

# TAPER: A Two-Step Approach for All-Strong-Pairs Correlation Query in Large Databases

Hui Xiong, *Member, IEEE*, Shashi Shekhar, *Fellow, IEEE*,  
Pang-Ning Tan, *Member, IEEE*, and Vipin Kumar, *Fellow, IEEE*

**Abstract**—Given a user-specified minimum correlation threshold  $\theta$  and a market-basket database with  $N$  items and  $T$  transactions, an all-strong-pairs correlation query finds all item pairs with correlations above the threshold  $\theta$ . However, when the number of items and transactions are large, the computation cost of this query can be very high. The goal of this paper is to provide computationally efficient algorithms to answer the all-strong-pairs correlation query. Indeed, we identify an upper bound of Pearson's correlation coefficient for binary variables. This upper bound is not only much cheaper to compute than Pearson's correlation coefficient, but also exhibits special monotone properties which allow pruning of many item pairs even without computing their upper bounds. A Two-step All-strong-Pairs correlation query (TAPER) algorithm is proposed to exploit these properties in a filter-and-refine manner. Furthermore, we provide an algebraic cost model which shows that the computation savings from pruning is independent of or improves when the number of items is increased in data sets with Zipf-like or linear rank-support distributions. Experimental results from synthetic and real-world data sets exhibit similar trends and show that the TAPER algorithm can be an order of magnitude faster than brute-force alternatives. Finally, we demonstrate that the algorithmic ideas developed in the TAPER algorithm can be extended to efficiently compute negative correlation and uncentered Pearson's correlation coefficient.

**Index Terms**—Association analysis, data mining, Pearson's correlation coefficient, statistical computing.

## 1 INTRODUCTION

GIVEN a large set of items and observation data sets about cooccurring items, association analysis is concerned with identification of strongly related (e.g., as measured by Pearson's correlation coefficient [22]) subsets of items. Association analysis is a core problem in data mining and databases. It plays an important role in many application domains such as market-basket analysis [2], climate studies [25], public health [8], and bioinformatics [17]. For instance, association analysis in market-basket study can reveal how the sales of a product are related to the sales of other products. This type of information can be useful for sales promotions, catalog design, and store layout.

The focus of this paper is on computing an *all-strong-pairs correlation query* that returns pairs of high positively correlated items (or binary attributes). The all-strong-pairs correlation query problem can be formalized as follows: Given a user-specified minimum correlation threshold  $\theta$  and a market-basket database with  $N$  items and  $T$  transactions, an all-strong-pairs correlation query finds all item pairs with correlations above the minimum correlation threshold  $\theta$ .

However, as the number of items and transactions in the data set increases, the computation cost for an all-strong-

pairs correlation query becomes prohibitively expensive. For example, consider a database of  $10^6$  items, which may represent the collection of books available at an e-commerce Web site. Answering the all-strong-pairs correlation query from such a massive database requires computing the correlations of  $\binom{10^6}{2} \approx 0.5 \times 10^{12}$  possible item pairs. Thus, it may not be computationally feasible to apply a brute-force approach to compute correlations for all half-trillion pairs, particularly when the number of transactions in the data set is also very large.

### 1.1 Related Work

Jermaine [13] investigated the implication of incorporating chi-square ( $\chi^2$ ) [22] based queries to data-cube computations. He showed that finding the subcubes that satisfy statistical tests such as  $\chi^2$  is inherently NP-hard, but can be made more tractable using approximation schemes. Jermaine [14] also presented an iterative procedure for high-dimensional correlation analysis by shaving off part of the database via feedback from human experts. Finally, Brin et al. [4] proposed a  $\chi^2$ -based correlation rule mining strategy. However,  $\chi^2$  does not possess a desired upward closure property for exploiting efficient computation [9].

This paper focuses on the efficient computation of statistical correlation for all pairs of items. Given  $n$  items, a traditional brute-force approach computes Pearson's correlation coefficient for all  $\binom{n}{2} = \frac{n(n-1)}{2}$  item pairs. This approach is often implemented using matrix algebra in a statistical software package as the "correlation matrix" [15] function, which computes Pearson's correlation coefficient for all pairs of columns. This approach is applicable to but not efficient for the case of Boolean matrices, which can model market-basket-type data sets. Recently, Ilyas et al. [12] proposed a more efficient method for identifying correlated

• H. Xiong is with the Management Science and Information Systems Department, Rutgers University, 180 University Avenue, Newark, NJ 07102. E-mail: hui@rbs.rutgers.edu.

• S. Shekhar and V. Kumar are with the Department of Computer Science and Engineering, University of Minnesota, 200 Union Street SE, Minneapolis, MN 55455. E-mail: {shekhar, kumar}@cs.umn.edu.

• P.-N. Tan is with the Department of Computer Science and Engineering, Michigan State University, East Lansing, MI 48824-1226. E-mail: ptan@cse.msu.edu.

Manuscript received 6 Nov. 2004; revised 15 June 2005; accepted 23 Aug. 2005; published online 17 Feb. 2006.

For information on obtaining reprints of this article, please send e-mail to: tkde@computer.org, and reference IEEECS Log Number TKDE-0394-1104.

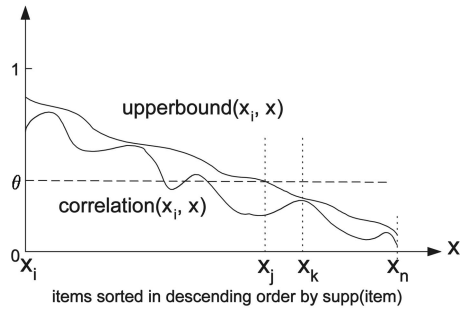


Fig. 1. Illustration of the filtering techniques. (The curves are only used for illustration purposes.)

pairs. In this method, the sampling techniques are applied to exploit efficient computation. As a result, this method cannot avoid false-positive and false-negative correlations.

In contrast, unlike the correlation-matrix approach, the method proposed in this paper does not need to compute all  $\binom{n}{2}$  pairs. In particular, for market-basket-type data sets with a Zipf-like rank-support distribution, we show that only a small portion of the item pairs needs to be examined. In the real world, Zipf-like distributions have been observed in a variety of application domains, including commercial retail data, Web click streams, and telecommunication data. Also, we show that our method is complete and correct since we do not apply any approximation schemes, such as sampling techniques.

## 1.2 Contributions

In our preliminary work [28], we provide an upper bound of Pearson's correlation coefficient for binary variables. The computation of this upper bound is much cheaper than the computation of the exact correlation, since this upper bound can be computed as a function of the support of individual items. Furthermore, we show that this upper bound has a special 1D monotone property that allows elimination of many item pairs even without computing their upper bounds, as shown in Fig. 1. The  $x$ -axis in the figure represents the set of items having a lower level of support than the support for item  $x_i$ . These items are sorted from left to right in decreasing order of their individual support values. The  $y$ -axis indicates the correlation between each item  $x$  and item  $x_i$ .  $Upperbound(x_i, x)$  represents the upper bound of  $correlation(x_i, x)$  and has a monotone decreasing behavior. This behavior guarantees that an item pair  $(x_i, x_k)$  can be pruned if there exists an item  $x_j$  such that  $upperbound(x_i, x_j) < \theta$  and  $supp(x_k) < supp(x_j)$ .

A Two-step All-strong-Pairs correlation query (TAPER) algorithm is proposed to exploit this 1D monotone property in a filter-and-refine manner, which consists of two steps: filtering and refinement. In the filtering step, many item pairs are filtered out using the easy-to-compute  $upperbound(x_i, x)$  and its monotone property. In the refinement step, the exact correlation is computed for remaining pairs to determine the final query results. In addition, we have proved the completeness and correctness of TAPER and provided an algebraic cost model to quantify the computational savings. As demonstrated by our experiments on both real and synthetic data sets, TAPER can be an order of magnitude faster than brute-force alternatives and the computational savings by TAPER are independent or improve when the number of items is increased in data sets with common Zipf [29] or linear rank-support distributions.

In this paper, we identify that the upper bound of Pearson's correlation coefficient for binary variables has special 2D monotone properties. Indeed, besides the 1D monotone property mentioned above, the upper bound has another 1-D monotone property, as shown in Fig. 2a. The  $x$ -axis in the figure represents a sorted item list in ascending order of their individual support values. The  $y$ -axis indicates the correlation between item  $x_i$  and  $x$ . In this case, the item  $x_i$  is fixed and the upper bound of  $correlation(x, x_i)$  is monotonically decreasing with the increase of the support value of item  $x$ .

Fig. 2 shows 2D monotone properties of the upper bound. In the figure, for an item list  $\{1, 2, 3, 4, 5, 6\}$ , which is sorted by item support in nonincreasing order, the upper bound of item pairs is decreasing following the arrow direction. For instance, the upper bound of item pair  $\{5, 6\}$  is greater than that of item pair  $\{4, 6\}$ . Also, the upper bound of item pair  $\{1, 2\}$  is greater than that of item pair  $\{1, 3\}$ .

With 2D monotone properties of the upper bound, we can further refine the TAPER algorithm by reducing the upper bound computation in the coarse-filter step. We show that the number of upper bound computations is reduced from  $\frac{n(n-1)}{2}$  to  $2n - 3$  for the worst case. In addition, we present experimental results to show this computational improvement.

Finally, the method proposed in this paper is not limited to finding all pairs of high positively correlated pairs. We show that the algorithmic ideas developed in the TAPER algorithm can also be extended for identifying pairs of high negatively correlated pairs and for efficiently computing uncentered Pearson's correlation coefficient. To this end, we provide the theoretical basis, algorithmic ideas, and experimental results for such an extension.

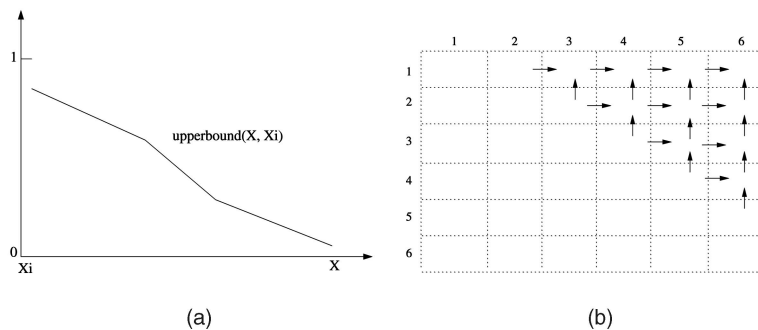


Fig. 2. (a) A monotone property of the upper bound. (b) Illustration of 2D monotone properties of the upper bound.

TABLE 1  
A Two-Way Table of Item A and Item B

		B		Row Total
		0	1	
A	0	$P_{(00)}$	$P_{(01)}$	$P_{(0+)}$
	1	$P_{(10)}$	$P_{(11)}$	$P_{(1+)}$
Column Total		$P_{(+0)}$	$P_{(+1)}$	$N$

### 1.3 Scope and Outline

The scope of the all-strong-pairs correlation query problem proposed in this paper is restricted to market-basket databases with binary variables and the correlation computational form is Pearson's correlation coefficient for binary variables, which is also called the  $\phi$  correlation coefficient.

