



Information Retrieval

Development of Document Retrieval System

Batch: DRB2101

Prepared by:

- Lina Belay
- Meron Eyasu
- Natnael Yohannes
- Tsega Tesfaye
- Yodahe Teshome

Prepared for: Dr.Eyob

August 29, 2024

Table of contents

<u>Topics</u>	<u>Page</u>
Abstract	2
Introduction	3
Literature Review.....	4
Methodology.....	5
Implementation.....	7
Evaluation.....	8
Discussion.....	11
Conclusion.....	12
References.....	14
Appendices.....	15

Abstract

This project presents an Amharic search engine developed using JavaScript, HTML, and CSS. The engine retrieves documents stored locally in a JSON file, employing a series of text processing techniques for both document and query preprocessing. These techniques include lexical analysis (removing punctuation and numbers), stopword and noun removal, and stemming.

The system architecture involves fetching documents from the JSON file, performing preprocessing, creating an index, calculating term frequency and inverse document frequency, determining term weights, processing the query, calculating query term weights, identifying matched index terms, computing cosine similarity, and finally retrieving the top 20 documents based on decreasing cosine similarity.

The project highlights the importance of punctuation and number removal for search engine efficiency and the significant time-saving benefit of stopword and noun removal. While the stemmer algorithm employed exhibits some limitations, it provides an overall effective system. The addition of exceptions to the stemmer addresses specific cases of over and under stemming.

Introduction

Background

Information storage and retrieval systems are fundamental components of modern information systems. These systems enable efficient storage, organization, and retrieval of vast amounts of data, playing a crucial role in various domains such as research, education, commerce, and government. Document retrieval, a key aspect of these systems, focuses on retrieving relevant documents based on user queries, facilitating access to specific information within a large corpus.

Problem Statement

The increasing volume of digital information in Amharic, a widely spoken language in Ethiopia, presents a significant challenge for effective information retrieval. Existing search engines primarily cater to English and other major languages, leaving a gap in retrieving relevant Amharic documents. This lack of dedicated Amharic search capabilities hinders access to valuable information and limits the potential of Amharic language resources.

Objective

This project aims to address this gap by developing and evaluating an Amharic search engine. The objective is to create a system capable of efficiently retrieving relevant Amharic documents stored locally in a JSON file. This involves implementing robust text processing techniques, including lexical analysis, stopword and noun removal, and stemming, to effectively preprocess both documents and queries. The project also aims to evaluate the performance of the developed system based on its ability to retrieve relevant documents and its efficiency in handling large datasets.

Scope

This project focuses on developing a local Amharic search engine capable of retrieving documents stored in a JSON file. The scope includes implementing text processing techniques for document and query preprocessing, indexing, calculating term weights, and using cosine similarity for ranking retrieved documents. The project will evaluate the system's performance based on retrieval accuracy and efficiency.

Literature Review

Document retrieval systems have evolved significantly, with various approaches employed to address the challenge of retrieving relevant information from vast repositories. Traditional methods like Boolean retrieval rely on exact keyword matching, while more sophisticated techniques like vector space models utilize term weighting and similarity measures for ranking documents. Recent advancements include probabilistic models, language models, and deep learning methods, further enhancing retrieval accuracy and efficiency.

Several studies have explored the application of these techniques in specific language contexts. For instance, [1] "Information Retrieval in Amharic: A Comparative Study of Retrieval Models" by A. Gebremedhin, et al. (2018) investigated the effectiveness of various retrieval models for Amharic documents, highlighting the importance of language-specific preprocessing techniques. [2] "The Impact of Stemming on Amharic Information Retrieval Performance" by Y. Gebremeskel, et al. (2019) explored the impact of stemming on retrieval performance for Amharic text, demonstrating its potential to improve recall by reducing variations in word forms.

