

Efficient Keyword-Based Search for Top-K Cells in Text Cube

Bolin Ding, Bo Zhao, Cindy Xide Lin, Jiawei Han, *Fellow, IEEE*, Chengxiang Zhai, Ashok Srivastava, *Senior Member, IEEE*, Nikunj C. Oza

Abstract—Previous studies on supporting free-form keyword queries over RDBMSs provide users with linked-structures (e.g., a set of joined tuples) that are relevant to a given keyword query. Most of them focus on ranking individual tuples from one table or joins of multiple tables containing a set of keywords.

In this paper, we study the problem of keyword search in a data cube with text-rich dimension(s) (so-called *text cube*). The text cube is built on a multidimensional text database, where each row is associated with some text data (a document) and other structural dimensions (attributes). A cell in the text cube aggregates a set of documents with matching attribute values in a subset of dimensions. We define a keyword-based query language and an IR-style relevance model for scoring/ranking cells in the text cube. Given a keyword query, our goal is to find the top- k most relevant cells. We propose four approaches, *inverted-index one-scan*, *document sorted-scan*, *bottom-up dynamic programming*, and *search-space ordering*. The *search-space ordering* algorithm explores only a small portion of the text cube for finding the top- k answers, and enables early termination. Extensive experimental studies are conducted to verify the effectiveness and efficiency of the proposed approaches.

Index Terms—Keyword search, multidimensional text data, data cube.

1 INTRODUCTION

THE boom of Internet has given rise to an ever increasing amount of text data associated with multiple dimensions (attributes), for example, customer reviews in shopping websites (e.g., Amazon) are always associated with attributes like price, model, and rate. A traditional OLAP *data cube* can be naturally extended to summarize and navigate structured data together with unstructured text data. Such a cube model is called *text cube* [1]. A *cell* in the text cube aggregates a set of documents/tuples with matching attribute values in a subset of dimensions.

Keyword query, one of the most popular and easy-to-use ways to retrieve useful information from a collection of plain documents, is being extended to RDBMSs to retrieve information from text-rich attributes [2], [3], [4], [5], [6], [7], [8], [9], [10], [11], [12], [13], [14], [15], [16]. Given a set of keywords, existing methods aim to find relevant tuples or joins of tuples (e.g., linked by foreign keys) that contain all or some of the keywords.

In this paper, we study supporting keyword-based search in text cube. Unlike in plain documents nor RDBMSs, the goal of keyword search in text cube is to find the top- k most *relevant cells* for a given keyword query.

Example 1.1 (Motivation): Table 1 shows a small sample database of customer reviews about laptops. *Customer*

TABLE 1
Motivation Example

| Brand | Model | CPU | OS | Customer Review |
|-------|--------|--------|-------|---|
| Acer | AOA110 | 1.6GHz | Linux | light weight as little as 2.2 lb, fun and powerful computing features |
| Acer | AOA110 | 1.8GHz | XP | weight just over two pounds, with powerful Intel Atom Processor |
| Asus | EEE PC | 1.8GHz | XP | comes in pearl white style, images are sharp, disk is large |

Review is the text attribute, whereas *Brand*, *Model*, *CPU*, and *OS* are structural attributes.

Jim, a marketing analyst, wants to find which laptops are commented as *light weight* and *powerful performance*. He types a set of keywords: {"light", "weight", "powerful"}. Using IR techniques, the system can rank all the customer reviews and output the most relevant ones. However, when there are many customer reviews relevant to the query, Jim has to browse through a lot of reviews and summarize different opinions by himself. As multidimensional attributes are associated with each review, is it more desirable that a system provides users with "aggregated information", such as "*Acer AOA110 laptops are usually lightweight and have powerful performance*", than returning individual reviews? This is our intention to study such a new mechanism.

A cell in the text cube is in the form of (Brand = Acer, Model = AOA110, CPU = *, OS = *), which aggregates all the customer reviews (the first two in Table 1) for Acer AOA110 laptops. Another cell (Brand = *, Model = *, CPU = 1.8GHz, OS = XP) aggregates the second two reviews. It

- B. Ding, B. Zhao, C. Lin, J. Han, and C. Zhai are with the Department of Computer Science, University of Illinois at Urbana-Champaign, 201 North Goodwin Avenue, Urbana, IL 61801, USA.
E-mail: {bding3, bozhao3, xidelin2, hanj, czhai}@uiuc.edu
- A. Srivastava and N. C. Oza are with Intelligent Systems Division, NASA Ames Research Center, Moffett Field, CA 94035, USA.
E-mail: ashok@email.arc.nasa.gov, nikunj.c.oza@nasa.gov

can be seen that the first cell is more relevant to Jim’s query (Acer AOA110 laptops are lightweight and powerful) than the second. The goal of our system is to provide Jim with such aggregated information (i.e., the cell “Acer AOA110 laptops” that is relevant to his query), instead of a ranked list of individual customer reviews.

If Jim is interested in gaming computers, he can types {“fast”, “gaming”} into the system, and high-performance computer models with good customer reviewers are expected to be output, e.g., (Brand = *, Model = Xtreme XT, CPU = Intel Core i7, OS = *). □

Traditional IR techniques can be used to rank documents according to the relevance. In a large text database, however, the number of relevant documents to a query could be large, and a user has to spend much time reading them.

If a document is associated with attribute information, in a data cube model (a multidimensional space induced by the attributes), e.g., the text cube, a *cell* aggregates the documents with matching values in a subset of attributes. Such a collection of documents is associated with each cell, corresponding to an object, e.g., “Acer AOA110 laptops” or “1.8GHz XP laptops” (in the above example), that can be directly recommended to the user for the given query.

This paper studies the problem of *keyword-based top-k search in text cube*, i.e., *given a keyword query, find the top-k most relevant cells in a text cube*. When users want to retrieve information from a text cube using keyword queries, we believe that relevant cells, rather than relevant documents, are preferred as the answers, because: (i) relevant cells are easy for users to browse; and (ii) relevant cells provide users insights about the relationship between the values of relational attributes and the text data.

Overview of Model and Techniques. Following is an overview of our work, with important issues highlighted.

Relevance Score: Given a keyword query, the first question is how to compute the “relevance” of a cell in a text cube for ranking. Note that a cell corresponds to a collection of documents. We employ the *average model* studied in [17] to assign the relevance scores of the cells. Specifically, ANY IR scoring function (e.g., Okapi) can be used to compute the relevance score of each single document, and the *relevance score of a cell* (a document collection) is the (weighted) average of relevance scores of documents in this cell. An alternative is the *big document model*: documents in a cell are concatenated into a *big document*, and the relevance of the cell is the relevance of this big document. [17] shows the average model is better than the big document model.

Support Requirement: We allow a user to specify a *support threshold* minsup, and we only output the relevant cells with no less than minsup documents. Users may not prefer a cell with too few documents, since this cell is not “popular” (e.g., if there are too few reviews on a computer model, the user is unlikely to choose this one); users may feel more confident browsing a cell with larger number of documents (reviews). Also, when high-level information is preferred, minsup can be set to a relatively large value.

Challenges: The major computational challenge of this problem is the huge number of cells in a text cube. It

increases exponentially w.r.t. the dimensionality and is much larger than the number of documents in the database.

Efficient Algorithms: We design four approaches for the problem of *keyword-based top-k search in text cube*. The first one (*inverted-index one-scan*) scans the inverted index of each keyword in the query once to compute the relevance scores of all cells in the text cube. The second one (*document sorted-scan*) scans the documents in the order of relevance, and at the same time, updates the relevance scores of cells containing each document being visited; early stop for top-k is possible because upper bounds of relevance scores can be estimated. The third one (*bottom-up dynamic programming*) computes relevance scores of cells from bottom to top in a level-by-level manner. The fourth one (*search-space ordering*), based on the third one, explores cells in the order of relevance and prunes the search space by estimating relevance score upper bound in subspaces, so as to explore as few cells in the text cube as possible before outputting the top-k answers.

Contributions. In this paper, we propose and study the problem of *keyword-based top-k search in text cube* (or multidimensional text data): *finding the top-k cells relevant to a user-given keyword query*. Flexible keyword-based query language and relevance scoring formula of cells (aggregation of text data) are developed. We analyze the computational challenges and propose four approaches, *inverted-index one-scan*, *document sorted-scan*, *bottom-up dynamic programming*, and *search-space ordering*, to support the query language in text cube. We compare the four approaches and study their efficiency and effectiveness experimentally using both real and synthetic datasets.

Organization. Section 2 introduces the text cube model of multidimensional text data, defines the keyword search problem and keyword-based query language in text cube, and analyzes computational challenges. The four approaches for finding the top-k most relevant cells for a given keyword query are introduced in Sections 3 and 4. Experimental study is reported in Section 5, followed by discussion on how to extend our approach in Section 6 and related work in Section 7. Section 8 concludes this paper.

2 KEYWORD QUERIES IN TEXT CUBE

In this section, we introduce our data cube model and the keyword search problem. Preprocessing for our system and challenges for query processing are also discussed.

2.1 A Data Cube Model for Text Data

A set of documents \mathcal{D} is stored in an n -dimensional database $DB = (\mathcal{A}_1, \mathcal{A}_2, \dots, \mathcal{A}_n, \mathcal{D})$. Each row of DB is in the form of $r = (a_1, a_2, \dots, a_n, d)$: let $r[\mathcal{A}_i] = a_i \in \mathcal{A}_i$ be the *value* of *attribute* (or *dimension*) \mathcal{A}_i , and $r[\mathcal{D}] = d$ be the *document* in this row.

A document d is a multi(sub)set of the *term set* $\mathcal{W} = \{w_1, \dots, w_M\}$: a term w_i may appear multiple times in d .

The *data cube* model extended to the above multidimensional text database is called *text cube* [1]. Several important concepts are introduced as follows.

TABLE 2
A 4-Dimensional Text Database \mathcal{DB}

| Dimensions | | | | Text Data |
|--------------|--------------|-------------|--------------|--------------------------------|
| M (Model) | P (Price) | T (Time) | S (Score) | d (Document) |
| m1 | p1 | t1 | s1 | $d_1 = \{w1, w1, w2, w2, w8\}$ |
| m1 | p2 | t1 | s2 | $d_2 = \{w2, w2, w5, w6, w7\}$ |
| m1 | p3 | t2 | s2 | $d_3 = \{w3, w4, w3, w6, w7\}$ |
| m2 | p1 | t2 | s2 | $d_4 = \{w1, w1, w1, w1, w4\}$ |
| m2 | p2 | t2 | s2 | $d_5 = \{w5, w6, w7, w8, w8\}$ |
| m2 | p3 | t1 | s1 | $d_6 = \{w4, w5, w8, w8, w9\}$ |

TABLE 3
Some Cells in the Text Cube

| Cell | M | P | T | S | D |
|----------|----|----|----|----|------------------------------------|
| C_0 | * | * | * | * | $\{d_1, d_2, d_3, d_4, d_5, d_6\}$ |
| C_1 | m1 | * | * | * | $\{d_1, d_2, d_3\}$ |
| C_2 | m2 | * | * | * | $\{d_4, d_5, d_6\}$ |
| C_3 | m1 | * | t1 | * | $\{d_1, d_2\}$ |
| C_4 | m1 | * | t2 | * | $\{d_3\}$ |
| C_5 | m1 | * | * | s1 | $\{d_1\}$ |
| C_6 | m1 | * | * | s2 | $\{d_2, d_3\}$ |
| C_7 | m2 | * | t1 | * | $\{d_6\}$ |
| C_8 | m2 | * | t2 | * | $\{d_4, d_5\}$ |
| C_9 | m2 | * | t2 | s2 | $\{d_4, d_5\}$ |
| C_{10} | m1 | * | t1 | s1 | $\{d_1\}$ |
| C_{11} | m1 | p1 | t1 | s1 | $\{d_1\}$ |
| C_{12} | m1 | p2 | t1 | s2 | $\{d_2\}$ |
| C_{13} | * | p2 | * | s1 | \emptyset (empty cell) |
| C_{14} | m1 | * | t2 | s1 | \emptyset (empty cell) |

Cell and Aggregated Document. In the text cube built on \mathcal{DB} , a cell is in the form of $C = (v_1, v_2, \dots, v_n : D)$, where either $v_i \in \mathcal{A}_i$ (a value of dimension \mathcal{A}_i) or $v_i = *$ (the dimension \mathcal{A}_i is aggregated in C), and D (the aggregated document in C) is the set of documents in the rows of \mathcal{DB} having the same dimension values as C on all the non-* dimensions. Formally, for a cell $C = (v_1, v_2, \dots, v_n : D)$,

$$D = \{r[\mathcal{D}] \mid \text{for } r \in \mathcal{DB} \text{ s.t. } r[\mathcal{A}_i] = v_i \text{ if } v_i \neq *\}.$$

We use $C[\mathcal{A}_i]$ to denote the value v_i of dimension \mathcal{A}_i in the cell C , and $C[\mathcal{D}]$ to denote the aggregated document D of the cell C . Define the *support* of a cell C , denoted by $|C|$, to be the number of documents in $C[\mathcal{D}]$, i.e., $|C| = |C[\mathcal{D}]|$. A cell is said to be *empty* if $C[\mathcal{D}] = \emptyset$. For simplicity, a cell is also written as $C = (v_1, v_2, \dots, v_n)$.

