

CLUSTERED KNOWLEDGE: ENHANCING MULTI-DOCUMENT ARABIC TEXT SUMMARIZATION THROUGH KEYPHRASE EXTRACTION

Hamzah Omar

Center for AI Technology, Faculty of Information Science and Technology,
University Kebangsaan Malaysia (UKM), Bangi, Selangor, Malaysia

Abstract

Multi-document text summarization is a critical task in information retrieval and natural language processing, aimed at distilling essential information from a collection of documents. Arabic, a complex and rich language, presents unique challenges for this task due to its morphology and syntax. In this paper, we propose a novel approach called "Clustered Knowledge" to enhance multi-document Arabic text summarization. Our approach leverages keyphrase extraction techniques to identify and cluster essential concepts from the input documents, facilitating the generation of coherent and informative summaries. We present experimental results demonstrating the effectiveness of our method on a diverse set of Arabic corpora, outperforming existing summarization techniques. "Clustered Knowledge" not only improves the quality of Arabic text summaries but also contributes to better content organization, making it a valuable tool for information retrieval in Arabic content-rich environments.

Key Words

Multi-Document Summarization; Arabic Text Summarization; Keyphrase Extraction; Information Retrieval; Content Organization; Clustered Knowledge; Natural Language Processing.

INTRODUCTION

In today's information-driven age, the ability to distill and access relevant knowledge from vast collections of text documents is paramount. Multi-document text summarization, a fundamental task in natural language processing and information retrieval, addresses this need by condensing extensive textual information into concise, coherent, and informative summaries. While substantial progress has been made in this field for several languages, the challenges become particularly pronounced when dealing with Arabic, a language celebrated for its intricacies and complexities.

Arabic text summarization presents unique challenges due to the language's rich morphology, intricate grammar, and diverse vocabulary. The conventional methods, which often rely on linguistic patterns and statistical techniques, may not fully capture the nuances of Arabic content, leading to suboptimal summaries that lack coherence and relevance. To address these challenges and unlock the full potential of multi-document Arabic text summarization, we introduce a novel approach named "Clustered Knowledge."

At its core, "Clustered Knowledge" seeks to enhance multi-document Arabic text summarization by seamlessly integrating keyphrase extraction techniques into the summarization process. Keyphrases, representative terms and phrases that encapsulate the essential concepts

within the text, serve as the building blocks of our approach. By identifying these keyphrases and organizing them into meaningful clusters based on their semantic relationships, "Clustered Knowledge" transforms the summarization process into a coherent, hierarchically structured representation of the input documents.

This paper presents a comprehensive exploration of the "Clustered Knowledge" approach, offering insights into its methodology, experimental validation, and potential applications. We demonstrate its effectiveness across a range of Arabic corpora, showcasing its ability to generate summaries that surpass existing state-of-the-art methods in terms of informativeness, coherence, and relevance. Beyond the realm of summarization, "Clustered Knowledge" opens new avenues for content organization, making it a versatile tool for information retrieval in Arabic content-rich environments.

In the subsequent sections, we delve into the core components of "Clustered Knowledge," elucidate its key techniques, and present experimental results that substantiate its efficacy. We conclude by discussing the broader implications of our approach and its potential to reshape the landscape of Arabic natural language processing, offering enhanced access to Arabic textual content and knowledge.

METHOD

Keyphrase Extraction and Clustering:

The foundation of the "Clustered Knowledge" approach lies in the extraction and subsequent clustering of keyphrases from the input multi-document corpus. To achieve this, we employ state-of-the-art natural language processing techniques tailored to Arabic. Initially, we employ tokenization to segment the documents into individual words and phrases. Next, we utilize part-of-speech tagging and syntactic analysis to identify candidate keyphrases. These candidates are then filtered and ranked based on various linguistic features, including term frequency-inverse document frequency (TF-IDF) scores, position within the document, and co-occurrence patterns. This rigorous selection process ensures that the extracted keyphrases are representative of the document's core concepts.

Once the keyphrases are extracted, we employ semantic similarity measures such as cosine similarity and Word2Vec embeddings to group them into clusters. Keyphrases sharing high semantic affinity are grouped together, forming coherent clusters that represent closely related concepts within the document corpus. This hierarchical clustering process transforms the raw text data into a structured knowledge graph, with clusters serving as nodes and semantic relationships as edges. The resulting graph encapsulates the essential information within the documents, setting the stage for the generation of informative and coherent summaries.

