

An Efficient and Effective Index Structure for Query Evaluation in Search Engines

W

Yangjun Chen

University of Winnipeg, Canada

INTRODUCTION

Indexing the Web for fast keyword search is among the most challenging applications for scalable data management. In the past several decades, different indexing methods have been developed to speed up text search, such as inverted files, signature files and signature trees for indexing texts (Anh and Moffat, 2005; Chen et al., 2004; Chen et al. 2006; Faloutsos, 1985; Faloutsos et al., 1988); and suffix trees and tries (Knuth, 1975) for string matching. Especially, different variants of inverted files have been used by the Web search engines to find pages satisfying a query (Arasu, 2001; Lemple et al., 2003).

A text database can be roughly viewed as a collection of documents and each document is stored as a list of words. Over the documents, there are two kinds of Boolean queries, that is, queries that can be constructed from query terms by conjunction (\wedge) or disjunction (\vee). A document D is an answer to a conjunctive query $w_1 \wedge w_2 \wedge \dots \wedge w_k$ if it contains every w_i for $1 \leq i \leq k$ while D is an answer to a disjunctive query $w_1 \vee w_2 \vee \dots \vee w_l$ if it contains any w_i for $1 \leq i \leq l$. Conjunction and disjunction can be nested to arbitrary depth, but can always be transformed to a conjunctive normal form:

$$(w_{11} \dots \vee w_{1l_1} \dots) \wedge (w_{k1} \dots \vee w_{kl_k} \dots)$$

In this chapter, we discuss a new method to evaluate both conjunctive and disjunctive queries by decomposing an inverted list into a collection of disjoint sub-lists. The decomposition is based

on the construction of a trie structure T over documents and then associating each document word with an interval sequence generated by labeling T by using a kind of tree encoding.

With this method, we can improve the efficiency of traditional methods by an order of magnitude or more.

BACKGROUND

In order to efficiently evaluate such queries, indexes need to be established. It is well known that English texts typically contain many different variants of basic words, by using variant word endings such as ‘ing’, ‘ed’, ‘ses’, and ‘ation’. All the variants of a word should be regarded as a match and therefore it is efficient for an index only include these basic words, or say, stems. Different algorithms have been developed to extract stems from documents. Among them, the algorithm proposed by Lovins (1968) is widely used.

By the signature file, a word is hashed to a bit string (called a signature) and all the words’ signatures of a document are superimposed (bit-wise OR operation) into a document signature. When a query arrives, its signature will be created using the same hash function and the document signatures are scanned and many nonqualifying documents are discarded. The rest are either checked (so that the ‘false drops’ are removed) or they are returned to the user as they are. The main disadvantage of this method is the false drop (Kitagawa et al., 1997), which needs extra time to check. The signature file is greatly improved by the so-called signature tree (Chen et al. 2006),

by which a set of signatures is organized into a binary tree structure and a sequential search of signatures is replaced with a search of binary trees. However, signature-based methods can be used only for evaluating conjunctive queries. For disjunctive queries, they are not efficient.

As pointed out by many researchers (Anh et al., 2005; Ao et al., 2011; Zobel et al., 2006), the inverted file is a more competitive indexing method than signature-based approaches. It is extensively used by different web search engines due to its efficiency and simplicity. Structurally, it contains two parts: a search structure or vocabulary, containing all the distinct words to be indexed, and a set of inverted lists with each constructed for a distinct word w , storing the identifiers of all those documents containing w . Queries are evaluated by fetching the inverted lists for the query terms, and then intersecting them for conjunctive queries, or merging them (by a set union operation) for disjunctive queries. According to (Zobel et al., 1998), the inverted file is superior to the signature file in almost every respect, including functionality, query time, and space overhead.

