Chapter 4

# Methodologies and Technologies to Retrieve Information From Text Sources

**Anu Singha**
*South Asian University, India*

**Phub Namgay**
*Sherubtse College, Royal University of Bhutan, Bhutan*

## ABSTRACT

*A tool which algorithmically traces the effectiveness of the text files would be helpful in determining whether the text file have all the characteristic of important concepts. Every text source is build up on key phrases, and these paramount phrases follow a certain grammatical linguistic pattern widely used. An enormous amount of information can be derived from these key concepts for the further analysis such as their dispersion, relationship among the concepts etc. The relationship among the key concepts can be used to draw a concept graphs. So, this chapter presents a detailed methodologies and technologies which evaluate the effectiveness of the extracted information from text files.*

## INTRODUCTION

Before the advent of internet, text files formed the only source to get the content knowledge. Text files are read as a part of academics for a university student or for leisure, who fancy reading. With the rise of modern technology, paper based text files is being replaced by electronic version and is easily available online. Though

the readers embrace comfort and portability provided by electronic version of the text files, paper based text files are still popular. Unlike decades back, text files in any field of study is readily available now. With the demand for text files rising, the qualities of the text files are compromised. With the rate reading habits on decline, impatient readers hardly spent some time to evaluate the content of the file. A technique to analyze the effectiveness by evaluating the text files in detail will be of great importance. This ultimately would help the readers in finding the text file most suitable to his / her needs and comprehension capability.

A tool which algorithmically traces the effectiveness the text files would be helpful in determining whether the text file have all the characteristic of a good source. Every file is build up on key concepts, and these key concepts form the foundation of a good source. The text sources that contain concepts that share some common properties and semantically related are more lucid and intelligible than those text sources which contain many unrelated concepts. These paramount phrases follow a certain grammatical linguistic pattern widely used. An enormous amount of information can be derived from these key concepts for the further analysis such as their dispersion across the file, relationship among the concepts, etc. Such analysis will help in better assessment of the text file. The relationship among the key concepts can be used to draw a concept graphs. Since we live in an increasingly visual society, pictorial representation of the key concepts as a graph would help the readers in easily judging the text source and their content.

The goal of this chapter is to confer the methodologies on the key concepts retrieval from the text files. The authors investigate the techniques for examining the key concepts in the text files. This chapter also presents some of the different tools used in natural language processing. Their uses and implementation methods are explicitly discussed. The key concepts in our context correspond to the terminological noun phrases. These extracted set of phrases can be further analyzed to check the credibility of the text file in conveying the required set of information to the readers. It is based on the intuition that a source which contains right set of related key concepts is more beneficial and comprehensible. The set of key concept which form the cornerstone of a text file can be further used to draw concept graphs. The noun phrases from the candidate set of extracted phrases form the nodes of the graph and the relationship that exists between the nodes can be denoted by a link between the concept pairs. The 'in' degree and 'out' degree of each vertex of the graph i.e., the noun phrases can be used to determine the most important key concepts. Such representation of the source in a visual form helps readers in easily judging and grasping the key concepts. This ultimately serves as a preface to the text files and reduces the cognitive burden of the readers.

23 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/methodologies-and-technologies-to-retrieve-information-from-text-sources/185980

## Related Content

Social Media Mining: A New Framework and Literature Review
Vipul Guptaand Mayank Gupta (2016). *Big Data: Concepts, Methodologies, Tools, and Applications (pp. 2401-2414).*
www.irma-international.org/chapter/social-media-mining/150271

Finding Non-Coincidental Sporadic Rules Using Apriori-Inverse
Yun Sing Koh, Nathan Rountreeand Richard O'Keefe (2006). *International Journal of Data Warehousing and Mining (pp. 38-54).*
www.irma-international.org/article/finding-non-coincidental-sporadic-rules/1765

A Comprehensive Workflow for Enhancing Business Bankruptcy Prediction
Rui Sarmento, Luís Trigoand Liliana Fonseca (2015). *Integration of Data Mining in Business Intelligence Systems (pp. 216-238).*
www.irma-international.org/chapter/a-comprehensive-workflow-for-enhancing-business-bankruptcy-prediction/116817

Discovering Higher Level Correlations from XML Data
Luca Cagliero, Tania Cerquitelliand Paolo Garza (2012). *XML Data Mining: Models, Methods, and Applications (pp. 288-315).*
www.irma-international.org/chapter/discovering-higher-level-correlations-xml/60914

Adding Electric Vehicle Modeling Capability to an Agent-Based Transport Simulation
Rashid A. Waraich, Gil Georges, Matthias D. Galusand Kay W. Axhausen (2014). *Data Science and Simulation in Transportation Research (pp. 282-318).*
www.irma-international.org/chapter/adding-electric-vehicle-modeling-capability-to-an-agent-based-transport-simulation/90076