

An Analytical Study of Text Summarization Techniques

By Bavrabi Ghosh, Aritra Ghosh, Subhojit Ghosh, Anupam Mondal

INSTITUTE OF ENGINEERING & MANAGEMENT, KOLKATA

Introduction

Text summarization is crucial for condensing large information into concise overviews, enabling efficient comprehension across various domains like journalism, research, business, and social media. However, existing summarization models often fail to capture the nuances and details required by users. To address these limitations and ease research efforts, we aimed to create a comprehensive document on summarization techniques and applications.

Input Article

Marseille, France (CNN) The French prosecutor leading an investigation into the crash of Germanwings Flight 9525 insisted Wednesday that he was not aware of any video footage from on board the plane. Marseille prosecutor Brice Robin told CNN that "so far no videos were used in the crash investigation." He added, "A person who has such a video needs to immediately give it to the investigators." Robin's comments follow claims by two magazines, German daily Bild and French Paris Match, of a cell phone video showing the harrowing final seconds from on board Germanwings Flight 9525 as it crashed into the French Alps. All 150 on board were killed. Paris Match and Bild reported that the video was recovered from a phone at the wreckage site. ...

Abstractive summarization

Text Summarization Models

Extractive summarization

Generated summary

Prosecutor : " So far no videos were used in the crash investigation "

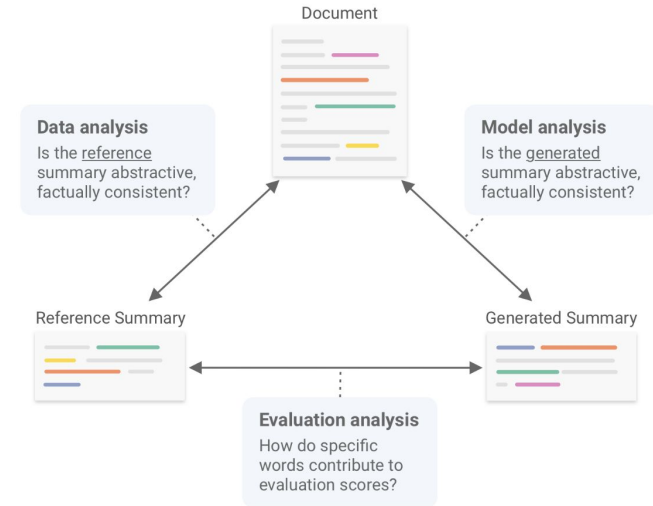
Extractive summary

marseille prosecutor brice robin told cnn that " so far no videos were used in the crash investigation . " robin 's comments follow claims by two magazines , german daily bild and french paris match , of a cell phone video showing the harrowing final seconds from on board germanwings flight 9525 as it crashed into the french alps . paris match and bild reported that the video was recovered from a phone at the wreckage site .



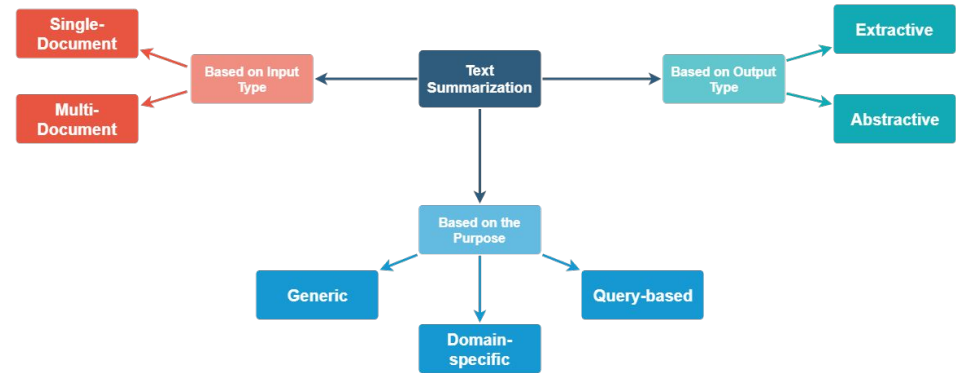
Motive

We explored various text summarization approaches and techniques to gain insights into existing methods' complexities and identify factors behind suboptimal performance. An analysis of research papers revealed extractive and abstractive summarization as prevalent methods. Extractive summarization utilizes ranking algorithms to extract salient sentences, while abstractive summarization attempts to generate new sentences capturing the meaning. We delved into technical details of early systems like the cueing, title, location methods and the Trainable Document Summarizer.



Background

We explored various text summarization approaches and techniques to gain insights into existing methods' complexities and identify factors behind suboptimal performance. An analysis of research papers revealed extractive and abstractive summarization as prevalent methods. Extractive summarization utilizes ranking algorithms to extract salient sentences, while abstractive summarization attempts to generate new sentences capturing the meaning. We delved into technical details of early systems like the cueing, title, location methods and the Trainable Document Summarizer.



Summarization Methods



Extractive Method

- Inverse Document Frequency method
- Cluster Method
- Graph Based Approach
- Latent Semantic Analysis Method

Abstractive Summarization

- Seq2Seq (Sequence-to-Sequence) Models
- Pointer-Generator Networks
- BERTSUM

Deep Learning Summarization

- Bidirectional and Auto-Regressive Transformers
- Generative Pre-trained Transformer
- Text-To-Text Transfer Transformer
- Pegasus



Overview of Important Features

Text representation models in NLP transform words into numerical forms for pattern detection.

