

Review Summarization :

1. Collected the dataset from the site given in the question.
2. Imported all the packages needed to preprocess reviews like nltk, numpy and pd etc.
3. Lowered texts.
4. Tokenized texts.
5. Removed stop words.
6. Corrected the spacing for the output generated and then removed punctuation.
7. Finally, performed stemming and lemmatization.
8. I have applied the preprocessing to Text and Summary columns.
9. Then saved the result to a new csv file '`final_csv_after_preprocessing.csv`'.
10. Then took the first 6667 rows of the dataset as the training was taking a lot of time.
11. Split the data into train and test set and saved both dataframes to csv files train.csv and test.csv respectively.

Custom dataset class:

- SummaryDataset()
 - I have initialized the block size.
 - I have initialized a data frame after conversion from csv file and then two other data frames for Text and Summary respectively.
 - Initialized a tokenizer.
 - Defined `len` function.
 - Defined `__get__` function.
 - In which I have first converted a specific row data into encodings.
 - Then returned the Input IDs and attention masks and the embedding of respective labels.
- Training of Model.
 - Fine tuned the model on the food review dataset.
 - In the train function:
 - I have loaded the model and the tokenizer and The DataCollator as there might be a chance that inputs are not padded to the same length, so to pad them we need a collator.
 - Then passed the training arguments like number of epoch, batch size and at which checkpoint should we save the model, etc.
 - Trained the model and saved the trained model in '`final_model_aditya`' directory as this directory changes as we tune the model.
- Load and Generation of output
 - Then I loaded the saved model and the tokenizer and created a new function named "generate_text" and then generated the output from the saved model using generated ids from tokenizer.
 - I have formed a list '`summaries_and_generated`' in which I am saving the test label given and their corresponding generated summary.
- Evaluation.

- Calculated the Rouge score of unigram, bigram and longest common subsequence(for entire testing set).

Hyper Parameter Tuning:

- Code flow is the same.
- Model and loading and generation of text is the same.
- 2nd model I have changed the block size to 512 and the number of epochs to 3.
- 3rd model I have again changed the block size to 128 and increased the number of epochs to 8.

Result Analysis:

Result for model 3 was the best. The reason I think is because of increasing the number of epochs the model ran for more number of times so the loss decreased, hence outputting the best result.

