

Fourth: Communicate with Stakeholders

I am addressing this email to a product manager named **Ashton**.

Subject: Data Quality Issues, Concerns, and Plans

Dear Ashton,

I hope this message finds you well. I am sharing some insights and observations from my recent data analysis work that could significantly impact our decision-making process and help improve our business outcomes.

I encountered several interesting findings and data quality issues that are worth addressing:

1) Data Quality Issues:

The data contains lots of missing values in the instances of brands and receipts. This could impact our understanding of the brand recognition and receipt validation process and affect the accuracy of our analysis. One of the most serious concerns is in the receiptItemList instance of the receipts data as for most of the keys a value is missing which could affect our tracking of items.

I also noticed some inconsistencies in data which is needed to be addressed for making our analysis efficient and effective.

2) Data Quality Check and Cleaning:

I used Python and Pandas, a popular data manipulation library, to explore and clean the data. I identified the data quality issues by checking for missing values, inconsistencies, and anomalies.

3) Resolution of Data Quality Issues:

To resolve the data inconsistency quality issues, we need to know about the data sources, revisit the data collection process and put some constraints to avoid inconsistencies. And, to resolve the missing data issue we need to impute the values based on the data available or we can use standard imputing techniques like mean imputation, substitution, or regression imputation.

4) Information Required for Optimization:

It would be helpful to know more about the data collection process, potential data sources, and the business rules that govern certain fields. This context would aid in making more accurate data-related decisions. Understanding the criteria for 'Accepted' and 'Rejected' receipts would enable us to interpret the impact of these categories on our analysis. There are some fields in the item list which are mostly empty, so knowing how important each piece of data is would help us in efficient data processing.

5) Performance and Scaling Considerations:

As we scale up our data analysis, we might encounter challenges related to processing time and memory usage. We can consider optimizing our code for efficiency and utilizing distributed computing frameworks like Spark, and Kafka. Additionally, ensuring that our data storage infrastructure is robust and capable of handling increasing volumes of data will be crucial for maintaining performance. This could be achieved by creating ETL pipelines to extract, transform, and load data as per requirements. Also, the collaboration with domain experts to validate data assumptions and refine data definitions will enhance the quality and utility of our data assets.

I'd greatly appreciate the opportunity to discuss these findings further and explore potential strategies for addressing these concerns. Your expertise and input will be invaluable in charting the best path forward.

Please let me know a convenient time for us to connect and delve into these observations. Looking forward to our productive discussion.

Best regards,
Ashish Patidar