

The Use of Structured Text Retrieval to Search for Scientific Publications

Betreuer:

Dipl.-Ing. Dr.techn. Roman Kern

Thomas Mauerhofer (1031957)

Graz, May 19, 2017

Contents

| | | |
|----------|------------------------------|----------|
| 1 | Introduction | 3 |
| 2 | Related Work | 3 |
| 3 | Structure | 4 |
| 4 | Outline | 4 |
| 5 | Selected Bibliography | 5 |
| 6 | Schedule | 6 |

1 Introduction

In modern society searching for information via the Internet and specifically via Google's search engine has become a regular part of day-to-day life. While Google registered 10,000 searches a day in 1998, the same number of queries was sent per second in 2006 [1]. This shows that searching for data becomes more and more important. On a daily basis, new websites are created, articles are written and scientific papers are published. In order to manage this amount of data and to distinguish between relevant and non-relevant sources various search engines are used. Therefore, searches should be simple and precise.

Looking at search engines for scientific publications and research there is a broad selection of possibilities. One of the best-known engines is Google Scholar. It uses a simple input interface and lists the results with regards to their relevance. These listings contain data in various formats from different years covering a variety of topics which do not necessarily reflect the search criteria. However, especially while working on a scientific paper it is vital to obtain precise results.

This paper is concerned with improving the quality of search queries and their results of different scientific publications in the Portable Document Format (PDF). A common problem while searching for a specific author, for example, is that the most popular search engines often do not only list all of his or her articles and books but also sources that cite these publications. The objective of this paper is to provide measures which enhance the stability of the search terms and optimize the usability in order to deliver more concise results. To do so, simple search query structures, an intuitive front end, spell checks and post-processing in the back end will be implemented.

2 Related Work

As stated in the **Introduction** this master's thesis is concerned with improving the quality of search queries and their results for scientific publications in PDF. According to [3] all academic publications consist of similar structures. These structures are divided into chapters, sections, subsections and so on. The Structured Text Retrieval Model, as proposed in [5], can be applied to these documents perfectly. The model describes the handling of search queries and their corresponding results as well as the post-processing in the back end using structural meta information. In order to sort the documents with regards to their relevance, various ranking strategies such as the Jelenik-Mercer smoothing are used. Through extensions such as the contextualization strategy and aggregation strategy the ranking will be refined further.

To ensure simple and stable search query structures the syntax will be constructed as described in [2]. Thereby the queries consist of multiple terms which are structured according to *label:keyword* format. Furthermore, the terms can be prefixed with *+* to increase their priority. Due to this structure, the user interface will be extended with an advanced search in order to improve usability and therefore create an intuitive front end.

The final usability improvement is achieved through spell checks. These can be implemented as described in [6], where the Levenshtein Distance is used to verify that a word is spelled correctly.

3 Structure

In the first stage the fundamental structure for the project will be implemented. The micro framework Flask requires the front end to be coded in JavaScript and the back end to be in Python. Via simple requests inputs from the user side queries are sent to the server side. There the query is processed and a response is generated. These components form a solid ground work for the Text Retrieval Model.

Python provides a native library to connect to MySQL databases. In the next Stage it will be used to generate a database as proposed in [7]. To fill up the database with data a Util to classify the document structure, as described in [3], is used. Then a ranking system, according to [4, 5] can be implemented.

Thereafter spell checking of the search terms, as proposed in [6], is applied in order to enhance the usability.

Finally, an advanced search will be added that allows the user to formulate searches without prior knowledge regarding the syntax of the queries, which will result in an intuitive and simple user interface.

4 Outline

1. Introduction
2. Related Work
 - a) Unsupervised Document Structure Analysis
 - b) Structured Text Retrieval
 - c) Structured-Text Retrieval System with an Object-Oriented Database System
 - d) Search Queries
 - e) Real-Word Error Detection and Correction
3. Implementation
 - a) Fundamental Structure
 - b) Setup of the Database
 - c) User Interface
 - d) Satisfaction of Search Queries
 - e) Ranking system
 - f) Spell Checks

4. Results

5. Conclusion

5 Selected Bibliography

References

- [1] Google search statistics. <http://www.internetlivestats.com/google-search-statistics/>. Accessed: 2017-05-15.
- [2] Sara Cohen, Jonathan Mamou, Yaron Kanza, and Yehoshua Sagiv. Xsearch: A semantic search engine for xml. In Johann Christoph Freytag, Peter C. Lockemann, Serge Abiteboul, Michael J. Carey, Patricia G. Selinger, and Andreas Heuer, editors, *VLDB*, pages 45–56. Morgan Kaufmann, 2003.
- [3] Stefan Klampfl, Michael Granitzer, Kris Jack, and Roman Kern. Unsupervised document structure analysis of digital scientific articles. *Int. J. on Digital Libraries*, 14(3-4):83–99, 2014.
- [4] Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schütze. *Introduction to Information Retrieval*. Cambridge University Press, 2008.
- [5] Berthier Ribeiro-Neto and Ricardo Baeza-Yates. *Modern Information Retrieval*. ACM Press / Addison-Wesley, 1999.
- [6] Pratip Samanta and Bidyut B. Chaudhuri. A simple real-word error detection and correction using local word bigram and trigram. In *ROCLING*. Association for Computational Linguistics and Chinese Language Processing (ACLCLP), Taiwan, 2013.
- [7] Tak W. Yan and Jurgen Annevelink. Integrating a structured-text retrieval system with an object-oriented database system. In Jorge B. Bocca, Matthias Jarke, and Carlo Zaniolo, editors, *VLDB*, pages 740–749. Morgan Kaufmann, 1994.

6 Schedule

| Step | Task in Detail | Deadline |
|-------------------------------|---------------------------------------------|------------|
| Implementation Process | Creation of the Fundamental Structure | 2017-05-31 |
| | Database | 2017-06-30 |
| | Import of PDFs via PDF-Extractor | 2017-06-30 |
| | Satisfaction of Search Queries | 2017-07-31 |
| | Simple Ranking System | 2017-07-31 |
| | Improved Ranking System | 2017-08-31 |
| | Implement Spell Checks | 2017-09-30 |
| | Improve User Interface | 2017-10-31 |
| Writing Process | Introduction | 2017-10-31 |
| | Related Work | 2017-11-30 |
| | Unsupervised Document Structure Analysis | 2017-11-30 |
| | Structured Text Retrieval | 2017-11-30 |
| | STR with an Object-Oriented Database System | 2017-12-31 |
| | Search Queries | 2017-12-31 |
| | Spelling Checks | 2018-01-31 |
| | Implementation | 2018-01-31 |
| | Fundamental Structure | 2018-02-28 |
| | Setup of the Database | 2018-02-28 |
| | User Interface | 2018-03-31 |
| | Satisfaction of Search Queries | 2018-03-31 |
| | Ranking System | 2018-03-31 |
| | Spell Checks | 2018-04-30 |
| | Results | 2018-04-30 |
| | Conclusion | 2018-04-30 |
| | Correction | 2018-05-31 |