Note that the all-strong-pairs correlation query problem is different from the standard association rule mining problem [1], [3], [4], [5], [6], [7], [10], [11], [19], [20], [21], [26]. Given a set of transactions, the objective of association-rule mining is to extract all subsets of items that satisfy a minimum support threshold. Support measures the fraction of transactions that contain a particular subset of items. The notions of support and correlation may not necessarily agree with each other. This is because item pairs with high support may be poorly correlated, while those that are highly correlated may have very low support. For instance, suppose we have an item pair  $\{A, B\}$ , where  $supp(A) = supp(B) = 0.8$  and  $supp(A, B) = 0.64$ . Both items are uncorrelated because  $supp(A, B) = supp(A)supp(B)$ . In contrast, an item pair  $\{A, B\}$  with  $supp(A) = supp(B) = supp(A, B) = 0.001$  is perfectly correlated despite its low support. Patterns with low support but high correlation are useful for capturing interesting associations among rare anomalous events or rare but expensive items such as gold necklaces and earrings.

The remainder of this paper is organized as follows: Section 2 presents basic concepts. In Section 3, we introduce the upper bound of Pearson's correlation coefficient for binary variables. Section 4 proposes the TAPER algorithm. In Section 5, we analyze the TAPER algorithm in the areas of completeness, correctness, and computation gain. Section 6 discusses how to generalize our method. In Section 7, we present the experimental results. Finally, we draw conclusions and suggest future work in Section 8.

## 2 PEARSON'S CORRELATION COEFFICIENT

In statistics, a measure of association is a numerical index which describes the strength or magnitude of a relationship among variables. Although literally dozens of measures exist, they can be categorized into two broad groups: ordinal and nominal. Relationships among ordinal variables can be analyzed with ordinal measures of association such as Kendall's tau [16] and Spearman's rank correlation coefficient [18]. In contrast, relationships among nominal variables can be analyzed with nominal measures of association such as Pearson's correlation coefficient and measures based on chi-square [22].

The  $\phi$  correlation coefficient [22] is the computation form of Pearson's correlation coefficient for binary variables. In this section, we describe the  $\phi$  correlation coefficient and show how it can be computed using the support measure of association-rule mining [1].

TID	Items
1	1, 2, 3
2	1, 2, 3
3	1, 3
4	1, 2
5	1, 2
6	1, 2
7	1, 2, 3, 4, 5, 6
8	1, 2, 4, 5
9	1, 2, 4
10	3

(a)

Pair	Correlation	Support
{1, 2}	0.667	0.8
{1, 3}	-0.333	0.4
{1, 4}	0.218	0.3
{1, 5}	0.167	0.2
{1, 6}	0.111	0.1
{2, 3}	-0.5	0.3
{2, 4}	0.327	0.3
{2, 5}	0.25	0.2
{2, 6}	0.167	0.1
{3, 4}	-0.218	0.1
{3, 5}	0	0.1
{3, 6}	0.333	0.1
{4, 5}	0.764	0.2
{4, 6}	0.509	0.1
{5, 6}	0.667	0.1

(b)

Fig. 3. An example database.

In a  $2 \times 2$  two-way table shown in Table 1, the calculation of the  $\phi$  correlation coefficient reduces to

$$\phi = \frac{P_{(00)}P_{(11)} - P_{(01)}P_{(10)}}{\sqrt{P_{(0+)}P_{(1+)}P_{(+0)}P_{(+1)}}, \quad (1)$$

where  $P_{(ij)}$ , for  $i = 0, 1$  and  $j = 0, 1$ , denotes the number of samples which are classified in the  $i$ th row and  $j$ th column of the table. Furthermore, we let  $P_{(i+)}$  denote the total number of samples classified in the  $i$ th row and we let  $P_{(+j)}$  denote the total number of samples classified in the  $j$ th column. Thus,  $P_{(i+)} = \sum_{j=0}^1 P_{(ij)}$  and  $P_{(+j)} = \sum_{i=0}^1 P_{(ij)}$ . In the two-way table, given that  $N$  is the total number of samples, we can transform (1) as follows:

$$\begin{aligned} \phi &= \frac{(N - P_{(01)} - P_{(10)} - P_{(11)})P_{(11)} - P_{(01)}P_{(10)}}{\sqrt{P_{(0+)}P_{(1+)}P_{(+0)}P_{(+1)}}} \\ &= \frac{NP_{(11)} - (P_{(11)} + P_{(10)})(P_{(01)} + P_{(11)})}{\sqrt{P_{(0+)}P_{(1+)}P_{(+0)}P_{(+1)}}} \\ &= \frac{\frac{P_{(11)}}{N} - \frac{P_{(1+)}P_{(+1)}}{N}}{\sqrt{\frac{P_{(0+)}P_{(1+)}P_{(+0)}P_{(+1)}}{N}}} \end{aligned}$$

Hence, when adopting the support measure of association-rule mining [1], for two items  $A$  and  $B$  in a market-basket database, we have  $supp(A) = P_{(1+)}/N$ ,  $supp(B) = P_{(+1)}/N$ , and  $supp(A, B) = P_{(11)}/N$ . With support notations and the above new derivations of (1), we can derive the support form of the  $\phi$  correlation coefficient as shown below in (2):

$$\phi = \frac{supp(A, B) - supp(A)supp(B)}{\sqrt{supp(A)supp(B)(1 - supp(A))(1 - supp(B))}} \quad (2)$$

**Example 1.** Consider the market-basket database shown in Fig. 3. For item 1 and item 2 in the database, we can construct a two-way table as shown in Table 2. Then, by (1), we get  $\phi = \frac{1 \times 8 - 1 \times 0}{\sqrt{1 \times 9 \times 2 \times 8}} = \frac{2}{3}$ . Also, since  $supp(1, 2) = 0.8$ ,  $supp(1) = 0.9$ , and  $supp(2) = 0.8$ , by (2), we get  $\phi = \frac{0.8 - 0.8 \times 0.9}{\sqrt{0.8 \times 0.9 \times 0.1 \times 0.2}} = \frac{0.08}{0.12} = \frac{2}{3}$ . As can be seen, the results from the two equations are identical. Finally,  $\phi$  correlation coefficients for other item pairs can be computed similarly. Fig. 3b shows  $\phi$  correlation coefficients for all item pairs.

TABLE 2  
A Two-Way Table of Item 1 and Item 2

		{2}		Row Total
		0	1	
{1}	0	1	0	1
	1	1	8	9
Column Total		2	8	10

### 3 PROPERTIES OF THE $\phi$ CORRELATION COEFFICIENT

In this section, we present some properties of the  $\phi$  correlation coefficient. These properties are useful for the efficient computation of all-strong-pairs correlation queries.

#### 3.1 An Upper Bound

In this subsection, we propose an upper bound of the  $\phi$  correlation coefficient for a given pair  $\{A, B\}$  in terms of the support value of item  $A$  and the support value of item  $B$ .

**Lemma 1.** *Given an item pair  $\{A, B\}$ , the support value  $supp(A)$  for item  $A$ , and the support value  $supp(B)$  for item  $B$ , without loss of generality, let  $supp(A) \geq supp(B)$ . The upper bound  $upper(\phi_{\{A,B\}})$  of the  $\phi$  correlation coefficient for  $\{A, B\}$  can be obtained when  $supp(A, B) = supp(B)$  and*

$$upper(\phi_{\{A,B\}}) = \sqrt{\frac{supp(B)}{supp(A)}} \sqrt{\frac{1 - supp(A)}{1 - supp(B)}}. \quad (3)$$

**Proof.** According to (2), for an item pair  $\{A, B\}$ ,

$$\phi_{\{A,B\}} = \frac{supp(A, B) - supp(A)supp(B)}{\sqrt{supp(A)supp(B)(1 - supp(A))(1 - supp(B))}}.$$

When the support values  $supp(A)$  and  $supp(B)$  are fixed,  $\phi_{\{A,B\}}$  is monotone increasing with the increase of the support value  $supp(A, B)$ . By the given condition  $supp(A) \geq supp(B)$  and the antimonotone property of the support measure, we get the maximum possible value of  $supp(A, B)$  is  $supp(B)$ . As a result, the upper bound  $upper(\phi_{\{A,B\}})$  of the  $\phi$  correlation coefficient for an item pair  $\{A, B\}$  can be obtained when  $supp(A, B) = supp(B)$ . Hence,

$$upper(\phi_{\{A,B\}}) = \sqrt{\frac{supp(B)}{supp(A)}} \sqrt{\frac{1 - supp(A)}{1 - supp(B)}}.$$

□

As can be seen in (3), the upper bound of the  $\phi$  correlation coefficient for an item pair  $\{A, B\}$  relies only on the support value of item  $A$  and the support value of item  $B$ . In other words, there is no requirement to get the support value  $supp(A, B)$  of an item pair  $\{A, B\}$  for the calculation of this upper bound. As already noted, when the number of items  $N$  becomes very large, it is difficult to store the support of every item pair in the memory, since  $N(N-1)/2$  is a huge number. However, it is possible to store the support of individual items in the main memory. As a result, this upper bound can serve as a coarse filter to filter out item pairs which are of no interest, thus saving I/O cost

Pair	UPPER ( $\Phi$ )	Correlation
{1, 2}	0.667	0.667
{1, 3}	0.333	-0.333
{1, 4}	0.218	0.218
{1, 5}	0.167	0.167
{1, 6}	0.111	0.111
{2, 3}	0.5	-0.5
{2, 4}	0.327	0.327
{2, 5}	0.25	0.25
{2, 6}	0.167	0.167
{3, 4}	0.655	-0.218
{3, 5}	0.5	0
{3, 6}	0.333	0.333
{4, 5}	0.764	0.764
{4, 6}	0.509	0.509
{5, 6}	0.667	0.667

Fig. 4. All pairs in the example database.

by reducing the computation of the support values of those pruned pairs.

**Example 2.** Fig. 4 shows all item pairs with their upper bound values for the example data set shown in Fig. 3. If we consider item pair  $\{1, 2\}$ , then, by (3), we have

$$\begin{aligned} upper(\phi_{\{1,2\}}) &= \sqrt{\frac{supp(2)}{supp(1)}} \sqrt{\frac{1 - supp(1)}{1 - supp(2)}} = \sqrt{\frac{0.8}{0.9}} \sqrt{\frac{1 - 0.9}{1 - 0.8}} \\ &= 0.667, \end{aligned}$$

as shown in Fig. 4. Finally, upper bounds of the  $\phi$  correlation coefficient for other item pairs can be computed in a similar manner. Fig. 4 shows the upper bounds and  $\phi$  correlation coefficients for all item pairs.

#### 3.2 2D Monotone Properties

In this section, we present a conditional monotone property of the upper bound of the  $\phi$  correlation coefficient as shown below in Lemma 2.

**Lemma 2.** *For an item pair  $\{A, B\}$ , if we let  $supp(A) > supp(B)$  and fix item  $A$ , the upper bound  $upper(\phi_{\{A,B\}})$  of  $\{A, B\}$  is monotone decreasing with the decrease of the support value of item  $B$ .*

**Proof.** By Lemma 1, we get

$$upper(\phi_{\{A,B\}}) = \sqrt{\frac{supp(B)}{supp(A)}} \sqrt{\frac{1 - supp(A)}{1 - supp(B)}}.$$

For any given two items  $B_1$  and  $B_2$  with  $supp(A) > supp(B_1) > supp(B_2)$ , we need to prove  $upper(\phi_{\{A,B_1\}}) > upper(\phi_{\{A,B_2\}})$ . This claim can be proved as follows:

$$\frac{upper(\phi_{\{A,B_1\}})}{upper(\phi_{\{A,B_2\}})} = \frac{\sqrt{\frac{supp(B_1)}{supp(A)}} \sqrt{\frac{1 - supp(A)}{1 - supp(B_1)}}}{\sqrt{\frac{supp(B_2)}{supp(A)}} \sqrt{\frac{1 - supp(A)}{1 - supp(B_2)}}} > 1.$$

This follows the given condition that  $supp(B_1) > supp(B_2)$  and  $(1 - supp(B_1)) < (1 - supp(B_2))$ . □

**TAPER ALGORITHM**

Input:  $S'$ : an item list sorted by item supports in non-increasing order.  
 $\theta$ : a user-specified minimum correlation threshold.  
Output:  $P$ : the result of all-strong-pairs correlation query.  
Variables:  $n$ : the size of item set  $S'$ .  
 $A$ : the item with larger support.  
 $B$ : the item with smaller support.

