# Chapter 3.4
# Full–Text Search Engines for Databases

**László Kovács**
*University of Miskolc, Hungary*

**Domonkos Tikk**
*Budapest University of Technology and Economics, Hungary*

## INTRODUCTION

Current databases are able to store several Tbytes of free-text documents. The main purpose of a database from the user's viewpoint is the efficient information retrieval. In the case of textual data, information retrieval mostly concerns the selection and the ranking of documents. The selection criteria can contain elements that apply to the content or the grammar of the language. In the traditional database management systems (DBMS), text manipulation is restricted to the usual string manipulation facilities, i.e. the exact matching of substrings. Although the new SQL1999 standard enables the usage of more powerful regular expressions, this traditional approach has some major drawbacks. The traditional string-level operations are very costly for large documents as they work without task-oriented index structures.

The required full-text management operations belong to text mining, an interdisciplinary field of natural language processing and data mining. As the traditional DBMS engine is inefficient for these operations, database management systems are usually extended with a special full-text search (FTS) engine module. We present here the particular solution of Oracle; there for making the full-text querying more efficient, a special engine was developed that performs the preparation of full-text queries and provides a set of language and semantic specific query operators.

## BACKGROUND

Traditional DBMS engines are not adequate to meet the users' requirements on the management of free-text data as they handles the whole text field as an atom (Codd, 1985). A special extension to the DBMS engine is needed for the efficient implementation of text manipulating operations. There is a significant demand on the market on the

usage of free text and text mining operations, since information is often stored as free text. Typical application areas are, e.g., text analysis in medical systems, analysis of customer feedbacks, and bibliographic databases. In these cases, a simple character-level string matching would retrieve only a fraction of related documents, thus an FST engine is required that can identify the semantic similarities between terms.

There are several alternatives for implementing an FTS engine. In some DBMS products, such as Oracle, Microsoft SQLServer, Postgres, and mySQL, a built-in FTS engine module is implemented. Some other DBMS vendors extended the DBMS configuration with a DBMS-independent FTS engine. In this segment the main vendors are: SPSS LexiQuest (SPSS, 2007), SAS Text Miner (SAS, 2007), dtSearch (dtSearch, 2007), and Statistica Text Miner (Statsoft, 2007).

The market of FTS engines is very promising since the amount of textual information stored in databases rises steadily. According to the study of Meryll Lynch (Blumberg & Arte, 2003), 85% of business information are text documents – e-mails, business and research reports, memos, presentations, advertisements, news, etc. – and their proportion still increases. In 2006, there were more than 20 billion documents available on the Internet (Chang, 2006). The estimated size of the pool increases to 550 billion documents when the documents of the hidden (or deep) web

– which are e.g. dynamically generated ones – are also considered.
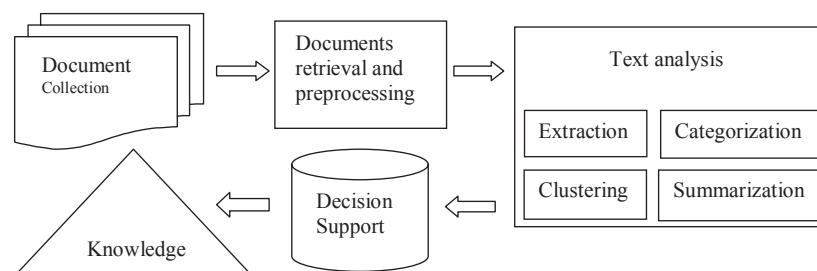
## TEXT MINING

The subfield of document management that aims at processing, searching, and analyzing text documents is *text mining*. The goal of text mining is to discover the non-trivial or hidden characteristics of individual documents or document collections. Text mining is an application oriented interdisciplinary field of machine learning which exploits tools and resources from computational linguistics, natural language processing, information retrieval, and data mining.

The general application schema of text mining is depicted in Figure 1 (Fan, Wallace, Rich & Zhang, 2006). For giving a brief summary of text mining, four main areas are presented here: information extraction, text categorization/classification, document clustering, and summarization.

## Information Extraction

The goal of information extraction (IE) is to collect the text fragments (facts, places, people, etc.) from documents relevant to the given application. The extracted information can be stored in structured databases. IE is typically applied in such

*Figure 1. The text mining module*

7 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage:
www.igi-global.com/chapter/full-text-search-engines-databases/7950

# Related Content

### Self-Tuning Database Management Systems
Camilo Porto Nunes, Cláudio de Souza Baptistaand Marcus Costa Sampaio (2009). *Handbook of Research on Innovations in Database Technologies and Applications: Current and Future Trends  (pp. 753-761).*
www.irma-international.org/chapter/self-tuning-database-management-systems/20761

### Integrity Constraints Checking in a Distributed Database
Hamidah Ibrahim (2010). *Soft Computing Applications for Database Technologies: Techniques and Issues (pp. 153-169).*
www.irma-international.org/chapter/integrity-constraints-checking-distributed-database/44387

### Representing Classes of Things and Properties in General in Conceptual Modelling: An Empirical Evaluation
Graeme Shanks, Daniel Moody, Jasmina Nuredini, Daniel Tobinand Ron Weber (2010). *Journal of Database Management (pp. 1-25).*
www.irma-international.org/article/representing-classes-things-properties-general/42083

### XML Integration and Toolkit for B2B Applications
Christophe Nicolle, Kokou Yetongnonand Jean-Claude Simon (2003). *Journal of Database Management (pp. 33-58).*
www.irma-international.org/article/xml-integration-toolkit-b2b-applications/3302

### Assuring Database Integrity
David Stemple, Eric Simon, Subhasish Mazumdarand Matthias Jarke (1990). *Journal of Database Administration (pp. 12-27).*
www.irma-international.org/article/assuring-database-integrity/51075