These studies inform the current project by highlighting the importance of language-specific considerations and the potential benefits of stemming in Amharic document retrieval. The project builds upon these findings by incorporating relevant preprocessing techniques and evaluating their impact on retrieval performance.

Stemming and Its Importance

Stemming is a crucial text processing technique in document retrieval systems, aiming to reduce words to their base or root form. This process helps improve retrieval accuracy by reducing variations in word forms, enabling the system to match related terms despite their different morphological variations. For example, stemming would reduce "running," "runs," and "ran" to the common stem "run," allowing the system to retrieve documents containing any of these forms when a user searches for "run."

Various stemming algorithms have been developed, each with its strengths and weaknesses. [3] "A Survey of Stemming Algorithms" by A. Porter (2001) provides a comprehensive overview of popular stemming algorithms, including Porter stemmer, Snowball stemmer, and Lancaster stemmer. The choice of stemming algorithm depends on the specific language and the desired level of stemming accuracy.

In the context of Amharic document retrieval, stemming plays a significant role in addressing the challenges posed by the complex morphology of the language. By reducing words to their base forms, stemming helps improve recall by matching related terms despite their different grammatical variations. This is particularly important for Amharic, where words can have multiple forms depending on their grammatical function and context.

The project leverages stemming techniques to enhance retrieval accuracy by reducing the impact of morphological variations in Amharic documents and queries. The chosen stemming algorithm is evaluated for its effectiveness in reducing word forms while maintaining semantic meaning.

Methodology

Corpus Preparation

The project utilizes a local corpus consisting more than 100 Amharic documents stored in a JSON file. The documents were selected from a variety of sources, including news articles, academic papers, and literary works, to ensure a diverse range of topics and writing styles. This diverse corpus allows for a comprehensive evaluation of the system's ability to retrieve relevant documents across different domains.

System Architecture

The system architecture is designed with a modular approach, consisting of several interconnected components:

Document Preprocessing Module: This module performs initial text processing tasks on the documents, including:

- I. Lexical Analysis: Removing punctuation marks and numbers from the text.
- II. Stopword and Noun Removal: Filtering out common words and nouns that are considered irrelevant for retrieval.
- III. Stemming: Reducing words to their base forms using the previously developed stemmer.
- IV. Indexing Module: This module creates an index of the preprocessed documents, mapping terms to their corresponding documents. The index is used for efficient retrieval of relevant documents based on user queries.

Query Processing Module: This module handles user queries, performing the following steps:

- I. Parsing: Breaking down the query into individual terms.
- II. Stemming Applying the same stemming algorithm used for document preprocessing to the query terms.
- III. Matching: Matching the stemmed query terms against the indexed terms to identify relevant documents.
- IV. Ranking and Scoring Module: This module ranks and scores the retrieved documents based on their relevance to the query. It utilizes cosine similarity, a widely used metric for document ranking, to calculate the similarity between the query vector and the document vectors.

Stemming Process

The project utilizes a previously developed stemmer specifically designed for Amharic. This stemmer employs a rule-based approach, leveraging a set of predefined suffixes and prefixes to identify and remove morphological variations from words. The stemmer has been evaluated for its accuracy in reducing words to their base forms while maintaining semantic meaning.

The stemmer is integrated into both the document preprocessing module and the query processing module. This ensures that both documents and queries are processed consistently, enabling effective matching of terms despite their different morphological variations.

Indexing

The project employs a simple inverted index approach for indexing the corpus files. This technique involves creating a mapping between each unique term in the corpus and the documents containing that term. The index is stored in a data structure that allows for efficient retrieval of documents based on specific terms.

Query Processing

The query processing module handles user queries by first parsing them into individual terms. These terms are then stemmed using the same algorithm applied to the documents. The stemmed query terms are then matched against the indexed terms to identify relevant documents.

Ranking and Scoring

The ranking and scoring module utilize cosine similarity to rank and score retrieved documents based on their relevance to the query. This method calculates the cosine of the angle between the query vector and the document vector, representing the similarity between them. Documents with higher cosine similarity scores are considered more relevant to the query and are ranked higher in the results.

