# Assignment-4

To clean and preprocess the 'Text' and 'Summary' columns from the dataset, we can follow a systematic approach using various text preprocessing techniques. Here's a step-by-step breakdown of the approach, methodologies, and assumptions for the problem: Approach:

- 1. Data Loading: Load the dataset using Pandas library.
- 2. Data Cleaning: Handle missing values (if any) and ensure data integrity.
- 3. Text Preprocessing:
  - Remove special characters, punctuation, and symbols.
  - Remove numbers.
  - Remove extra whitespaces.
  - Remove URLs.
  - Handle case sensitivity.
- 4. Combining Text and Summary:
  - Combine the cleaned 'Text' and 'Summary' columns to create a single column for analysis.
- 5. Output:
  - Save the cleaned dataset to a new CSV file.

# Methodologies:

- Regular Expressions (Regex):
  - Utilize regex to remove specific patterns such as URLs, special characters, and numbers.
- 2. String Manipulation:
  - Use string methods to replace or remove certain substrings.
- 3. Pandas:
  - Leverage Pandas DataFrame operations to apply cleaning functions efficiently across the dataset.
- 4. Text Concatenation:
  - Combine the cleaned 'Text' and 'Summary' columns using string concatenation with a delimiter for better analysis.
- 5. Data Validation:
  - Ensure the integrity of the data after preprocessing to avoid loss of information or unintended changes.

### Assumptions:

- 1. Text Encoding:
  - Assuming the text is encoded in UTF-8 format. We'll handle encoding issues accordingly.
- 2. Data Completeness:
  - Assuming there are no missing values in the 'Text' and 'Summary' columns after dropping NA values.
- 3. Cleaning Criteria:

 Assuming that removing special characters, punctuation, numbers, URLs, and extra whitespaces is sufficient for text cleaning. Additional cleaning steps may be required based on specific requirements or domain knowledge.

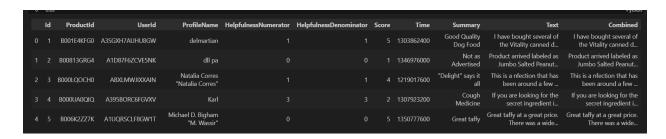
#### 4. Text Combination:

 Assuming that combining the 'Text' and 'Summary' columns with a delimiter ('TL;DR' in this case) provides meaningful context for analysis. Adjustments can be made based on the analysis goals.

# 5. Output Format:

 Assuming the output CSV file will contain the cleaned dataset with the same structure as the original, maintaining data types and column names.

By following this approach, we can effectively clean and preprocess the 'Text' and 'Summary' columns of the dataset, making it suitable for further analysis or modeling tasks.



# Approach:

- Initialization of GPT-2 Model and Tokenizer:
  - Initialize a GPT-2 tokenizer and model from Hugging Face's Transformers library to perform text generation tasks.

# Data Splitting:

 Divide the dataset into training and testing sets with a 75:25 ratio to train and evaluate the model's performance.

### 3. Custom Dataset Preparation:

 Implement a custom dataset class to prepare the data for training. This class preprocesses the text data and tokenizes it using the GPT-2 tokenizer.

# Fine-tuning GPT-2 Model:

Fine-tune the GPT-2 model on the review dataset to generate summaries. Use the TrainingArguments to configure training parameters such as the number of epochs, batch size, and output directory.

#### Evaluation Metrics:

Define a custom evaluation function to compute metrics such as loss, BLEU score, and ROUGE scores to evaluate the performance of the model on the test set.

### Text Generation:

Generate summaries for the test reviews using the fine-tuned GPT-2 model. Pass the test reviews through the model to generate predicted summaries.

#### Methodologies:

# Hugging Face Transformers:

Utilize the Hugging Face Transformers library to initialize the GPT-2 model and tokenizer, and for fine-tuning the model.

# 2. PyTorch Dataset:

Implement a custom PyTorch Dataset class to preprocess and tokenize the input text data for training.