**CoarseFilter**( $S', \theta$ ) //The filtering Step

1DFiltering( $S', \theta$ ) or 2DFiltering( $S', \theta$ )

**Refine**( $A, B, \theta$ ) //The Refinement Step

1. Get the support  $\text{supp}(A, B)$  of item set  $\{A, B\}$
2.  $\phi = \frac{\text{supp}(A, B) - \text{supp}(A)\text{supp}(B)}{\sqrt{\text{supp}(A)\text{supp}(B)(1 - \text{supp}(A))(1 - \text{supp}(B))}}$
3. **if**  $\phi < \theta$  **then**
4.     **return**  $\emptyset$  //return NULL
5.     **else**
6.         **return**  $\{\{A, B\}, \phi\}$

Fig. 5. The TAPER algorithm.

Along the same line of Lemma 2, we can also derive another conditional monotone property of the upper bound of the  $\phi$  correlation coefficient as follows:

**Lemma 3.** For a pair of items  $\{A, B\}$ , let  $\text{supp}(A) > \text{supp}(B)$  and fix the item  $B$ , the upper bound  $\text{upper}(\phi_{\{A, B\}})$  of  $\{A, B\}$  is monotone increasing with the decreasing of the support value of item  $A$ .

**Example 3.** This example illustrates Lemma 2 and Lemma 3. As shown in Fig. 4, if item 1 is fixed, the upper bounds of item pairs  $\{1, 2\}$ ,  $\{1, 3\}$ ,  $\{1, 4\}$ ,  $\{1, 5\}$ , and  $\{1, 6\}$  are in a decreasing order, which follows Lemma 2. Also, the upper bounds of item pairs  $\{1, 6\}$ ,  $\{2, 6\}$ ,  $\{3, 6\}$ ,  $\{4, 6\}$ ,  $\{5, 6\}$  are in an increasing order. This follows Lemma 3.

Lemma 2 and Lemma 3 are 2D monotone properties of the upper bound. These two lemmas allow us to push the upper bound into the search algorithm, thus efficiently pruning the search space. In the following section, we will introduce how this pruning works.

## 4 ALGORITHM DESCRIPTIONS

Here, we first present the Two-step All-strong-Pairs correlation query (TAPER) algorithm to exploit the proposed upper bound and its monotone properties.

### 4.1 Overview

The TAPER algorithm is a filter-and-refine query processing strategy which consists of two steps: filtering and refinement.

**The Filtering Step.** In this step, the TAPER algorithm applies two pruning techniques. The first technique uses the upper bound of the  $\phi$  correlation coefficient as a coarse filter. In other words, if the upper bound of the  $\phi$  correlation coefficient for an item pair is less than the user-specified correlation threshold, we can prune this item pair right away. The second pruning technique prunes item pairs based on special monotone properties of the upper bound of the  $\phi$  correlation coefficient.

**The Refinement Step.** In the refinement step, the TAPER algorithm computes the exact correlation for each surviving pair from the filtering step and retrieves the pairs with correlations above the user-specified minimum correlation threshold as the query results.

Fig. 5 shows the pseudocode of the TAPER algorithm, which includes the *CoarseFilter* and *Refine* procedures. In

**1DFiltering**( $S', \theta$ ) //The Filtering Step

1.  $n = \text{size}(S')$ ,  $P = \emptyset$
2.     **for**  $i$  from 1 to  $n-1$
3.          $A = S'[i]$
4.         **for**  $j$  from  $i+1$  to  $n$
5.              $B = S'[j]$
6.              $\text{upper}(\phi) = \sqrt{\frac{\text{supp}(B)}{\text{supp}(A)}} \sqrt{\frac{1 - \text{supp}(A)}{1 - \text{supp}(B)}}$
7.             **if**( $\text{upper}(\phi) < \theta$ ) **then** //Pruning by the monotone property
8.                 **break** from inner loop
9.             **else**
10.                  $P = P \cup \text{Refine}(A, B, \theta)$

Fig. 6. TAPER 1D: TAPER with a 1D filter.

the *CoarseFilter* procedure, monotone properties of the upper bound of the  $\phi$  correlation coefficient can be used in two different ways, namely, 1D filtering and 2D filtering. The details of these two filtering techniques are introduced in the following sections.

Procedure *Refine* works as follows: Line 1 gets the support for the item pair  $\{A, B\}$ . Note that the I/O cost can be very expensive for line 1 when the number of items is large since we cannot store the support of all item pairs in the memory. Line 2 calculates the exact correlation coefficient of this item pair. If the correlation is greater than the user-specified minimum correlation threshold, this item pair is returned as a query result in line 6. Otherwise, the procedure returns NULL in line 4.

### 4.2 TAPER 1D: TAPER with 1D Filter

In this section, we illustrate TAPER with 1D filter, denoted as TAPER 1D. The working mechanism of 1D filtering is illustrated as the following corollary.

**Corollary 1.** When searching for all pairs of items with correlations above a user-specified threshold  $\theta$ , if an item list  $\{i_1, i_2, \dots, i_m\}$  is sorted by item supports in nonincreasing order, an item pair  $\{i_a, i_c\}$  with  $\text{supp}(i_a) > \text{supp}(i_c)$  can be pruned if  $\text{upper}(\phi_{\{i_a, i_b\}}) < \theta$  and  $\text{supp}(i_c) \leq \text{supp}(i_b)$ .

**Proof.** First, when  $\text{supp}(i_c) = \text{supp}(i_b)$ , we get  $\text{upper}(\phi_{\{i_a, i_c\}}) = \text{upper}(\phi_{\{i_a, i_b\}}) < \theta$  according to (3) and the given condition  $\text{upper}(\phi_{\{i_a, i_b\}}) < \theta$ ; then we can prune the item pair  $\{i_a, i_c\}$ . Next, we consider  $\text{supp}(i_c) < \text{supp}(i_b)$ . Since  $\text{supp}(i_a) > \text{supp}(i_b) > \text{supp}(i_c)$ , by Lemma 2, we get  $\text{upper}(\phi_{\{i_a, i_c\}}) < \text{upper}(\phi_{\{i_a, i_b\}}) < \theta$ . Hence, the pair  $\{i_a, i_c\}$  is pruned.  $\square$

Fig. 6 shows the pseudocode of 1D filtering, which works as follows: Line 1 initializes the variables and creates an empty query result set  $P$ . Lines 2-10 iteratively enumerate candidate pairs and filter out item pairs whose correlations are obviously less than the user-specified correlation threshold  $\theta$ . This is implemented by two loops. Line 2 starts an outer loop. Line 3 specifies reference item  $A$  and line 4 starts a search within each branch. Line 5 specifies the target item  $B$ , and line 6 computes the upper bound of the  $\phi$  correlation coefficient for item pair  $\{A, B\}$ . In line 7, if this upper bound is less than the user-specified correlation threshold  $\theta$ , the search within this loop can stop by exiting from the inner loop, as shown in line 8. The reason is as follows: First, the reference item  $A$  is fixed in each branch and it has the maximum support value due to the way we construct the branch. Also, items within each branch are sorted based on their support in nonincreasing order. Then, by Lemma 2, the upper bound of the  $\phi$  correlation coefficient for the item pair  $\{A, B\}$  is monotone decreasing with the decrease of the support of item  $B$ . Hence, if we find

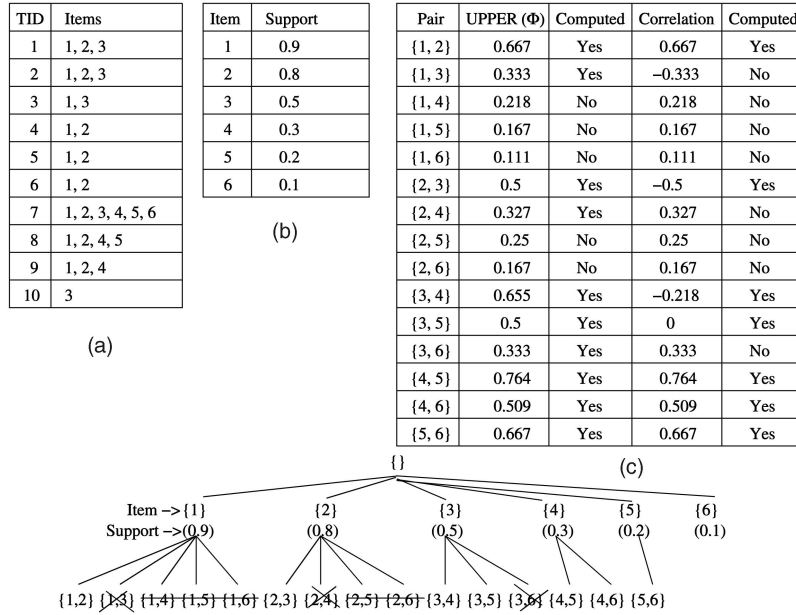


Fig. 7. Illustration of the filter-and-refine strategy of the TAPER algorithm with 1D filter.

the first target item B which results in an upper bound  $upper(\phi_{\{A,B\}})$  that is less than the user-specified correlation threshold  $\theta$ , we can stop the search. Line 10 calls the procedure *Refine* to compute the exact correlation for each surviving candidate pair and continues to check the next target item until no target item is left.

**Example 4.** To illustrate the TAPER algorithm with 1D filter, consider a database shown in Fig. 7. To simplify the discussion, we use an item list  $\{1, 2, 3, 4, 5, 6\}$  that is sorted by item support in nonincreasing order. For a given correlation threshold 0.36, we can use Rymon's generic set-enumeration tree search framework [24] to demonstrate how filter-and-refine query processing works. For instance, for the branch starting from item 1, we identify that the upper bound of the  $\phi$  correlation coefficient for the item pair  $\{1, 3\}$  is 0.333, which is less than the given correlation threshold 0.36. Hence, we can prune this item pair immediately. Also, since the item list  $\{1, 2, 3, 4, 5, 6\}$  is sorted by item supports in nonincreasing order, we can prune pairs  $\{1, 4\}$ ,  $\{1, 5\}$ , and  $\{1, 6\}$  by Lemma 2 without any further computation cost. In contrast, for the traditional filter-and-refine paradigm, the coarse filter can only prune the item pair  $\{1, 3\}$ . There is no technique to prune item pairs  $\{1, 4\}$ ,  $\{1, 5\}$ , and  $\{1, 6\}$ . Finally, in the refinement step, only seven item pairs are required to compute the exact correlation coefficients, as shown in Fig. 7c. More than half of the item pairs are pruned in the filter step even though the correlation threshold is as low as 0.36. Please note that the Rymon's set-enumeration tree is used for illustration purposes. In our algorithm, there is no requirement to construct such a tree structure.