All the cells with the same set of non-* dimensions form a *cuboid*. A cuboid with m non-* dimensions is an m -dim cuboid. The n -dim cuboid (with no aggregated dimension) is called the *base cuboid*. Cells in an m -dim cuboid are m -dim cells, and cells in a base cuboid are *base cells*. Note that each base cell may aggregate multiple rows/documents, since different rows may have the same values of attributes.

Ancestor and Descendant. Cell C' is an *ancestor* of C (or C is a *descendant* of C') iff $\forall i : C'[\mathcal{A}_i] \neq * \Rightarrow C'[\mathcal{A}_i] = C[\mathcal{A}_i]$. Cell C is an ancestor (or descendant) of itself. We use $\text{ans}(C)$ to denote the set of ancestors of a cell C , and $\text{des}(C)$ to denote the set of descendants of a cell C' .

Parent and Child. *Parents* and *children* are immediate ancestors and descendants of a cell, respectively. Cell C' is a *parent* of C (or C is a *child* of C') iff (i) $C' \in \text{ans}(C)$, and (ii) C' is an i -dim cell while C is an $(i+1)$ -dim cell. Let $\text{par}(C)$ ($\text{chd}(C')$) denote the set of parents of a cell C (children of a cell C'). Note that the 0-dim cell (i.e., the cell in the form of $(*, *, \dots, * : D)$) has no parent.

A-Parent and A-Child. The *A-parent* of a cell C is C' (or, C is an *A-child* of C'), denoted by $\text{par}_{\mathcal{A}}(C) = C'$, iff (i) C' is a parent of C , (ii) $C'[\mathcal{A}] \neq *$, and (iii) $C'[\mathcal{A}] = *$. A cell may have more than one *A-child*, and we use $\text{chd}_{\mathcal{A}}(C')$ to denote the set of all *A-children* of C' .

It is well-known all the cells in a data cube (or text cube) form a lattice, according to the parent-child relationship.

Example 2.1 (Text Cube): Table 2 shows a text database \mathcal{DB} , with four dimensions, M, P, T, and S. Term set $\mathcal{W} = \{w1, w2, \dots, w8\}$. A total of six documents are stored.

Table 3 shows some cells in the text cube generated from \mathcal{DB} . C_0 has support 6 and C_1 has support 3, i.e., $|C_0| = 6$ and $|C_1| = 3$. C_0 is the 0-dim cell and C_{11}, C_{12} are base cells. $C_3, C_4, C_5, C_6, C_{10}, C_{11}$, and C_{12} are descendants of C_1 , while only C_3, C_4, C_5 and C_6 are children of C_1 . Note C_1 has some other descendants and children not listed in this table. A cell may have more than one parent, e.g., both C_3 and C_5 are parents of C_{10} .

Cell C_1 is the T-parent of C_3 and C_4 , i.e., $\text{par}_T(C_3) = \text{par}_T(C_4) = C_1$. Cells C_3 and C_4 are the T-children of C_1 , i.e., $\text{chd}_T(C_1) = \{C_3, C_4\}$. Cells C_5 and C_6 are the S-children of C_1 , i.e., $\text{chd}_S(C_1) = \{C_5, C_6\}$. And similarly, $\text{par}_S(C_5) = \text{par}_S(C_6) = C_1$. \square

Multiple Documents in a Row. In a database where each row contains more than one document, to generalize our model, we can create a linked list of documents and keep the total number of documents in each row. By aggregating this count into cells on higher levels, we can easily calculate the support of each cell. Note that this modification will not affect our algorithms introduced later.

2.2 Keyword Search Problem in Text Cube

In traditional data cubes, operations like *drill-down* and *roll-up* suffice for users to explore multidimensional data. In text cube, however, a large portion of data is text. Since keyword query is an effective way for users to explore text data, we propose the keyword search problem in text cube.

2.2.1 Keyword Search Problem

A *keyword query* is a set of terms, i.e., $q = \{t_1, t_2, \dots, t_{|q|}\} \subseteq \mathcal{W}$. Given a keyword query q and a minimum support requirement minsup , the goal is to *find the top-k cells C 's s.t. supports $|C| \geq \text{minsup}$ with the top-k highest relevance scores in the text cube of \mathcal{DB} .*

Note that a cell relevant to the query q may contain all or some of the terms $t_1, \dots, t_{|q|}$. The *relevance score* of a cell C w.r.t. the query q is defined as a function $\text{rel}(q, C[\mathcal{D}])$ of the aggregated document $C[\mathcal{D}]$ and the query q . For brevity, it is also written as $\text{rel}(q, C)$. Because the total number of

cells in the text cube could be huge and it is not possible for a user to browse all of them, we return the top- k cells in the non-increasing order of relevance scores, where k can be specified by a user.

Recall the support of a cell C is the number of documents aggregated in $C[\mathcal{D}]$. We allow users to specify the minimum support minsup , because they may be only interested in cells “popular” enough. Users feel more confident when browsing a cell containing a large number of documents (e.g., customer reviews) than a cell with only one or two.

2.2.2 Relevance Scoring Formula

To rank all the cells and find the top- k ones, we define the relevance scoring function $\text{rel}(q, C[\mathcal{D}])$ (or $\text{rel}(q, C)$ for brevity). Recall an aggregated document $C[\mathcal{D}] \subseteq \mathcal{D}$ is a set of documents aggregated in C . Here, we treat every document in $C[\mathcal{D}]$ equally when calculating relevance score of C w.r.t. a keyword query q . $\text{rel}(q, C[\mathcal{D}])$ is defined as the average of all the relevance scores of documents in $C[\mathcal{D}]$:

$$\text{rel}(q, C[\mathcal{D}]) = \frac{1}{|C[\mathcal{D}]|} \sum_{d \in C[\mathcal{D}]} s(q, d). \quad (1)$$

$s(q, d)$ is the relevance score of a document d w.r.t. q

$$s(q, d) = \sum_{t \in q} \ln \frac{N - \text{df}_t + 0.5}{\text{df}_t + 0.5} \frac{(k_1 + 1) \text{tf}_{t,d}}{k_1((1-b) + b \frac{\text{dl}_d}{\text{avdl}}) + \text{tf}_{t,d}} \frac{(k_3 + 1) \text{qtf}_{t,q}}{k_3 + \text{qtf}_{t,q}} \quad (2)$$

(Okapi weighting [18])

$$= \sum_{t \in q} w_{\text{idf}}(t) \cdot w_{\text{tf}}(t, d) \cdot w_{\text{qtf}}(t, q),$$

where $N = |\mathcal{DB}|$, $\text{tf}_{t,d}$ is the *term frequency* of term $t \in q$ in d (the number of times t appearing in d), df_t is the number of documents in \mathcal{DB} containing t , dl_d is the length of document d , avdl is the average length of documents in \mathcal{DB} , $\text{qtf}_{t,q}$ is the number of times t appearing in q , and, k_1, b, k_3 are constants.

Our algorithms introduced later can also handle other formulas $s(q, d)$ for scoring documents, say pivoted normalization weighting [19]. In general, our algorithms allow $s(q, d)$ to be in the form of $s(q, d) = s((\text{tf}_{t,d}), \Delta_q, \Delta_d)$, where Δ_q is the set of parameters about q , e.g., idf (inverted document frequency), and Δ_d is the set of parameters about d , e.g., document length (for normalization). This form covers most IR models for evaluating the relevance of a document w.r.t. a query.

[17] shows that the weighted version of (1) (document with higher weight in $C[\mathcal{D}]$ contributes more to the relevance score $\text{rel}(q, C)$) is effective for ranking document corporuses. In this paper, we use the unweighted version (i.e., (1)) for the ease of explanation. Our algorithms can be easily extended to the weighted version, if the weights of documents can be precomputed.

Note that, for query q , $\text{rel}(q, C)$ is the relevance score of a set (corpus) of documents (i.e., $C[\mathcal{D}]$), while $s(q, d)$ is the relevance score of one document. In the rest part of this paper, we call $\text{rel}(q, C)$ a *relevance scoring formula*, and $s(q, d)$ a *document scoring formula* to distinguish them.

2.2.3 Extended Form of Keyword Query

Users may want to retrieve answers from a certain part of the text cube, by specifying a subset dimensions of interests and/or values of some dimensions. Motivated by this, the simplest form of keyword queries q can be extended by adding *dimension-value constraints*.

In an n -dimensional text cube, an *extended keyword query* is in the form of $Q = (u_1, u_2, \dots, u_n : q)$, where $u_i \in \mathcal{A}_i \cup \{*, ?\}$. We also use $Q[\mathcal{A}_i]$ to denote u_i . $Q[\mathcal{A}_i] \in \mathcal{A}_i$ specifies the value of dimension \mathcal{A}_i in a cell C ; $Q[\mathcal{A}_i] = *$ means the dimension \mathcal{A}_i in a cell C must be aggregated; and $Q[\mathcal{A}_i] = ?$ (question mark) imposes no constraint on the dimension \mathcal{A}_i of a cell C . A cell C is said to be *feasible* w.r.t. the query Q iff

- (i) for dimension \mathcal{A}_i s.t. $Q[\mathcal{A}_i] = *$, $C[\mathcal{A}_i] = *$ (\mathcal{A}_i is aggregated in C);
- (ii) for \mathcal{A}_i s.t. $Q[\mathcal{A}_i] \in \mathcal{A}_i$, $C[\mathcal{A}_i] = Q[\mathcal{A}_i]$; and
- (iii) for \mathcal{A}_i s.t. $Q[\mathcal{A}_i] = ?$, no constraint on $C[\mathcal{A}_i]$.

Given an *extended keyword query* $Q = (u_1, u_2, \dots, u_n : q)$ and a minimum support minsup , our goal is to find the top- k *feasible* cells C 's s.t. $\text{support } |C| \geq \text{minsup}$ and *relevance scores* $\text{rel}(q, C[\mathcal{D}])$'s are the top- k highest.

Example 2.2 (Queries, Answers, and Relevance Scores): In the text cube of \mathcal{DB} in Table 2, consider a keyword query $q = \{w1, w2\}$ and $\text{minsup} = 2$.

Using Okapi (2) with $k_1 = k_3 = 1$ and $b = 0.5$: $s(q, d_1) = 1.6$, $s(q, d_2) = 0.8$, $s(q, d_4) = 0.9$. Base cells $(m1, p1, t1, s1, \{d_1\})$, $(m1, p2, t1, s2, \{d_2\})$, and $(m2, p1, t2, s2, \{d_4\})$ have relevance scores 1.6, 0.8, and 0.9, respectively. But they will not be output because their supports are all 1, less than the threshold minsup . Cell $(*, p1, *, *, \{d_1, d_4\})$ has relevance score $(1.6 + 0.9)/2 = 1.25$, and $(m1, *, t1, *, \{d_1, d_2\})$ has relevance score $(1.6 + 0.8)/2 = 1.2$. They are the top-2 cells in the output. So we may infer, when the price $P = p1$, documents in the text cube are highly relevant to the given query q .

An extended keyword query $Q = (?, *, ?, * : q)$ (i.e., no constraint on dimensions M and T, while dimensions P and S must be aggregated). C_0 - C_4 in Table 3 are feasible cells, but C_5 and C_6 are not, since $C_5[S] \neq *$ and $C_6[S] \neq *$.

Another query is $Q' = (?, *, ?, s2 : q)$, which imposes an additional constraint that a cell must value $s2$ on dimension S, i.e., $C[S] = s2$. Then, C_9 in Table 3 is a feasible cell ($C_9[S] = s2$); but C_{10} is not, for $C_{10}[S] = s1 \neq s2$. \square

2.3 Preprocessing

In the algorithms introduced later, some parameters are assumed to be precomputed to speedup the query processing.

First, the support of each cell, i.e., the number of documents aggregated in this cell, is precomputed and stored, since supports are query-independent. This only requires additional $O(1)$ space for each cell. Later we can use this to efficiently check whether the support of a cell is above the threshold, and whether a cell is *fully-aggregated*.

Second, note that (2) can be rewritten as a more general form: $s(q, d) = \sum_{t \in q} w_{\text{idf}}(t) \cdot w_{\text{tf}}(t, d) \cdot w_{\text{qtf}}(t, q)$. Both $w_{\text{idf}}(t)$ and $w_{\text{tf}}(t, d)$ are precomputed for all terms and

documents, and stored in inverted indexes. So as an online query q comes, $s(q, d)$ can be computed conveniently.

Note that the above two steps can be done at the same time, as documents are scanned one by one. All algorithms presented later can benefit from this precomputation.

2.4 Challenges of Query Processing

There are two challenges of this keyword search problem: First, the size of a text cube could be huge; and second, it is impractical to precompute relevance scores for all queries.

2.4.1 Size of Text Cube

There is an n -dimensional database $\mathcal{DB} = (\mathcal{A}_1, \mathcal{A}_2, \dots, \mathcal{A}_n, \mathcal{D})$ with N rows, s.t. the non-empty cells in the text cube of \mathcal{DB} is $\Omega(N \cdot 2^n)$ or $\Omega(\prod_{i=1}^n (|\mathcal{A}_i| + 1))$, where $|\mathcal{A}_i|$ is the number of different values in dimension \mathcal{A}_i .

