

**Subject:** Identifying and Addressing Data Quality Issues in Receipts, Users, and Brands Data

Hi Mike,

I hope this email finds you well. After a detailed exploratory analysis of our Receipts, Users, and Brands datasets, I've identified several critical data quality issues that I believe are essential to address:

1. **Missing Data:**

- **finishedDate:** 49% of receipts lack the date indicating when they became invalid.
- **pointsEarned:** 45% of records are missing data on points earned.
- **purchasedItemCount:** Significant gaps in this data may affect the assessment of special offers and bonus points eligibility.
- **totalSpent:** Missing values in these fields impact our ability to track transaction amounts and items purchased, which in turn affects the accuracy of pointsEarned data.
- **topBrand:** Incomplete data on whether a brand should be featured as a 'top brand'.

2. **Anomalous Values:**

- For **pointsEarned**, **purchasedItemCount**, and **totalSpent**, there are numerous entries with values significantly higher than the norm. These anomalies need investigation to ensure they are not the result of errors in our app processes.

3. **Duplicate Records:**

- Over half of the records in the Users dataset are duplicates. This issue needs immediate attention to remove redundancies and prevent future occurrences.

4. **Inconsistent Date Formats:**

- Dates are recorded in varying formats, deviating from standard formats like MM/DD/YYYY. We need to standardize date formats across our database.

**Questions About the Data:**

- What processes are currently in place for capturing and validating the data in these fields?
- Are there any known issues or events that could have led to the high number of missing or anomalous values?
- What steps have been taken previously to address similar data quality issues?

**Information Needed to Resolve Issues:**

- Detailed documentation on the data collection and processing workflows.
- Access to raw data logs and user activity records for further investigation.
- Input from the development team regarding the handling and storage of these data fields.

**Additional Information for Optimization:**

- Insights into user behavior and transaction patterns to better understand the context of the data.
- Historical data trends to identify persistent issues and their impact.
- Feedback from stakeholders on key metrics and priorities for data quality.

**Performance and Scaling Concerns:**

- Potential increased load on the database during data cleaning and validation processes.
- Ensuring real-time data updates do not compromise performance.
- Strategies to manage larger datasets effectively as we scale, such as indexing and partitioning.

I have developed a plan to address these issues and would like to discuss it with you in detail. Please let me know a convenient time for us to meet and go over the proposed solutions.

Best Regards,