### 4.3 TAPER 2D: TAPER with 2D Filter

The coarse filter step of TAPER 1D can be improved by reducing the number of upper bounds to be computed; this leads to TAPER 2D. Indeed, we can reduce the computation of upper bounds in each inner loop and produce an improved coarse filter step as shown in Fig. 8. The key difference between TAPER 1D and TAPER 2D is that

TAPER 2D records the break point in the last inner loop and starts computing upper bounds from the recorded point instead of going through every candidate pair. The correctness of this additional filtering is guaranteed by the following corollary:

**Corollary 2.** *Given an item list  $\{i_1, i_2, \dots, i_n\}$ , which is sorted by item supports in nonincreasing order, the upper bound computation of an item pair  $\{i_a, i_c\}$  with  $supp(\{i_a\}) > supp(\{i_c\})$  can be saved if  $upper(\phi_{\{i_{a-1}, i_c\}}) > \theta$ , where  $supp(\{i_{a-1}\}) > supp(\{i_c\})$ .*

**Proof.** First, we have  $supp(\{i_{a-1}\}) > supp(\{i_a\})$ , since items are sorted by support in nonincreasing order. Also, by Lemma 3, we have  $upper(\phi_{\{i_{a-1}, i_c\}}) < upper(\phi_{\{i_a, i_c\}})$ . Therefore,  $upper(\phi_{\{i_a, i_c\}}) > \theta$ . As a result, we do not need to compute the upper bound of  $\{i_a, i_c\}$ .  $\square$

TAPER 2D is an improvement over TAPER 1D. The following Lemma 4 shows that the number of upper bound computations is reduced from  $\frac{n(n-1)}{2}$  in TAPER 1D to  $2n-3$  in TAPER 2D for the worst case.

**Lemma 4.** *The number of upper bounds required to compute in the coarse filter step of TAPER 2D is  $2n-3$ , where  $n$  is the number of objects in the data set.*

**Proof.** Fig. 8 shows that the number of upper bounds to be computed is determined by the value  $k$ . For each outer loop, there is an inner loop where the upper bound is computed.  $k$  starts from 2 and ends at  $n$ . In each outer loop  $i$ , there are at most  $k_i - k_{i-1} + 1$  number of upper bound computations. In total, there are  $\sum_{i=1}^{n-1} (k_i - k_{i-1} + 1) = 2n - 3$  number of upper bounds to be computed.  $\square$

**Example 5.** This example illustrates the coarse-filter steps of the TAPER algorithm with 1D and 2D filtering. Here, we use the data set shown in Fig. 7. For the item list  $\{1, 2, 3, 4, 5, 6\}$ , which is sorted by item support in nonincreasing order, if the correlation threshold is 0.2, the coarse filter step of TAPER 1D is illustrated in Fig. 9a. In the figure, the dot line indicates the item pairs whose upper bounds need to be computed. As can be seen, there are 14 upper-bound computations. In contrast, as shown in Fig. 9b,

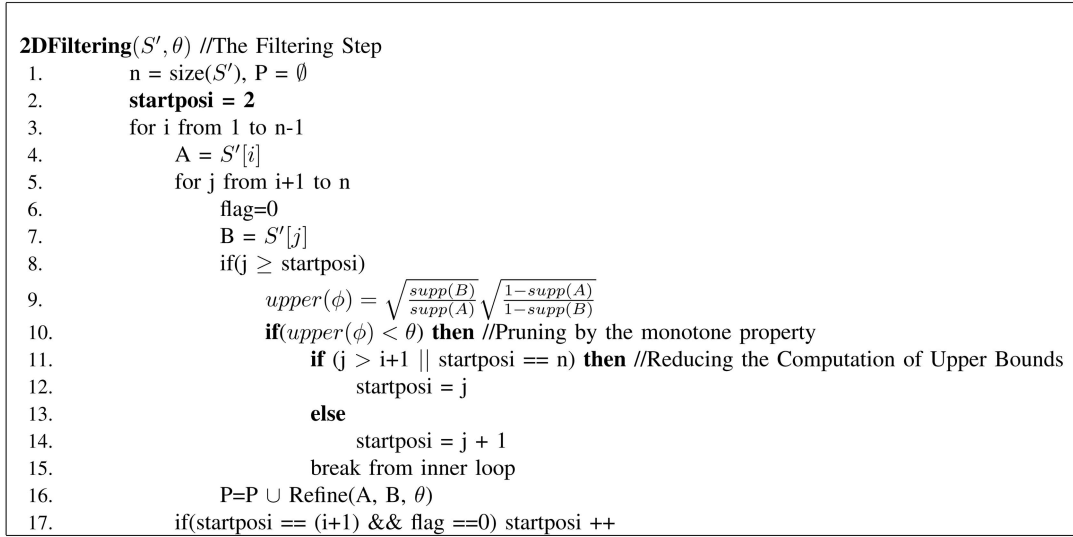


Fig. 8. TAPER 2D: TAPER with a 2D filter.

TAPER 2D computes upper bounds for  $2n - 3 = 2 \times 6 - 3 = 9$  number of item pairs, which are  $\{1, 2\}, \{1, 3\}, \{1, 4\}, \{1, 5\}, \{2, 5\}, \{2, 6\}, \{3, 6\}, \{4, 6\},$  and  $\{5, 6\}$ .

## 5 ANALYSIS OF THE TAPER ALGORITHM

In this section, we analyze TAPER in the areas of completeness, correctness, and computation savings. Note that TAPER here will stand for both TAPER 1D and TAPER 2D to simplify the discussion.

### 5.1 Completeness and Correctness

**Lemma 5.** *The TAPER algorithm is complete. In other words, this algorithm finds all pairs which have correlations above a user-specified minimum correlation threshold.*

**Proof.** The completeness of the TAPER algorithm can be shown by the following two facts: First, all item pairs in the database have the opportunity to be checked during the iteration process. Second, the filtering step only prunes item pairs if the upper bounds of the  $\phi$  correlation coefficient for these pairs are less than the user-specified correlation threshold. This is guaranteed by Corollary 1 as well as Corollary 2. Also, the refinement step only

prunes item pairs whose correlations are less than the user-specified correlation threshold. The second fact guarantees that all the pruned item pairs cannot have correlations above the user-specified minimum correlation threshold.  $\square$

**Lemma 6.** *The TAPER algorithm is correct. In other words, every pair that both algorithms find has a correlation above a user-specified minimum correlation threshold.*

**Proof.** The correctness of the TAPER algorithm can be guaranteed by the refinement step of these two algorithms since the exact correlation of each candidate pair is calculated in the refinement step and every pair with a correlation lower than the user-specified correlation threshold will be pruned.  $\square$

### 5.2 Quantifying the Computation Savings

This section presents analytical results for the amount of computational savings obtained by TAPER. First, we illustrate the relationship between the choices of the minimum correlation threshold and the size of the reduced search space (after performing the filtering step). Knowing the relationship gives us an idea of the amount of pruning achieved using the upper bound function of correlation.

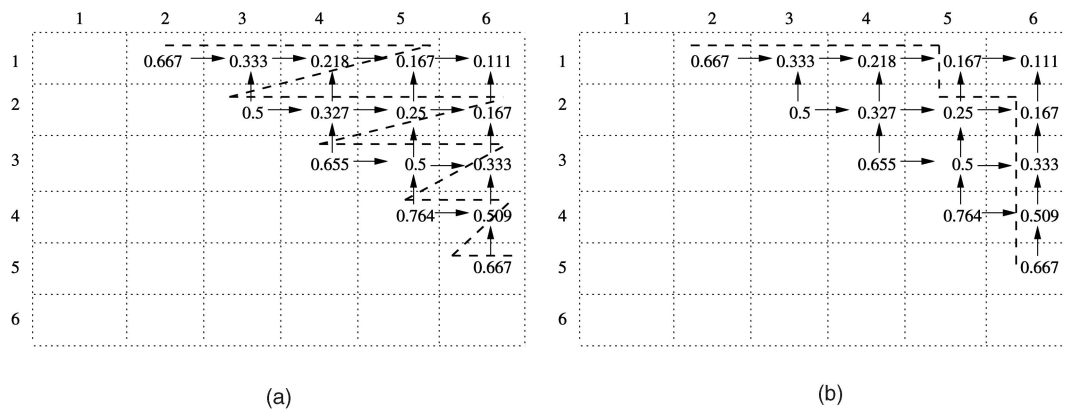


Fig. 9. (a) TAPER 1D versus (b) TAPER 2D.

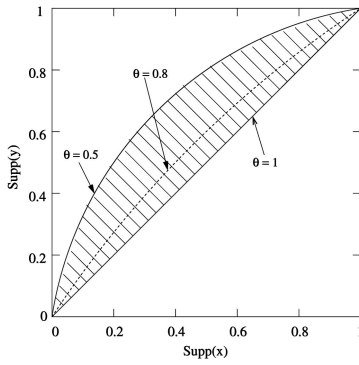


Fig. 10. An illustration of the reduced search space for the refinement step of the TAPER algorithm. Only item pairs within the shaded region must be computed for their correlation.

Fig. 10 illustrates a two-dimensional plot for every possible combination of support pairs,  $supp(x)$  and  $supp(y)$ . If we impose the constraint that  $supp(x) \leq supp(y)$ , then all item pairs must be projected to the upper left triangle since the diagonal line represents the condition  $supp(x) = supp(y)$ .

To determine the size of the reduced search space, let us start from the upper bound of the correlation:

$$\begin{aligned} upper(\phi_{x,y}) &= \sqrt{\frac{supp(x)}{supp(y)}} \sqrt{\frac{1 - supp(x)}{1 - supp(y)}} < \theta \\ \implies supp(x)(1 - supp(y)) &< \theta^2 supp(y)(1 - supp(x)) \quad (4) \\ \implies supp(y) &> \frac{supp(x)}{\theta^2 + (1 - \theta^2)supp(x)}. \end{aligned}$$

The above inequality provides a lower bound on  $supp(y)$  such that any item pair involving  $x$  and  $y$  can be pruned using the conditional monotone property of the upper bound function. In other words, any surviving item pair that undergoes the refinement step must violate the condition given in (4). These item pairs are indicated by the shaded region shown in Fig. 10. During the refinement step, TAPER has to compute the exact correlation for all item pairs that fall in the shaded region between the diagonal and the polyline drawn by (5).

$$supp(y) = \frac{supp(x)}{\theta^2 + (1 - \theta^2)supp(x)} \quad (5)$$

As can be seen from Fig. 10, the size of the reduced search space depends on the choice of minimum correlation threshold. If we increase the threshold from 0.5 to 0.8, the search space for the refinement step is reduced substantially. When the correlation threshold is 1.0, the polyline from (5) overlaps with the diagonal line. In this limit, the search space for the refinement step becomes zero.

The above analysis shows only the size of the reduced search space that must be explored during the refinement step of the TAPER algorithm. The actual amount of pruning achieved by TAPER depends on the support distribution of items in the database. To facilitate our discussion, we first introduce the definitions of several concepts used in the remainder of this section.

**Definition 1.** The pruning ratio of the TAPER algorithm is defined by the following equation:

$$\gamma(\theta) = \frac{S(\theta)}{T}, \quad (6)$$

where  $\theta$  is the minimum correlation threshold,  $S(\theta)$  is the number of item pairs which are pruned before computing their exact correlations at the correlation threshold  $\theta$ , and  $T$  is the total number of item pairs in the database. For a given database,  $T$  is a fixed number and is equal to  $\binom{n}{2} = \frac{n(n-1)}{2}$ , where  $n$  is the number of items.