Even when $\text{minsup} (> 0)$ is nontrivial, the number of cells to be considered (those with support $\geq \text{minsup}$) is still huge, since a data cube is “fat” in the middle. For example, suppose minsup is large enough s.t. only d -dim cells with $d \leq n/2$ have support $\geq \text{minsup}$, we may still need to consider at least $N \cdot \binom{n}{n/2} = \Omega(N \cdot 2^{n/2})$ cells to select the top- k relevant ones w.r.t. a keyword query.

So *only when* the number of dimensions is small (2 to 4), we can compute the relevance scores of all cells and then sort them to find the top- k cells efficiently.

2.4.2 Limits of Preprocessing

In text cube, relevance scores $\text{rel}(q, C)$ of cells have to be computed online because they are query-dependent, although $w_{\text{idf}}(t)$ and $w_{\text{tf}}(t, d)$ can be precomputed (as in Section 2.3). Also, the ranking of cells according to $\text{rel}(q, C)$ cannot be precomputed either, since for different queries, the rankings could be different. Online computation of $\text{rel}(q, C)$ and ranking of cells dominate the cost of processing a query; so, we aim to explore only a small portion of the text cube for the top- k answers.

3 DOCUMENT LINEAR-SCAN APPROACHES

Two approaches for the keyword-based top- k search problem based on linear scans of rows are introduced in Section 3.1 and 3.2. The first one scans all rows to compute the relevance scores of all cells, and then output the top- k cells. The second one scans the rows in the descending order of relevance scores; in this process, cells containing each document are updated with upper bounds of relevance scores estimated to enable the early stop for top- k . We first focus on simple keyword query $q = \{t_1, \dots, t_{|q|}\}$ with support requirement minsup and show how to handle extended keyword queries in Section 3.3.

3.1 Inverted-Index One-Scan Approach

We construct a *row-based inverted index* for each term. Recall $\text{tf}_{t,d}$ denotes the term frequency of t in the document d . From the *row-based inverted index* $\text{IV}(t)$ of term t , we can efficiently retrieve all the rows $r = (a_1, \dots, a_n, d)$'s in

\mathcal{DB} with term frequencies $\text{tf}_{t,d} > 0$. The size of row-based inverted indices for all terms in \mathcal{W} is bounded by the total length of all documents in \mathcal{DB} .

Given a query q against the text cube of \mathcal{DB} and minimum support minsup , the *inverted-index one-scan* approach is outlined in Algorithm 1. We first compute the score $s(q, d)$ of each document in a row that contains at least one term in q (since otherwise $s(q, d)$ is 0) using the inverted indexes (line 1). Then, for each row r , we can efficiently identify all the cells C 's in the text cube s.t. $r[\mathcal{D}] \in C[\mathcal{D}]$ and support $|C[\mathcal{D}]| \geq \text{minsup}$ (lines 2-3: note $|C[\mathcal{D}]|$ can be precomputed and materialized). These cells are the ancestors of the base cell corresponding to row r . We use $\text{ans}(r)$ to denote this set of cells, and note that $|\text{ans}(r)| \leq 2^n$ in an n -dimensional text cube. We update $\text{rel}(q, C)$ using $s(q, r[\mathcal{D}])$ (line 4). Finally, after all rows are scanned, each $\text{rel}(q, C)$ is correctly computed as the average of $s(q, r[\mathcal{D}])$'s with r contained in C (some rows with $s(q, r[\mathcal{D}]) = 0$ are not scanned, but we know the total number of documents in C , so the average still can be computed), and the top- k ones are output (line 5).

Algorithm 1 Inverted-Index One-Scan Algorithm

Input: keyword query q against \mathcal{DB} , parameter k , minsup

- 1: Compute relevance score $s(q, r[\mathcal{D}])$ for each row r ;
 - 2: **for each** row $r \in \mathcal{DB}$ with $s(q, r[\mathcal{D}]) > 0$ **do**
 - 3: **for each** cell $C \in \text{ans}(r)$ and $|C[\mathcal{D}]| \geq \text{minsup}$ **do**
 - 4: Update $\text{rel}(q, C)$ using $s(q, r[\mathcal{D}])$;
 - 5: Output the top- k relevant cells with support $\geq \text{minsup}$ (for cells with the same relevance scores, output them in the descending order of support).
-

Example 3.1: Consider a row $r = (m1, p1, t1, s1, d1)$ in the text database in Table 2(a). A cell C with $d_1 \in C[\mathcal{D}]$ must be one of the 2^4 ancestors of the base cell $(m1, p1, t1, s1)$, e.g., $(m1, p1, *, s1)$, $(m1, *, t1, *)$, and $(*, *, t1, *)$, which are considered in lines 3-4. \square

Theorem 1: Algorithm 1 uses $O(N \cdot |q| + 2^n \cdot N + k \cdot \log M)$ time and $O(n \cdot M)$ space, where $N = |\mathcal{DB}|$ is the total number of rows, n is the dimensionality, and M is the total number of non-empty cells.

Proof: Algorithm 1 scans the row-based inverted index of each query term only once. The size of each row-based inverted index is bounded by $N = |\mathcal{DB}|$. Also, the 2^n ancestors of a row r can be enumerated in linear time. So Algorithm 1 requires $O(N \cdot |q| + 2^n \cdot N)$ time to compute the relevance scores of all non-empty cells. After that, it constructs a (Fibonacci) heap of all non-empty cells (ordered by relevance) in linear time $O(M) \leq O(2^n \cdot N)$ [20] and outputs the top- k in $O(k \cdot \log M)$ time. Totally, we need $O(N \cdot |q| + 2^n \cdot N + k \cdot \log M)$ time.

Since there are at most $M \leq 2^n \cdot N$ nonempty cells, and for each cell we store its dimension values, relevance score, and support, we need totally $O(n \cdot M)$ space (besides the storage of inverted indexes and text cube). \square

Algorithm 2 Document Sorted-Scan Algorithm

Q: a priority queue of cells (sorted in the non-increasing order of relevance scores)

top(Q): the top (first) cell of a priority queue **Q**

Input: keyword query q against \mathcal{DB} , parameter k , minsup

```

1: Rows are scanned in the descending order of relevance:
    $r_1, r_2, \dots, r_N$  such that  $s(q, r_i[\mathcal{D}]) \geq s(q, r_{i+1}[\mathcal{D}])$ ;
2: for  $i = 1$  to  $N$  do
3:   for each cell  $C \in \text{ans}(r_i)$  and  $|C[\mathcal{D}]| \geq \text{minsup}$  do
4:     Aggregate  $r_i$  into  $C$ : compute  $\text{rel}(q, C)_i$  as
       the average of relevance scores of documents in
       both  $r_1, \dots, r_i$  and  $C$ ;
5:     if  $C$  contains NO doc in  $\{r_{i+1}, \dots, r_N\}$  then
6:       Insert  $C$  into Q ( $\text{rel}(q, C) = \text{rel}(q, C)_i$ );
       (for other cells, we have  $\text{rel}(q, C)_i = \text{rel}(q, C)_{i-1}$ )
7:     while (i)  $\text{rel}(q, \text{top}(\mathbf{Q})) \geq \max_C \{\text{rel}(q, C)_i\}$  and
       (ii)  $\text{rel}(q, \text{top}(\mathbf{Q})) \geq s(q, r_{i+1}[\mathcal{D}])$  do
8:       Output top(Q), and delete it from Q;
9:        $\text{cnt} \leftarrow \text{cnt} + 1$ ; if  $\text{cnt} = k$  then terminate.

```

3.2 Document Sorted-Scan Approach

Our second approach is outlined in Algorithm 2. Compared to Algorithm 1, it uses an additional data structure, a priority queue **Q**, to keep cells (candidates) in the descending order of relevance. Given a query q , initially, **Q** is empty. All rows in \mathcal{DB} are scanned in the descending order of relevance (line 1): r_1, r_2, \dots, r_N such that $s(q, r_i[\mathcal{D}]) \geq s(q, r_{i+1}[\mathcal{D}])$. To obtain this ranked list, we simply compute scores of all rows first and then construct a heap to keep them in the descending order of scores; a row is popped if it needs to be scanned.¹ Recall $\text{ans}(r_i)$ is the set of ancestors of r_i , i.e., the set of cells containing document $r[\mathcal{D}]$. When r_i is scanned, for each $C \in \text{ans}(r_i)$, we compute $\text{rel}(q, C)_i$ (line 4) as the average relevance score of documents among $r_1[\mathcal{D}], \dots, r_i[\mathcal{D}]$ that are contained in C (recall that $\text{rel}(q, C)$ is the average relevance score of all documents that are contained in C). If C contains NO document in $\{r_{i+1}, \dots, r_N\}$, we can conclude $\text{rel}(q, C) = \text{rel}(q, C)_i$ and insert C into the priority queue **Q** (lines 5-6). For a cell containing some rows that are already scanned, we have $\text{rel}(q, C) \leq \text{rel}(q, C)_i$, as rows are scanned in the descending order of relevance; for a cell C not touched yet (containing no row that is scanned), $\text{rel}(q, C)$ is bounded by $s(q, r_{i+1}[\mathcal{D}])$, the maximum relevance score among the rest individual documents. So for the cell with max relevance score in **Q**, $C' = \text{top}(\mathbf{Q})$, once $\text{rel}(q, C') \geq \text{rel}(q, C)_i$ for all cells C (line 7(i)), and $\text{rel}(q, C') \geq s(q, r_{i+1}[\mathcal{D}])$ (line 7(ii)), C' can be output and popped from **Q** (line 8).

Theorem 2: Algorithm 2 outputs the top- k cells with the highest relevance scores and support $\geq \text{minsup}$ in the non-increasing order of $\text{rel}(q, C)$.

Proof: We only need to prove that, in any iteration of lines 3-9, if $C' = \text{top}(\mathbf{Q})$ is output, then $\text{rel}(q, C')$ is the

maximum among all cells that are not output yet.

For any cell $C \in \mathbf{Q}$, by the definition of **Q**, $\text{rel}(q, C') \geq \text{rel}(q, C)$. For any other cell C , we may have either

$$\text{rel}(q, C) \leq \text{rel}(q, C)_i \quad \text{or} \quad \text{rel}(q, C) \leq s(q, r_{i+1}[\mathcal{D}]).$$

So from line 7(i)-(ii), we have $\text{rel}(q, C') \geq \text{rel}(q, C)$. \square

Implementation Issues. On line 4, $\text{rel}(q, C)_i$ can be efficiently computed from $\text{rel}(q, C)_j$ (r_j is the last row aggregated into C). On line 5, if the number of documents in a cell C is precomputed, we can test whether C contains NO document in $\{r_{i+1}, \dots, r_N\}$ by counting the number of documents that are already aggregated into C . On line 7, $\max_C \{\text{rel}(q, C)_i\}$ can be efficiently fetched if $\text{rel}(q, C)_i$'s are maintained in another priority queue.

Tie Breaker. Ties (for cells with the same relevance scores) break arbitrarily in Algorithm 2. If for cells with the same relevance scores, we want to output them in the descending order of support, then we modify Algorithm 2 as follows. First, the priority queue **Q** keeps cells (candidates) in the descending order of relevance, and keeps cells with the same relevance scores in the descending order of support. Second, “ \geq ” should be changed into “ $>$ ” on line 7, so that, when a cell is about to be output, all the cells (not-output-yet) with the same relevance scores are in **Q**.

Theorem 3: Algorithm 2 uses $O(N \cdot (|q| + \log N) + 2^n \cdot N \cdot \log M)$ time and $O(n \cdot M)$ space, where $N = |\mathcal{DB}|$ is the total number of rows, n is the dimensionality, and M is the total number of non-empty cells.

Proof: In the worst case, it needs to scan all rows with relevance score computed, in the descending order of relevance, which takes $O(N \cdot (|q| + \log N))$ time. Lines 4-6 repeat at most $O(2^n \cdot N)$ times, and each iteration takes $O(\log M)$ time since we need to maintain a priority queue **Q** and another queue to maintain $\max_C \{\text{rel}(q, C)_i\}$. Lines 7-9 also repeat at most $O(M) \leq O(2^n \cdot N)$ times and each iteration takes $O(\log M)$ time. So the time complexity follows. The space complexity is the same as Algorithm 1, as additional priority queues need linear space. \square

The complexity of Algorithm 2 is worse than Algorithm 1, but may terminate earlier before scanning all rows.

3.3 Handling Extended Keyword Queries

Algorithm 1 and 2 can be easily adapted to handle extended keyword query $Q = (u_1, \dots, u_n : q)$. With the dimension-value constraints u_1, \dots, u_n , we actually restrict our attention to a subspace of the text cube. Lines 2-3 of Algorithm 1 are modified a bit: we update $\text{rel}(q, C)$ only for the *feasible* cells in $\text{ans}(r)$. For example, if $Q = (?, ?, t1, * : q)$ over the text cube of \mathcal{DB} in Table 2 is given, then for row $r = (m1, p1, t1, s1)$, we only update relevance scores of its ancestors $(m1, p1, t1, *)$, $(*, p1, t1, *)$, $(m1, *, t1, *)$, and $(*, *, t1, *)$ in lines 3-4. Similar modification can be also applied to lines 2-3 of Algorithm 2.

