

titulo

Remigio Hurtado¹

University Politecnica Salesiana Rhurtadoo@ups.edu.ec
ups.edu.ec

Abstract.

Keywords:

1 Introduction

This paper presents a comprehensive methodology for the development of survey recommendation systems using Natural Language Processing (NLP) and Machine Learning (ML) techniques. The approach ensures systematic analysis and generation of knowledge representations from scientific literature.

2 Methodology

Our methodology is designed to ensure a comprehensive and systematic analysis of scientific literature. The methodology used in this study builds on our previous work, including a methodology that leverages machine learning and natural language processing techniques [5], and a novel method for predicting the importance of scientific articles on topics of interest using natural language processing and recurrent neural networks [6]. The process is structured into three phases: data preparation, topic modeling, and the generation and integration of knowledge representations. Each phase is essential for transforming raw text data into meaningful insights, and the detailed parameters and algorithm are explained below. The table 1 describes the parameters for understanding the overall process and algorithm. The high-level process is presented in Fig. 1, and the detailed algorithm is outlined in Table 2.

In the data preparation phase, we focus on extracting and cleaning text from scientific documents using natural language processing (NLP) techniques. This phase involves several steps to ensure that the text data is ready for analysis. First, text is extracted from the documents, and non-alphabetic characters that do not add value to the analysis are removed. Next, the text is converted to lowercase, and stopwords (common words that do not contribute much meaning) are removed. We then apply lemmatization, which transforms words to their base form (e.g., "running" becomes "run"). Each document is tokenized (split into individual words or terms), and n-grams (combinations of words) are identified to

find common terms. We generate a unified set of common terms, denoted as TE , which includes both the terms extracted from the documents and basic terms relevant to any field of study, such as [Fundamentals, Evaluation of Solutions, Trends]. Finally, each document is vectorized with respect to TE , resulting in the Document-Term Matrix (DTM).

The DTM is a crucial component for topic modeling. It is a matrix where the rows represent the documents in the corpus, and the columns represent the terms (words or n-grams) extracted from the corpus. Each cell in the matrix contains a value indicating the presence or frequency of a term in a document. This structured representation of the text data allows us to apply machine learning techniques to uncover hidden patterns.

In the topic modeling phase, we use Latent Dirichlet Allocation (LDA), a popular machine learning technique for identifying topics within a set of documents. By applying LDA to the DTM, we transform the matrix into a space of topics. Specifically, LDA provides us with two key matrices: the Topic-Term Matrix (β) and the Document-Topic Matrix (θ). The Topic-Term Matrix (β) indicates the probability that a term is associated with a specific topic, while the Document-Topic Matrix (θ) indicates the probability that a document belongs to a specific topic. For each topic in β , we assign a name by combining the most probable terms associated with that topic. This process results in the set T , which contains the most relevant topics, and the set K , which contains the most relevant terms. Additionally, we generate a graph G_{tk} to represent the relationships between topics and terms. Using a language generation model, we refine the names of the topics to ensure they are concise and meaningful, with a maximum of W_t words.

The final phase involves the generation and integration of knowledge representations, which include summaries, keywords, and interactive visualizations. For each topic t in T , we identify the N most relevant documents—those with the highest probabilities in θ . This set, D_t , represents the documents most closely related to each topic. We then generate summaries of these documents, each with a maximum of W_b words, using a language generation model. These summaries include references to the most relevant documents, which are added to the set R_d if they are not already included. We also integrate all the summaries and generate J suggested keywords using a language generation model. Additionally, we create interactive graphs from the G_{tk} graph, showcasing nodes and relationships between the topics and terms, and highlighting significant connections.

This methodology provides a clear, structured approach to analyzing scientific literature, leveraging advanced NLP and machine learning techniques to generate useful and comprehensible knowledge representations. This is the methodology used in the development of this study, which presents the fundamentals, solution evaluation techniques, trends, and other topics of interest

within this field of study. This structured approach ensures a thorough review and synthesis of the current state of knowledge, providing valuable insights and a solid foundation for future research.