Since it was first proposed in mid-1960s, the inverted file has been adopted in information retrieval, database systems, distributed systems (Büttcher et al., 2005; Camel et al., 2001), and different search engines. Also, much effort has been spent on the improvement of its performance by using integer coding (Golomb, 1966), bitmap compression (Apaydin et al., 2006; Bjørklund et al., 2009), caching (Saraiva et al., 2001), and parallelism (Ao et al., 2011). For a static environment, the bitmap compression is most efficient. However, it is obviously not suitable for a dynamical environment like the Web. So, different set intersection algorithms have been proposed to directly manipulate sorted arrays (Barbay et al., 2009), with caching (Barbay et al., 2009; Saraiva et al., 2001), parallelism (Ao et al., 2011), interpolation (Demaine et al., 2004), and even hardware characteristics (Tsirogiannis et al., 2009) being utilized to enhance performance. The method discussed in (Ding et al., 2011) is an in-

memory algorithm for set intersection, by which an inverted list is partitioned into a collection of sublists. Then, generate a hash image for each of them, by which a number in an inverted list is mapped to a single bit in an image. In this way, a set intersection is done by a series of hash value intersections.

INDEX STRUCTURES

In this section, we mainly discuss our index structure, by which each word with high frequency will be assigned an interval sequence. We will then associate intervals, instead of words, with inverted sub-lists. To clarify this mechanism, we will first discuss interval sequences for words in Subsection 1. Then, in Subsection 2, how to associate inverted lists with intervals will be addressed.

1. Interval Sequences Assigned to Words

Let $D = \{D_1, \dots, D_n\}$ be a set of documents. Let $W_i = \{w_{i1}, \dots, w_{in_i}\} (i = 1, \dots, n)$ be all of the words appearing in D_i , to be indexed. Denote $W = \bigcup_{i=1}^n W_i$, called the vocabulary. We define the word appearance frequency by the following formula:

$$f(w) = \frac{\text{num. of documents containing } w}{\text{num. of documents}}, (w \in W).$$

We then define a frequency threshold ζ . For any word w with $f(w) < \zeta$, we will associate it with an inverted list in a normal way, denoted as $\delta(w)$, exactly as in the method of inverted files. However, for all those with $f(w) \geq \zeta$, we will create a new index. For this, we will represent each D_i as a sequence containing all those words w with $f(w) \geq \zeta$, decreasingly sorted by $f(w)$. That is, in such a sequence, a word w precedes another w' if w is more frequent than w' in all documents.

9 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/an-efficient-and-effective-index-structure-for-query-evaluation-in-search-engines/184495

Related Content

Gene Expression Analysis based on Ant Colony Optimisation Classification

Gerald Schaefer (2016). *International Journal of Rough Sets and Data Analysis* (pp. 51-59).

www.irma-international.org/article/gene-expression-analysis-based-on-ant-colony-optimisation-classification/156478

Information-As-System in Information Systems: A Systems Thinking Perspective

Tuan M. Nguyen and Huy V. Vo (2008). *International Journal of Information Technologies and Systems Approach* (pp. 1-19).

www.irma-international.org/article/information-system-information-systems/2536

Piracy and Intellectual Property Theft in the Internet Era

Shun-Yung Kevin Wang and Jeremy J. McDaniel (2018). *Encyclopedia of Information Science and Technology, Fourth Edition* (pp. 1656-1666).

www.irma-international.org/chapter/piracy-and-intellectual-property-theft-in-the-internet-era/183881

Design of an Integrated Project Management Information System for Large Scale Public Projects: Iranian Case Study

Mona Taghavi, Ahmed Patel and Hamed Taghavi (2013). *Interdisciplinary Advances in Information Technology Research* (pp. 150-164).

www.irma-international.org/chapter/design-integrated-project-management-information/74539

How Mobile Technologies Are Leading to Economic Development in Sub-Saharan Africa

Nigel McKelvey, Adam Crossan and Kevin Curran (2021). *Encyclopedia of Information Science and Technology, Fifth Edition* (pp. 1719-1726).

www.irma-international.org/chapter/how-mobile-technologies-are-leading-to-economic-development-in-sub-saharan-africa/260300