

Chapter 8

System Analysis and Design for Document Classification

ABSTRACT

The text-mining process starts with a keyword search in text collections. Current text processing technology allows a search technique beyond simple Boolean searches by using natural language queries. Since search engines can recognize any of thousands of keywords and phrases but not the concepts behind the text, it is necessary for researchers to construct an automatic keyword extractor to generate the “Keyword List” for each document. Later, this list can act as the knowledge base to associate unorganized documents to meaningful classes. Failures in identifying the keywords for a certain concept will result in missing values or data for that specific concept.

FACT FINDING

Technology is nearly at a point that Text Classification can apply automated **Text Mining** approaches to develop strategic information. Much of the relevant information is contained on Textual documents and is freely available if one can gain access to it. **Text Mining** applications are expensive and relatively crude, but as interest grows in this, prices will diminish and functionality will improve. Domain intelligence plays an important role when **Text Mining** tools make strategic and operational decisions.

One particular focus for IT has been on using Data Mining techniques to extract meaningful patterns and build predictive customer relationship

DOI: 10.4018/978-1-7998-3772-5.ch008

models from textual data. Although widely used, data mining is currently widely available only to structured, numeric databases. However, a majority of business information exists in the form of unstructured or semi structured text documents or in Web based data sources. The traditional way of processing text information involves human actions in information gathering, analysis, and dissemination. This requires substantial investment of money, time, and human resources.

Moreover, it is difficult to combine qualitative text data with quantitative numeric data in business analyses. Therefore, there is a pressing need to develop a method that can accurately extract business intelligence from large text collections and integrate the fragmented information into business intelligence databases.

The text-mining process starts with a keyword search in text collections. Current text processing technology allows a search technique beyond simple Boolean searches by using natural language queries. Since search engines can recognize any of thousands of keywords and phrases but not the concepts behind the text, it is necessary for researchers to construct an automatic keyword extractor to generate the “Keyword List” for each Document. Later, this list can act as the knowledge base to associate unorganized Documents to meaningful Classes. Failures in identifying the keywords for a certain concept will result in missing values or data for that specific concept.

6 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/system-analysis-and-design-for-document-classification/268467

Related Content

Dimensionality Reduction with Unsupervised Feature Selection and Applying Non-Euclidean Norms for Classification Accuracy

Amit Saxena and John Wang (2010). *International Journal of Data Warehousing and Mining* (pp. 22-40).

www.irma-international.org/article/dimensionality-reduction-unsupervised-feature-selection/42150

Reuse of Excess Research Data for New Researches

Amiram Porath (2013). *Ethical Data Mining Applications for Socio-Economic Development* (pp. 40-52).

www.irma-international.org/chapter/reuse-excess-research-data-new/76256

Analyzing the Impact of e-WOM Text on Overall Hotel Performances: A Text Analytics Approach

Aakash Aakash, Anu G. Aggarwal and Sanchita Aggarwal (2022). *Research Anthology on Implementing Sentiment Analysis Across Multiple Disciplines* (pp. 1805-1830).

www.irma-international.org/chapter/analyzing-the-impact-of-e-wom-text-on-overall-hotel-performances/308576

A Hybrid Method for High-Utility Itemsets Mining in Large High-Dimensional Data

Guangzhu Yu, Shihuang Shao, Bin Luo and Xianhui Zeng (2009). *International Journal of Data Warehousing and Mining* (pp. 57-73).

www.irma-international.org/article/hybrid-method-high-utility-itemsets/1823

Efficient Top-k Keyword Search Over Multidimensional Databases

Ziqiang Yu, Xiaohui Yu and Yang Liu (2013). *International Journal of Data Warehousing and Mining* (pp. 1-21).

www.irma-international.org/article/efficient-top-keyword-search-over/78373