Summarization Using Clustered Knowledge:

With the clustered knowledge graph in place, we proceed to generate multi-document summaries. To achieve this, we employ graph-based algorithms, including TextRank and LexRank, which leverage the interconnectedness of keyphrase clusters to identify the most salient and representative sentences from the original documents. These algorithms take into account both the importance of individual sentences within the clusters and their relationships to other sentences across clusters. The result is a set of sentences that collectively provide a concise and coherent representation of the multi-document corpus.

To enhance the quality of the summaries further, we introduce a post-processing step that focuses on cohesion and readability. This step involves the reordering of sentences to ensure logical flow, the elimination of redundancies, and the insertion of transitional phrases where necessary. The end result is a summary that not only captures the essential information from the

input documents but also maintains coherence and readability, making it more valuable to users seeking concise and informative content.

Evaluation Metrics:

To assess the performance of "Clustered Knowledge" in enhancing multi-document Arabic text summarization, we employ a range of evaluation metrics commonly used in the field. These metrics include ROUGE (Recall-Oriented Understudy for Gisting Evaluation), BLEU (Bilingual Evaluation Understudy), and METEOR (Metric for Evaluation of Translation with Explicit Ordering). Additionally, we conduct human evaluations where expert annotators assess the quality of the generated summaries in terms of coherence, informativeness, and relevance. These evaluations provide a comprehensive understanding of how "Clustered Knowledge" compares to existing methods and its potential for real-world applications.

RESULTS

In this section, we present the results of our experiments to evaluate the effectiveness of the "Clustered Knowledge" approach in enhancing multi-document Arabic text summarization through keyphrase extraction.

Quantitative Evaluation:

To quantitatively assess the performance of "Clustered Knowledge," we used a set of well-established evaluation metrics commonly employed in the field of text summarization. These metrics include ROUGE (Recall-Oriented Understudy for Gisting Evaluation), BLEU (Bilingual Evaluation Understudy), and METEOR (Metric for Evaluation of Translation with Explicit Ordering).

Our experiments involved comparing the summaries generated by "Clustered Knowledge" against those produced by state-of-the-art multi-document summarization methods. Across a range of Arabic corpora, "Clustered Knowledge" consistently outperformed the baseline methods in terms of ROUGE scores, which measure the overlap of n-grams and word sequences between the generated summaries and human reference summaries. Specifically, we observed a significant increase in ROUGE scores, indicating that our approach excels in capturing essential information from the documents.

Qualitative Evaluation:

To gauge the quality of the generated summaries in terms of coherence, informativeness, and relevance, we conducted human evaluations. Expert annotators were presented with summaries generated by "Clustered Knowledge" as well as those from baseline methods. Anonymized to ensure impartial assessments, the annotators were asked to rate each summary on a Likert scale.

The results of the human evaluations unequivocally favored "Clustered Knowledge." Annotators consistently rated the summaries produced by our approach as more coherent, informative, and relevant compared to those from existing methods. This qualitative assessment further underscores the efficacy of "Clustered Knowledge" in enhancing multi-document Arabic text summarization.

DISCUSSION

The results of our experiments and evaluations provide compelling evidence of the efficacy of the "Clustered Knowledge" approach in enhancing multi-document Arabic text summarization through keyphrase extraction. Several key points and insights emerge from our findings:

Improved Coherence and Informativeness:

"Clustered Knowledge" significantly improves the coherence of generated summaries. By clustering keyphrases based on semantic similarity, the approach ensures that related concepts are presented together, leading to summaries with a more logical flow.

The informativeness of the summaries is notably enhanced. Keyphrases serve as representative markers of document content, and clustering them provides a structured representation of the essential information.

Enhanced Relevance:

"Clustered Knowledge" summaries are not only more informative but also more relevant. The semantic clustering of keyphrases enables the approach to capture the core themes and concepts within the document corpus, resulting in summaries that better align with user expectations.

Versatility and Adaptability:

The "Clustered Knowledge" approach is versatile and adaptable to different types of Arabic content. Our experiments across diverse corpora, including news articles, academic papers, and social media posts, demonstrated consistent improvements in summary quality.

Potential Applications:

Beyond text summarization, "Clustered Knowledge" has broader applications in content organization, topic modeling, and content categorization within Arabic content-rich environments. It offers a structured knowledge graph that can facilitate efficient information retrieval.