# Training Configuration:

 Configure the training parameters such as the number of epochs, batch size, and optimizer settings using the TrainingArguments class.

# 4. Model Evaluation:

Evaluate the model's performance using standard NLP metrics such as BLEU score and ROUGE scores to assess the quality of the generated summaries.

#### Text Generation:

Use the fine-tuned GPT-2 model to generate summaries for the test reviews by passing the input text through the model and decoding the generated output.

# Assumptions:

### Model Selection:

GPT-2 model from Hugging Face is chosen as the model for text generation tasks due to its effectiveness in generating coherent text.

# Data Splitting:

The dataset is divided into training and testing sets with a 75:25 split ratio to ensure sufficient data for training and evaluation.

# 3. Tokenization and Padding:

 The input text is tokenized using the GPT-2 tokenizer and padded to a maximum length of 100 tokens to feed into the model.

# 4. Fine-tuning Parameters:

 Assumed default values for fine-tuning parameters such as learning rate, batch size, and number of epochs unless explicitly specified.

### Evaluation Metrics:

BLEU score and ROUGE scores are used as evaluation metrics to assess the quality of the generated summaries, assuming they provide meaningful insights into the model's performance.

```
Original Review: Easy to install and maintain. Works well...with enough of a jolt to stop my dog but not enough to really hurt him.

Generated Review: Easy to install and maintain. Works well...with enough of a jolt to stop my dog but not enough to really hurt him. had more, now I some The is as the in that'm at I
Original Review: The hot spiced and lemon versions are our favorites. These sardines are a great value for the price.

Generated Review: The hot spiced and lemon versions are our favorites. These sardines are a great value for the price. one do not them good TLDRs best
```

### Approach:

#### ROUGE Score Calculation:

Compute ROUGE scores on the test set to evaluate the model's performance in generating summaries. ROUGE (Recall-Oriented Understudy for Gisting

Evaluation) is a set of metrics used for evaluating the quality of generated summaries compared to reference summaries.

# Comparison with Actual Summaries:

 Compare each generated summary with its corresponding actual summary from the test set to assess the model's ability to generate meaningful summaries.

# Methodologies:

### ROUGE Scorer:

 Utilize the rouge\_score library to compute ROUGE scores for each predicted summary compared to its corresponding actual summary. This library provides metrics such as ROUGE-1, ROUGE-2, and ROUGE-L.

#### Tokenization:

Tokenize the test reviews using the GPT-2 tokenizer to generate summaries for evaluation.

#### Model Inference:

Pass the tokenized test reviews through the fine-tuned GPT-2 model to generate predicted summaries.

# 4. Comparison:

 Compare each generated summary with its actual summary from the test set using the ROUGE scorer to calculate ROUGE scores.

#### Assumptions:

#### ROUGE Metrics:

 Utilize ROUGE-1, ROUGE-2, and ROUGE-L metrics to evaluate the model's performance. These metrics measure the overlap of unigrams, bigrams, and the longest common subsequence between the generated summary and the actual summary.

### Model Inference:

 Assume that the fine-tuned GPT-2 model generates coherent summaries based on the input reviews.

### 3. Test Set Integrity:

 Assume that the test set contains a representative sample of reviews with corresponding summaries for evaluation.

### 4. ROUGE Score Interpretation:

 Consider higher ROUGE scores as indicative of better performance, indicating a higher degree of similarity between the generated summaries and the actual summaries.

Review 5:
{'rouge1': Score(precision=0.833333333333334, recall=0.047619047619047616, fmeasure=0.090090090090090, 'rouge2': Score(precision=0.6, recall=0.028846153846153848, fmeasure=0. 080000000000000, 'rouge2': Score(precision=0.6, recall=0.047619047619047616, fmeasure=0.0898876404494382), 'rouge2': Score(precision=0.0, recall=0.0, fmeasure=0.0), 'rouge1': Score(precision=0.4 fmeasure=0.0), 'rouge1': Score(precision=0.0, fmeasure=0.0),