Table 1. Description of parameters of the proposed methodology

Parameter	Description
T	Set of topics to be generated with the LDA model.
$\#T$	Cardinality of T . That is, the number of topics (dimensions).
K	Set of common terms with the highest probability in topic t used to label that topic.
$\#K$	Cardinality of K .
W_t	Maximum number of words for combining the common terms of topic t into a new topic name.
N	Number of the most relevant documents to generate a summary of topic t .
W_b	Maximum number of words to generate a summary of the list of N most related documents to topic t .
J	Number of suggested keywords to be generated as knowledge representation.

3 Conclusion

This paper presents a robust methodology for creating survey recommendation systems. By integrating NLP and ML techniques, we ensure a systematic and comprehensive analysis, leading to high-quality knowledge representations and user-friendly reports.

References

1. Hurtado, Remigio; Picón, Cristian; Muñoz, Arantxa; Hurtado, Juan. "Survey of Intent-Based Networks and a Methodology Based on Machine Learning and Natural Language Processing." In Proceedings of Eighth International Congress on Information and Communication Technology. Springer Nature Singapore, Singapore, 2024.
2. Park, Keunheung; Kim, Jinmi; Lee, Jiwoong. "Visual Field Prediction using Recurrent Neural Network," *Scientific Reports*, vol. 9, no. 1, p. 8385, 2019, <https://doi.org/10.1038/s41598-019-44852-6>.
3. Xu, M.; Du, J.; Guan, Z.; Xue, Z.; Kou, F.; Shi, L.; Xu, X.; Li, A. "A Multi-RNN Research Topic Prediction Model Based on Spatial Attention and Semantic Consistency-Based Scientific Influence Modeling," *Comput Intell Neurosci*, vol. 2021, 2021, p. 1766743, doi: 10.1155/2021/1766743.
4. Kreutz, Christin; Schenkel, Ralf. "Scientific Paper Recommendation Systems: a Literature Review of recent Publications," 2022/01/03.
5. Hurtado, R., et al. "Survey of Intent-Based Networks and a Methodology Based on Machine Learning and Natural Language Processing." International Congress on Information and Communication Technology. Singapore: Springer Nature Singapore, 2023.