3.4 Deficiencies of Linear Scan Approaches

Algorithm 1 is efficient only when the number of dimensions is small (from 2 to 4). Because, it scores all the non-

1. There are other top- k algorithms to obtain this ranked list, e.g., [21]. But the cost of Algorithm 2 is dominated by the cost of ranking cells, instead of ranking documents, so we just apply this simple method.

empty cells but the number of non-empty cells increases exponentially w.r.t. the dimensionality (see Section 2.4).

Algorithm 2 aims to improve Algorithm 1 by enabling early-stop. However, once a row is scanned, all the cells containing it are explored. There are totally 2^n cells containing a document in a n -dim cube. So the numbers of candidate cells and explored cells increase very quickly.

In the next section, we will introduce a more delicately-designed top- k algorithm, *search-space ordering*, which carries out cell-based search and explores as small number of cells in the cube as possible to find the top- k answers.

4 SEARCH-SPACE ORDERING APPROACH

Unlike the two algorithms in Section 3 which compute the relevance score of a cell from documents/rows in the database, the algorithms introduced in this section compute the score from its children cells, in a dynamic-programming manner. We introduce the basic *bottom-up dynamic programming* algorithm in Section 4.1, followed by the design and analysis of a more efficient approach, *search-space ordering*. We start with simple keyword queries and discuss how to extend our algorithm in Section 4.4.

4.1 Bottom-up Dynamic Programming

Consider a cell C' and all of its \mathcal{A} -children $\text{chd}_{\mathcal{A}}(C')$ for some aggregated dimension \mathcal{A} ($C'[\mathcal{A}] = *$ but $C[\mathcal{A}] \neq *$ for $C \in \text{chd}_{\mathcal{A}}(C')$). It is clear that the sets of documents aggregated in different cells in $\text{chd}_{\mathcal{A}}(C')$ are disjoint and their union is equal to $C'[\mathcal{D}]$, i.e., $C'[\mathcal{D}] = \bigcup_{C \in \text{chd}_{\mathcal{A}}(C')} C[\mathcal{D}]$. From the definition of relevance (1), we have recursion:

$$\text{rel}(q, C') = \frac{\sum_{C \in \text{chd}_{\mathcal{A}}(C')} \text{rel}(q, C) |C[\mathcal{D}]|}{\sum_{C \in \text{chd}_{\mathcal{A}}(C')} |C[\mathcal{D}]|}, \quad (3)$$

where $|C[\mathcal{D}]|$ is the number of documents in $C[\mathcal{D}]$.

Naive Dynamic Programming Algorithm

A naive dynamic programming algorithm that directly utilizes (3) is as follows. We first compute the relevance scores $\text{rel}(q, C)$ for all n -dim (base) cells, as in (1). Then we compute the relevance scores $\text{rel}(q, C')$ of $(n-1)$ -dim cells, $(n-2)$ -dim cells, \dots , and 0-dim cell, in the decreasing order of dimensionality, using the recursion (3); this is because for a d -dim cell C' , the right-hand side of (3) involves only $d+1$ cells. Finally, after the relevance scores of all cells are obtained, output the top- k relevant ones with supports no less than the threshold minsup .

Bottom-up Aggregation:

From Partially-aggregated to Fully-aggregated

To compute $\text{rel}(q, C')$, a direct implementation of (3) could be inefficient if $\min_{\mathcal{A}} |\text{chd}_{\mathcal{A}}(C')| = O(\min_{\mathcal{A}} |\mathcal{A}|)$ is large. We will introduce a bottom-up aggregation approach whose complexity is independent of $\min_{\mathcal{A}} |\mathcal{A}|$, but only dependent on n (the total number of dimensions). It is also an important basis of our search-space ordering algorithm which will be introduced in Section 4.3.

For each aggregated dimension \mathcal{A} of a cell C' ($C'[\mathcal{A}] = *$), we maintain a *partial relevant score* $\overline{\text{rel}}_{\mathcal{A}}(q, C')$ based on the \mathcal{A} -children that are already aggregated into C' . Let $L_{\mathcal{A}}(C')$ be the number of documents contained by the \mathcal{A} -children that are already aggregated into C' . A non-base cell C' is said to be *fully-aggregated* iff for some dimension \mathcal{A} , $L_{\mathcal{A}}(C') = |C'[\mathcal{D}]|$ and $\overline{\text{rel}}_{\mathcal{A}}(q, C') = \text{rel}(q, C')$; otherwise, C' is said to be *partially-aggregated*.

Initially, all cells are partially-aggregated and $\text{rel}(q, C)$'s are *unknown* for all cells; and we set $\overline{\text{rel}}_{\mathcal{A}}(q, C) = \infty$ and $L_{\mathcal{A}}(C) = 0$ for each cell C and each of its aggregated dimension \mathcal{A} . We compute $\text{rel}(q, C)$ for all base cells from the database. After that, *all base cells are fully-aggregated*.

For each fully-aggregated cell C and each of its non-aggregated dimensions \mathcal{A} ($C[\mathcal{A}] \neq *$), we aggregate C to its \mathcal{A} -parent C' as follows (define $\infty \times 0 = 0$):

Aggregate(C, C'): $C \in \text{chd}_{\mathcal{A}}(C')$

$$1 : \overline{\text{rel}}_{\mathcal{A}}(q, C') \leftarrow \frac{\overline{\text{rel}}_{\mathcal{A}}(q, C) L_{\mathcal{A}}(C) + \text{rel}(q, C) |C[\mathcal{D}]|}{L_{\mathcal{A}}(C') + |C[\mathcal{D}]|}; \quad (4)$$

$$2 : L_{\mathcal{A}}(C') \leftarrow L_{\mathcal{A}}(C') + |C[\mathcal{D}]|; \quad (5)$$

$$3 : \text{if } L_{\mathcal{A}}(C') = |C'[\mathcal{D}]| \text{ then } \text{rel}(q, C') \leftarrow \overline{\text{rel}}_{\mathcal{A}}(q, C'). \quad (6)$$

When $L_{\mathcal{A}}(C') = |C'[\mathcal{D}]| = \sum_{C \in \text{chd}_{\mathcal{A}}(C')} |C[\mathcal{D}]|$, i.e., all documents in $C'[\mathcal{D}]$ have been aggregated into C' through its \mathcal{A} -children, from (3), we have $\text{rel}(q, C') = \overline{\text{rel}}_{\mathcal{A}}(q, C')$ obtained and thus C' becomes fully-aggregated. We repeat this procedure until all cells are fully-aggregated. Note that from the definition of “fully-aggregated”, we only need to store $\overline{\text{rel}}_{\mathcal{A}}(q, C')$ and $L_{\mathcal{A}}(C')$ on one aggregated dimension \mathcal{A} (e.g., the one with the smallest index) for each cell.

Algorithm 3 Bottom-up Dynamic Programming Algorithm

Input: keyword query q against \mathcal{DB} , parameter k , minsup

- 1: Compute $\text{rel}(q, C)$'s for all non-empty base cells;
 - 2: **for** $d = n$ **downto** 1 **do**
 - 3: **for each** non-empty d -dim cell C **do**
 - 4: **for each** dimension \mathcal{A} **do**
 - 5: Let C' be C 's \mathcal{A} -parent;
 - 6: Call **Aggregate**(C, C');
 - 7: Output the top- k relevant cells with support $\geq \text{minsup}$ (for cells with the same relevance scores, output them in the descending order of support).
-

Algorithm 3 describes this approach. In each iteration of line 2-6, all of non-empty d -dim cells are fully-aggregated, because all of non-empty n -dim, $(n-1)$ -dim, \dots , $(d+1)$ -dim cells are scanned. We keep all non-empty cells in a queue so that we do not need to touch any empty cell. Note that a non-empty cell C has $C[\mathcal{D}] \neq \emptyset$ but possibly $\text{rel}(q, C) = 0$. From the recursion (3) how and $\overline{\text{rel}}_{\mathcal{A}}(q, C')$ is updated in (4)-(5), we can prove Algorithm 3 is correct.

Example 4.1: We shows how $\overline{\text{rel}}(q, C')$ is computed in Algorithm 3. Suppose there are five documents d_1, \dots, d_5 with document scores $s(q, d_1) = 8$, $s(q, d_2) = 4$, $s(q, d_3) = 4$, $s(q, d_4) = 6$, and $s(q, d_5) = 2$. Also, suppose $C' = (a1, *, *)$ has three B-children and three C-children,

TABLE 4
B-children and C-children of $C' = (a1, *, *)$

| chd _B (C') | | | | | chd _C (C') | | | | |
|-----------------------|----|---|------------------------------------|-----|-----------------------|---|----|------------------------------------|-----|
| A | B | C | D | rel | A | B | C | D | rel |
| a1 | b1 | * | {d ₁ , d ₂ } | 6 | a1 | * | c1 | {d ₁ , d ₂ } | 6 |
| a1 | b2 | * | {d ₃ , d ₄ } | 5 | a1 | * | c2 | {d ₃ } | 4 |
| a1 | b3 | * | {d ₅ } | 2 | a1 | * | c3 | {d ₄ , d ₅ } | 4 |

which are all fully-aggregated. Documents in C' 's children and cell relevance scores are shown in Table 4.

When $C = (a1, b1, *)$ and $(a1, *, c1)$ are aggregated into C' , using **Aggregate**(C, C'), we have $L_B(C') = L_C(C') = 2$, and partial scores $\overline{\text{rel}}_B(q, C') = \overline{\text{rel}}_C(q, C') = 6$.

After B-child $(a1, b2, *)$ is aggregated into C' , we have $\overline{\text{rel}}_B(q, C') = (6 \times 2 + 5 \times 2)/4 = 5.5$, and $L_B(C') = 4$.

After C-child $(a1, *, c2)$ is aggregated into C' , we have $\overline{\text{rel}}_C(q, C') = (6 \times 2 + 4 \times 1)/3 = 5.33$, and $L_C(C') = 3$.

Later when at least one of $(a1, b3, *)$ and $(a1, *, c3)$ is aggregated into C' , C' is fully-aggregated. But actually we only need to keep either $(\overline{\text{rel}}_B, L_B)$ or $(\overline{\text{rel}}_C, L_C)$. \square

Theorem 4: Algorithm 3 uses $O(N \cdot |q| + n \cdot M)$ time and $O(n \cdot M)$ space, where $N = |\mathcal{DB}|$, M is the total number of non-empty cells, and n is the number of dimensions.

Proof: Non-empty base cells are computed in $O(N \cdot |q|)$ time. All non-empty cells are kept in queue s.t. the ordered access of each takes $O(1)$ time. Lines 4-6 need $O(n)$ time. So we need $O(n \cdot M)$ time to compute $\text{rel}(q, C)$ for all non-empty cells. We use Hoare's selection algorithm to select the top- k from all non-empty cells in $O(M)$ time.

A simple implementation needs $O(n \cdot M)$ space: for each cell we keep $\overline{\text{rel}}_A(q, C)$ and $L_A(C)$ for every dimension A . But we actually only need to keep them on one aggregated dimension (e.g., the one with the smallest index). \square

In the rest of this paper, for each cell C we store $\overline{\text{rel}}_A(C)$ and $L_A(C)$ for only one dimension. So we will also write $\overline{\text{rel}}_A(C)$ as $\overline{\text{rel}}(C)$, if not otherwise specified.

4.2 Upper/Lower Bounds of Relevance Score

We first introduce a pair of upper/lower bounds of relevance scores, which will be used to prune search space of our algorithm in Section 4.3. The main idea is: the relevance score of the union of two document sets lies between the relevance scores of these two document sets.

Lemma 1: For any query q and any two disjoint document sets D_1, D_2 s.t. $\text{rel}(q, D_1) \leq \text{rel}(q, D_2)$, we have $\text{rel}(q, D_1) \leq \text{rel}(q, D_1 \cup D_2) \leq \text{rel}(q, D_2)$.

Proof: From the definition in Equation (1),

$$\text{rel}(q, D_1 \cup D_2) = \frac{|D_1| \cdot \text{rel}(q, D_1) + |D_2| \cdot \text{rel}(q, D_2)}{|D_1| + |D_2|}.$$

Since $\text{rel}(q, D_1) \leq \text{rel}(q, D_2)$,

$$\text{rel}(q, D_1 \cup D_2) \leq \frac{|D_1| + |D_2|}{|D_1| + |D_2|} \cdot \text{rel}(q, D_2) = \text{rel}(q, D_2),$$

and similarly, $\text{rel}(q, D_1 \cup D_2) \geq \text{rel}(q, D_1)$. \square

Example 4.2: To verify Lemma 1, consider $D_1 = \{d_1, d_2\}$ and $D_2 = \{d_3\}$. Suppose for some query q , we

have $s(q, d_1) = 1, s(q, d_2) = 3, s(q, d_3) = 8$. Then, $\text{rel}(q, D_1) = 2, \text{rel}(q, D_2) = 8$, and $\text{rel}(q, D_1 \cup D_2) = 4$. \square

Lemma 2: (*Two-side bound property*) For any non-base cell C in text cube and any query q , there exist two children C_1 and C_2 of C s.t. $\text{rel}(q, C_1) \leq \text{rel}(q, C) \leq \text{rel}(q, C_2)$.

