**Informatica Data Quality Project Report**

---

**Project Overview**

This project focuses on implementing data quality rules using Informatica Data Quality (IDQ) to cleanse and standardize customer data. The processed data is subsequently exported into Informatica PowerCenter for further integration. The project ensures data accuracy, completeness, and consistency by applying specific rules and validations.

---

**Project Objectives**

1. **Ensure Customer ID Consistency:**

   o Validate that Cust_id contains only digits and has a fixed length of 8.

   o Null values in Cust_id are rejected and moved to the rejection section.

2. **Name Splitting:**

   o Concatenate FNAME and LNAME attributes and split them into separate fields.

3. **Customer Tier Standardization:**

   o Replace invalid values in the cust_tier attribute with predefined valid values.

4. **State Code Standardization:**

   o Ensure that the state field contains exactly 2 characters (e.g., "NW").

5. **ZIP Code Validation:**

   o Ensure zip contains exactly 5 numeric digits.

6. **City, State, and Country Splitting:**

   o Separate concatenated city, state, and country fields into individual columns.

7. **Record Rejection Criteria:**

   o Reject records with null Cust_id.

   o Reject records with null city, state, and country fields.

8. **Review Section Handling:**

   o If any attributes are null (except those leading to rejection), send records to the review section.

9. **Matching Process:**

   o Records not meeting rejection criteria are processed for matching.

   o Extract the first 3 characters from fname, lname, city, and state for matching.

   o Records with a match score of 80% or higher are auto-consolidated and loaded into the standard output.

   o Records with a match score between 60% and 80% are sent to the duplicate review section.

---

## Data Profiling

Data profiling was conducted to assess the quality of the input dataset, identifying key issues such as:

- Null values in Address2, Address3.

- Inconsistent values in cust_tier.

- Length and format inconsistencies in zip and Cust_id.

In the Developer :

**Column Profiling**

All 2499 rows. Last run on:Jan 15, 2025 3:49:28 PM EET

| Column | Distinct Values | % Distinct | Nulls | % Null | Datatype | Documented Datatype | Max Value | Min Value | Last Profiled |
|---|---|---|---|---|---|---|---|---|---|
| ⊟ All_Customers | | | | | | | | | |
| CustomerID | 2495 | 99.83 | 4 | 0.16 | Integer(8) [100.00] | decimal(8) | 19134145 | 10110095 | Jan 15, 2025 3:49:28 PM... |
| CustomerTier | 9 | 0.36 | 171 | 6.84 | String(7) [100.00] | string(7) | Silver | 1 | Jan 15, 2025 3:49:28 PM... |
| Lastname | 274 | 10.96 | - | - | String(11) [100.00] | string(11) | Zack | Abedini | Jan 15, 2025 3:49:28 PM... |
| Firstname | 207 | 8.28 | - | - | String(9) [100.00] | string(9) | Zakir | Aaron | Jan 15, 2025 3:49:28 PM... |
| Company | 2157 | 86.31 | - | - | String(51) [100.00] | string(51) | ZURICH... | #INPUT... | Jan 15, 2025 3:49:28 PM... |
| Address1 | 2251 | 90.07 | - | - | String(60) [100.00] | string(60) | WEST LO... | #1 TOM... | Jan 15, 2025 3:49:28 PM... |
| Address2 | 0 | - | 2499 | 100.0 | | string(1) | | | Jan 15, 2025 3:49:28 PM... |
| Address3 | 0 | - | 2499 | 100.0 | | string(1) | | | Jan 15, 2025 3:49:28 PM... |
| City | 553 | 22.12 | 7 | 0.28 | String(20) [100.00] | string(20) | WYCKOFF | ** SYSTE... | Jan 15, 2025 3:49:28 PM... |
| City2 | 3 | 0.12 | 2495 | 99.83 | Fixed Length String(2) [1... | string(2) | TX | TX | Jan 15, 2025 3:49:28 PM... |
| State | 64 | 2.56 | 107 | 4.28 | Fixed Length String(2) [9... | string(13) | WY | AK | Jan 15, 2025 3:49:28 PM... |
| Zip | 1053 | 42.13 | - | - | Integer(5) [96.43] | string(10) | 99701 | 00906-7... | Jan 15, 2025 3:49:28 PM... |
| PrefDelivMethod | 12 | 0.48 | - | - | String(35) [100.00] | string(35) | United P... | Courier | Jan 15, 2025 3:49:28 PM... |
| Country | 4 | 0.16 | 58 | 2.32 | Fixed Length String(2) [9... | string(6) | USA | U.S. | Jan 15, 2025 3:49:28 PM... |
| CreateDate | 661 | 26.45 | - | - | Date [100.00] | string(10) | 9/9/1990 | 1/10/1983 | Jan 15, 2025 3:49:28 PM... |
| Currency | 5 | 0.20 | 81 | 3.24 | Fixed Length String(3) [9... | string(3) | USD | #in | Jan 15, 2025 3:49:28 PM... |
| OrderAmount | 2179 | 87.19 | 1 | 0.04 | Integer(6) [100.00] | decimal(5) | 277698 | -47064 | Jan 15, 2025 3:49:28 PM... |
| Status | 6 | 0.24 | - | - | Fixed Length String(4) [9... | string(8) | pending | A | Jan 15, 2025 3:49:28 PM... |
| MiscDate | 1 | 0.04 | - | - | Date [100.00] | string(8) | 31/12/99 | 31/12/99 | Jan 15, 2025 3:49:28 PM... |