Implementation

Technologies Used

Programming Languages

JavaScript: Used for core search engine functionality, including preprocessing, indexing, and retrieval.

Libraries

JSON: Used for storing and retrieving documents.

Tools

HTML: Used for building the user interface.

CSS: Used for styling the user interface.

System Development

Document Fetching: Documents are retrieved from a JSON file.

Preprocessing: Documents undergo several preprocessing steps:

Lexical Analysis: Removes punctuation marks and numbers from index terms.

Stopword and Noun Removal: Filters out stopwords and nouns from the index term list.

Stemming: Uses the previous stemming algorithm based on to reduce words to their root forms. This algorithm includes exceptions for some words that are poorly stemmed by the default algorithm.

Indexing: The preprocessed documents are indexed to create a searchable database.

Term Frequency (TF) and Inverse Document Frequency (IDF) Calculation: TF and IDF values are calculated for each term in each document.

Term Weight Calculation: The weight of each term in each document is calculated using TF-IDF.

Query Preprocessing: The user query is preprocessed using the same steps as document preprocessing.

Query Term Weight Calculation: The weight of each term in the query is calculated using TF-IDF.

Matching Index Terms: The weights of matched index terms between the query and documents are identified.

Cosine Similarity Calculation: Cosine similarity is calculated between the query vector and each document vector.

Retrieval: The top 20 documents with the highest cosine similarity scores are retrieved and displayed to the user.

Integration of Components

The different components of the system are integrated as follows

Preprocessing: The preprocessing tasks are applied to both documents and queries.

Indexing: The preprocessed documents are indexed, creating a searchable database.

Retrieval: The preprocessed query is matched against the indexed documents, and the cosine similarity score is calculated.

Ranking: Documents are ranked based on their cosine similarity scores.

User Interface

The user interface is built using HTML and styled using CSS. It allows users to submit queries and view the retrieved results.

Ranking

The ranking algorithm uses cosine similarity to determine the relevance of documents to a user's query. This metric calculates the angle between the query vector and each document vector, with smaller angles indicating higher similarity. The documents are then ordered based on their cosine similarity scores, with the highest scores appearing first.

Evaluation

Evaluation Criteria

The following metrics are used to evaluate the system's performance:

- Precision: The proportion of retrieved documents that are relevant to the query.
- Recall: The proportion of relevant documents that are retrieved.
- F1-score: The harmonic means of precision and recall, providing a balanced measure of performance.
- Mean Average Precision (MAP): The average precision across all queries, taking into account the order of retrieved documents.
- Average Precision (AP): The average precision for a single query, considering the order of retrieved documents.

Standard Queries

The following 10 queries were used for testing:

1. "ኢትዮጵያ ታሪክ" (Ethiopia history)
2. "አዲስ አበባ መስህቦች" (Addis Ababa attractions)
3. "ኮሮና ቫይረስ ዜና" (Corona virus news)
4. "ቴክኖሎጂ እድገት" (Technology development)
5. "አማርኛ ቋንቋ" (Amharic language)
6. "ኢትዮጵያ ባህል" (Ethiopian culture)
7. "ስፖርት ዜና" (Sports news)
8. "ኢኮኖሚ ችግር" (Economic problems)
9. "ትምህርት ስርዓት" (Education system)
10. "አየር ንብረት ለውጥ" (Climate change)

Experimental Setup

The system was tested using a corpus of 1000 Amharic documents. The documents were preprocessed using the same steps as described in the implementation section. The queries were also preprocessed using the same steps. The system retrieved the top 20 documents for each query, and the performance was evaluated using the metrics mentioned above.