Proof: Consider \mathcal{A} -children of C : $\bigcup_{C' \in \text{chd}_A(C)} C'[\mathcal{D}] = C[\mathcal{D}]$, where $C'[\mathcal{D}]$'s are disjoint. So from Lemma 1, the proof can be completed by induction on $|\text{chd}_A(C)|$. \square

Remark. Lemmas 1 and 2 are still true for different document scoring formulas $s(q, d)$ and the weighted version of $\text{rel}(q, C)$ (i.e., $\text{rel}(q, C)$ is a weighted average of $s(q, d)$ over all documents d in $C[\mathcal{D}]$, where the weight of a document w_d is fixed and positive).

4.3 Search-Space Ordering Algorithm

Our *search-space ordering* approach is based on the following idea. Although the search space is huge (especially if the dimensionality is reasonably large, say, ≥ 10), when only the top- k answers are interesting to users, only a small portion of the text cube is ‘‘closely’’ relevant to the given query q . So, (i) we estimate the upper bound of relevance scores in a subspace of the text cube, and (ii) we *order* and *prune* subspaces so that we can follow a shortcut to the top- k relevant cells with support $\geq \text{minsup}$.

4.3.1 Algorithm Framework

The main idea is: *order* and *prune* the search space of cells using (the upper bound of) relevance scores.

Bottom-up Aggregation: Our algorithm works in a similar manner as the bottom-up dynamic programming algorithm (Algorithm 3). Relevance scores are computed from the base cells (whose relevance scores can be computed with the row-based inverted indexes efficiently) to higher levels.

Search-Space Ordering: But unlike line 3 in Algorithm 3 (where the next cell C to be aggregated to its parents is selected arbitrarily), we pick the most relevant *fully-aggregated cell* in each iteration and aggregate it to its parents. This is because a highly-relevant cell is likely to have highly-relevant children. Note that a fully-aggregated cell needs to be aggregated to its parents at most once.

Search-Space Pruning: The relevance score of a fully-aggregated cell is known. For a partially-aggregated cell, we have the upper bound from Lemma 2. When a *fully-aggregated cell* has its relevance score no lower than the upper bound of relevance scores of all partially-aggregated cells, it is output as one of the top- k (if support $\geq \text{minsup}$).

We will discuss the details of this algorithm and upper bound estimation for relevance scores in Section 4.3.2.

4.3.2 Algorithm Description

Three structures are used in our algorithm: Initially, all base cells are put into *aggregation queue* (Q_A), and the base cells with support no less than minsup are put into *candidate queue* (Q_C), both maintained in the non-increasing order of relevance scores. The *partially-aggregated pool* (P_A) is empty. More cells will be inserted into Q_A and Q_C later. If

TABLE 5
Base Cuboid of a 2-Dimensional Text Database

| M | P | $C[\mathcal{D}]$ | $\text{rel}(q, C)$ |
|----|----|------------------------------------|--------------------|
| m1 | p1 | {d ₁ } | 8 |
| m1 | p2 | {d ₂ } | 6 |
| m1 | p3 | {d ₃ , d ₇ } | 4 |
| m2 | p1 | {d ₄ } | 8 |
| m2 | p2 | {d ₅ } | 6 |
| m2 | p3 | {d ₆ , d ₈ } | 4 |

an ancestor has the same relevance score as its descendant, the ancestor has higher priority in these two queues.

- **Aggregation Queue (Q_A):** In each iteration, the cell C with the highest score $\text{rel}(q, C)$ in Q_A is popped out and aggregated into each of its parents C' using **Aggregate**(C, C'), i.e., (4)-(6), with the partial relevance score $\overline{\text{rel}}(q, C')$ computed/updated (since a highly-relevant cell likely has highly-relevant children). After that, if C' is partially-aggregated, it is put into the partially-aggregated pool P_A if C' is not in P_A ; if C' becomes fully-aggregated, then: (i) C' is deleted from P_A and inserted into Q_A , and (ii) if its support $\geq \text{minsup}$, C' is inserted into Q_C .
- **Partially-aggregated Pool (P_A):** It maintains partially-aggregated cells (and partial relevance scores $\overline{\text{rel}}(q, \cdot)$), each of which has at least one child aggregated into the cell. Note that it is possible no child of a cell C is aggregated into C , and thus C (still called partially-aggregated cell) is not in P_A .
- **Candidate Queue (Q_C):** All fully-aggregated cells are maintained in Q_C if support $\geq \text{minsup}$. When a cell in Q_C has a relevance score higher than the upper bound of relevance scores of all partially-aggregated cells, it can be output as one of the top- k . We will prove such an upper bound is: $\max\{\text{maximum partial relevance score in } P_A, \text{maximum relevance score in } Q_A\}$.

The framework of this algorithm is outlined in Algorithm 4. Given a keyword query q , the relevance scores of base cells can be computed using inverted indexes (line 1). Initially, P_A is empty, all base cells are inserted into Q_A , and the ones with support $\geq \text{minsup}$ are inserted into Q_C . In each iteration of lines 3-14, the top cell C in Q_A is popped (deleted from Q_A) for aggregation (line 4). For each partially-aggregated parent C' of C , we insert it into P_A if necessary (line 6); and we can compute/update the partial relevance score $\overline{\text{rel}}(q, C')$ using **Aggregate**(C, C') (line 7), as in (4)-(6). If the cell C' is *fully-aggregated* (for some dimension \mathcal{A} , all of its \mathcal{A} -children have already been aggregated into C'), we can compute its exact relevance score $\text{rel}(q, C')$, delete it from P_A , and add it into Q_A ; it is also added into Q_C if its support $\geq \text{minsup}$ (lines 8-11). At the end of each iteration (lines 12-14), the top cells in Q_C can be output (also deleted from Q_C) if *their scores are no less than the maximum partial relevance score in P_A and the maximum relevance score in Q_A (i.e., the upper bound of relevance scores of partially-aggregated cells)*.

TABLE 6
Running Example

| | |
|---|--|
| <i>Initialization:</i> All base cells are put into Q_A . Two cells with support $\geq \text{minsup}$ are put into Q_C . P_A is empty. | |
| $P_A(0)$ | \emptyset |
| $Q_A(0)$ | (m1,p1):8, (m2,p1):8, (m1,p2):6, (m2,p2):6, (m1,p3):4, (m2,p3):4 |
| $Q_C(0)$ | (m1,p3):4, (m2,p3):4 |

| | |
|--|---|
| <i>1st iteration:</i> Cell (m1,p1) is popped from $Q_A(0)$, and aggregated into its two parents, (m1,*) and (*,p1) (inserted into $P_A(1)$). The partial relevance scores of (m1,*) and (*,p1) are computed as 8. | |
| $P_A(1)$ | (m1,*):8, (*,p1):8 |
| $Q_A(1)$ | (m2,p1):8, (m1,p2):6, (m2,p2):6, (m1,p3):4, (m2,p3):4 |
| $Q_C(1)$ | (m1,p3):4, (m2,p3):4 |

| | |
|--|--|
| <i>2nd iteration:</i> (m2,p1) is popped from $Q_A(1)$, and aggregated into the two parents, (m2,*) and (*,p1): (m2,*) is inserted into $P_A(2)$ with the partial relevance score computed as 8; (*,p1) are already in $P_A(1)$. Then, (*,p1) is <i>fully aggregated</i> , since it has only two M-children, (m1,p1) and (m2,p1). So (*,p1) is deleted from $P_A(1)$ and inserted into $Q_A(2)$. It is also inserted into $Q_C(2)$ for its support $\geq \text{minsup}$. Finally, the top cell in $Q_C(2)$, (*,p1), has relevance score no less than the maximum partial relevance score in $P_A(2)$ and the maximum relevance score in $Q_A(2)$, so (*,p1) is output. | |
| $P_A(2)$ | (m1,*):8, (m2,*):8 |
| $Q_A(2)$ | (*,p1):8, (m1,p2):6, (m2,p2):6, (m1,p3):4, (m2,p3):4 |
| $Q_C(2)$ | (*,p1):8, (m1,p3):4, (m2,p3):4 |

| | |
|--|--|
| <i>3rd iteration:</i> Cell (*,p1) is popped from $Q_A(2)$, and its parent (*,*) is added into $P_A(3)$. | |
| $P_A(3)$ | (m1,*):8, (m2,*):8, (*,*):8 |
| $Q_A(3)$ | (m1,p2):6, (m2,p2):6, (m1,p3):4, (m2,p3):4 |
| $Q_C(3)$ | (m1,p3):4, (m2,p3):4 |

| | |
|--|---------------------------------------|
| <i>4th iteration:</i> (m1,p2) is popped from $Q_A(3)$, and aggregated into its two parents, (m1,*) and (*,p2). (m1,*) is already in $P_A(3)$, and its partial relevance score is updated as 7 in $P_A(4)$, since we have seen two children of (m1,*) and (m1,p2), with score 8 and 6, respectively. (*,p2) is inserted into $P_A(4)$ with partial score 6. | |
| $P_A(4)$ | (m2,*):8, (*,*):8, (m1,*):7, (*,p2):6 |
| $Q_A(4)$ | (m2,p2):6, (m1,p3):4, (m2,p3):4 |
| $Q_C(4)$ | (m1,p3):4, (m2,p3):4 |

We give some intuitions on why the above upper bound is valid. In our algorithm, a loop invariant is: for any partially-aggregated cell C (no matter whether it is in P_A), the union of P_A and Q_A forms a cut between C and the base cells in the cube lattice. So from the two-side bound property (Lemma 2), the relevance score of C is upper bounded by the (partial) relevance score of a descendant cell in $P_A \cup Q_A$. Formal proofs will be given in Section 4.3.3.

Example 4.3: This example shows the basic idea of our algorithm, using a 2-dimensional text database. Table 5 shows the base cuboid of this database—each base cell may contain more than one document. Consider a keyword query $q = \{w1, w2\}$ and $\text{minsup} = 2$. Suppose the relevance score $\text{rel}(q, C)$ (defined in Equations (1)-(2)) of each base cell is computed as in the fourth column of Table 5.

Initially, all base cells are put into Q_A , and the ones with support no less than minsup are put into Q_C , both in the non-decreasing order of relevance scores. P_A is empty.

We show how to maintain Q_A, Q_C, P_A in Table 6. Number in bracket in the first column is the number of iterations, e.g., $P_A(2)$ means the status of P_A at the end of the 2nd

Algorithm 4 Search-Space Ordering Algorithm

Q_A and Q_C : priority queues of cells (sorted in the non-increasing order of relevance scores)
 $\text{top}(Q)$: the top (first) cell of a priority queue Q
 P_A : the set of partially-aggregated cells

Input: keyword query q against \mathcal{DB} , parameter k , minsup

- 1: Compute $\text{rel}(q, C)$'s for all non-empty base cells;
- 2: $\text{cnt} \leftarrow 0$, $Q_A \leftarrow \{\text{all non-empty base cells}\}$,
 $Q_C \leftarrow \{\text{non-empty base cells with support} \geq \text{minsup}\}$;
- 3: **while** $\text{cnt} < k$ **do**
- 4: Let C be the top cell in Q_A , and delete C from Q_A ;
- 5: **for** each partially-aggregated parent $C' \in \text{par}(C)$ **do**
- 6: **if** $C' \notin P_A$ **then** insert C' into P_A ;
- 7: Aggregate C into C' using **Aggregate**(C, C'):
Update partial relevance score $\overline{\text{rel}}(q, C')$;
- 8: **if** C' is fully-aggregated (all documents in $C'[\mathcal{D}]$
are aggregated into C' on some dimension) **then**
- 9: Compute $\text{rel}(q, C')$;
- 10: Delete C' from P_A ;
- 11: Insert C' into Q_A ;
- 12: Insert C' into Q_C **if** $\text{support } |C'| \geq \text{minsup}$;
- 13: **while** $\text{rel}(q, \text{top}(Q_C)) \geq \max_{C' \in P_A} \{\overline{\text{rel}}(q, C')\}$ and
 $\text{rel}(q, \text{top}(Q_C)) \geq \text{rel}(q, \text{top}(Q_A))$ **do**
- 14: Output the top cell of Q_C , and delete it from Q_C ;
- 15: $\text{cnt} \leftarrow \text{cnt} + 1$;

iteration; the number after each cell in the second column is the relevance score (or partial relevance score in P_A). \square

Enlarging k at Running Time without Recomputation.

With k unspecified, repeating the iteration of lines 3-15, cells with support $\geq \text{minsup}$ will be output in the non-increasing order of relevance scores. Users can stop it at any time when they feel the results are enough.

Tie Breaker. Similar to Algorithm 2, ties (for cells with the same relevance scores) break arbitrarily in Algorithm 4. If we want to rank cells with the same relevance scores in the descending order of their support, Algorithm 4 is revised as follows. First, the candidate queue Q_C keeps cells in the descending order of relevance, and keeps cells with the same relevance scores in the descending order of support. Second, " \geq " should be changed into " $>$ " on line 12, so that, when a cell is about to be output, all the cells (not-output-yet) with the same relevance scores are in Q_C .

