

## **EMAIL TO STAKEHOLDER**

Subject: Addressing Data Quality Issues in the Data

Dear Stakeholder,

I hope this email finds you well. I wanted to provide you with an update on my recent assessment of the data quality for the Data.

Here are some key points from assessment:

### **Data Quality Issues:**

I identified several data quality issues in the 3 data files, including:

1. Invalid JSON format: Some JSON files were found to be invalid, which could potentially cause errors during data processing. I fixed these issues to ensure the data is in a valid format.
2. Incorrect UUID format: I found instances where UUIDs were not in the correct format. I have corrected these UUIDs to ensure data consistency and integrity.
3. Missing Values: There were missing values in some columns of the data. I have addressed these missing values by replacing them with default values.

### **Discovery of Data Quality Issues:**

I discovered these data quality issues through a systematic analysis of the Data(preprocessing.py). I developed scripts to check for JSON validity, UUID format, and missing values in the data.

### **Resolving Data Quality Issues:**

To resolve the identified data quality issues, I implemented automated scripts(DataSolutions.py) to fix invalid JSON formats and incorrect UUIDs. For missing values, I replaced them with default values to ensure completeness of the data.

### **Additional Information Needed:**

Moving forward, it would be helpful to have documentation or guidelines on data entry standards to prevent future occurrences of data quality issues. I also need access to any data dictionaries or metadata that describe the structure and semantics of the data. In order to optimize the data assets, I would require insights into the business requirements and use cases for the data. Understanding the end goals will help me design and structure the data in a way that best supports your needs.

**Performance and Scaling Concerns:**

As we scale our data processing pipelines, I anticipate potential performance bottlenecks and scalability challenges. I plan to address these concerns by optimizing query performance, implementing data partitioning strategies, and leveraging cloud-based infrastructure for scalability.

Please feel free to reach out if you have any questions or require further clarification on any aspect of the data quality assessment. I am committed to ensuring the integrity and reliability of our data assets.

Best regards,

Siva Rama Pavan Kumar Buddi