In the Analyst :



## Data Quality Rules Implementation

Key transformations and rules implemented include:
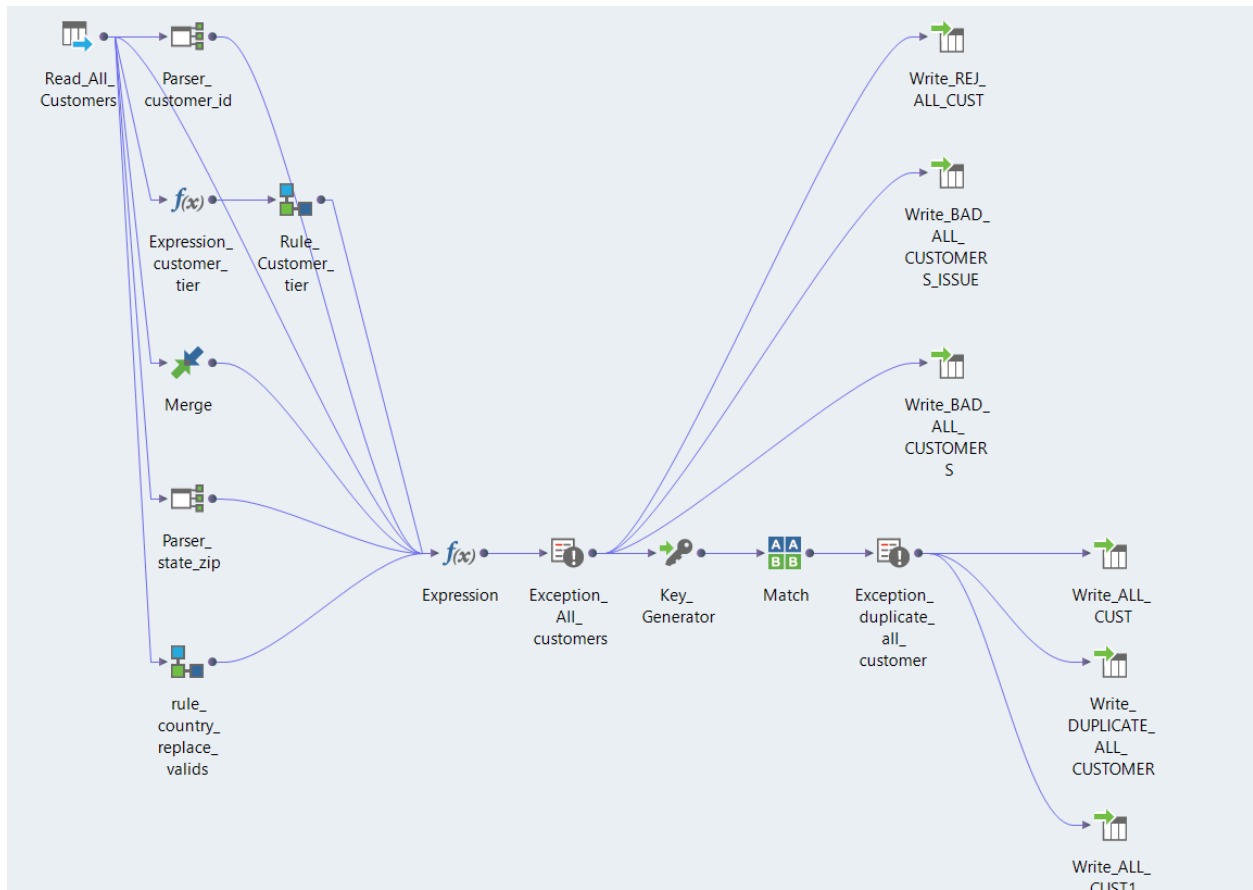
- **Parsing and Validation:**
  - Cust_id parsed and checked for numeric values.
  - zip length validated to ensure 5 digits.

