**Talend Preparation**

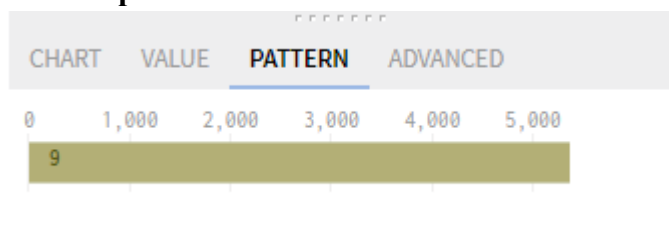**Dataset1: Customer**

**Before**

| | CustomerID fr_postal_code | Aging integer | Gender gender | Location(CityTier) integer | Churn integer | PreferedOrderCat text | MembershipLevel last_name | MaritalStatus text | Complain integer | SatisfactionScore integer |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 50001 | 41 | Female | 3 | 1 | Laptop & Accessory | Medium | Married | 1 | 2 |
| 2 | 50002 | 35 | Male | 1 | 1 | Mobile | Medium | Single | 1 | 3 |
| 3 | 50003 | 41 | Male | 1 | 1 | Mobile | Critical | Single | 1 | 3 |
| 4 | 50004 | 40 | Male | 3 | 1 | Laptop & Accessory | High | Single | 0 | 5 |
| 5 | 50005 | 42 | Male | 1 | 1 | Mobile | Critical | Single | 0 | 5 |
| 6 | 50006 | 41 | Female | 1 | 1 | Mobile Phone | Critical | Single | 1 | 5 |
| 7 | 50007 | 34 | Male | 3 | 1 | Laptop & Accessory | High | Divorced | 0 | 2 |
| 8 | 50008 | 40 | Male | 1 | 1 | Mobile | Critical | Divorced | 1 | 2 |
| 9 | 50009 | 40 | Male | 3 | 1 | Mobile | Critical | Divorced | 1 | 3 |
| 10 | 50010 | 43 | Male | 1 | 1 | Mobile | Critical | Single | 0 | 3 |
| 11 | 50011 | 43 | Female | 1 | 1 | Others | Critical | Divorced | 0 | 3 |
| 12 | 50012 | 34 | Male | 1 | 1 | Fashion | High | Single | 1 | 3 |
| 13 | 50013 | 40 | Male | 1 | 1 | Mobile | Critical | Single | 1 | 3 |

- **Operation1:**
  Number of missing values for each column:
  Gender:2

- **Operation 2:** Check "Parten" to make sure each column is formatted uniformly



- **Operation 3:** Modified the "MembershipLevel" column, adjusting the original levels "Medium, High, Critical" to numbers "1, 2, 3" to facilitate subsequent processing.



- **Operation 4:** Modify the data type of "Customer ID"

## Dataset2: Order

**Before**



| | CustomerID fr_postal_code | OrderCount integer | LastPurchaseDate date | Sales en_money_amount | DaySinceLastOr... integer |
|---|---|---|---|---|---|
| 1 | 50001 | 1 | 9/15/2011 | $140.00 | 5 |
| 2 | 50002 | 1 | 6/30/15 | $211.00 | 0 |
| 3 | 50003 | 1 | 5/15/2012 | $117.00 | 3 |
| 4 | 50004 | 1 | 9/15/2005 | $118.00 | 3 |
| 5 | 50005 | 1 | 9/15/2007 | $250.00 | 3 |
| 6 | 50006 | 6 | 2/25/15 | $72.00 | 7 |
| 7 | 50007 | 1 | 9/15/2004 | $54.00 | 0 |
| 8 | 50008 | 2 | 3/30/15 | $114.00 | 0 |
| 9 | 50009 | 1 | 9/15/2002 | $231.00 | 2 |
| 10 | 50010 | 1 | 4/21/15 | $140.00 | 1 |

- **Operation1:**

  Number of missing values for each column:

  OrderCount:245

  DaySinceLastOrder:282

- **Operation2:** Unified date format



- **Operation3:** Format the "Sales" column



- **Operation4:** Modify the data type of "Customer ID"



| | CustomerID integer | OrderCount integer | LastPurchaseDate date | Sales en_money_amount | DaySinceLastOr... integer |
|---|---|---|---|---|---|
| 1 | 50001 | 1 | 2011-09-15 | $140.00 | 5 |
| 2 | 50002 | 1 | 2015-06-30 | $211.00 | 0 |
| 3 | 50003 | 1 | 2012-05-15 | $117.00 | 3 |
| 4 | 50004 | 1 | 2005-09-15 | $118.00 | 3 |
| 5 | 50005 | 1 | 2007-09-15 | $250.00 | 3 |
| 6 | 50006 | 6 | 2015-02-25 | $72.00 | 7 |
| 7 | 50007 | 1 | 2004-09-15 | $54.00 | 0 |
| 8 | 50008 | 2 | 2015-03-30 | $114.00 | 0 |
| 9 | 50009 | 1 | 2002-09-15 | $231.00 | 2 |

- **Operation5:** Remove extra symbols for "sales"

| MembershipLevel | MaritalStatus | Complain | SatisfactionScore | OrderCount | LastPurchaseDate | Sales | DaySinceLastOr... |
|---|---|---|---|---|---|---|---|
| integer | text | integer | integer | integer | date | decimal | integer |
| 1 | Married | 1 | 2 | 1 | 15-Sep-2011 | 140.00 | 5 |
| 1 | Single | 1 | 3 | 1 | 30-Jun-2015 | 211.00 | 0 |
| 3 | Single | 1 | 3 | 1 | 15-May-2012 | 117.00 | 3 |
| 2 | Single | 0 | 5 | 1 | 15-Sep-2005 | 118.00 | 3 |
| 3 | Single | 0 | 5 | 1 | 15-Sep-2007 | 250.00 | 3 |
| 3 | Single | 1 | 5 | 6 | 25-Feb-2015 | 72.00 | 7 |
| 2 | Divorced | 0 | 2 | 1 | 15-Sep-2004 | 54.00 | 0 |
| 3 | Divorced | 1 | 2 | 2 | 30-Mar-2015 | 114.00 | 0 |
| 3 | Divorced | 1 | 3 | 1 | 15-Sep-2002 | 231.00 | 2 |
| 3 | Single | 0 | 3 | 1 | 21-Apr-2015 | 140.00 | 1 |

| MembershipLevel | MaritalStatus | Complain | SatisfactionScore | OrderCount | LastPurchaseDate | Sales | DaySinceLastOr... |
|---|---|---|---|---|---|---|---|
| integer | text | integer | integer | integer | date | decimal | integer |
| 1 | Married | 1 | 2 | 1 | 15-Sep-2011 | 140.00 | 5 |