**Definition 2.** For a sorted item list, the rank-support function  $f(k)$  is a discrete function which presents the support in terms of the rank  $k$ .

For a given database, let  $I = \{A_1, A_2, \dots, A_n\}$  be an item list sorted by item supports in nonincreasing order. Then item  $A_1$  has the maximum support and the rank-support function  $f(k) = supp(A_k)$ ,  $\forall 1 \leq k \leq n$ , which is monotone decreasing with the increase of the rank  $k$ . To quantify the computation savings for a given item  $A_j$  ( $1 \leq j < n$ ) at the threshold  $\theta$ , we need to find only the first item  $A_l$  ( $j < l \leq n$ ) such that  $upper(\phi_{\{A_j, A_l\}}) < \theta$ . By Lemma 2, if  $upper(\phi_{\{A_j, A_l\}}) < \theta$ , we can guarantee that  $upper(\phi_{\{A_j, A_i\}}) < \theta$ , where  $l \leq i \leq n$ , is less than the correlation threshold  $\theta$ . In other words, all these  $n - l + 1$  pairs can be pruned without a further computation requirement. According to Lemma 1, we get

$$\begin{aligned} upper(\phi_{\{A_j, A_l\}}) &= \sqrt{\frac{supp(A_l)}{supp(A_j)}} \sqrt{\frac{1 - supp(A_j)}{1 - supp(A_l)}} < \sqrt{\frac{supp(A_l)}{supp(A_j)}} \\ &= \sqrt{\frac{f(l)}{f(j)}} < \theta \implies \frac{f(l)}{f(j)} < \theta^2. \end{aligned}$$

Since the rank-support function  $f(k)$  is monotone decreasing with the increase of the rank  $k$ , we get  $l > f^{-1}(\theta^2 f(j))$ . To make the computation simple, we let  $l = f^{-1}(\theta^2 f(j)) + 1$ . Therefore, for a given item  $A_j$  ( $1 < j \leq n$ ), the computation cost for  $(n - f^{-1}(\theta^2 f(j)))$  item pairs can be saved. As a result, the total computation savings of the TAPER algorithm is shown below in (7). Note that the computation savings shown in (7) is an underestimated value of the achieved computation savings.

$$S(\theta) = \sum_{j=2}^n \{n - f^{-1}(\theta^2 f(j))\} \quad (7)$$

Finally, we conduct computation savings analysis on the data sets with some special rank-support distributions. Specifically, we consider three special rank-support distributions: a uniform distribution, a linear distribution, and a generalized Zipf distribution [29], as shown in the following three cases.

### 5.2.1 CASE I: A Uniform Distribution

In this case, the rank-support function  $f(k) = C$ , where  $C$  is a constant. According to (3), the upper bound of the  $\phi$  correlation coefficient for any item pair is 1, which is the maximum possible value for the correlation. Hence, for any given item  $A_j$ , we cannot find an item  $A_l$  ( $j < l \leq n$ ) such that  $upper(\phi_{\{A_j, A_l\}}) < \theta$ , where  $\theta \leq 1$ . As a result, the total computation savings  $S(\theta)$  is zero.

### 5.2.2 CASE II: A Linear Distribution

In this case, the rank-support function has a linear distribution and  $f(k) = a - mk$ , where  $m$  is the absolute value of the slope and  $a$  is the intercept.



**Lemma 7.** When a database has a linear rank-support distribution  $f(k)$  and  $f(k) = a - mk$  ( $a > 0$ ,  $m > 0$ ), for a user-specified minimum correlation threshold  $\theta$ , the pruning ratio of the TAPER algorithm increases with the decrease of the ratio  $a/m$ , the increase of the correlation threshold  $\theta$ , and the increase of the number of items, where  $0 < \theta \leq 1$ .

**Proof.** For the given database, let  $I = \{A_1, A_2, \dots, A_n\}$  be the item list sorted by item support in nonincreasing order. Then, item  $A_1$  has the maximum support. Also, let the rank-support function  $f(k) = a - mk$ , where  $m$  is the absolute value of the slope and  $a$  is the intercept. From the rank-support function  $f(k)$ , we can derive the inverse function  $f^{-1}(y) = \frac{a-y}{m}$ . Thus,  $f^{-1}(\theta^2 f(j)) = \frac{a-\theta^2(a-mj)}{m} = \frac{a}{m}(1-\theta^2) + j\theta^2$ . According to (7), we can get

$$\begin{aligned} S(\theta) &= \sum_{j=2}^n \{n - f^{-1}(\theta^2 f(j))\} \\ &= n(n-1) - \sum_{j=2}^n \frac{a}{m}(1-\theta^2) - \sum_{j=2}^n j\theta^2 \\ &= n(n-1) - \frac{a}{m}(n-1)(1-\theta^2) - \frac{(n-1)(n+2)}{2}\theta^2 \\ &= (n-1) \left( \frac{n-2}{2} - \frac{a}{m} - \frac{n+2}{2} \right) (1-\theta^2). \end{aligned}$$

Since the pruning ratio  $\gamma(\theta) = \frac{S(\theta)}{T}$  and  $T = \frac{n(n-1)}{2}$ ,

$$\Rightarrow \gamma(\theta) = \frac{(n-2) - (2\frac{a}{m} - (n+2))(1-\theta^2)}{n}.$$

Also, we know  $\text{supp}(A_n) = f(n) = a - mn > 0 \Rightarrow 2\frac{a}{m} > 2n \geq (n+2)$ , when  $n \geq 2$ .

Thus, we can derive three rules as follows:

$$\begin{aligned} \text{rule 1: } \theta \nearrow &\Rightarrow (1-\theta^2) \searrow \Rightarrow \gamma(\theta) \nearrow, \\ \text{rule 2: } a/m \searrow &\Rightarrow (2\frac{a}{m} - (n+2)) \searrow \Rightarrow \gamma(\theta) \nearrow, \\ \text{rule 3: } n \nearrow &\Rightarrow (2\frac{a}{m} - (n+2))/n \searrow \Rightarrow \gamma(\theta) \nearrow. \end{aligned}$$

Therefore, the claim that the pruning ratio of the TAPER algorithm is increased with the decrease of the ratio  $a/m$ , the increase of the correlation threshold  $\theta$ , and the increase of the number of items holds.  $\square$

### 5.2.3 CASE III: A Generalized Zipf Distribution

In this case, the rank-support function has a generalized Zipf distribution and  $f(k) = \frac{c}{k^p}$ , where  $c$  and  $p$  are constants and  $p \geq 1$ . When  $p$  is equal to 1, the rank-support function has a Zipf distribution.

**Lemma 8.** When a database has a generalized Zipf rank-support distribution  $f(k)$  and  $f(k) = \frac{c}{k^p}$ , for a user-specified minimum correlation threshold  $\theta$ , the pruning ratio of the TAPER algorithm increases with the increase of  $p$  and the correlation threshold  $\theta$ , where  $0 < \theta \leq 1$ . Furthermore, the pruning ratio is independent when the number of items is increased.

**Proof.** Since the rank-support function  $f(k) = \frac{c}{k^p}$ , the inverse function  $f^{-1}(y) = (\frac{c}{y})^{\frac{1}{p}}$ . Accordingly,

$$f^{-1}(\theta^2 f(j)) = (\frac{c}{\theta^2 \frac{c}{j^p}})^{\frac{1}{p}} = \frac{j}{(\theta^2)^{\frac{1}{p}}}$$

Applying (7), we get

$$\begin{aligned} S(\theta) &= \sum_{j=2}^n \{n - f^{-1}(\theta^2 f(j))\} \\ &= n(n-1) - \sum_{j=2}^n \frac{j}{(\theta^2)^{\frac{1}{p}}} \\ &= n(n-1) - \frac{(n-1)(n+2)}{2} \frac{1}{\theta^{\frac{2}{p}}}. \end{aligned}$$

Since the pruning ratio  $\gamma(\theta) = \frac{S(\theta)}{T}$  and  $T = \frac{n(n-1)}{2}$ ,

$$\Rightarrow \gamma(\theta) = 2 - \frac{n+2}{n} \frac{1}{\theta^{\frac{2}{p}}}.$$

Thus, we can derive three rules as follows:

$$\begin{aligned} \text{rule 1: } \theta \nearrow &\Rightarrow \frac{n+2}{n} \frac{1}{\theta^{\frac{2}{p}}} \searrow \Rightarrow \gamma(\theta) \nearrow, \\ \text{rule 2: } p \nearrow &\Rightarrow \frac{n+2}{n} \frac{1}{\theta^{\frac{2}{p}}} \searrow \Rightarrow \gamma(\theta) \nearrow, \\ \text{rule 3: } n \rightarrow \infty &\Rightarrow \lim_{n \rightarrow \infty} \frac{n+2}{n} \frac{1}{\theta^{\frac{2}{p}}} = \frac{1}{\theta^{\frac{2}{p}}}. \end{aligned}$$

Therefore, the claim that the pruning ratio of the TAPER algorithm increases with the increase of  $p$  and the correlation threshold  $\theta$  holds. Also, Rule 3 indicates that the pruning ratio is independent when the number of items is increased in data sets with Zipf distributions.  $\square$

### 5.3 Dominance Zone Analysis

Here, we provide simple algebraic cost models for the computational cost of the brute-force algorithm and the TAPER algorithm. We assume that the total number of objects in the data set is  $n$ .

The main cost of the brute-force algorithm is the cost of computing  $\frac{n(n-1)}{2}$  number of exact correlation coefficients. Let  $C_{corr}$  indicate the cost of computing the correlation coefficient for an item pair. The cost model of the brute-force algorithm is given as follows:

$$Cost_{Brute} = O(n^2) * C_{corr}. \quad (8)$$

The main cost of the TAPER algorithm consists of three parts: the sorting cost denoted by  $C_{sort}$ , the cost of computing upper bounds denoted by  $C_{tu}$ , and the cost of computing exact correlation coefficients denoted by  $C_{tc}$ . The cost model of the TAPER algorithm is given as follows:

$$\begin{aligned} Cost_{TAPER} &= C_{sort} + C_{tu} + C_{tc} \\ &= O(n \log n) + O(n) * C_{upper} + (1 - \gamma(\theta))O(n^2) \\ &\quad * C_{corr}. \end{aligned}$$

Where  $C_{upper}$  indicates the cost of computing the upper bound for an item pair,

$$\begin{aligned} Cost_{Brute} - Cost_{TAPER} &= \gamma(\theta)O(n^2) * C_{corr} - O(n \log n) - O(n) * C_{tu}. \end{aligned}$$

With the increase of  $n$ , the computation savings can be more dramatic for the TAPER algorithm, particularly for data sets with nonuniform rank-support distributions, such as Zipf distributions. Note that  $C_{corr}$  is much larger

than  $C_{upper}$  since there is a very expensive computation cost of finding support values for item pairs when computing exact correlation coefficients. Indeed, when the number of objects is large, we cannot store the support values of all item pairs in the memory, that is, we may need to scan the whole data set once in order to find the support value of every item pair.

## 6 DISCUSSION

In this section, we extend our algorithm for finding item pairs with strong negative correlations and demonstrate that the algorithmic ideas developed here can also be applied to some other association measures, such as uncentered Pearson's correlation coefficients.