Theorem 5: Algorithm 4 needs $O(N \cdot |q| + n \cdot T \cdot \log T)$ time and $O(n \cdot T)$ space, where $N = |\mathcal{DB}|$, n is the number of dimensions, and T is the number of cells explored by this algorithm (i.e., the ones used to be added into Q_A, Q_C, P_A).

Proof: Again, non-empty base cells are computed in $O(N \cdot |q|)$ time. After a cell explored is fully-aggregated, on lines 5-11, it is aggregated into its (at most n) parents (using $O(n)$ time), and one operation on Q_A, Q_C, P_A (using $O(\log T)$ time) is invoked for each parent.

The space needed is $O(n \cdot T)$. This is because, for each cell explored, we keep its n -dim identifier and $\overline{\text{rel}}_{\mathcal{A}}(q, C)$ & $L_{\mathcal{A}}(C)$ for one dimension \mathcal{A} . And, a pointer to this cell is kept in Q_A, Q_C, P_A using $O(1)$ space. \square

The total number of non-empty cells, M (recall Theorem 4), is an upper bound of T . However, T is determined by both data distribution in the database and the query q . So we cannot get a preciser estimation of it beforehand. We will report the actual values of T in our experiments, and show that Algorithm 4 explores a much smaller portion of cells than other algorithms.

4.3.3 Correctness

We prove the correctness of Algorithm 4 in the following part ($s(\cdot, \cdot)$ not restricted to the form of Equation (2)).

Lemma 3 (Estimating Upper Bound of Relevance Scores):

For any query q , on line 12 of Algorithm 4, if

(a) $\text{rel}(q, \text{top}(Q_C)) \geq \max_{C' \in P_A} \{\overline{\text{rel}}(q, C')\}$ and

(b) $\text{rel}(q, \text{top}(Q_C)) \geq \text{rel}(q, \text{top}(Q_A))$,

then cell $\text{top}(Q_C)$ has a relevance score no less than the scores of all partially-aggregated cells.

Proof: Let $\Delta = \text{rel}(q, \text{top}(Q_C))$. For the purpose of contradiction, suppose there is a cell C'' that is *partially-aggregated* and the relevance score $\text{rel}(q, C'') > \Delta$.

We focus on \mathcal{A} -children of C'' for its selected aggregated dimension \mathcal{A} , s.t. $\overline{\text{rel}}(q, C'') = \overline{\text{rel}}_{\mathcal{A}}(q, C'')$.

Let $P(C'')$ denote the set of C'' 's \mathcal{A} -children that are partially-aggregated, and $F(C'')$ denote the set of C'' 's \mathcal{A} -children that are fully-aggregated but not aggregated into C yet. Recall $L_{\mathcal{A}}(C'')$ is the number of documents that are already aggregated into C'' from its \mathcal{A} -children.

If $L_{\mathcal{A}}(C'') = 0$, we know that no child is aggregated into C'' . Then, $P(C'') \cup F(C'')$ covers all of C'' 's \mathcal{A} -children. Note $\text{rel}(q, C'') > \Delta$. From Lemma 2, there exists $C \in P(C'') \cup F(C'')$ s.t. $\text{rel}(q, C) \geq \text{rel}(q, C'') > \Delta$.

If $L_{\mathcal{A}}(C'') \neq 0$, $C'' \in P_A$. From condition (a) in the lemma, $\overline{\text{rel}}(q, C'') \leq \max_{C' \in P_A} \overline{\text{rel}}(q, C') \leq \Delta$. Again note $\text{rel}(q, C'') > \Delta$. From Lemma 1-2, there exists a child $C \in P(C'') \cup F(C'')$ s.t. $\text{rel}(q, C) > \Delta$.

For the selected children C of C'' , there are two cases:

- (i) $C \in F(C'')$: We must have $C \in Q_A$, since otherwise, C is already aggregated into C'' when C is popped from Q_A . Then, $\text{rel}(q, \text{top}(Q_A)) \geq \text{rel}(q, C) > \Delta$, which contradicts with $\Delta \geq \text{rel}(q, \text{top}(Q_A))$.
- (ii) $C \in P(C'')$: C is partially aggregated. Repeat the above argument about C'' for C (or induction on the number of aggregated dimensions in C'').

Since the text cube is finite, that means, if $\text{rel}(q, C'') > \Delta$ for some C'' , we will always terminate with a descendant C of C'' with contradiction. So the proof is completed. \square

The above proof holds if we choose $\overline{\text{rel}}(q, C'') = \overline{\text{rel}}_{\mathcal{A}}(q, C'')$ for any aggregated dimension \mathcal{A} . So if choosing $\overline{\text{rel}}(q, C'') = \min_{\mathcal{A}} \overline{\text{rel}}_{\mathcal{A}}(q, C'')$, our algorithm has a tighter upper bound (it may terminate earlier but incurs more space to store $\overline{\text{rel}}_{\mathcal{A}}$ for each dimension).

Theorem 6: Algorithm 4 outputs the top- k cells with the highest relevance scores and support $\geq \text{minsup}$ in the non-increasing order of $\text{rel}(q, C)$.

Proof: For any cell C with support $\geq \text{minsup}$, before it is output, it is either partially-aggregated or in Q_C . From Lemma 3, when the cell $\text{top}(Q_C)$ is output, its relevance

TABLE 7
Laptop Review Database: Dimensions and Statistics

| Dimension | Description | Cardinality |
|------------------|-------------------------|-------------|
| 1. Laptop Brand | brand of the laptop | 20 |
| 2. Laptop Family | family of the laptop | 26 |
| 3. Screen | size of the screen | 26 |
| 4. Memory | size of the main memory | 9 |
| 5. Disk | size of the hard drive | 28 |
| 6. CPU Brand | brand of the processor | 5 |
| 7. CPU Model | model of the processor | 148 |
| 8. CPU Speed | speed of the processor | 33 |
| 9. OS | pre-installed os | 12 |
| 10. Color | color of the laptop | 10 |

score is no less than the relevance score of any partially-aggregated cell. So, from the way that Q_C is ordered (non-increasing order of relevance scores), for any $cnt = k$ on lines 13-14 of Algorithm 4, $\text{top}(Q_C)$ is the cell with the top- k highest relevance score among all cells with support $\geq \text{minsup}$; thus, top- k cells are output in the non-increasing order of relevance scores. \square

4.4 Handling Extended Keyword Queries

We modify Algorithm 4 as follows to handle extended keyword queries $Q = (u_1, u_2, \dots, u_n : q)$.

- (i) On lines 1-2, instead of inserting base cells into Q_A and Q_C , we insert the cells C 's with (a) $C[\mathcal{A}_i] = *$ if $u_i = *$ (dimension \mathcal{A}_i must be aggregated), (b) $C[\mathcal{A}_i] = u_i$ if $u_i \in \mathcal{A}_i$ is a specific value, and (c) $C[\mathcal{A}_i] \in \mathcal{A}_i$ is any specific value if $u_i = ?$. One scan of inverted indexes of terms in q suffices to identify these cells and insert them into Q_A and Q_C .
- (ii) On lines 5-14, we consider $C' \in \text{par}(C)$, only if C' is a feasible cell, which satisfies the constraint in Q .

Algorithm 4 is even *faster* when handling an extended keyword query than when handling a simple query with the same set of keywords, because more constraints on dimensions actually reduce the size of search space.

5 EXPERIMENTAL STUDY

5.1 Datasets and Environment Setup

A real dataset of multidimensional text data is used in the experiments. We crawled laptop reviews with multidimensional attributes from Amazon.com, and then select 10 attributes as the structural dimensions in our database (Table 7). Each review (*document*) together with its categorical attributes (*dimensions*) forms a row in our database. We have 13,930 reviews in sum, and each review has 73.4 keywords with stop words removed on average. The total length of documents in this text database is 1,022,462.

Synthetic datasets are generated by adding randomly-generated documents into the real dataset. To this end, we estimate the overall distributions of terms, document length, and attribute values for different dimensions. Five synthetic datasets, DB_1, DB_2, \dots, DB_5 are used in Section 5.2 to test the scalability of our algorithms w.r.t. the size of text database. The numbers of rows are: $|DB_1| = 28K$,

$|DB_2| = 56K$, $|DB_3| = 84K$, $|DB_4| = 112K$, and $|DB_5| = 140K$. DB_i is a subset of DB_{i+1} for $i = 1, 2, 3, 4$.

All the experiments were conducted on a PC running the Microsoft Windows XP SP2 Professional OS, with a 2.5 GHz Intel Core 2 Duo T9300 CPU, 4 GB of RAM, and 150 GB hard drive. Our algorithms were implemented in C/C++ and compiled on Microsoft Visual Studio 2008.

We use Okapi weighting (2) for $s(q, d)$ in experiments.

5.2 Performance Study on Efficiency

We test our four algorithms, **One-Scan** (Algorithm 1), **Sorted-Scan** (Algorithm 2), **BU-DP** (Algorithm 3) and **SS-Ordering** (Algorithm 4), to demonstrate their scalability on both the real dataset (laptop reviews) and the synthetic datasets. We report the time and memory used by these four algorithms for finding top- k answers to given keyword queries. Time is measured in terms of milliseconds (MSec), and memory is measured in terms of KBs. We also report the number of cells that are explored by these algorithms, to demonstrate the pruning power of our algorithms.

We report only the *running-time* memory, *i.e.*, the amount of memory consumed by data structures for query processing (excluding the inverted indexes). **One-Scan** and **BU-DP** consume almost the same amount memory and explore the same number of cells. So we only report **One-Scan** for memory and the number of cells.

In the real dataset, we vary the number of dimensions (Figure 1), the parameter of k (Figure 4), the number of keywords $|q|$ (Figure 5), and the minimum support requirement minsup (Figure 6). Also, the five synthetic datasets are used to test the scalability of these three algorithms w.r.t. the sizes of databases (Figure 2).

All the results reported below are the averages for 20 typical keyword queries (with no dimension-value constraint) for laptops. Each query has 4 keywords on average.

We omit the results of efficiency experiments for extended keyword queries, because as discussed in Sections 3.3 and 4.4, adding such constraints will not slow down our algorithms. In fact, to handle an extended keyword query $Q = (u_1, u_2, \dots, u_n : q)$ with m dimensions \mathcal{A}_i 's s.t. $Q[\mathcal{A}_i] = u_i = ?$ (dimension \mathcal{A}_i has no constraint), it is equivalent for our algorithms to handle a simple keyword query q in an m -dim text cube.

Preprocessing: As discussed in Section 2.3, our algorithms benefit from preprocessing (*i.e.*, precomputing the support of each cell and constructing inverted indexes). In our implementation, it takes four and a half minutes to preprocess the real dataset described above, which is nontrivial but affordable since we only need to preprocess a dataset once.

Experiment I: Number of Dimensions (Figure 1) We use the real dataset (with about 14K rows). Set $k = 10$ and $\text{minsup} = 1$. The database is projected into the first n dimensions, for $n = 2, 4, 6, 8, 10$ (refer to Table 7).

SS-Ordering is consistently at least 10 times faster than the others, and consumes much less memory (except when $n = 2$). BU-DP is the second fastest one when $n \leq 6$ but becomes worst later on.

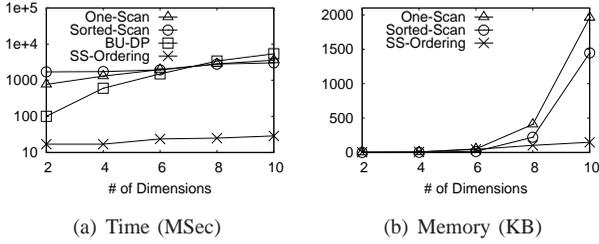
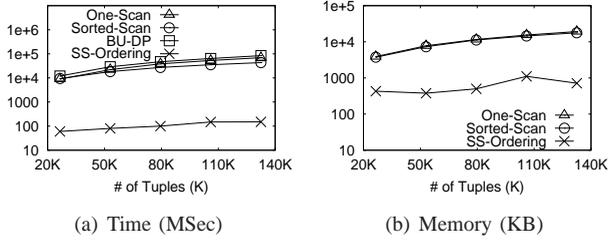


Fig. 1. Varying the Number of Dimensions

Fig. 2. Varying the Size of Database $|DB|$

Sorted-Scan has an improvement of 20% over One-Scan for 10-dim cube, but is slower than One-Scan for 4-dim cube. This is because they have the same worst-case complexity and Sorted-Scan needs to maintain additional data structures; so sometimes when its pruning power is not strong, it needs more time than One-Scan.

When n increases, the number of non-empty cells increases exponentially, and thus all algorithms use more time and memory. SS-Ordering is less sensitive than the rest three, because SS-Ordering only explores a small portion of the cells, while One-Scan and BU-DP always explore all the non-empty cells. Sorted-Scan explores less cells than One-Scan, but for the additional operation required, the improvement over One-Scan on efficiency is not significant.

The experimental result of testing extended keyword queries $Q = (u_1, u_2, \dots, u_n : q)$ with m no-constraint dimensions \mathcal{A}_i 's (i.e. $Q[\mathcal{A}_i] = u_i = ?$) is similar to the above, when we vary m from 2 to 10.

