

Topic Aware Text Summarization

Using:- NLP, Topic Modeling, Text summarization

09.09.2024

Team 24

- 1. Jhansi Laxmi Polagani 11710412
- 2. Sahithi Padigela 11702556
- 3. Shivam Jadhav 11673635
- 4. Siddhartha Alwala 11661116

GitHub Repository Link:-

 $\underline{https://github.com/shivamjadhav2000/Topic-Aware-Text-Summarization.git}$

Motivation

Why do we need text summarization in the new world of exponentially growing data? It's estimated that 2.5 quintillion bytes of data are created daily, which presents a challenge to efficiently study any given text corpus and extract knowledge from it. This is especially important in academic and research areas where time is limited, but data loss cannot be compromised. Summarizing the text makes it more helpful to subdivide the corpus into sections or individual topics because a summary may overly focus on a single **Topic**, leading to the loss of crucial information in subsections of the text. This study/project aims to solve this problem

Significance

The two-model approach allows each model to specialize. Topic modeling identifies key themes, while summarization focuses on condensing content based on those themes. This ensures more focused and relevant summaries. Additionally, it offers more modularity and flexibility, enabling fine-tuning or swapping of models independently, thereby improving performance and interpretability. In the field of Natural Language Processing (NLP), generating summaries that retain key essence elements enhances information retrieval and improves user understanding. This project can have implications for various sectors, such as academia (research paper summarization) and business (summarizing customer reviews by product aspects).

Objectives

In this project, we will focus on building a two-phase model. In the first phase, we will train the model to extract topics, and in the second phase, the model will focus on text summarization.

Specifically, we aim to:

- Implement a topic modeling system (LDA, NMF,BERT) to extract key topics from a corpus.
- Fine-tune a pre-trained text summarization model (LSA,T5) to condition its output on the identified topics.
- Evaluate the quality of the summaries in terms of topic coherence and information relevance.

Features

I. Neural Topic Modeling

The first stage of the pipeline will involve discovering topics from a large corpus using neural topic models like LDA, and NMF Topic, which provide more traditional Topic Modeling techniques, additionally we shall explore the more advanced models such as transformers BERT.

II. Summarization Model

We will use a text summarization model (such as LSA, T5, BART) to generate summaries. The model will be fine-tuned to focus on the most important topics identified by the topic modeling stage.

III. Evaluation

We will use standard NLP evaluation metrics such as ROUGE and BLEU to assess the quality of the summaries, while also incorporating topic-coherence metrics to gauge relevance.

Dataset

For this project, we will use a dataset like the Quora Questions dataset provided in the project ideas file, which contains Questions on Quora on many topics. This dataset is suitable for extractive and abstractive summarization tasks and will help in training and evaluating the performance of our model. Additional datasets such as the Amazon Product Reviews dataset may be used for experiments with topic-specific summarization in the review domain.

- Size: The Quora Questions dataset contains over 400,000 Questions.
- Type: The dataset consists of Many different with human-written question on various topic summaries, making it an ideal choice for both topic modeling and summarization.
- Preprocessing: Before analysis, the text will be cleaned by removing unnecessary characters, and tokenization will be applied. For topic modeling, we will also experiment with stop-word removal and lemmatization to improve the quality of the topics extracted.

• Sentence Segmentation: we are dealing with a text corpus which has multiple questions or data which is in multiple sentences, segregating those sentences is a challenge given the vividness of the data.

Visualization

Workflow diagram



