

Text Summarization using Python

Shreya Vishwakarma

Prasoon

Parul Goel

+91 9696908131

+91 7739581188

+91 8924047448

vishreya363@gmail.com

spyprasoon@gmail.com

goelparul2002@gmail.com

Abstract: With the advent of the digital age, huge amounts of text are created daily, and it is becoming ever harder to process and comprehend long content manually. This project introduces a Python-based text summarization system that automates the reduction of long documents with the retention of their crucial details. Using Natural Language Processing (NLP) methods, the system is able to offer both extractive and abstractive summarization. Libraries like NLTK, spaCy, and Transformers are used to analyze language structures, extract key sentences, and create meaningful summaries. The solution is scalable, lightweight, and can be used in domains like news summarization, academic research, and information retrieval.

1. Introduction:

As the internet expands exponentially, individuals and institutions are increasingly confronted with the problem of processing huge amounts of text information. From research papers and news articles to legal documents and customer feedback, being able to comprehend important information at a quick rate is more vital than ever before. Manual summarization of these texts takes time and is also subject to inconsistency, particularly when handling repetitive work.

Text summarization is the operation of generating a concise and logical summary of a longer text. It is used to boost productivity and decision-making by allowing users to quickly comprehend key content. In this project, we hope to create a Python-based tool for text summarization that utilizes Natural Language Processing methods to extract or generate automatically summaries from provided input text.

It targets both abstractive and extractive summarization techniques. Abstractive summarization paraphrases and rewrites the text using deep learning models, while extractive summarization picks the most important sentences from the original text. Python provides a strong environment with extensive libraries such as NLTK, spaCy, and Hugging Face Transformers to implement and test these techniques. Our tool is meant to be easy, effective, and flexible in many domains where quick information processing is needed.

2. Literature Review

2.1 Types of Text Summarization

Extractive Summarization: Selects key sentences from the original text without altering them.

Abstractive Summarization: Generates new phrases and sentences to convey the original meaning.

2.2 NLP Techniques and Tools

Tokenization, part-of-speech tagging, named entity recognition

Pre-trained transformer models (e.g., BART, T5)

2.3 Prior Research Early work by Luhn (1958) focused on word frequency-based methods.

Rush et al. (2015) introduced neural attention models for summarization. More recently, Transformer-based architectures like BERTSUM and PEGASUS have shown state-of-the-art performance in both extractive and abstractive summarization tasks.

3. Methodology

3.1 Data Sources

News articles from CNN/DailyMail dataset

Scientific articles from arXiv dataset

3.2 Preprocessing Steps

Text normalization (lowercasing, punctuation removal)

Tokenization and stopword removal

Named entity and sentence segmentation

3.3 Summarization Techniques

Extractive: Ranking sentences using TF-IDF based scoring and frequency matrices

Abstractive: Generating new text using pre-trained Transformer models (e.g., BART, T5)

3.4 Evaluation Metrics

ROUGE (Recall-Oriented Understudy for Gisting Evaluation) scores

Human evaluation for coherence and fluency

4. Implementation

4.1 Tools and Libraries Used

Python: Main programming language

spaCy: Linguistic processing and sentence segmentation

Transformers (Hugging Face): Abstractive summarization models

Custom TF-IDF Pipeline: For extractive summarization

4.2 The text is first preprocessed to eliminate irrelevant information. The text is tokenized into words and sentences. A term frequency matrix then computes how frequently each word occurs. Subsequently, inverse document frequency (IDF) gives weight to infrequent words in sentences. With TF and IDF, a TF-IDF score is then calculated for each sentence. Sentences with greater scores are ranked as more important. Lastly, the most appropriate sentences are chosen to come up with a brief summary that encapsulates the key ideas of the paper.

4.3 Implementation Challenges

1. Efficaciously handling variability in text structure and language delicacies.
2. Balancing brevity of summary with maintenance of major facts.
3. Working with large datasets without compromising on processing speed and accuracy.
4. Handling ambiguity and loss of context in automated summarization.
5. Choosing proper algorithms for extractive and abstractive summarization.

5. Results

5.1 Performance Metrics

Extractive Summarization (TF-IDF): ROUGE-1: 0.43, ROUGE-L: 0.40

Abstractive Summarization (BART): ROUGE-1: 0.47, ROUGE-L: 0.44

5.2 User Feedback

A usability survey of 20 participants showed that:

90% found the summaries accurate

85% appreciated the clean and responsive UI

75% preferred abstractive summarization for readability.

6. Limitations

Large transformer models require significant memory and computation

Summaries may lose factual correctness in abstractive mode

Limited support for multi-document summarization

6.1 Suggested Improvements:

Use model distillation for lighter models

Add factual consistency checks

Implement multi-document input support

7. Future Improvements

7.1 Features to be implemented

Real-time extraction from web links or PDF files

Multilingual support for summarization

Platform integration with smartphones

7.2 Directions Of Research

Researching reinforcement learning for improving quality of summaries

Adding user feedback loops to enhance relevance

8. Ethical Issues

8.1 Privacy and Consent

Ensure user information entered into the summarizer tool is not retained or abused.

Give users proper consent and data removal policies.

8.2 Bias and Misinformation

Use care with language models that have the potential to inject bias or produce disinformation. Ongoing review is critical.

8.3 Intellectual Property

Summarizing copyrighted content needs to adhere to fair use policies. Always properly cite sources when using third-party content.

9. Conclusion

This project illustrates the implementation of a proficient text summarization tool using Python and NLP libraries. Through the implementation of extractive and abstractive approaches, the system offers versatility and accuracy in diverse content types. In spite of present drawbacks, future improvements and ethical practices can further enhance its abilities. The tool is poised for extensive usage in academia, journalism, customer service, and other areas.

Abbreviations

- NLP: Natural Language Processing
- ROUGE: Recall-Oriented Understudy for Gisting Evaluation
- UI: User Interface
- API: Application Programming Interface

Glossary

- **Tokenization:** Splitting text into words or sentences
- **ROUGE Score:** A metric to evaluate summary quality based on content overlap
- **Transformer:** A deep learning model architecture for sequence-to-sequence tasks

References

- [1] H.P. Luhn, "The automatic creation of literature abstracts," IBM Journal, 1958.
- [2] A.M. Rush, S. Chopra, J. Weston, "A Neural Attention Model for Abstractive Sentence Summarization," EMNLP 2015.
- [3] Y. Liu, M. Lapata, "Text Summarization with Pretrained Encoders," EMNLP 2019.
- [4] Hugging Face Transformers. [Online]. Available: <https://huggingface.co/transformers>
- [5] Streamlit Documentation. [Online]. Available: <https://docs.streamlit.io>