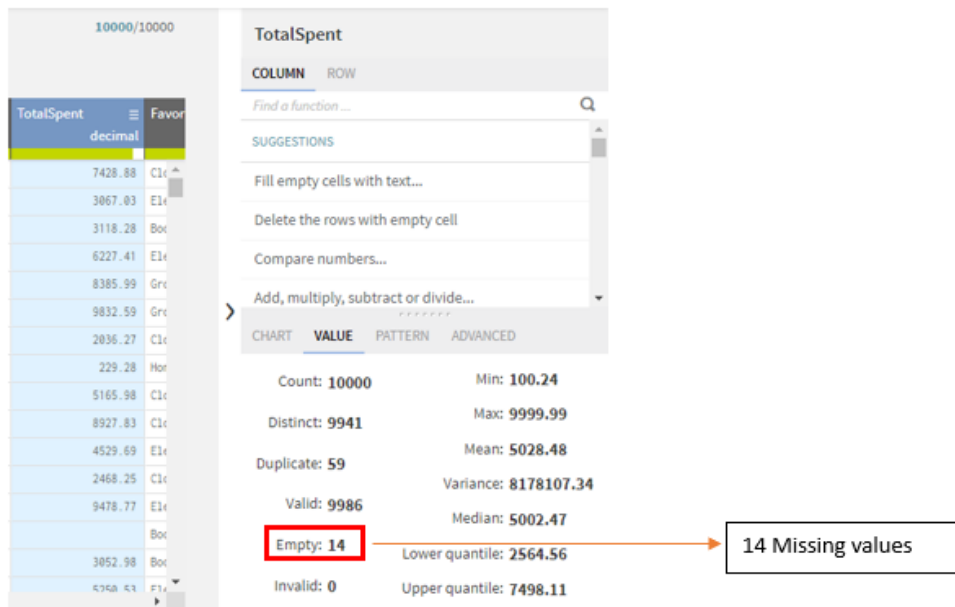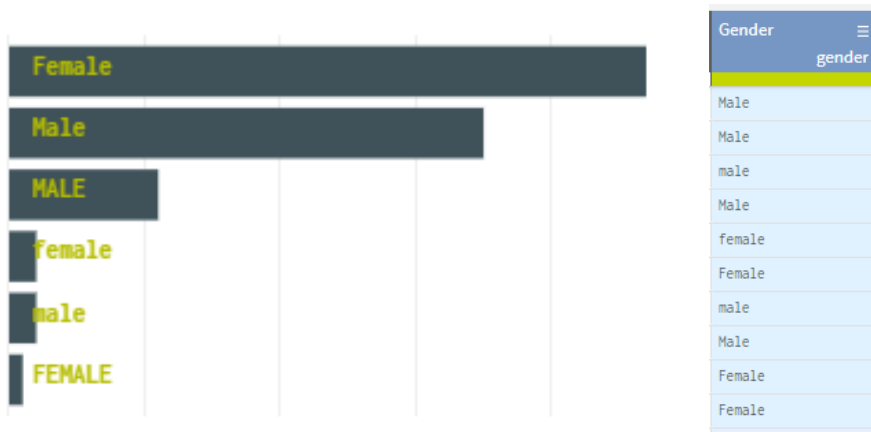


**Data Cleaning**

The data cleaning process involves individually clicking on each column. During this, I noticed that the 'TotalSpend' column contains 14 missing values.

**TotalSpent**

COLUMN    ROW

Find a function ...

SUGGESTIONS

Fill empty cells with text...

Delete the rows with empty cell

Compare numbers...

Add, multiply, subtract or divide...

CHART    **VALUE**    PATTERN    ADVANCED

| | |
|---|---|
| Count: 10000 | Min: 100.24 |
| Distinct: 9941 | Max: 9999.99 |
| | Mean: 5028.48 |
| Duplicate: 59 | Variance: 8178107.34 |
| Valid: 9986 | Median: 5002.47 |
| Empty: 14 | Lower quantile: 2564.56 |
| Invalid: 0 | Upper quantile: 7498.11 |

Empty: 14 → 14 Missing values

TotalSpent / decimal column values: 7428.88, 3067.03, 3118.28, 6227.41, 8385.99, 9832.59, 2036.27, 229.28, 5165.98, 8927.83, 4529.69, 2468.25, 9478.77, 3052.98, 5250.53

**Customer_Transactions PREPARATION** ✎

**1  Fill empty cells with text** on column
TotalSpent

Use with:

Value ▾

Value:

