



Chapter XII

Mining Free Text for Structure

Vladimir A. Kulyukin
Utah State University, USA

Robin Burke
DePaul University, USA

ABSTRACT

Knowledge of the structural organization of information in documents can be of significant assistance to information systems that use documents as their knowledge bases. In particular, such knowledge is of use to information retrieval systems that retrieve documents in response to user queries. This chapter presents an approach to mining free-text documents for structure that is qualitative in nature. It complements the statistical and machine-learning approaches, inasmuch as the structural organization of information in documents is discovered through mining free text for content markers left behind by document writers. The ultimate objective is to find scalable data mining (DM) solutions for free-text documents in exchange for modest knowledge-engineering requirements. The problem of mining free text for structure is addressed in the context of finding structural components of files of frequently asked questions (FAQs) associated with many USENET newsgroups. The chapter describes a system that mines FAQs for structural components. The chapter concludes with an outline of possible future trends in the structural mining of free text.

INTRODUCTION

When the manager of a mutual fund sits down to write an update of the fund's prospectus, he does not start his job from scratch. He knows what the fund's sharehold-

ers expect to see in the document and arranges the information accordingly. An inventor, ready to register his idea with the Patent and Trademark Office of the U.S. Department of Commerce, writes it up in accordance with the rules specifying the format of patent submissions. A researcher who wants to submit a paper to a scientific conference must be aware of the format specifications set up by the conference committee. Each of these examples suggests that domains of human activity that produce numerous documents are likely to have standards specifying how information must be presented in them.

Such standards, or presentation patterns, are a matter of economic necessity; documents whose visual structure reflects their logical organization are much easier to mine for information than unconstrained text. The ability to find the needed content in the document by taking advantage of its structural organization allows the readers to deal with large quantities of data efficiently. For example, when one needs to find out if a person's name is mentioned in a book, one does not have to read it from cover to cover; going to the index section is a more sensible solution.

Knowledge of the structural organization of information in documents¹ can be of significant assistance to information systems that use documents as their knowledge bases. In particular, such knowledge is of use to information retrieval systems (Salton & McGill, 1983) that retrieve documents in response to user queries. For example, an information retrieval system can match a query against the structural components of a document, e.g., sections of an article, and make a retrieval decision based on some combination of matches. More generally, knowledge of the structural organization of information in documents makes it easier to mine those documents for information.

The advent of the World Wide Web and the Internet have resulted in the creation of millions of documents containing unstructured, structured, and semi-structured data. Consequently, research on the automated discovery of structural organization of information in documents has come to the forefront of both information retrieval and natural language processing (Freitag, 1998; Hammer, Garcia-Molina, Cho, Aranha, & Crespo, 1997; Hsu & Chang, 1999; Jacquemin & Bush, 2000; Kushmerick, Weld, & Doorenbos, 1997). Most researchers adhere to numerical approaches of machine learning and information retrieval. Information retrieval approaches view texts as sets of terms, each of which exhibits some form of frequency distribution. By tracking the frequency distributions of terms, one can attempt to partition the document into smaller chunks, thus claiming to have discovered a structural organization of information in a given document. Machine-learning approaches view texts as objects with features whose combinations can be automatically learned by inductive methods.

Powerful as they are, these approaches to mining documents for structure have two major drawbacks. First, statistical computations are based on the idea of statistical significance (Moore & McCabe, 1993). Achieving statistical significance requires large quantities of data. The same is true for machine-learning approaches that require large training sets to reliably learn needed regularities. Since many documents are small in size, the reliable discovery of their structural components using numerical methods alone is problematic. Second, numerical approaches ignore the fact that document writers leave explicit markers of content structure in document texts. The presence of these markers in document texts helps the reader digest the information contained in the document. If these markers are ignored, document texts become much harder to navigate and understand.

21 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/mining-free-text-structure/7605

Related Content

A Porter Framework for Understanding the Strategic Potential of Data Mining for the Australian Banking Industry

Kate A. Smith and Mark S. Dale (2004). *Organizational Data Mining: Leveraging Enterprise Data Resources for Optimal Performance* (pp. 25-45).

www.irma-international.org/chapter/porter-framework-understanding-strategic-potential/27906

A Method of Sanitizing Privacy-Sensitive Sequence Pattern Networks Mined From Trajectories Released

Haitao Zhang and Yunhong Zhu (2019). *International Journal of Data Warehousing and Mining* (pp. 63-89).

www.irma-international.org/article/a-method-of-sanitizing-privacy-sensitive-sequence-pattern-networks-mined-from-trajectories-released/228938

Uncertainty-Based Clustering Algorithms for Large Data Sets

B. K. Tripathy, Hari Seetha and M. N. Murty (2018). *Modern Technologies for Big Data Classification and Clustering* (pp. 1-33).

www.irma-international.org/chapter/uncertainty-based-clustering-algorithms-for-large-data-sets/185977

Multidimensional Model Design using Data Mining: A Rapid Prototyping Methodology

Sandro Bimonte, Lucile Sautot, Ludovic Journaux and Bruno Faivre (2017). *International Journal of Data Warehousing and Mining* (pp. 1-35).

www.irma-international.org/article/multidimensional-model-design-using-data-mining/173704

Analysis on the Allocation of Control Right of Intergenerational Inheritance of Family Enterprises in the New Era

Zhanzhong Wang and Jiajun Li (2023). *International Journal of Data Warehousing and Mining* (pp. 1-20).

www.irma-international.org/article/analysis-on-the-allocation-of-control-right-of-intergenerational-inheritance-of-family-enterprises-in-the-new-era/319966