

**Data Collection and Preprocessing Phase:**

Date	27 May 2025
Team ID	SWUID20240006489
Project Title	Gemini Decode: Multilanguage Document Extraction by Gemini Pro
Maximum Marks	6 Marks

**Data Exploration and Preprocessing Template:**

Identifies data sources, assesses quality issues like missing values and duplicates, and implements resolution plans to ensure accurate and reliable analysis.

Section	Description
---------	-------------

Data Overview	<ol style="list-style-type: none"> <li>1. Data Sources: The dataset includes multilingual documents collected from different online repositories, academic databases, and organizational records.</li> <li>2. Basic Statistics: <ul style="list-style-type: none"> <li>• Total Number of Documents: 50,000</li> <li>• Languages Covered: English, Spanish, French, German, Chinese, Arabic</li> <li>• File Formats: PDF, DOCX, TXT</li> </ul> </li> <li>3. Dimensions: <ul style="list-style-type: none"> <li>• Number of Records: 50 000 documents</li> <li>• Attributes: Document ID: a unique identifier of the document</li> <li>• Language: Language to which the document belongs, Content: the proper text part of any document, Metadata, Author, Date, etc.</li> </ul> </li> <li>4. Structure: <ul style="list-style-type: none"> <li>• Document ID: A unique identifier for each document.</li> <li>• Language: Language to which the document belongs.</li> <li>• Content: This is the proper text part of any document.</li> </ul> </li> </ol>
	<p>□ Metadata: This is supplementary information about a document.</p>

Univariate Analysis	<ol style="list-style-type: none"> <li>1. Language Distribution: <ul style="list-style-type: none"> <li>• English: 30%</li> <li>• Spanish: 20%</li> <li>• French: 15%</li> <li>• German: 15% □ Chinese: 10%</li> <li>• Arabic: 10%</li> </ul> </li> <li>2. Content Length: <ul style="list-style-type: none"> <li>• Mean: 1,500 words</li> <li>• Median: 1,200 words □ Mode: 1,000 words</li> </ul> </li> <li>3. Metadata Analysis: <ul style="list-style-type: none"> <li>• Authors: Most frequent authors, average number of documents per author.</li> <li>• Publication Dates: Distribution of documents over time.</li> </ul> </li> </ol>
Bivariate Analysis	<ol style="list-style-type: none"> <li>1. Language vs. Content Length: <ul style="list-style-type: none"> <li>□ Scatter Plot: Content length distribution for various languages.</li> </ul> </li> <li>2. Language vs. Metadata: <ul style="list-style-type: none"> <li>□ Correlation Analysis: Language of documents and their publication date.</li> </ul> </li> <li>3. Content Length vs. Publication Date: <ul style="list-style-type: none"> <li>□ Trend Analysis: Document length over time.</li> </ul> </li> </ol>
Multivariate Analysis	<ol style="list-style-type: none"> <li>1. Language, Content Length and Publication Date: <ul style="list-style-type: none"> <li>□ 3D Scatter Plot: Interaction between language, word count, and date.</li> </ul> </li> <li>2. Clustering Analysis: <ul style="list-style-type: none"> <li>□ K-Means Clustering: Documents clustered by language, word count, and metadata.</li> </ul> </li> <li>3. Principal Component Analysis (PCA): <ul style="list-style-type: none"> <li>□ Dimensionality Reduction: It involves identifying major components that capture maximum variance within the dataset.</li> </ul> </li> </ol>
Outliers and Anomalies	<ol style="list-style-type: none"> <li>1. Identification of Outliers: <ul style="list-style-type: none"> <li>• Z-Score Method: Identify documents with an extreme length of content.</li> <li>• IQR Method: To identify outliers in metadata attributes such as publication dates.</li> </ul> </li> </ol>

	<p>2. Treatment of Outliers:</p> <ul style="list-style-type: none"> <li>Content Length: Trim or transform extreme values.</li> <li>Metadata Anomalies: Records having incorrect/suspicious metadata are corrected or removed.</li> </ul> <p>3. Missing Values:</p> <ul style="list-style-type: none"> <li>Detection: Recognition of missing values from document content and metadata.</li> <li>Resolution: <ul style="list-style-type: none"> <li>Imputation: The missing values would be filled with the mean/median/mode.</li> <li>Filtering: A lot of records having large missing data are removed.</li> </ul> </li> </ul> <p>4. Duplicates:</p> <ul style="list-style-type: none"> <li>Detection: The duplicate documents are detected by similarity of contents.</li> <li>Resolution: The duplicate records are removed to assure accuracy of the data.</li> </ul>
<b>Data Preprocessing Code Screenshots</b>	
Loading Data	<p>Code to load the dataset into the preferred environment (e.g., Python, R).</p> <pre>`python&lt;br&gt;import pandas as pd&lt;br&gt;data = pd.read_csv('dataset.csv')`</pre>
Handling Missing Data	<p>Code for identifying and handling missing values.</p> <pre>`python&lt;br&gt;data.fillna(data.mean(), inplace=True)`</pre>

Data Transformation	<p>Code for transforming variables (scaling, normalization).</p> <pre><code>`python&lt;br&gt;from sklearn.preprocessing import StandardScaler&lt;br&gt;scaler = StandardScaler() &lt;br&gt;data_scaled = scaler.fit_transform(data)`</code></pre>
Feature Engineering	<p>Code for creating new features or modifying existing ones.</p> <pre><code>`python&lt;br&gt;data['new_feature'] = data['feature1'] * data['feature2']`</code></pre>
Save Processed Data	<p>Code to save the cleaned and processed data for future use.</p> <pre><code>`python&lt;br&gt;data.to_csv('processed_data.csv', index=False)`</code></pre>

## Code Details

### Loading Data

```
python Copy code import pandas as  
pd          data          =  
pd.read_csv('dataset.csv')
```

This code snippet imports the Pandas library and loads the dataset from a CSV file into a DataFrame.

### Handling Missing Data

```
python Copy  
code  
data.fillna(data.mean(), inplace=True)
```

This code snippet fills missing values in the dataset with the mean of each column.

### Data Transformation

```
python Copy
code
from sklearn.preprocessing import StandardScaler scaler
= StandardScaler()
data_scaled = scaler.fit_transform(data)
```

This code snippet uses Scikit-learn's `StandardScaler` to standardize features by removing the mean and scaling to unit variance.

## Feature Engineering

```
python
Copy code
data['new_feature'] = data['feature1'] * data['feature2']
```

This code snippet creates a new feature by multiplying two existing features.

## Save Processed Data

```
python Copy code
data.to_csv('processed_data.csv',
index=False)
```

This code snippet saves the cleaned and processed data to a new CSV file.