### 6.1 Negative Correlations

In this paper, our focus is to find all pairs of high positively correlated items. However, in some application domains, there may be interest in knowing pairs of high negatively correlated items [27]. In the following, we present a lower bound of the  $\phi$  correlation coefficient.

**Lemma 9.** *Given a pair of items  $\{A, B\}$ , without loss of generality, let  $supp(A) \geq supp(B)$ . The lower bound,  $lower(\phi_{\{A,B\}})$ , of the  $\phi$  correlation coefficient is equal to*

$$\begin{cases} -\frac{\sqrt{supp(A)supp(B)}}{\sqrt{(1-supp(A))(1-supp(B))}} & \text{if } supp(A) + supp(B) \leq 1 \\ -\frac{\sqrt{(1-supp(A))(1-supp(B))}}{\sqrt{supp(A)supp(B)}} & \text{if } supp(A) + supp(B) > 1 \end{cases}$$

**Proof.** According to Equation 2, for an item pair  $\{A, B\}$

$$\phi_{\{A,B\}} = \frac{supp(A, B) - supp(A)supp(B)}{\sqrt{supp(A)supp(B)(1-supp(A))(1-supp(B))}}$$

When the support values  $supp(A)$  and  $supp(B)$  are fixed,  $\phi_{\{A,B\}}$  is monotone decreasing with the decrease of  $supp(A, B)$ . Let us consider the following two cases:

**CASE 1:** *if  $supp(A) + supp(B) \leq 1$ . The minimum possible value of  $supp(A, B)$  is zero. Hence,*

$$\begin{aligned} lower(\phi_{\{A,B\}}) &= -\frac{supp(A)supp(B)}{\sqrt{supp(A)supp(B)(1-supp(A))(1-supp(B))}} \\ &= -\frac{\sqrt{supp(A)supp(B)}}{\sqrt{(1-supp(A))(1-supp(B))}} \end{aligned}$$

**CASE 2:** *if  $supp(A) + supp(B) > 1$ . The minimum possible value of  $supp(A, B)$  is equal to  $supp(A) + supp(B) - 1$ . Hence,*

$$\begin{aligned} lower(\phi_{\{A,B\}}) &= \frac{supp(A) + supp(B) - 1 - supp(A)supp(B)}{\sqrt{supp(A)supp(B)(1-supp(A))(1-supp(B))}} \\ &= -\frac{\sqrt{(1-supp(A))(1-supp(B))}}{\sqrt{supp(A)supp(B)}} \end{aligned}$$

From the above, this lemma holds.  $\square$

We also present a conditional monotone property of the lower bound of the  $\phi$  correlation coefficient.

**Lemma 10.** *For a pair  $\{A, B\}$ , let  $supp(A) > supp(B)$ . We have the following two cases: 1) If  $supp(A) + supp(B) \leq 1$  and  $supp(A)$  is fixed, the  $lower(\phi_{\{A,B\}})$  is monotone increasing with the decrease of  $supp(B)$ . 2) If  $supp(A) + supp(B) > 1$  and  $supp(B)$  is fixed, the  $lower(\phi_{\{A,B\}})$  is monotone increasing with the increase of  $supp(A)$ .*

**Proof.** Let us consider the following two cases:

**CASE 1:**  *$supp(A) + supp(B) \leq 1$  and  $supp(A)$  is fixed. By Lemma 9, we have*

$$lower(\phi_{\{A,B\}}) = -\frac{\sqrt{supp(A)supp(B)}}{\sqrt{(1-supp(A))(1-supp(B))}}$$

For any given two items  $B_1$  and  $B_2$  with  $supp(A) > supp(B_1) > supp(B_2)$ , we need to prove  $lower(\phi_{\{A,B_1\}}) < lower(\phi_{\{A,B_2\}})$ . This claim can be proven as follows:

$$\frac{|lower(\phi_{\{A,B_1\}})|}{|lower(\phi_{\{A,B_2\}})|} = \sqrt{\frac{supp(B_1)}{supp(B_2)}} \sqrt{\frac{1-supp(B_2)}{1-supp(B_1)}} > 1.$$

The above follows the given condition that  $supp(B_1) > supp(B_2)$  and  $(1-supp(B_1)) < (1-supp(B_2))$ . Since  $lower(\phi_{\{A,B\}}) < 0$ , we have

$$lower(\phi_{\{A,B_1\}}) < lower(\phi_{\{A,B_2\}}).$$

**CASE 2:**  *$supp(A) + supp(B) > 1$  and  $supp(B)$  is fixed. By Lemma 9, we have*

$$lower(\phi_{\{A,B\}}) = -\frac{\sqrt{(1-supp(A))(1-supp(B))}}{\sqrt{supp(A)supp(B)}}$$

For any given two items  $A_1$  and  $A_2$  with  $supp(A_1) > supp(A_2) > supp(B)$ , we need to prove  $lower(\phi_{\{A_1,B\}}) > lower(\phi_{\{A_2,B\}})$ . This claim can be proven as follows:

$$\frac{|lower(\phi_{\{A_1,B\}})|}{|lower(\phi_{\{A_2,B\}})|} = \sqrt{\frac{supp(A_2)}{supp(A_1)}} \sqrt{\frac{1-supp(A_1)}{1-supp(A_2)}} < 1.$$

The above follows the conditions that  $supp(A_1) > supp(A_2)$  and  $(1-supp(A_1)) < (1-supp(A_2))$ . Since  $lower(\phi_{\{A,B\}}) < 0$ , we get  $lower(\phi_{\{A_1,B\}}) > lower(\phi_{\{A_2,B\}})$ .  $\square$

Based on Lemmas 9 and 10, we can extend the TAPER algorithm to find all pairs of high negatively correlated items as well. Note that Pearson's correlation coefficient has some limitations when capturing negative correlation between items with low support. However, in this paper, our focus is on the computational perspective of Pearson's correlation coefficient.

### 6.2 An Extension to Uncentered Pearson's Correlation Coefficient

Here, we present how to extend the algorithmic ideas developed in this paper for uncentered Pearson's correlation coefficient, also known as the cosine measure [23]. Using the support notation, uncentered Pearson's correlation coefficient is defined as the following equation:

$$\langle \{A, B\} \rangle = \frac{supp(\{A, B\})}{\sqrt{supp(A)supp(B)}} \quad (9)$$

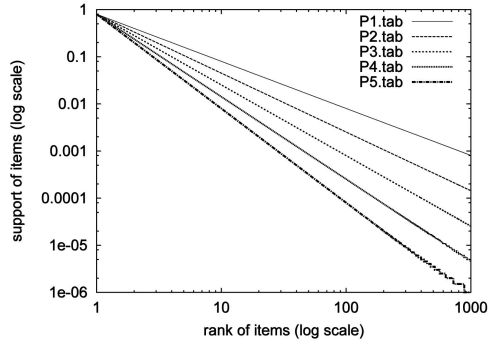


Fig. 11. The plot of the Zipf rank-support distributions of synthetic data sets in log-log scale.

Indeed, similar to the  $\phi$  correlation coefficient, uncentered Pearson's correlation coefficient has an upper bound and this upper bound has a conditional monotone property as shown in the following.

**Lemma 11.** *Given an item pair  $\{A, B\}$ , the support value  $\text{supp}(A)$  for item  $A$ , and the support value  $\text{supp}(B)$  for item  $B$ , without loss of generality, let  $\text{supp}(A) \geq \text{supp}(B)$ . The upper bound  $\text{upper}(\text{uncentered}(\{A, B\}))$  of the uncentered Pearson's correlation coefficient for an item pair  $\{A, B\}$  can be obtained when  $\text{supp}(A, B) = \text{supp}(B)$  and*

$$\text{upper}(\text{uncentered}(\{A, B\})) = \sqrt{\frac{\text{supp}(B)}{\text{supp}(A)}}.$$

**Proof.** According to (9), for an item pair  $\{A, B\}$ ,

$$\text{uncentered}(\{A, B\}) = \frac{\text{supp}(\{A, B\})}{\sqrt{\text{supp}(A)\text{supp}(B)}}.$$

By the given condition  $\text{supp}(A) \geq \text{supp}(B)$  and the antimonotone property of the support measure, the upper bound  $\text{upper}(\text{uncentered}(\{A, B\}))$  of uncentered Pearson's correlation coefficient can be obtained when  $\text{supp}(A, B) = \text{supp}(B)$ . Hence,

$$\text{upper}(\text{uncentered}(\{A, B\})) = \sqrt{\frac{\text{supp}(B)}{\text{supp}(A)}}. \quad \square$$

**Lemma 12.** *For a pair of items  $\{A, B\}$ , if we let  $\text{supp}(A) > \text{supp}(B)$  and fix item  $A$ , the upper bound of uncentered Pearson's correlation coefficient  $\text{upper}(\text{uncentered}(\{A, B\}))$*

TABLE 3  
Parameters of the Synthetic Data Sets

Data set name	T	N	C	P
P1.tab	2000000	1000	0.8	1
P2.tab	2000000	1000	0.8	1.25
P3.tab	2000000	1000	0.8	1.5
P4.tab	2000000	1000	0.8	1.75
P5.tab	2000000	1000	0.8	2

TABLE 4  
Real-Life Data Set Characteristics

Data set	#Item	#Transaction	Source
Pumsb	2113	49046	IBM Almaden
Pumsb*	2089	49046	IBM Almaden
Chess	75	3196	UCI Repository
Mushroom	119	8124	UCI Repository
Connect	127	67557	UCI Repository
LA1	29704	3204	TREC-5
Retail	14462	57671	Retail Store

for item pair  $\{A, B\}$  is monotone decreasing with the decrease of the support value of item  $B$ .

**Proof.** Since

$$\text{upper}(\text{uncentered}(\{A, B\})) = \sqrt{\frac{\text{supp}(B)}{\text{supp}(A)}},$$

if item  $A$  is fixed, the  $\text{upper}(\text{uncentered}(\{A, B\}))$  becomes one variable function with  $\text{supp}(B)$  as the variable. It can be seen that this function is monotone decreasing with the decrease of  $\text{supp}(B)$ .  $\square$

Lemma 11 and Lemma 12 are analogous to Lemma 1 and Lemma 2, which form the basis of the TAPER algorithm. Therefore, the algorithmic ideas in the TAPER algorithm can also be applied to efficiently compute uncentered Pearson's correlation coefficient.

## 7 EXPERIMENTAL RESULTS

In this section, we present the results of experiments to evaluate the performance of the TAPER algorithm. Specifically, we demonstrate

1. a performance comparison between TAPER 1D and a brute-force approach,
2. the effectiveness of the proposed algebraic cost model,
3. the scalability of the TAPER algorithm,
4. a performance evaluation of the choices between 1-D and 2-D monotone properties in the coarse-filter step of TAPER, and
5. the extension of algorithmic ideas developed in TAPER for computing negative correlation as well as the uncentered Pearson's correlation coefficient.

### 7.1 The Experimental Setup

Our experiments were performed on both real-life and synthetic data sets. Synthetic data sets were generated such that the rank-support distributions follow Zipf's law, as shown in Fig. 11. Note that, in log-log scales, the rank-support plot of a Zipf distribution will be a straight line with a slope equal to the exponent  $P$  in the Zipf distribution. A summary of the parameter settings used to generate the synthetic data sets is presented in Table 3, where  $T$  is the number of transactions,  $N$  is the number of items,  $C$  is the constant of a generalized Zipf distribution, and  $P$  is the exponent of a generalized Zipf distribution.