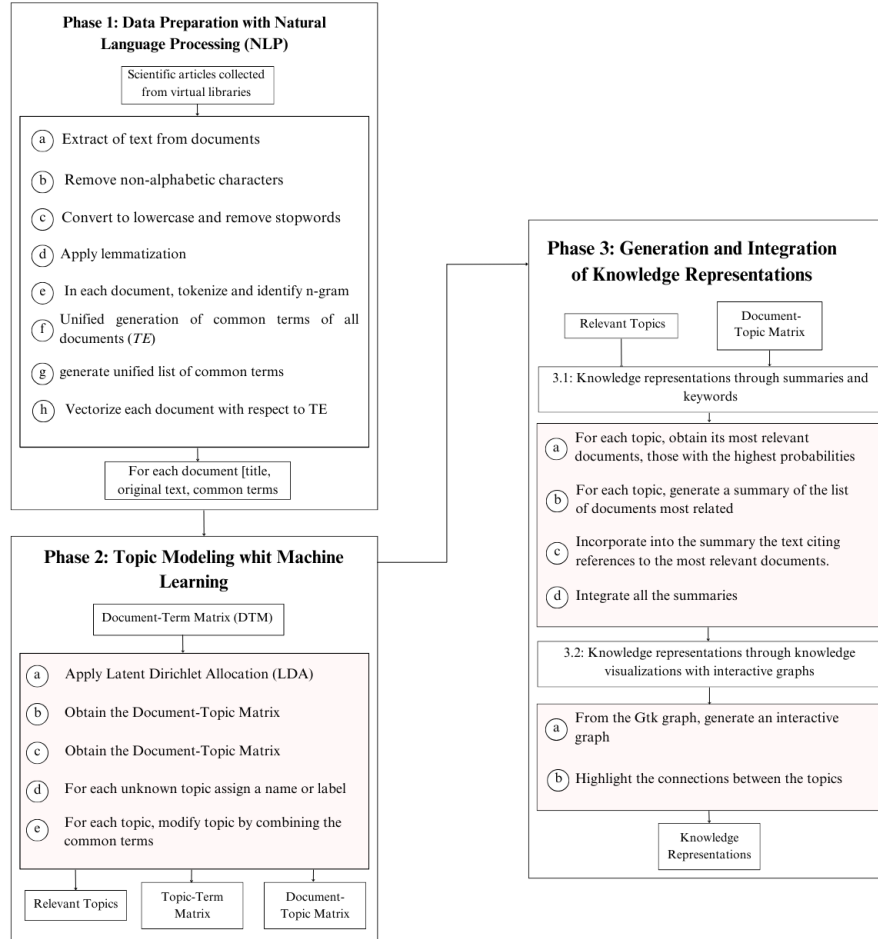


Fig. 1. Methodology for the generation of knowledge representations

6. Lopez, A., Dutan, D., Hurtado, R. "A New Method for Predicting the Importance of Scientific Articles on Topics of Interest Using Natural Language Processing and Recurrent Neural Networks." In: Yang, X.S., Sherratt, S., Dey, N., Joshi, A. (eds) Proceedings of Ninth International Congress on Information and Communication Technology. ICICT 2024 2024. Lecture Notes in Networks and Systems, vol 1013.

General Algorithm

Input: Field of study, Scientific articles collected from virtual libraries, $\#T$, $\#K$, W_t , N , W_b , J

Phase 1: Data Preparation with Natural Language Processing (NLP)

- a. Extract of text from documents.
- b. Remove non-alphabetic characters that do not add value to the analysis.
- c. Convert to lowercase and remove stopwords.
- d. Apply lemmatization to transform words to their base form.
- e. In each document, tokenize and identify n-grams to identify common terms (words or n-grams).
- f. Unified generation of common terms of all documents. Where TE is the unified set of common terms.
- g. Add in TE the basic terms for any field of study, such as: [Fundamentals, Evaluation of Solutions, Trends].
- h. Vectorize each document with respect to TE and generate Document-Term Matrix (DTM).

Output: For each document [title, original text, common terms], and Document-Term Matrix (DTM)

Phase 2: Topic Modeling whit Machine Learning

Input: Document-Term Matrix (DTM)

- a. Apply Latent Dirichlet Allocation (LDA) to transform the DTM Matrix into a space of $\#T$ topics (dimensions).
- b. Obtain the Topic-Term Matrix (β) that indicates the probability that a term is generated by a specific topic.
- c. Obtain the Document-Topic Matrix (θ) that indicates the probability that a document belongs to a specific topic.
- d. For each unknown t topic in β , assign a name or label to the t topic by combining the $\#K$ highest probability common terms in β associated with that t topic. This generates:
 The T Set with the $\#T$ most relevant topics.
 The K Set with the $\#K$ most relevant terms.
 The G_{tk} Graph of the relationships between the most relevant topics (T) and the most relevant terms (K).
- e. For each topic t in T , modify topic t by combining the common terms of t into a new topic name with at most W_t words using a language generation model.

Output: Relevant Topics (T), Topic-Term Matrix (β), Document-Topic Matrix (θ)

Phase 3: Generation and Integration of Knowledge Representations

Input: Relevant Topics T , Document-Topic Matrix (θ)

3.1: Knowledge representations through summaries and keywords

- a. For each t topic in T , obtain its N most relevant documents, i.e., those with the highest probabilities in θ .
 Thus, D_t represents the set of documents most related to each topic t .
- b. For each topic t in T , and from D_t generate a summary of the list of documents most related to that topic t with at most W_b words using a language generation model.
- c. Incorporate into the summary the text citing references to the most relevant documents. Add these references to the set R_d , which will contain all cited references, including new references if they have not been previously included.
- d. Integrate all the summaries and from them generate J suggested keywords using a language generation model.

3.2: Knowledge representations through knowledge visualizations with interactive graphs

- a. From the G_{tk} graph, generate an interactive graph with nodes and relationships between the topics T and the K most relevant terms.
- b. Highlight the connections between the topics T and the K most relevant terms.

3.3: Integration of Knowledge Representations

Output: Knowledge Representations

Fig. 2. General algorithm of the methodology incorporating natural language processing, machine learning techniques and language generation models