- **N-grams** - Group words into N components, offering reasonable vector sizes but computational challenges.
- **Bag of Words** - It disregards word order but captures multiplicity.
- **TF-IDF** - Determines term importance relative to corpus but assumes independent term counts.
- **Word Embedding** - maps words to vectors, placing similar words close together, capturing semantic and syntactic relationships through algorithms like FastText, GloVe, and Word2Vec.



Overview of Findings

Algorithms			Limitations
Extractive	Unsupervised	Fuzzy logic	Post-processing should remove redundancies to improve the quality of summarization.
		Concept-based	The summary should use similarity measures to reduce redundancy, which can affect quality.
		Latent-Semantic	LSA-generated summaries take a long time.
	Supervised	Machine Learning	To make good summaries, it has to be trained and improved on a large set of data.
		Neural Network	Both the training phase and the application phase are quite slow with neural networks. Training data also requires human interruption.
		Conditional Random Fields	Linguistic features are not taken into account in the use of CRF. It also needs an external domain specific corpus.
Abstractive	Structural	Trees	The text ignores context and important phrases in the text, resulting in a failure to recognize the relationships between sentences. Another issue is that it consistently emphasises syntax rather than meaning.
		Template based	As the templates are pre-defined using this technique, the summaries lack variation.
		Based on Rules	It takes a long time to create regulations. It is also difficult to manually write the rules.
		Ontology method	The process of creating a suitable ontology is time-consuming and limited to a single domain.
	Semantic	Multi-model semantic	The framework must be automatically analysed because humans now manually evaluate it.
		Information item	Generating grammatical and meaningful sentences from the material is difficult. The linguistic quality of summaries is low due to incorrect parses.
		Based on Semantic Graph	Limited to single document abstractive summarization.

Challenges

1. Evaluation
2. Important Sentence Selection
3. Anaphora Problem
4. Predefined Template
5. Long Sentences and Jargon
6. Interpretability
7. Cataphora Problem

Future Scope

1. Multilingual Summarization
2. Customizable Summaries
3. Real-Time Summarization
4. Video Summarization



References

1. Aker, A., Cohn, T., Gaizauskas, R.: Multi-document summarization using a* search and discriminative learning. In: Proceedings of the 2010 conference on empirical methods in natural language processing. pp. 482–491 (2010)
2. AL-Banna,A.A.,AL-Mashhadany,A.K.:Natural language processing for automatic text summarization [datasets]-survey. Wasit Journal of Computer and Mathematics Science 1(4), 156–170 (2022)
3. Bala, A., Mitra, R., Mondal, A.: Recommendation system to predict best academic program. In: 2023 7th International Conference on Electronics, Materials Engineering & Nano-Technology (IEMENTech). pp. 1–6. IEEE (2023)
4. Dey,M.,Mondal,A.,Das,D.:Ntcir-12mobileclick:Sense-based ranking and summarization of english queries. In: Ntcir (2016)
5. Edmundson,H.P.:New methods in automatic extracting. Journal of the ACM (JACM) 16(2), 264–285 (1969)
6. Ferreira, R., de Souza Cabral, L., Lins, R.D., e Silva, G.P., Freitas, F., Cavalcanti, G.D., Lima, R., Simske, S.J., Favaro, L.: Assessing sentence scoring techniques for extractive text summarization. Expert systems with applications 40(14), 5755–5764 (2013)
7. Garcí a-Herná ndez, R.A., Ledeneva, Y.: Word sequence models for single text summarization. In: 2009 Second International Conferences on Advances in Computer-Human Interactions. pp. 44–48. IEEE (2009)
8. Kupiec,J.,Pedersen,J.,Chen,F.:A trainable document summarizer. In: Proceedings of the 18th annual international ACM SIGIR conference on Research and development in information retrieval. pp. 68–73 (1995)
9. Mahata,S.K.,Mondal,A.,Dey,M.,Sarkar,D.:Sentiment analysis using machine translation. In: Applications of Machine intelligence in Engineering. pp. 371–377. CRC Press (2022)
10. Mondal, A., Cambria, E., Dey, M.: An annotation system of a medical corpus using sentiment-based models for summarization applications. In: Computational Intelligence Applications for Text and Sentiment Data Analysis, pp. 163–178. Elsevier (2023)
11. Mondal, A., Dey, M., Das, D., Nagpal, S., Garda, K.: Chatbot: An automated conversation system for the educational domain. In: 2018 International Joint Symposium on Artificial Intelligence and Natural Language Processing (iSAI-NLP). pp. 1–5. IEEE (2018)
12. Mondal,A.,Dey,M.,Mahata,S.K.,Sarkar,D.:An automatic summarization system to understand the impact of covid-19 on education. In: Applications of Machine intelligence in Engineering. pp. 379–386. CRC Press (2022)
13. Mridha, M.F., Lima, A.A., Nur, K., Das, S.C., Hasan, M., Kabir, M.M.: A survey of automatic text summarization: Progress, process and challenges. IEEE Access 9, 156043–156070 (2021)
14. Sinha,S.,Mandal,S.,Mondal,A.:Question answering system-based chatbot for healthcare. In: Proceedings of the Global AI Congress 2019. pp. 71–80. Springer (2020)
15. Syed,A.A.,Gaol,F.L.,Matsuo,T.:A survey of the state-of-the-art models in neural abstract text summarization. IEEE Access 9, 13248–13265 (2021)
16. Tas,O.,Kiyani,F.:A survey of automatic text summarization. Press Academia Procedia 5(1), 205–213 (2007)
17. Zhang, Y., Zircir-Heywood, N., Milios, E.: Narrative text classification for automatic key phrase extraction in web document corpora. In: Proceedings of the 7th annual ACM international workshop on Web information and data management. pp. 51–58 (2005)

Thank you