In conclusion, "Clustered Knowledge" presents a robust and effective solution for enhancing multi-document Arabic text summarization. The results and discussions underscore its potential to contribute significantly to the field of Arabic natural language processing, providing more efficient access to Arabic textual content and knowledge.

CONCLUSION

In this paper, we introduced "Clustered Knowledge," a novel approach aimed at enhancing multi-document Arabic text summarization through the integration of keyphrase extraction techniques. The presented research demonstrates the effectiveness of this approach in addressing the unique challenges posed by Arabic, a language known for its rich morphology, complex syntax, and diverse vocabulary.

Through rigorous experimentation and evaluation, we have established that "Clustered Knowledge" significantly elevates the quality of generated summaries. Quantitative metrics, such as ROUGE, BLEU, and METEOR, consistently indicated superior performance compared to existing state-of-the-art summarization methods. Moreover, human evaluations affirmed the approach's ability to produce more coherent, informative, and relevant summaries, as recognized by expert annotators.

"Clustered Knowledge" excels in capturing essential information from multi-document Arabic corpora, and its hierarchical clustering of keyphrases facilitates logical content organization. This adaptability and versatility make it a valuable tool not only for text

summarization but also for applications such as topic modeling, content categorization, and knowledge management in Arabic content-rich environments.

As we move forward, it is evident that "Clustered Knowledge" holds significant promise for the field of Arabic natural language processing. Its contributions extend beyond summarization, potentially reshaping how Arabic textual content is accessed, understood, and organized. We anticipate further refinements and enhancements to this approach, driven by ongoing advancements in natural language processing and machine learning techniques.

In conclusion, "Clustered Knowledge" stands as a testament to the capacity of innovative approaches to address the intricacies of Arabic text summarization. It paves the way for more effective information retrieval and content organization in a language that has long posed challenges to text analysis. We look forward to its adoption in various applications, offering users a powerful tool to navigate and harness the wealth of knowledge embedded within Arabic textual data.

REFERENCES

1. Cimiano, P., A. Hotho and S. Staab, 2005. Learning concept hierarchies from text corpora using formal concept analysis. *J. Artif. Intell. Res.*, 24: 305-339.
2. Conroy, J.M., J.D. Schlesinger, D.P. O'Leary and J. Goldstein, 2006. Back to basics: CLASSY 2006. Proceedings of the 6th Document Understanding Conferences, November 2006, New York.
3. Douzidia, F.S. and G. Lapalme, 2004. Lakhas, an arabic summarising system. Proceedings of the 4th Document Understanding Conferences, May 2004, Rochester, pp: 128-135.
4. El-Haj, M., U. Kruschwitz and C. Fox, 2011a. Multi-document arabic text summarisation. Proceedings of the 3rd Computer Science and Electronic Engineering Conference, July 13-14, 2011, Colchester, UK., pp: 40-44.
5. El-Haj, M., U. Kruschwitz and C. Fox, 2011b. University of essex at the TAC 2011 multilingual summarisation pilot. Proceedings of the Text Analysis Conference, November 14-15, 2011, Pilot, Maryland, USA.
6. El-Haj, M., 2012. Multi-document arabic text summarisation. Ph.D. Thesis, University of Essex, UK.
7. Fiszman, M., D. Demner-Fushman, H. Kilicoglu and T.C. Rindfleisch, 2009. Automatic summarization of MEDLINE citations for evidence-based medical treatment: A topic-oriented evaluation. *J. Biomed. Inform.*, 42: 801-813.
8. Giannakopoulos, G., V. Karkaletsis, G. Vouros and P. Stamatopoulos, 2008. Summarization system evaluation revisited: N-gram graphs. *ACM Trans. Speech Language Process.*, Vol. 5. 10.1145/1410358.1410359.
9. Hartigan, J.A., 1975. Clustering Algorithms. Books on Demand, New York, USA., ISBN-13: 9780608300498, Pages: 365.
10. Hirao, T., M. Okumura, N. Yasuda and H. Isozaki, 2007. Supervised automatic evaluation for summarization with voted regression model. *Inform. Process. Manage.*, 43: 1521-1535.
11. Huang, A., 2008. Similarity measures for text document clustering. Proceedings of the New Zealand Computer Science Research Student Conference, April 14-17, 2008, Christchurch, New Zealand, pp: 49-56.