Experiment II: Size of Database (Figure 2) Synthetic datasets DB_1, \dots, DB_5 are used. These datasets have 10 dimensions and their sizes vary from 28K to 140K. Top-10 answers are output with $\text{minsup} = 1$. Sorted-Scan is always faster than One-Scan and BU-DP by 10%-40%. Again, SS-Ordering is much faster and consumes less memory than the other three for its power of search space pruning.

Experiment III: Number of Cells Explored (Figure 3) We report the number of cells explored by our algorithms in Experiments I-II. SS-Ordering explores much less cells than others thanks to its pruning techniques, especially for datasets that have more dimensions or more tuples (that is why it needs much less time than other algorithms). Both of One-Scan and BU-DP (not in the figure) explore all the non-empty cells which contain at least one keyword, and Sorted-Scan explores a bit less than them but still much more and SS-Ordering in most cases.

We will report the running-time memory consumed by

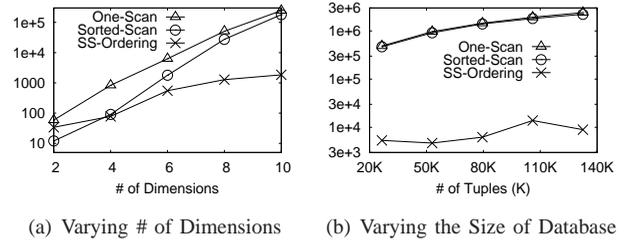
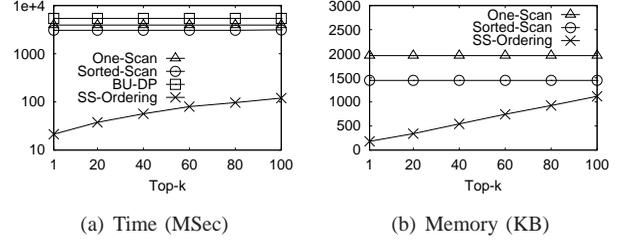


Fig. 3. Number of Cells Explored

Fig. 4. Varying the Parameter k

our algorithms for a fair comparison, since different algorithms consume different amounts of memory per cell.

Experiment IV: Parameter k (Figure 4) We use the real dataset with 10 dimensions and set $\text{minsup} = 1$. The time and memory used in One-Scan and BU-DP nearly unchange because they always explore all non-empty cells in the text cube. For Sorted-Scan, the time and memory nearly unchange, because the early stop condition (lines 7-9 of Algorithm 2) tends to output a batch of cells at certain time. Again, SS-Ordering consumes much less time and memory than others. Sorted-Scan has an improvement of 20%-40% over One-Scan and BU-DP in terms of the running time.

It can be also noticed that SS-Ordering uses more time and memory when k increases, because when more answers are required, more cells need to be explored in SS-Ordering. However, the increment is not too much, because, SS-Ordering algorithm follows a “short-cut” to the top- k cells, by first aggregating highly-relevant cells into their parents.

Experiment V: Number of Keywords (Figure 5) We study how the number of keywords $|q|$ affects the performance. We set $k = 10$ and $\text{minsup} = 1$, using the 10-dimensional laptop review dataset. Results are reported in Figure 5. SS-Ordering is always faster and consumes less memory than the other three, as in other experiments.

Note that the time and memory used by Sorted-Scan and SS-Ordering vibrate irregularly. We explain this as follows.

When the number of keywords increases, there are two factors which affect the performance of the four algorithms: (i) at the beginning, more inverted indexes need to be scanned to compute the relevance of rows/cells; (ii) the variance of relevance scores of different cells could be larger, so it is easier for Sorted-Scan to satisfy the early stop condition, and for the SS-Ordering algorithm to reach the highly-relevant cells (for it always seeks for a cell with high relevance score to aggregate it into its parent).

One-Scan and BU-DP are only affected by (i), so its running time and memory increase consistently as the num-

model (i.e., the relevance of a cell w.r.t. a query is the min/max relevance of documents in this cell).

When the relevance scoring formula *rel* does *not* satisfy the property in Lemma 1, how to efficiently answer the top-*k* keyword query is still an open question.

6.2 About Support Requirement

As discussed in Section 1, the purpose of specifying the support requirement *minsup* is to help users filter out cells/objects that are *not* “popular” enough (i.e., containing too few documents). However, when the user has little knowledge about the database (e.g., the size of this database), it is undesirable for her/him to provide such a parameter. We discuss two directions to handle this issue.

One way to address this issue is to integrate the support of a cell into its relevance score. But the resulting scoring formula may not satisfy the property in Lemma 1, and thus cannot be handled by Algorithm 4 (although still can be handled by less efficient Algorithm 1 and Algorithm 3).

Another way is more direct. Consider several typical support requirements: for example, “normal” (support $\geq \text{minsup} = x_1$), “somewhat popular” (support $\geq \text{minsup} = x_2$), “popular” (support $\geq \text{minsup} = x_3$), and “very popular” (support $\geq \text{minsup} = x_4$). Parameters $x_1 < \dots < x_4$ are specified by the database administrator, e.g., $x_i = 2^i$. For each incoming query (without specifying *minsup*), the system returns top-*k* answers for each of the four support requirements (i.e., running Algorithm 4 once for each value of *minsup*), so that the user can browse top-*k* cells in different levels (in terms of supports) at the same time.

7 RELATED WORK

Keyword Search in RDBMs. Although based on different applications and motivations, keyword search in text cube is related to keyword search in RDBMs, which has attracted a lot of attention recently [22], [23], [24]. Most previous studies on keyword search in RDBMs model the RDB as a graph (tuples/tables as nodes, and foreign-key links as edges) and focus on finding minimal connected tuple trees that contain all the keywords. They can be categorized into two types. The first type uses SQL to find the connected trees [2], [7], [6], [10], [12], [13]. The second type materializes the RDB graph and proposes algorithms to enumerate (top-*k*) subtrees in the graph [3], [4], [8], [14], [9].

Different from these two types of works, two recent studies [11] and [5] find single-center subgraphs from the RDB graph, and multi-center induced subgraphs, respectively.

OLAP on Multidimensional Text Data. The text cube model is firstly proposed in [1]. [1] mainly focuses on how to partially materialize inverted indexes and term frequency vectors in cells of text cube, and how to support OLAP queries (not keyword query) efficiently.

The topic cube model is proposed in [25]. Different from the text cube, the topic cube materializes the language model of the aggregated document in each cell. Efficient algorithms are proposed to compute this topic cube.

The techniques in [1] and [25] cannot be used directly to support keyword search, because the information materialized in text cube (term frequencies and inverted indexes) and in topic cube (language model) is query-independent.

Faceted Search. Several faceted search systems have been developed to help users navigate the search results if some meta-data about the results are available ([26], [27], [28], [29], [30], [31], [32], [33]). An initial faceted search systems [26] only uses the number of results in different facets to guide users. [30] introduces a model that allows dynamic constraints to be applied on the facets. [29], [32] build a decision tree with minimum height on the data so that the efforts users spend in order to reach the relevant tuples will be minimized.

Faceted search is different from our search problem on both the goal and the methodology. Faceted search systems aim to enhance the retrieval of individual documents or tuples by providing users exploration guide, and need a sequence of user inputs. With an orthogonal purpose, our work aims to rank aggregated objects/cells in a cube using only a one-time keyword query, i.e., the input is a set of keywords, and the output is a ranked list of objects/cells. For example, in Section 5.3, to find a stylish laptop using keywords, our system directly tells the user which models are commented to be stylish (e.g., (toshiba, *, red)), while faceted search tells the user which dimension(s) she/he should drill down to find relevant tuples/documents.

Keyword-based Search and OLAP in Data Cube. [34] studies answering keyword queries on RDB using minimal group-bys, which is the work most relevant to ours. For a keyword query against a multidimensional text database, it aims to find the minimal cells containing all (or some of) the query terms in the aggregated text data. “Minimal” here means there is no descendant of this cell containing more query terms. But, it is unnecessary that documents (cells) with more query terms are more relevant. And, [34] does not score or rank the answers. So when the number of returned answers is large (e.g., a thousand), it is difficult for the user to browse all the answers.

Another relevant work is keyword-driven analytical processing (KDAP) [35]. Motivated by a different scenario, KDAP supports interactive data exploration using keywords. Candidate subspaces are output to disambiguate the keyword terms. Efficiency is not a major concern in KDAP.

A keyword-driven OLAP system is proposed in [36], [37]. It populates static and dynamic dimensions, allowing for roll-up and drill-down operations based on the content and the link-structure of the dynamically selected document subset. [36] and [37] still focus on the ranking of rows in the databases, instead of aggregated cells in our work.

Ranking Cube. Another related work is Ranking Cube [38], which addresses top-*k* queries with multidimensional selection using semi-materialization and semi-online computation model. Different from our work, it does not support keyword queries. More importantly, it does not aim to find aggregations (cells) of rows. It can be applied for ranking base cells or rows (initialization part of our algorithms),

which, however, is not the major computational bottleneck in our problem of ranking cells in all levels.

8 CONCLUSIONS AND FUTURE WORK

We have proposed and studied the problem of keyword-based top- k search in text cube (*i.e.*, multidimensional text data). Flexible query language and relevance scoring formula are developed. Four efficient algorithms are designed for this problem. Among them, the most efficient one, search-space ordering approach, uses different search space pruning techniques for finding the top- k relevant cells by exploring only a small portion of the whole text cube (when k is small). Extensive performance studies are conducted to verify the efficiency and effectiveness of these approaches.

For future work, (i) it is interesting to compare the effectiveness of different scoring formulas experimentally in different real datasets; (ii) it is also interesting to design efficient algorithms for different forms of relevance scoring formula *rel*; and, (iii) while our approaches work in a bottom-up manner (starting from rows or base cells), it is interesting to study whether there is any efficient top-down approach (*e.g.*, starting from the 0-dim cell).

ACKNOWLEDGMENTS

The work was supported in part by the NASA Aviation Safety Program, Integrated Vehicle Health Management Project by NASA grant NNX08AC35A, the U.S. National Science Foundation grant IS-09-05215, an HP Research grant, Microsoft research Women's Scholarship, and the Network Science Collaborative Technology Alliance Program (NS-CTA/INARC) of U.S. Army Research Lab (ARL) under the contract number W911NF-09-2-0053. Any opinions, findings, and conclusions expressed here are those of the authors and do not necessarily reflect the views of the funding agencies. The U.S. Government is authorized to reproduce and distribute reprints for Government purposes notwithstanding any copyright notation here on.

This paper is an extended version of [39]. We thank the anonymous reviewers for their numerous insights and suggestions that immensely improved the paper.