Result 3 > Result 2 > Result 1

```
Pair 10: {'rouge-1': {'r': 0.5, 'p': 0.03896103896103896, 'f': 0.07228915528523736}, 'rouge-2': {'r': 0.0, 'p': 0.0, 'f': 0.0}, 'rouge-l': {'r': 0.5, 'p': 0.03896103896103896, 'f': 0.07228915528523736}}
Pair 11: {'rouge-1': {'r': 1.0, 'p': 0.04444444444444444, 'f': 0.08510638216387506}, 'rouge-2': {'r': 1.0, 'p': 0.017241379310344827, 'f': 0.0338983047515081}, 'rouge-l': {'r': 1.0, 'p': 0.017241379310344827, 'f': 0.0338983047515081}}
Pair 12: {'rouge-1': {'r': 0.5, 'p': 0.029850746268656716, 'f': 0.056338027105733005}, 'rouge-2': {'r': 0.3333333333333333, 'p': 0.01020408163265306, 'f': 0.020408163265306}, 'rouge-l': {'r': 0.5, 'p': 0.029850746268656716, 'f': 0.056338027105733005}}
Pair 13: {'rouge-1': {'r': 0.0, 'p': 0.0, 'f': 0.0}, 'rouge-2': {'r': 0.0, 'p': 0.0, 'f': 0.0}, 'rouge-l': {'r': 0.0, 'p': 0.0, 'f': 0.0}}
Pair 14: {'rouge-1': {'r': 0.6666666666666666, 'p': 0.0625, 'f': 0.11428571271836736}, 'rouge-2': {'r': 0.0, 'p': 0.0, 'f': 0.0}, 'rouge-l': {'r': 0.3333333333333333, 'p': 0.0625, 'f': 0.11428571271836736}}
Pair 15: {'rouge-1': {'r': 0.5, 'p': 0.01818181818181818, 'f': 0.0350877186211142}, 'rouge-2': {'r': 0.0, 'p': 0.0, 'f': 0.0}, 'rouge-l': {'r': 0.5, 'p': 0.01818181818181818, 'f': 0.0350877186211142}}
Pair 16: {'rouge-1': {'r': 0.0, 'p': 0.0, 'f': 0.0}, 'rouge-2': {'r': 0.0, 'p': 0.0, 'f': 0.0}, 'rouge-l': {'r': 0.0, 'p': 0.0, 'f': 0.0}}
Pair 17: {'rouge-1': {'r': 0.0, 'p': 0.0, 'f': 0.0}, 'rouge-2': {'r': 0.0, 'p': 0.0, 'f': 0.0}, 'rouge-l': {'r': 0.0, 'p': 0.0, 'f': 0.0}}
Pair 18: {'rouge-1': {'r': 1.0, 'p': 0.04166666666666666, 'f': 0.079999999232}, 'rouge-2': {'r': 0.5, 'p': 0.010416666666666666, 'f': 0.020408162865472725}, 'rouge-l': {'r': 1.0, 'p': 0.04166666666666666, 'f': 0.079999999232}}
Pair 19: {'rouge-1': {'r': 0.5, 'p': 0.021739130434782608, 'f': 0.041666665868055563}, 'rouge-2': {'r': 0.0, 'p': 0.0, 'f': 0.0}, 'rouge-l': {'r': 0.5, 'p': 0.021739130434782608, 'f': 0.041666665868055563}}
Pair 20: {'rouge-1': {'r': 0.25, 'p': 0.038461538461538464, 'f': 0.06666666435555564}, 'rouge-2': {'r': 0.0, 'p': 0.0, 'f': 0.0}, 'rouge-l': {'r': 0.25, 'p': 0.038461538461538464, 'f': 0.06666666435555564}}
Pair 21: {'rouge-1': {'r': 0.5, 'p': 0.03278688524590164, 'f': 0.06153846038343197}, 'rouge-2': {'r': 0.3333333333333333, 'p': 0.012048192771084338, 'f': 0.020408162865472725}, 'rouge-l': {'r': 0.5, 'p': 0.03278688524590164, 'f': 0.06153846038343197}}
Pair 22: {'rouge-1': {'r': 0.6, 'p': 0.07317073170731707, 'f': 0.13043478067107753}, 'rouge-2': {'r': 0.5, 'p': 0.034482758620689655, 'f': 0.06451612782518212}, 'rouge-l': {'r': 0.6, 'p': 0.07317073170731707, 'f': 0.13043478067107753}}
Pair 23: {'rouge-1': {'r': 0.4444444444444444, 'p': 0.04081632653061224, 'f': 0.07476635359944102}, 'rouge-2': {'r': 0.125, 'p': 0.007751937984496124, 'f': 0.020408162865472725}, 'rouge-l': {'r': 0.4444444444444444, 'p': 0.04081632653061224, 'f': 0.07476635359944102}}
Pair 24: {'rouge-1': {'r': 0.0, 'p': 0.0, 'f': 0.0}, 'rouge-2': {'r': 0.0, 'p': 0.0, 'f': 0.0}, 'rouge-l': {'r': 0.0, 'p': 0.0, 'f': 0.0}}
Pair 25: {'rouge-1': {'r': 0.0, 'p': 0.0, 'f': 0.0}, 'rouge-2': {'r': 0.0, 'p': 0.0, 'f': 0.0}, 'rouge-l': {'r': 0.0, 'p': 0.0, 'f': 0.0}}
...
Pair 1664: {'rouge-1': {'r': 1.0, 'p': 0.07142857142857142, 'f': 0.1333333320888889}, 'rouge-2': {'r': 0.25, 'p': 0.010309278350515464, 'f': 0.019801979437310}, 'rouge-l': {'r': 1.0, 'p': 0.07142857142857142, 'f': 0.1333333320888889}}
Pair 1665: {'rouge-1': {'r': 0.6666666666666666, 'p': 0.10526315789473684, 'f': 0.18181817946280993}, 'rouge-2': {'r': 0.5, 'p': 0.043478260869565216, 'f': 0.079999999232}, 'rouge-l': {'r': 0.6666666666666666, 'p': 0.10526315789473684, 'f': 0.18181817946280993}}
Pair 1666: {'rouge-1': {'r': 1.0, 'p': 0.04166666666666666, 'f': 0.079999999232}, 'rouge-2': {'r': 0.0, 'p': 0.0, 'f': 0.0}, 'rouge-l': {'r': 1.0, 'p': 0.04166666666666666, 'f': 0.079999999232}}
```