

Dear,

I hope this message finds you well. Following our recent discussions and the analysis of your dataset, I wanted to share some insights regarding the data quality issues we've identified, along with strategies to mitigate these issues effectively. Our assessment was guided by the Data Quality Framework Table and other relevant resources, focusing on key criteria and dimensions critical for maintaining high data quality.

### **Data Quality Issues Identified:**

1. **Completeness:** We noticed missing values in several columns, including 'address', 'postcode', and 'state'. This can impact the accuracy of analytics and decision-making processes.
2. **Consistency:** There are inconsistencies in the 'state' column, with variations in abbreviations and full names used for the same state. This could lead to discrepancies in data aggregation and analysis.
3. **Accuracy:** Some 'postcode' entries do not align with the corresponding 'state' or 'address', indicating potential inaccuracies in the data entry process.
4. **Uniqueness:** Duplicate records were found in the 'customer\_id' column, which could lead to skewed analysis and insights.
5. **Timeliness:** The dataset does not include timestamps for when the data was collected or last updated, making it challenging to assess its current relevance.

### **Strategies to Mitigate These Issues:**

1. **Data Cleaning:** Implement a routine data cleaning process to address missing values, either by removing records with missing data or imputing them based on other data points. Tools like Polymer's CSV Data Query Agent can automate parts of this process.
2. **Standardization:** Develop a data entry guideline that standardizes the format for entries in columns like 'state'. Utilize data validation rules to ensure consistency in future data entries.
3. **Data Validation:** Introduce a validation step in your data collection process to verify the accuracy of critical information, such as matching postcodes with the correct state and address.
4. **De-duplication:** Use data processing tools to identify and remove duplicate records, ensuring each 'customer\_id' is unique. Regular checks for duplicates should be part of your data maintenance routine.
5. **Timestamping:** Include timestamps for data collection and updates. This will help in assessing the timeliness of the data and in making informed decisions based on the most current information.

### **Details of the errors found on the dataset:**

Customer Demographic:

<b>FIELD NAME</b>	<b>ERRORS</b>
<u>DOB</u>	01 record 1843 87 records Blanks
<u>last_name</u>	125 records Blanks
<u>Gender</u>	88 records gender 'U' Values are not consistence M, Male, F, Female, Femal, U
<u>job_title</u>	506 records Blanks
<u>job_industry</u>	656 records mention 'N/A'
<u>Default</u>	3317 records value 'special characters' includes null and Blanks
<u>Tenure</u>	87 records Blanks

Transations:

<b>FIELD NAME</b>	<b>ERRORS</b>
<u>Online_order</u>	94 records Blanks
<u>brand</u>	48 records Blanks
<u>product_line</u>	48 records Blanks
<u>product_class</u>	48 records Blanks
<u>product_size</u>	48 records Blanks
<u>standard_cost</u>	48 records Blanks
<u>product_first_sold_date</u>	48 records Blanks

Implementing these strategies will significantly improve the quality of the data, enhancing the reliability of your analytics and the insights derived from them. Our team is here to assist you in this process, providing the necessary tools and expertise to ensure your data meets the highest standards of quality.

Please let us know if you have any questions or if there's anything specific, you'd like to discuss further. We're committed to supporting you in achieving and maintaining excellent data quality.

Best regards,

Malusi Msweli