The real-life data sets were obtained from several different application domains. Table 4 shows some characteristics of these data sets. The first five data sets in the table, i.e., pumsb, pumsb\*, chess, mushroom, and connect, are often used as benchmarks for evaluating

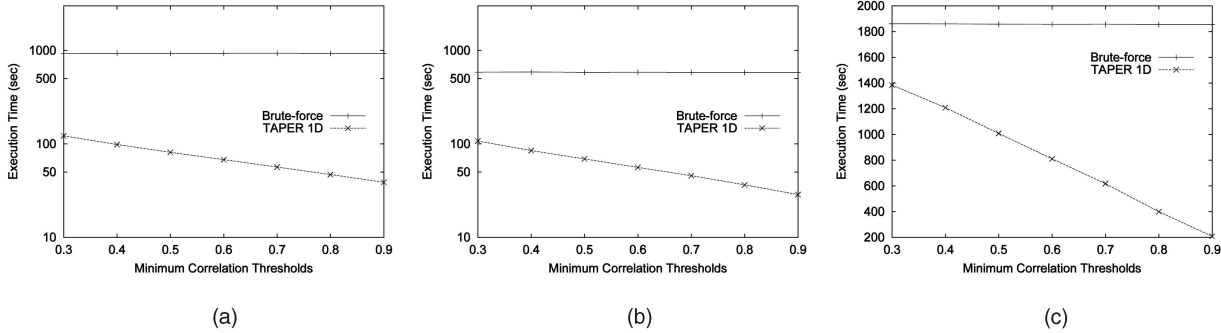


Fig. 12. TAPER 1D versus a brute-force approach on the (a) *pumbsb*, (b) *pumbsb\**, and (c) *retail* data sets.

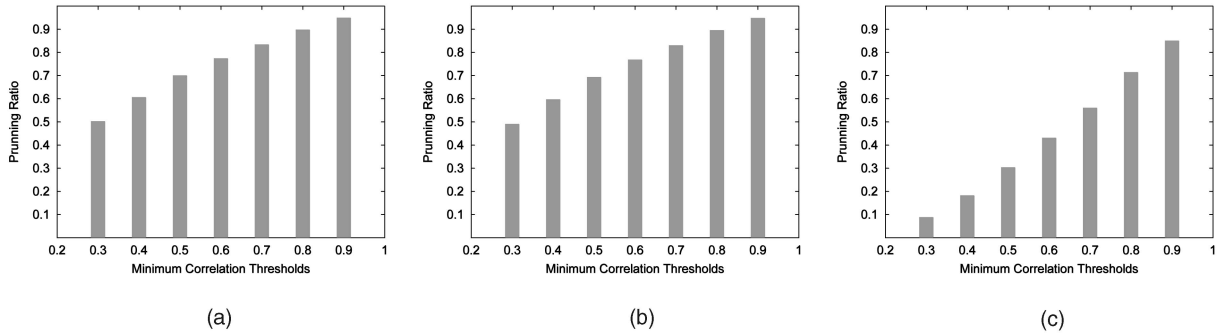


Fig. 13. The pruning effect of TAPER on (a) *pumbsb*, *pumbsb\**, and *retail* data sets.

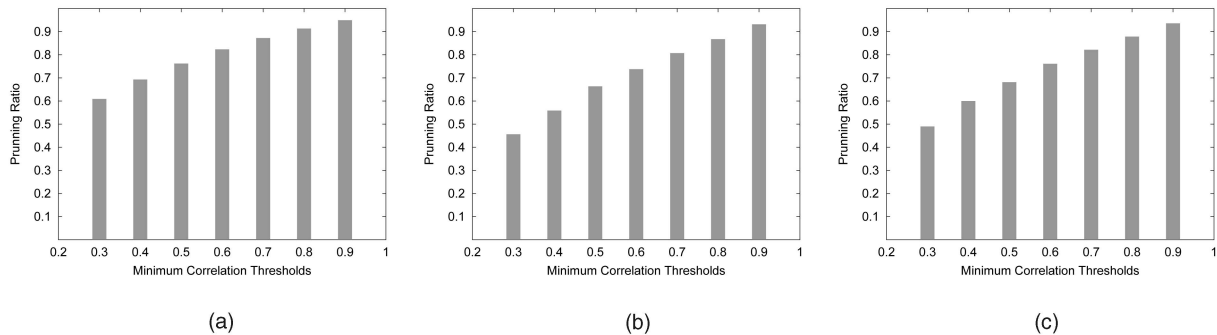


Fig. 14. The pruning effect of TAPER on UCI (a) *connect*, (b) *mushroom*, and (c) *chess* data sets.

the performance of association rule algorithms on dense data sets. The *pumbsb* and *pumbsb\** data sets correspond to binarized versions of a census data set from IBM.<sup>1</sup> The difference between them is that *pumbsb\** does not contain items with support greater than 80 percent. The *chess*, *mushroom*, and *connect* data sets are benchmark data sets from the UCI Machine Learning Repository.<sup>2</sup> The LA1 data set is part of the TREC-5 collection (<http://rec.nist.gov>) and contains news articles from the Los Angeles Times. Finally, *retail* is a masked data set obtained from a large mail-order company.

**Experimental Platform.** We implemented TAPER using C++ and all experiments were performed on a Sun Ultra 10 workstation with a 440 MHz CPU and 128 Mbytes of memory running the SunOS 5.7 operating system.

1. These data sets were obtained from IBM Almaden at <http://www.almaden.ibm.com/cs/quest/demos.html> (June 2005).

2. These data sets and data content descriptions are available at <http://www.ics.uci.edu/~mllearn/MLRepository.html> (June 2005).

## 7.2 TAPER 1-D versus the Brute-Force Approach

In this section, we present a performance comparison between TAPER 1D and a brute-force approach using several benchmark data sets from IBM, the UCI Machine Learning Repository, and some other sources, such as retail stores. The implementation of the brute-force approach is similar to that of TAPER 1D except that the filtering mechanism implemented in TAPER 1D is not included in the brute-force approach.

Fig. 12 shows the relative computation performance of TAPER 1-D and the brute-force approach on the *pumbsb*, *pumbsb\**, and *retail* data sets. As can be seen, the performance of the brute-force approach does not change much for any of the three data sets. However, the execution time of TAPER 1D can be an order of magnitude faster than the brute-force approach even if the minimum correlation threshold is low. For instance, as shown in Fig. 12a, the execution time of TAPER 1D on the *pumbsb* data set is one order of magnitude less than that of the brute-force approach at the correlation threshold 0.4. Also, when the minimum correlation threshold increases, the execution time of TAPER 1D dramatically decreases on the *pumbsb*

TABLE 5  
Groups of Items for the Retail Data Set

Group	I	II	III
# Items	4700	4700	4700
# Transactions	57671	57671	57671
$a/m$	10318	8149	4778

data set. Similar computation effects can also be observed on the `pumsb*` and `retail` data sets, although the computation savings on the `retail` data set is not as significant as it is on the other two data sets.

To better understand the above computation effects, we also present the pruning ratio of TAPER (both TAPER 1D and TAPER 2D) on these data sets. As illustrated in Fig. 13, the pruning ratio of TAPER on the `retail` data set is much smaller than that on the `pumsb` and `pumsb*` data sets. This smaller pruning ratio explains why the computation savings on `retail` is less than that on the other two data sets. Also, Fig. 14 shows the pruning ratio of TAPER on the UCI `connect`, `mushroom`, and `chess` data sets. The pruning ratios achieved on these data sets are comparable with the pruning ratio we obtained on the `pumsb` data set. This indicates that TAPER also achieves much better computation performance than the brute-force approach on UCI benchmark data sets.

### 7.3 The Effect of Correlation Thresholds

In this section, we present the effect of correlation thresholds on the computation savings of TAPER (both TAPER 1D and TAPER 2D). Recall that our algebraic cost model shows that the pruning ratio of TAPER increases with increases of the correlation thresholds for data sets with linear and Zipf-like distributions. Fig. 13 shows such an increasing trend of the pruning ratio on the `pumsb`, `pumsb*`, and `retail` data sets as correlation thresholds increase. Also, Fig. 14 shows a similar increasing trend of the pruning ratio on the UCI benchmark datasets including `mushroom`, `chess`, and `connect`.

One common feature of all the above data sets is the skewed nature of their rank-support distributions. As a result, these experimental results still exhibit a trend similar to that of the proposed algebraic cost model, although the rank-support distributions of these data sets do not follow Zipf's law exactly.

### 7.4 The Effect of the Slope $m$

Recall that the algebraic cost model for data sets with a linear rank-support distribution provides rules which indicate that the pruning ratio of TAPER (both TAPER 1D and TAPER 2D) increases with the decrease of the ratio  $a/m$  and the pruning ratio increases with the increase of the correlation threshold. In this subsection, we empirically evaluate the effect of the ratio  $a/m$  on the performance of the TAPER algorithm for data sets with a linear rank-support distribution.

First, we generated three groups of data from the retail data set by sorting all the items in the data set in non-decreasing order and then partitioning them into four groups. Each of the first three groups contains 4,700 items and the last group contains 362 items. The first three groups are the group data sets shown in Table 5. Fig. 15a shows the plot of the rank-support distribution of the retail data set and Figs. 15b, 15c, and 15d show the plots of the rank-support distributions of three groups of data generated from the retail data set. As can be seen, the rank-support distributions of the three groups approximately follow a

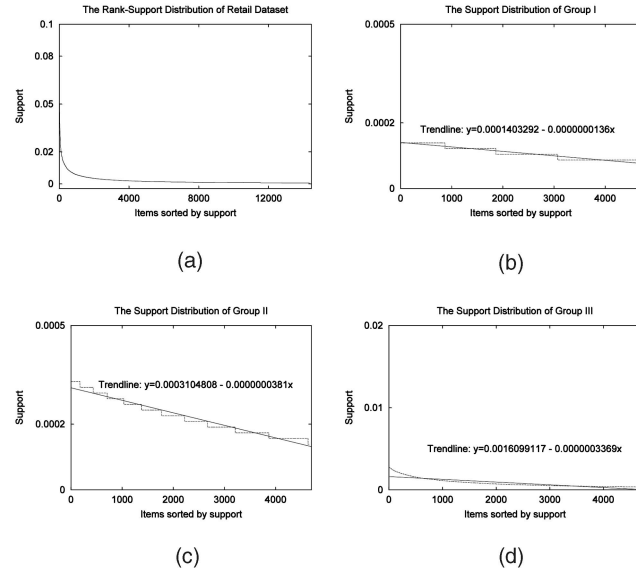


Fig. 15. The rank-support distributions of the retail data set and its three item groups with a linear regression fitting line (trendline). (a) Retail data set. (b) Group I. (c) Group II. (d) Group III.

linear distribution. Table 5 lists some of the characteristics of these data-set groups. Each group has the same number of items and transactions but a different  $a/m$  ratio. Group I has the highest  $a/m$  ratio and Group III has the lowest  $a/m$  ratio. Since the major difference among these three data-set groups is the ratio  $a/m$ , we can apply these data sets to show the impact of the  $a/m$  on the performance of the TAPER algorithm. Fig. 16 shows the pruning ratio of the TAPER algorithm on the data set with linear rank-support distributions. As expected, the pruning ratio increases as the  $a/m$  ratio decreases at different correlation thresholds. The pruning ratio also increases as correlation thresholds are increased. These experimental results confirm the trend exhibited by the cost model as shown in Lemma 7.