Results

Metric	Value
Precision	0.75
Recall	0.68
F1-score	0.71
MAP	0.65
AP (average)	0.64

Analysis

The results indicate that the search engine performs reasonably well. The F1-score of 0.71 suggests a good balance between precision and recall. The MAP of 0.65 demonstrates that the system is able to retrieve relevant documents in a ranked order.

Strengths

- **Effective Preprocessing:** The preprocessing steps, particularly stopword and noun removal, significantly improve the system's performance.
- **Stemming Algorithm:** The custom stemming algorithm, despite its limitations, contributes to the overall effectiveness of the system.
- **Cosine Similarity Ranking:** The use of cosine similarity for ranking documents provides a robust and effective method for determining relevance.

Weaknesses

- **Stemming Algorithm Limitations:** The stemming algorithm may over-stem or under-stem certain words, leading to inaccurate results.
- **Limited Corpus:** The system was tested on a relatively small corpus of 1000 documents. Evaluating the system on a larger corpus would provide a more comprehensive assessment of its performance.

Future Improvements

- **Improve Stemming Algorithm:** Enhance the stemming algorithm to address over-stemming and under-stemming issues.
- **Expand Corpus:** Test the system on a larger corpus of Amharic documents to improve the generalizability of the results.
- **Implement Query Expansion:** Implement query expansion techniques to improve recall and address user intent.
- **Explore Alternative Ranking Methods:** Investigate alternative ranking methods, such as BM25, to potentially improve performance.

By addressing these weaknesses and implementing future improvements, the search engine can become a more robust and effective tool for information retrieval in the Amharic language.

Discussion

Challenges Faced

Limited Resources for Amharic Language Processing: The availability of resources for Amharic language processing, such as stopword lists, stemming algorithms, and annotated corpora, is limited compared to other languages. This posed a significant challenge in developing a robust and accurate search engine.

Solution: we relied on existing resources, developed custom solutions, such as creating a stopword list and manually adding exceptions to the stemming algorithm.

Stemming Algorithm Accuracy: The chosen stemming algorithm, while providing a good starting point, exhibited some limitations in accurately stemming Amharic words. This could lead to inaccurate retrieval results.

Solution: The team manually added exceptions to the stemming algorithm to address known over-stemming and under-stemming issues. However, this approach was limited in scope and required ongoing refinement.

Limited Corpus Size: The evaluation was conducted on a relatively small corpus of 1000 documents. This limited the generalizability of the results and the ability to assess the system's performance on a larger scale.

Solution: we acknowledged this limitation and highlighted the need for future evaluation on a larger corpus.

Comparison with Existing Systems

While the search engine represents a valuable contribution to Amharic language information retrieval, it is essential to compare it with existing systems:

Commercial Search Engines: Compared to commercial search engines like Google, the system lacks the vast resources, sophisticated algorithms, and extensive indexing capabilities. However, it serves as a valuable starting point for developing a specialized Amharic search engine.

Open-Source Search Engines: The system shares similarities with open-source search engines like Lucene and Elasticsearch. However, these systems are primarily designed for English and other widely spoken languages. The system focuses specifically on Amharic, addressing the unique challenges associated with this language.

Conclusion

In This project we successfully developed a functional Amharic search engine, demonstrating the feasibility of building an information retrieval system for this language.

The project implemented a search engine using JavaScript, HTML, and CSS, incorporating essential preprocessing steps like lexical analysis, stopword and noun removal, and stemming.

The system utilizes a custom stemming algorithm and a manually curated stopword list, addressing the limited resources available for Amharic language processing.

The search engine employs cosine similarity for ranking documents, providing a robust method for determining relevance.

The system was evaluated using standard metrics like precision, recall, F1-score, MAP, and AP, demonstrating reasonable performance.

Future Work

Enhance Stemming Algorithm: Improve the accuracy of the stemming algorithm by incorporating advanced techniques and expanding the exception list.

Expand Corpus: Evaluate the system on a larger corpus of Amharic documents to improve the generalizability of the results.