5028.48

SUBMIT

CHART    **VALUE**    PATTERN    ADVANCED

| | |
|---|---|
| Count: 10000 | Min: 100.24 |
| Distinct: 9941 | Max: 9999.99 |
| Duplicate: 59 | Mean: 5028.48 |
| | Variance: 8166656.85 |
| Valid: 10000 | Median: 5008.35 |
| Empty: 0 | Lower quantile: 2569.14 |
| Invalid: 0 | Upper quantile: 7494.46 |

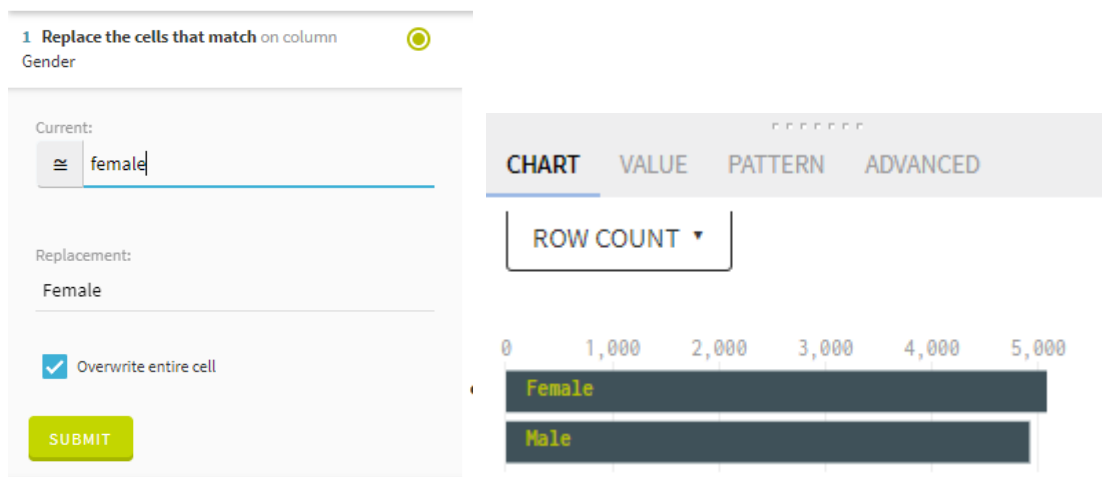For the "TotalSpend" column, I identified 14 missing values. These were replaced with the mean value of 5028.48 to ensure that the column is free from any empty entries.

**Data Standardization**

The data standardization process involves individually clicking on each column. During this, I noticed that the 'Gender' column contains varying letter cases, and the 'Date' column also lacks standardization.

The figure above shows the gender contains varying letter cases.



The figure above shows after standardized the gender.

After identifying the date columns, a function within Talend was selected to standardize the dates. The ISO 8601 format was specified within this function to transform all date entries to this uniform standard. Following the transformation, the data was reviewed to ensure that the dates were correctly standardized, with any anomalies or errors requiring further attention being identified and possibly corrected.
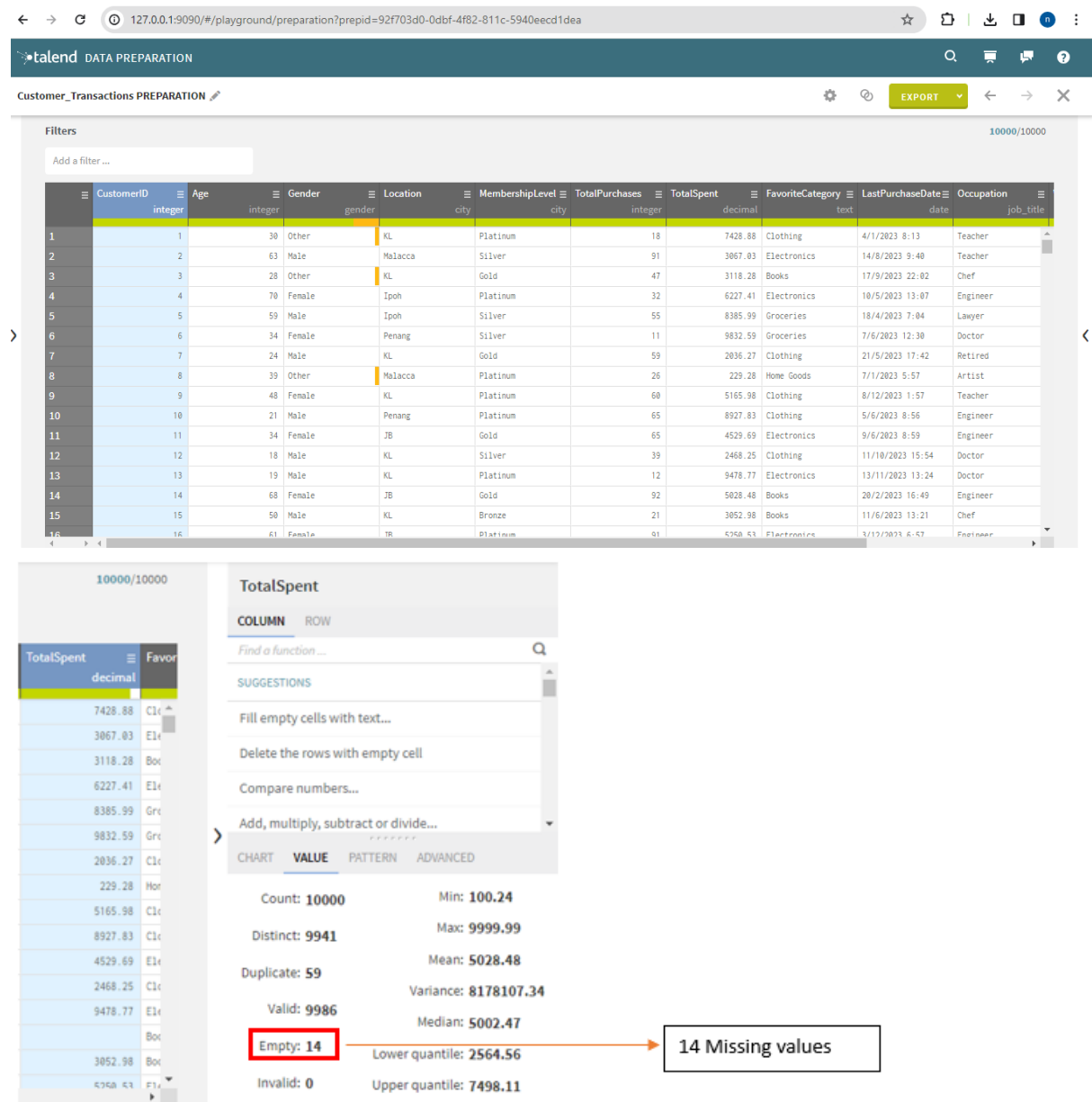


