

Note: I am writing an email to Hemant (fictional) Team Lead, informing about the data quality issues I found.

Subject: Data quality issues with Receipts, Users and Brands data

Hello Hemant,

After conducting a thorough Exploratory Analysis of the records in Receipts, Users and Brands, I came across the following data quality issues that I believe are important for you to know.

A considerable amount of data is missing for certain fields like-

finishedDate- for 49% (almost half) of the receipts we do not know when they become invalid (assuming that the date on which a receipt finishes processing is the date on which it becomes invalid)

pointsEarned- 45% of the values for the '*pointsEarned*' field are missing. This means that points were earned for certain receipts, but the data was not captured and that is why the large number of missing values.

purchasedItemCount- large number of missing values will pose problems for deciding if users who bought more than one unit of a product qualify for special offers/bonus points that require them to purchase certain number of products/brands.

totalSpent, *rewardsReceiptItemList*- Since data for the total amount spent on a receipt, and items shipped in a transaction is missing, it is natural that we do not have information about points earned (*pointsEarned* field) for those transactions.

topBrand (Boolean indicator for whether the brand should be featured as a 'top brand')

categoryCode (The category code that references a category of a brand)

For columns '*pointsEarned*', '*purchasedItemCount*', '*totalSpent*', there are a significant number of values that seem out of place (very large as compared to most values in respective fields). I would recommend investigating the processes in our App that produce these values, to determine if they are legit or result of something erroneous that is happening.

There are a lot of (more than half) duplicate records in the Users data. I strongly suggest going over our database, to eliminate redundant records and ensure there is no way such anomalies can happen again.

Lastly, I found the date formats to be inconsistent, against the usual MM/DD/YYYY or similar standard date formats. For this too, I'd consider going over through our database to ensure that date fields are being captured and stored in a consistent manner.

I have a plan to deal with the missing values issue and other problems as well, that I'd like to discuss with you in detail. Let me know what time works best for you so that we can set up a meeting.

Thank You

Regards,

Ankit Hemant Lade