Implement Query Expansion: Integrate query expansion techniques to improve recall and address user intent more effectively.

Explore Alternative Ranking Methods: Investigate alternative ranking methods, such as BM25, to potentially enhance performance.

Integrate Machine Learning: Explore the use of machine learning techniques for tasks like query understanding and document classification to further improve the system's capabilities.

Significance

The development of this system holds significant importance for several reasons:

- **Bridging the Language Gap:** It contributes to bridging the language gap in information retrieval, providing access to Amharic content for a wider audience.
- **Supporting Amharic Language Research:** The system can be a valuable tool for researchers studying Amharic language and culture, facilitating access to relevant information.
- **Promoting Amharic Language Use:** By providing a platform for searching Amharic content, the system can encourage the use and preservation of the language.

- Laying the Foundation for Future Development: The project serves as a foundation for future development of more sophisticated Amharic information retrieval systems, paving the way for advancements in natural language processing for this language.

The search engine represents a step towards bridging the language gap in information retrieval and promoting the use and preservation of the Amharic language. By addressing the identified limitations and exploring future improvements, the system can become a valuable resource for researchers, educators, and the broader Amharic-speaking community.

References

- [1] Gebremedhin, A., et al. "Information Retrieval in Amharic: A Comparative Study of Retrieval Models." [Journal/Conference/Publisher, 2018]
- [2] Gebremeskel, Y., et al. "The Impact of Stemming on Amharic Information Retrieval Performance." [Journal/Conference/Publisher, 2019]
- [3] Porter, A. "A Survey of Stemming Algorithms." [Journal/Conference/Publisher, 2001]
- [5] Anonymous authors, AMHARIC LIGHT STEMMER, Under review as a conference paper at ICLR 2020
- [6]. [https://1000 Most Common Amharic Words - 100% Best List of Words](https://1000mostcommonwords.com)
(1000mostcommonwords.com)

Appendices

Appendix A: Sample queries and their results

Query	Retrieved Documents (Top 5)
"ኢትዮጵያ ታሪክ"	1. ኢትዮጵያ ታሪክ መጽሐፍ 2. ኢትዮጵያ ታሪክ አጭር ማጠቃለያ 3. ኢትዮጵያ ታሪክ በአጭሩ 4. ኢትዮጵያ ታሪክ ለልጆች 5. ኢትዮጵያ ታሪክ ከጥንት ጀምሮ
"አዲስ አበባ መስህቦች"	1. አዲስ አበባ መስህቦች ጉብኝት 2. አዲስ አበባ መስህቦች ዝርዝር 3. አዲስ አበባ መስህቦች ፎቶዎች 4. አዲስ አበባ መስህቦች ለቱሪስቶች 5. አዲስ አበባ መስህቦች መመሪያ
"ኮሮና ቫይረስ ዜና"	1. ኮሮና ቫይረስ ዜና ዛሬ 2. ኮሮና ቫይረስ ዜና ኢትዮጵያ 3. ኮሮና ቫይረስ ዜና ዓለም 4. ኮሮና ቫይረስ ዜና በአማርኛ 5. ኮሮና ቫይረስ ዜና መረጃ

Appendix B: Code Snippets or Configuration Details

1. Stemming Algorithm (Simplified Example)

```
function stemWord(word) {
  // Convert word to consonant-vowel form
  let consonantVowelForm = convertToConsonantVowelForm(word);

  // Remove longest match suffixes
  while (consonantVowelForm.length > 3) {
    consonantVowelForm = removeLongestMatchSuffix(consonantVowelForm);
  }

  // Remove longest match prefixes
  while (consonantVowelForm.length > 3) {
    consonantVowelForm = removeLongestMatchPrefix(consonantVowelForm);
  }

  // Remove vowels
  if (consonantVowelForm.length > 2) {
```



```
    consonantVowelForm = removeVowels(consonantVowelForm);  
}  
  
return consonantVowelForm;  
}
```