The final step involves exporting the file from Talend Data Preparation in CSV format, ensuring that the delimiter is set to a comma.

# APPENDIX

**Customer_Transactions PREPARATION** ✏

**1 Fill empty cells with text** on column
TotalSpent ⚪

Use with:

Value ▼

Value:

5028.48

SUBMIT

CHART **VALUE** PATTERN ADVANCED

| | |
|---|---|
| Count: **10000** | Min: **100.24** |
| Distinct: **9941** | Max: **9999.99** |
| Duplicate: **59** | Mean: **5028.48** |
| | Variance: **8166656.85** |
| Valid: **10000** | Median: **5008.35** |
| Empty: **0** | Lower quantile: **2569.14** |
| Invalid: **0** | Upper quantile: **7494.46** |



| Gender | ≡ |
|---|---|
| | gender |
| Male | |
| Male | |
| male | |
| Male | |
| female | |
| Female | |
| male | |
| Male | |
| Female | |
| Female | |

**1 Replace the cells that match** on column
Gender ◎

Current:

≅   female

Replacement:

Female

☑ Overwrite entire cell

SUBMIT

CHART   VALUE   PATTERN   ADVANCED

ROW COUNT ▼

| 0 | 1,000 | 2,000 | 3,000 | 4,000 | 5,000 |
|---|---|---|---|---|---|

Female

Male

| LastPurchaseDate | Occupation | WebsiteVisitFreq... | Churn |
| --- | --- | --- | --- |
| date | job_title | text | |
| 2023-11-24 | Teacher | Rarely | |
| 2023-05-30 | Engineer | Frequently | |
| 2023-10-22 | Teacher | Regularly | |
| 2023-10-22 | Teacher | Regularly | |
| 2023-06-13 | Engineer | Occasionally | |
| 2023-01-16 | Doctor | Rarely | |
| 2023-04-11 | Teacher | Regularly | |
| 2023-07-18 | Engineer | Rarely | |
| 2023-07-08 | Doctor | Regularly | |
| 2023-06-07 | Student | Rarely | |
| 2023-12-23 | Doctor | Rarely | |
| 2023-10-22 | Student | Occasionally | |
| 2023-11-19 | Engineer | Occasionally | |
| 2023-04-28 | Teacher | Frequently | |
| 2023-07-05 | Doctor | Rarely | |
| 2023-05-20 | Artist | Frequently | |

### LastPurchaseDate

**COLUMN**    ROW

Find a function ...

New format:

ISO 8601 date ▼

**SUBMIT**    Learn more ...

## EXPORT TO CSV                                                    ✕

Delimiter:

Comma                                                               ▼

Filename:

Customer_Transactions PREPARATION

CANCEL    EXPORT