REFERENCES

- [1] C. X. Lin, B. Ding, J. Han, F. Zhu, and B. Zhao, "Text cube: Computing ir measures for multidimensional text database analysis," in *Proc. IEEE Int'l Conf. Data Mining (ICDM)*, 2008, pp. 905–910.
- [2] S. Agrawal, S. Chaudhuri, and G. Das, "Dbxplorer: A system for keyword-based search over relational databases," in *Proc. IEEE Int'l Conf. Data Eng. (ICDE)*, 2002, pp. 5–16.
- [3] G. Bhalotia, A. Hulgeri, C. Nakhe, S. Chakrabarti, and S. Sudarshan, "Keyword searching and browsing in databases using banks," in *Proc. IEEE Int'l Conf. Data Eng. (ICDE)*, 2002, pp. 431–440.
- [4] B. Ding, J. X. Yu, S. Wang, L. Qin, X. Zhang, and X. Lin, "Finding top- k min-cost connected trees in databases," in *Proc. IEEE Int'l Conf. Data Eng. (ICDE)*, 2007, pp. 836–845.
- [5] L. Qin, J. X. Yu, L. Chang, and Y. Tao, "Querying communities in relational databases," in *Proc. IEEE Int'l Conf. Data Eng. (ICDE)*, 2009, pp. 724–735.
- [6] V. Hristidis, L. Gravano, and Y. Papakonstantinou, "Efficient ir-style keyword search over relational databases," in *Proc. Int'l Conf. Very Large Data Bases (VLDB)*, 2003, pp. 850–861.
- [7] V. Hristidis and Y. Papakonstantinou, "Discover: Keyword search in relational databases," in *Proc. Int'l Conf. Very Large Data Bases (VLDB)*, 2002, pp. 670–681.
- [8] V. Kacholia, S. Pandit, S. Chakrabarti, S. Sudarshan, R. Desai, and H. Karambelkar, "Bidirectional expansion for keyword search on graph databases," in *Proc. Int'l Conf. Very Large Data Bases (VLDB)*, 2005, pp. 505–516.
- [9] K. Golenberg, B. Kimelfeld, and Y. Sagiv, "Keyword proximity search in complex data graphs," in *Proc. ACM SIGMOD*, 2008, pp. 927–940.
- [10] H. He, H. Wang, J. Yang, and P. S. Yu, "Blinks: ranked keyword searches on graphs," in *Proc. ACM SIGMOD*, 2007, pp. 305–316.
- [11] G. Li, B. C. Ooi, J. Feng, J. Wang, and L. Zhou, "Ease: an effective 3-in-1 keyword search method for unstructured, semi-structured and structured data," in *Proc. ACM SIGMOD*, 2008, pp. 903–914.
- [12] F. Liu, C. T. Yu, W. Meng, and A. Chowdhury, "Effective keyword search in relational databases," in *Proc. ACM SIGMOD*, 2006, pp. 563–574.
- [13] Y. Luo, X. Lin, W. Wang, and X. Zhou, "Spark: top- k keyword query in relational databases," in *Proc. ACM SIGMOD*, 2007, pp. 115–126.
- [14] B. Kimelfeld and Y. Sagiv, "Finding and approximating top- k answers in keyword proximity search," in *Proc. ACM Symp. Principles of Database Syst. (PODS)*, 2006, pp. 173–182.
- [15] —, "Efficient engines for keyword proximity search," in *Proc. Int'l Workshop on the Web and Databases (WebDB)*, 2005, pp. 67–72.
- [16] —, "Efficiently enumerating results of keyword search," in *Proc. Int'l Workshop on Database Programming Languages (DBPL)*, 2005, pp. 58–73.
- [17] J. L. Elsas, J. Arguello, J. Callan, and J. G. Carbonell, "Retrieval and feedback models for blog feed search," in *Proc. Int'l Conf. on Research and Development in Information Retrieval (SIGIR)*, 2008, pp. 347–354.
- [18] S. E. Robertson, S. Walker, and M. Hancock-Beaulieu, "Okapi at trec-7: Automatic ad hoc, filtering, vlc and interactive," in *Proc. Text REtrieval Conf. (TREC)*, 1998, pp. 199–210.
- [19] A. Singhal, J. Choi, D. Hindle, D. Lewis, and F. Pereira, "At&t at trec-7," in *Proc. Text REtrieval Conf. (TREC)*, 1998, pp. 239–252.
- [20] T. H. Cormen, C. E. Leiserson, R. L. Rivest, and C. Stein, *Introduction to Algorithm (2nd Edition)*. The MIT Press, 2001.
- [21] R. Fagin, A. Lotem, and M. Naor, "Optimal aggregation algorithms for middleware," *J. Comput. Syst. Sci.*, vol. 66, no. 4, pp. 614–656, 2003.
- [22] S. Chaudhuri, R. Ramakrishnan, and G. Weikum, "Integrating db and ir technologies: What is the sound of one hand clapping?" in *Proc. Conf. on Innovative Data Syst. Research (CIDR)*, 2005, pp. 1–12.
- [23] S. Amer-Yahia, P. Case, T. Rölleke, J. Shanmugasundaram, and G. Weikum, "Report on the db/ir panel at sigmod 2005," *SIGMOD Record*, vol. 34, no. 4, pp. 71–74, 2005.
- [24] G. Weikum, "Db&ir: both sides now," in *Proc. ACM SIGMOD*, 2007, pp. 25–30.
- [25] D. Zhang, C. Zhai, and J. Han, "Topic cube: Topic modeling for olap on multidimensional text databases," in *Proc. SIAM Int'l Conf. on Data Mining (SDM)*, 2009, pp. 1123–1134.
- [26] M. A. Hearst, A. Elliott, J. English, R. R. Sinha, K. Swearingen, and K.-P. Yee, "Finding the flow in web site search," *Commun. ACM*, vol. 45, no. 9, pp. 42–49, 2002.
- [27] V. Sinha and D. R. Karger, "Magnet: Supporting navigation in semistructured data environments," in *Proc. ACM SIGMOD*, 2005, pp. 97–106.
- [28] D. Dash, J. Rao, N. Megiddo, A. Ailamaki, and G. M. Lohman, "Dynamic faceted search for discovery-driven analysis," in *Proc. Int'l Conf. on Information and Knowledge Management (CIKM)*, 2008, pp. 3–12.
- [29] S. B. Roy, H. Wang, G. Das, U. Nambiar, and M. K. Mohania, "Minimum-effort driven dynamic faceted search in structured databases," in *Proc. Int'l Conf. on Information and Knowledge Management (CIKM)*, 2008, pp. 13–22.
- [30] O. Ben-Yitzhak, N. Golbandi, N. Har'El, R. Lempel, A. Neumann, S. Ofek-Koifman, D. Sheinwald, E. J. Shekita, B. Sznajder, and S. Yogev, "Beyond basic faceted search," in *Proc. Int'l Conf. on Web Search and Web Data Mining (WSDM)*, 2008, pp. 33–44.
- [31] J. Koren, Y. Z. 0001, and X. Liu, "Personalized interactive faceted search," in *Proc. Int'l World Wide Web Conf. (WWW)*, 2008, pp. 477–486.

- [32] S. B. Roy, H. Wang, U. Nambiar, G. Das, and M. K. Mohania, "Dynacet: Building dynamic faceted search systems over databases," in *Proc. IEEE Int'l Conf. Data Eng. (ICDE)*, 2009, pp. 1463–1466.
- [33] C. Li, N. Yan, S. B. Roy, L. Lisham, and G. Das, "Facetedpedia: dynamic generation of query-dependent faceted interfaces for wikipedia," in *Proc. Int'l World Wide Web Conf. (WWW)*, 2010, pp. 651–660.
- [34] B. Zhou and J. Pei, "Answering aggregate keyword queries on relational databases using minimal group-bys," in *Proc. Int'l Conf. Extending Database Technology (EDBT)*, 2009, pp. 108–119.
- [35] P. Wu, Y. Sismanis, and B. Reinwald, "Towards keyword-driven analytical processing," in *Proc. ACM SIGMOD*, 2007, pp. 617–628.
- [36] A. Simitis, A. Baid, Y. Sismanis, and B. Reinwald, "Multidimensional content exploration," *Proc. VLDB Endowment (PVLDB)*, vol. 1, no. 1, pp. 660–671, 2008.
- [37] A. Baid, A. Balmin, H. Hwang, E. Nijkamp, J. Rao, B. Reinwald, A. Simitis, Y. Sismanis, and F. van Ham, "Dbpubs: multidimensional exploration of database publications," *Proc. VLDB Endowment (PVLDB)*, vol. 1, no. 2, pp. 1456–1459, 2008.
- [38] D. Xin, J. Han, H. Cheng, and X. Li, "Answering top-k queries with multi-dimensional selections: The ranking cube approach," in *Proc. Int'l Conf. Very Large Data Bases (VLDB)*, 2006, pp. 463–475.
- [39] B. Ding, B. Zhao, C. X. Lin, J. Han, and C. Zhai, "Topcells: Keyword-based search of top-k aggregated documents in text cube," in *Proc. IEEE Int'l Conf. Data Eng. (ICDE)*, 2010, pp. 381–384.

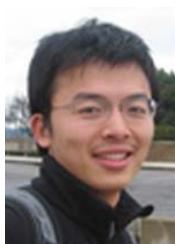


Jiawei Han is a Professor of Computer Science at the University of Illinois. He has served on the program committees of the major international conferences in the fields of data mining and database systems, and also served or is serving on the editorial boards for Data Mining and Knowledge Discovery, IEEE Transactions on Knowledge and Data Engineering, Journal of Computer Science and Technology, and Journal of Intelligent Information Systems. He is the founding Editor-in-Chief of ACM Transactions on Knowledge Discovery from Data (TKDD). Jiawei has received IBM Faculty Awards, HP Innovation Awards, ACM SIGKDD Innovation Award (2004), IEEE Computer Society Technical Achievement Award (2005), and IEEE W. Wallace McDowell Award (2009). He is a Fellow of ACM and IEEE. He is currently the Director of Information Network Academic Research Center (INARC) supported by the Network Science-Collaborative Technology Alliance (NS-CTA) program of U.S. Army Research Lab. His book "Data Mining: Concepts and Techniques" (Morgan Kaufmann) has been used worldwide as a textbook.



Bolin Ding received the BS degree in maths and applied mathematics from Renmin University of China (2005), and the MPhil degree in system engineering from the Chinese University of Hong Kong (2007). He is currently working toward the PhD degree in the Department of Computer Science, University of Illinois at Urbana-Champaign. His research interests include the mining and analysis of structured data, efficient and effective searching and indexing, and the design of randomized algorithms and approximation algorithms with applications in databases.

design of randomized algorithms and approximation algorithms with applications in databases.



Bo Zhao received the BS degree in computer science from Fudan University, Shanghai, in 2007. He is currently a PhD candidate in the Department of Computer Science, University of Illinois at Urbana-Champaign. His research interests include data mining, text mining, information retrieval and social networks.



Cindy Xide Lin received the BS degree from Tsinghua University in 2007 and the Master degree from University of Illinois at Urbana-Champaign in 2009, both in computer science. She is currently working toward her PhD degree in the Department of Computer Science, University of Illinois at Urbana-Champaign. Her research interests include data mining and information retrieval, especially topics related to web mining and social network. She has published more than

10 papers in international conferences, including SIGKDD, VLDB, WWW, ICDE, ICDM, EMNLP, CIKM, PAKDD and CIDU.



Chengxiang Zhai is an Associate Professor of Computer Science at the University of Illinois at Urbana-Champaign, where he also holds a joint appointment at the Institute for Genomic Biology, Statistics, and the Graduate School of Library and Information Science. He received a Ph.D. in Computer Science from Nanjing University in 1990, and a Ph.D. in Language and Information Technologies from Carnegie Mellon University in 2002. He worked at Clairvoyance Corp. as

a Research Scientist and a Senior Research Scientist from 1997 to 2000. His research interests include information retrieval, text mining, natural language processing, machine learning, and bioinformatics. He is an Associate Editor of ACM Transactions on Information Systems, and Information Processing and Management, and serves on the editorial board of Information Retrieval Journal. He is a program co-chair of ACM CIKM 2004, NAACL HLT 2007, and ACM SIGIR 2009. He is an ACM Distinguished Scientist, and received the 2004 Presidential Early Career Award for Scientists and Engineers (PECASE), the ACM SIGIR 2004 Best Paper Award, an Alfred P. Sloan Research Fellowship in 2008, and an IBM Faculty Award in 2009.



Ashok Srivastava, Ph.D. is the Principal Investigator for the Integrated Vehicle Health Management research project at NASA. His current research focuses on the development of data mining algorithms for anomaly detection in massive data streams, kernel methods in machine learning, and text mining algorithms. Dr. Srivastava is also the leader of the Intelligent Data Understanding group at NASA Ames Research Center. The group performs research and development of advanced machine learning and data mining algorithms in support of NASA missions. He performs data mining research in a number of areas in aviation safety and application domains such as earth sciences to study global climate processes and astrophysics to help characterize the large-scale structure of the universe. Dr. Srivastava is the author of many research articles in data mining, machine learning, and text mining, and has edited a book on Text Mining: Classification, Clustering, and Applications (with Mehran Sahami, 2009). He is currently editing two more books: Advances in Machine Learning and Data Mining for Astronomy (with Kamal Ali, Michael Way, and Jeff Scargle) and Data Mining in Systems Health Management (with Jiawei Han). Dr. Srivastava has given seminars at numerous international conferences. He has won numerous awards including the IEEE Computer Society Technical Achievement Award for "pioneering work in Intelligent Information Systems," the NASA Exceptional Achievement Medal for contributions to state-of-the-art data mining and analysis, the NASA Distinguished Performance Award, several NASA Group Achievement Awards, the IBM Golden Circle Award, and the Department of Education Merit Fellowship.

advanced machine learning and data mining algorithms in support of NASA missions. He performs data mining research in a number of areas in aviation safety and application domains such as earth sciences to study global climate processes and astrophysics to help characterize the large-scale structure of the universe. Dr. Srivastava is the author of many research articles in data mining, machine learning, and text mining, and has edited a book on Text Mining: Classification, Clustering, and Applications (with Mehran Sahami, 2009). He is currently editing two more books: Advances in Machine Learning and Data Mining for Astronomy (with Kamal Ali, Michael Way, and Jeff Scargle) and Data Mining in Systems Health Management (with Jiawei Han). Dr. Srivastava has given seminars at numerous international conferences. He has won numerous awards including the IEEE Computer Society Technical Achievement Award for "pioneering work in Intelligent Information Systems," the NASA Exceptional Achievement Medal for contributions to state-of-the-art data mining and analysis, the NASA Distinguished Performance Award, several NASA Group Achievement Awards, the IBM Golden Circle Award, and the Department of Education Merit Fellowship.



Nikunj C. Oza received his B.S. in Mathematics with Computer Science from the Massachusetts Institute of Technology (MIT) in 1994. He received his M.S. (in 1998) and Ph.D. (in 2001) in Computer Science from the University of California at Berkeley. He then joined NASA Ames Research Center as a research scientist, a member of the Intelligent Data Understanding group, and leads of a team applying data mining to the problem of health management for aerospace vehicles.

His 40+ research papers represent his research interests which include machine learning (especially ensemble learning and online learning), data mining, fault detection, combinations of machine learning with other areas of Artificial Intelligence, and their applications to Aeronautics and Earth Science. He was the Co-chair of the 2005 International Workshop on Multiple Classifier Systems (MCS), program committee member of several MCS and SIAM Data Mining (SDM) conferences, and serves as Associate Editor of the journal Information Fusion (Elsevier). He also received the 2007 Arch T. Colwell Award in 2007 for co-authoring one of the five most innovative technical papers selected out of over 3300 SAE technical papers published in 2005.