# Data Collection and Preprocessing Phase

| Date | 19 July 2024 |
|---|---|
| Team ID | SWTID1720099578 |
| Project Title | Inquisitive: A Multilingual AI Question Generator |
| 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 | The project uses user-provided text inputs in various languages to generate questions. The dataset includes:<br><br>• **Text**: Input text for question generation.<br>• **Language**: Detected language of the input text.<br>• **Translated Text**: English translation of the input text (if necessary).<br>• **Generated Questions**: Questions generated from the translated text.<br>**Dimensions:**<br><br>• Number of records: Dependent on user input.<br>• Number of features: 4 (Text, Language, Translated Text, Generated Questions) |

| | |
|---|---|
| Univariate Analysis | **Text**:<br><br>• Analyze text length.<br>• Common words.<br>• Language distribution.<br>**Language**:<br><br>• Frequency distribution of detected languages.<br>**Translated Text**:<br><br>• Analysis of text length after translation.<br>**Generated Questions**:<br><br>• Number and quality of generated questions. |
| Bivariate Analysis | **Text Length vs. Generated Questions**:<br><br>• Correlation between text length and the number or quality of generated questions.<br>**Language vs. Generated Questions**:<br><br>• Analysis to determine if certain languages produce more or better-quality questions. |
| Multivariate Analysis | **Language, Text Length, and Generated Questions**:<br><br>• Combined effect of language and text length on the generated questions.<br>**Sentiment of Text and Generated Questions**:<br><br>• Influence of input text sentiment on the generated questions. |
| Outliers and Anomalies | **Outliers in Text Length**:<br><br>• Identify unusually short or long texts.<br>**Anomalies in Generated Questions**:<br><br>• Detect nonsensical or irrelevant questions. |
| **Data Preprocessing Code Screenshots** | |

| Loading Data | ```python
import os
import google.generativeai as palm
from langdetect import detect
from googletrans import Translator
from dotenv import load_dotenv


load_dotenv()


api_key = os.getenv("API_KEY")
palm.configure(api_key=api_key)
translator = Translator()
``` |
|---|---|
| Handling Missing Data | ```python
if st.button("Generate Questions"):
    if user_text:
        questions = generate_questions(model_name, translated_text)
        if detected_language !='en':
            questions = translator.translate(questions, src="en", dest=detected_language).text
        st.subheader("Generated Questions:")
        st.write(questions)
    else:
        st.warning("Please enter some text.")
``` |
| Data Transformation | ```python
if user_text:
    detected_language = detect(user_text)
    if detected_language!='en':
        translated_text = translator.translate(user_text, src=detected_language, dest="en").text
    else:
        translated_text = user_text
``` |
| Feature Engineering | ```python
def generate_questions(model_name, text):
    response = palm.generate_text(
        model=model_name,
        prompt=f"Generate questions from the following text:\n\n{text}\n\nQuestions:",
        max_output_tokens=150
    )
    questions = response.result.strip() if response.result else "No questions generated."
    return questions
``` |
| Save Processed Data | The primary use case is for users to input text and receive questions immediately. This real-time interaction does not require storing the questions, as they are displayed directly to the user. |