

## Data Collection and Preprocessing Phase

Date	9 July 2024
Team ID	SWTID1719935665
Project Title	GeminiDecode: Multilanguage Document Extraction by Gemini Pro
Maximum Marks	2 Marks

## Data Quality Report

**Overview:** The Data Quality Report summarizes data quality issues from the document extraction process, including severity levels and resolution plans. It aims to systematically identify and rectify data discrepancies, ensuring high accuracy and efficiency in data extraction from multilingual documents.

Data Source	Data Quality Issue	Severity	Resolution Plan
<b>Extracted Invoices Dataset</b>	Missing text in specific fields due to OCR limitations	High	Enhance OCR model accuracy by training on a diverse dataset, and implement post-OCR validation checks.
<b>Extracted Legal Documents Dataset</b>	Misclassification of document types	Moderate	Improve the classification algorithm by incorporating additional features and refining the training data.
<b>Extracted Financial Statements Dataset</b>	Inconsistent numerical data formats	Low	Standardize numerical data formats using data preprocessing scripts before analysis.
<b>Extracted Medical Records Dataset</b>	Incorrect language detection for multilingual documents	High	Integrate advanced language detection algorithms and cross-validate with known language segments.

<b>Extracted Business Reports Dataset</b>	Duplicate data entries due to repeated scans	Moderate	Implement deduplication techniques to identify and merge duplicate entries based on unique identifiers.
<b>Extracted Research Papers Dataset</b>	Incomplete metadata extraction	Low	Automate metadata extraction processes and cross-reference with existing bibliographic databases.
<b>Extracted Contracts Dataset</b>	Data entry errors in extracted text	High	Use automated error detection tools to identify and correct data entry errors, and employ human verification for critical sections.
<b>Extracted Receipts Dataset</b>	Outliers in numerical data due to scanning artifacts	Moderate	Apply statistical methods to detect and correct outliers, such as z-score analysis or the IQR method.
<b>Extracted Surveys Dataset</b>	Missing responses in survey data	High	Use imputation techniques to fill missing responses, and ensure completeness by prompting users for mandatory fields during survey collection.
<b>Extracted Documents Dataset</b>	Data integrity issues due to format variations	Moderate	Establish data integrity constraints and conduct regular audits to ensure consistency across different formats.