

Data Collection and Preprocessing Phase:

Date	09 July 2024
Team ID	SWUID20240028082
Project Title	GeminiDecode: 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"> 1. Data Sources: The dataset includes multilingual documents collected from different online repositories, academic databases, and organizational records. 2. Basic Statistics: <ul style="list-style-type: none"> • Total Number of Documents: 50,000 • Languages Covered: English, Spanish, French, German, Chinese, Arabic • File Formats: PDF, DOCX, TXT 3. Dimensions: <ul style="list-style-type: none"> • Number of Records: 50 000 documents • Attributes: Document ID: a unique identifier of the document • Language: Language to which the document belongs, Content: the proper text part of any document, Metadata, Author, Date, etc. 4. Structure: <ul style="list-style-type: none"> • Document ID: A unique identifier for each document. • Language: Language to which the document belongs. • Content: This is the proper text part of any document.

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

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

Feature Engineering	Code for creating new features or modifying existing ones. <pre>`python
data['new_feature'] = data['feature1'] * data['feature2']`</pre>
Save Processed Data	Code to save the cleaned and processed data for future use. <pre>`python
data.to_csv('processed_data.csv', index=False)`</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.