- **Rule-Based Transformations:**
  - Standardizing cust_tier values.
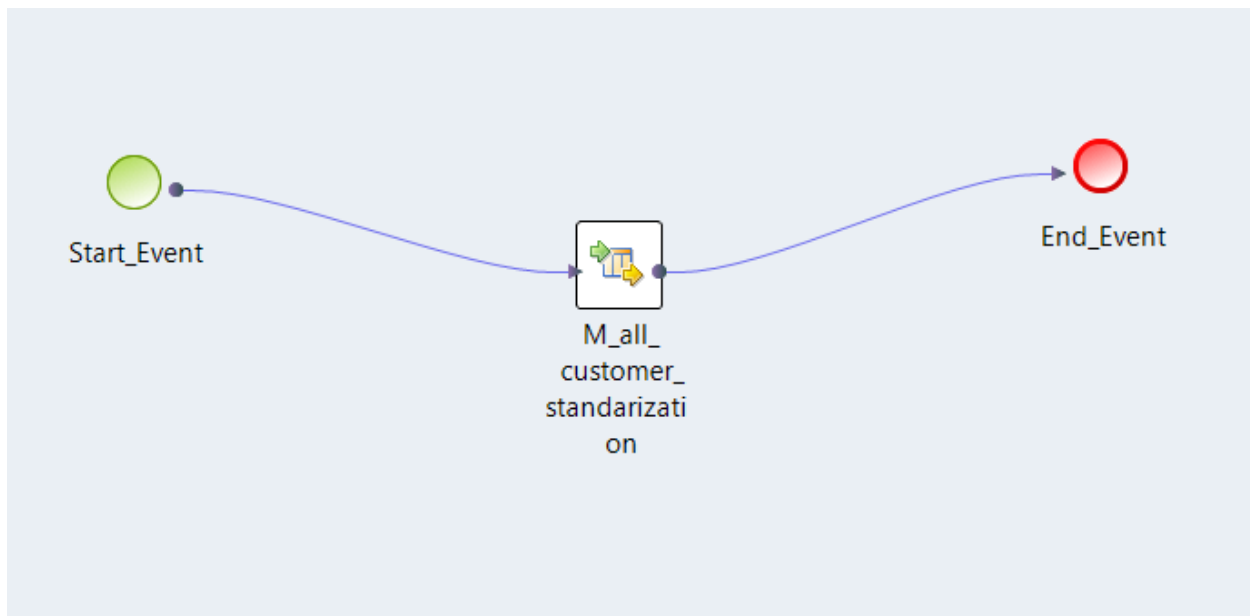  - Separating concatenated fields.



## Data Standardization Workflow

The workflow designed in Informatica Data Quality follows the steps below:

1. Reading customer data from the source.

2. Applying various data quality transformations.

3. Sending rejected records to the rejection table.

4. Consolidating valid records and exporting them to PowerCenter.

## Exporting to PowerCenter

After standardization, clean data is exported to Informatica PowerCenter for integration into downstream systems.

## Results and Observations

- Successfully processed 2499 records.

- Identified and rejected 5% of records due to null values.

- 10% of records were flagged for review due to partial completeness.

- The matching process identified duplicate records with a 70% match score.



## Conclusion

The Informatica Data Quality project successfully improved the accuracy and consistency of customer data. The project followed a structured approach to validate, standardize, and cleanse data, ensuring high-quality data for business processes.