### 7.5 Effect of the Exponent $p$

In this subsection, we examine the effect of the exponent  $P$  on the performance of TAPER (both TAPER 1D and TAPER 2D) for data sets with a generalized Zipf rank-support distribution. We used the synthetic data sets presented in Table 3 for this experiment. All the synthetic data sets in the table have the same number of transactions and items. The rank-support distributions of these data sets follow Zipf's law but with different exponent  $P$ . Fig. 17 displays the pruning ratio of the TAPER algorithm on data

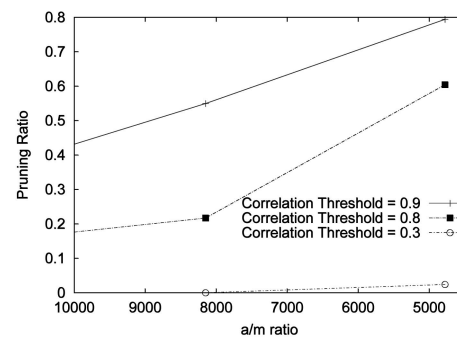


Fig. 16. Pruning ratios with the decrease of  $a/m$  for data sets with linear rank-support distribution.

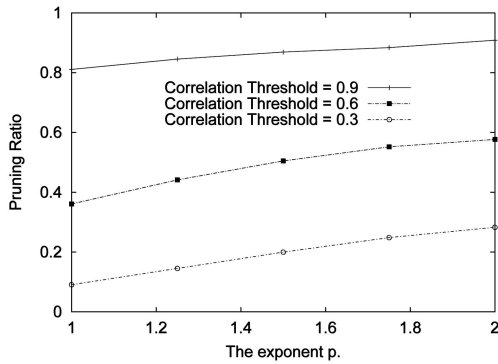


Fig. 17. The increase of pruning ratios with the increase of  $p$  for data sets with Zipf-like distribution.

sets with different exponent  $P$ . Again, the pruning ratios of the TAPER algorithm increase with the increase of the exponent  $P$  at different correlation thresholds. Also, we can observe that the pruning ratios of the TAPER algorithm increase with the increase of the correlation thresholds. Recall that the proposed algebraic cost model for data sets with a generalized Zipf distributions provides two rules which confirm the above two observations.

## 7.6 The Scalability of TAPER

In this subsection, we show the scalability of the TAPER algorithm with respect to database dimensions. Fig. 18a shows the plot of the rank-support distribution of the LA1 data set in log-log scale. Although this plot does not follow Zipf's law exactly, it does show Zipf-like behavior. In other words, the LA1 data set has an approximate Zipf-like distribution with the exponent  $P = 1.406$ . In this experiment, we generated three data sets, with 12,000, 18,000, and 24,000 items,

respectively, from the LA1 data set by random sampling on the item set. Due to the random sampling, the three data sets can have almost the same rank-support distributions as the LA1 data set. As a result, we used these three generated data sets and the LA1 data set for our scale-up experiments.

For data sets with Zipf-like rank-support distributions, Fig. 18b shows the effect of database dimensions on the performance of the TAPER algorithm. As can be seen, the pruning ratios of the TAPER algorithm show almost no change or slightly increase at different correlation thresholds. This indicates that the pruning ratios of the TAPER algorithm can be maintained when the number of items is increased. Recall that the proposed algebraic cost model for data sets with a generalized Zipf distribution exhibits a similar trend as the result of this experiment.

Finally, in Fig. 18c, we show that the execution time for our scale-up experiments increases linearly with the increase of the number of items at several different minimum correlation thresholds.

## 7.7 Evaluation of the Choices Between 1D and 2D Monotone Properties

Here, we present a performance evaluation of the choices between 1D and 2D monotone properties in the coarse-filter step of the TAPER algorithm. Fig. 19 shows the execution time of TAPER 1D and TAPER 2D on the UCI connect, mushroom, and chess data sets. In the figure, we can see that the execution time of TAPER 2D can be 10-15 percent less than that of TAPER 1D for the various correlation thresholds. This computation savings is due to the fact that TAPER 2D reduces the number of upper bounds that need to be computed as demonstrated by Lemma 4.

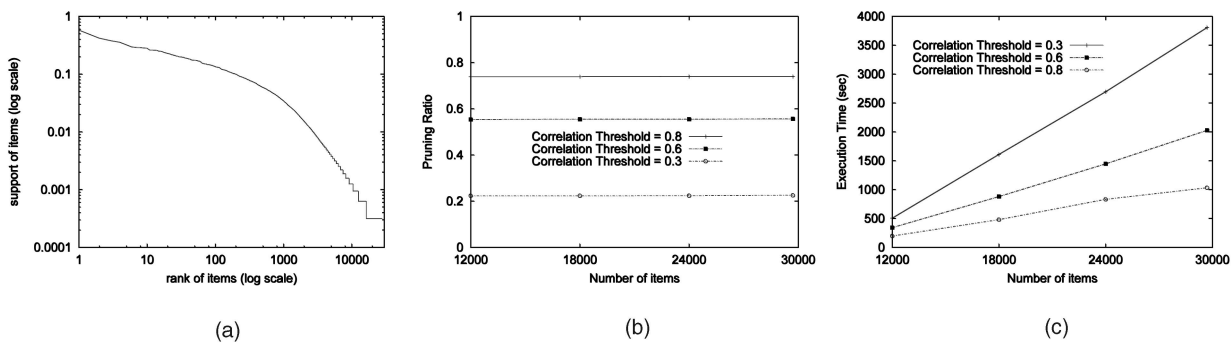


Fig. 18. (a) The plot of the rank-support distribution of the LA1 data set in log-log scale. (b) The effect of database dimensions on the pruning ratio for data sets with Zipf-like rank-support distributions. (c) The effect of database dimensions on the execution time for data sets with Zipf-like rank-support distributions.

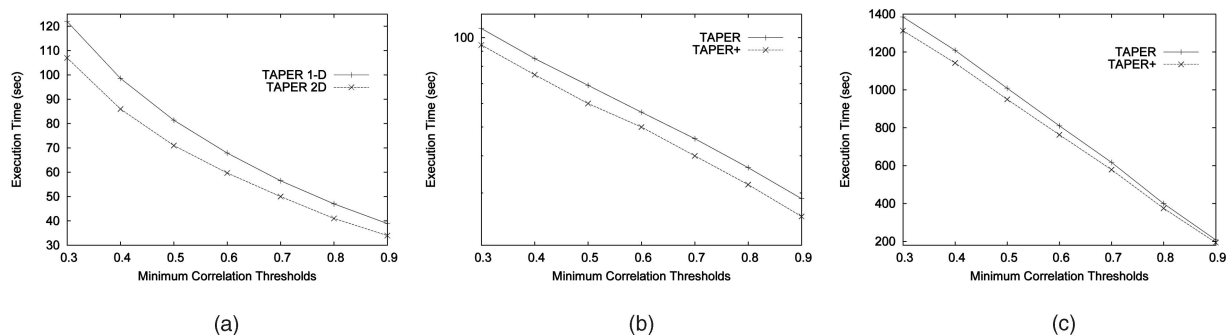


Fig. 19. TAPER 1D versus TAPER 2D on the *pumsb*, *pumsb\**, and *retail* data sets.

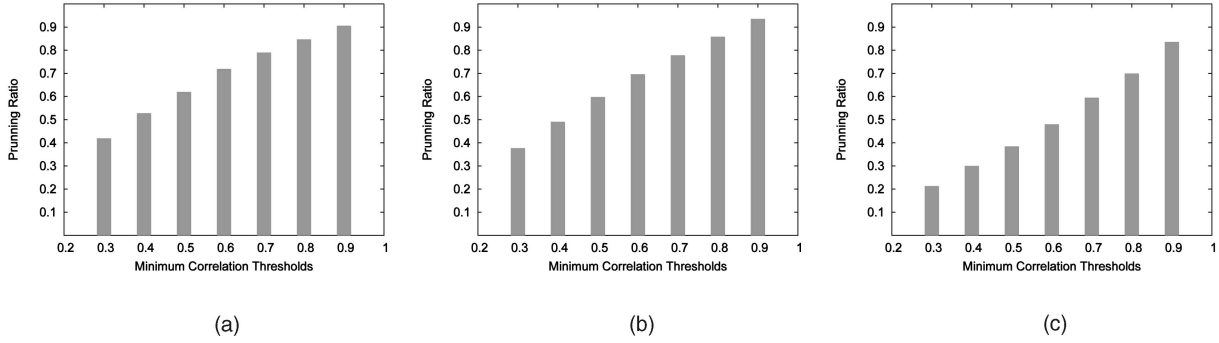


Fig. 20. The pruning effect for uncentered Pearson's correlation coefficient on UCI (a) *connect*, (b) *mushroom*, (c) *chess* data sets.

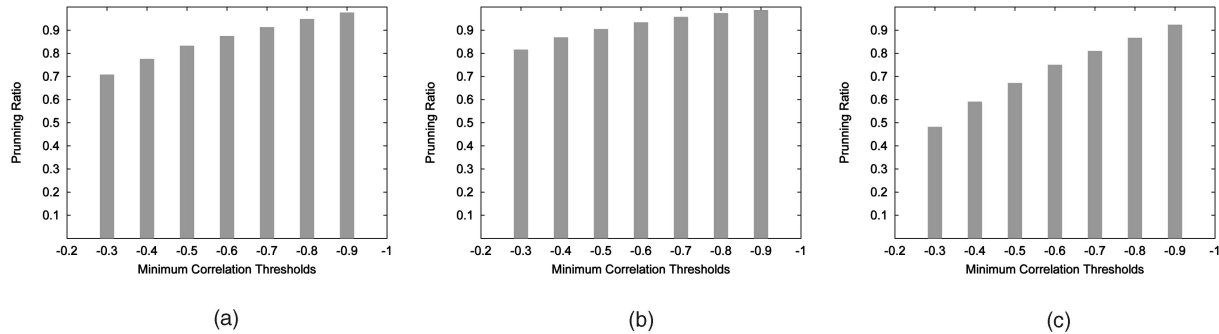


Fig. 21. The pruning effect for negative correlation computation on UCI (a) *connect*, (b) *mushroom*, (c) *chess* data sets.

## 7.8 The Extension of Algorithmic Ideas Developed in TAPER

In this subsection, we present experimental results to show the effectiveness of extending algorithmic ideas developed in the TAPER algorithm for computing the cosine measure as well as negative correlation coefficients. Fig. 20 shows the pruning ratios for computing the cosine measure (uncentered Pearson's correlation coefficient) on the UCI *connect*, *mushroom*, and *chess* data sets. Once again, the pruning ratios achieved on these data sets are quite significant. This indicates that the algorithmic ideas developed in TAPER can also achieve good computation performance for computing the cosine measure. Finally, Fig. 21 shows the pruning ratios for computing negative correlation on the UCI *connect*, *mushroom*, and *chess* data sets. The pruning ratios achieved in this case are even better than those for computing positive correlation.

## 8 CONCLUSIONS AND FUTURE WORK

In this paper, we designed an efficient two-step filter-and-refine algorithm, called TAPER, to search all the item pairs with correlations above a user-specified minimum correlation threshold. TAPER uses an upper bound of the  $\phi$  correlation coefficient, which shows 2D monotonic properties. In addition, we provided algebraic cost models to measure the computation performance of the TAPER algorithms. As demonstrated by our experimental results on both real and synthetic data sets, the pruning ratio of TAPER can be maintained or even increases with the increase of database dimensions, and the performance of TAPER confirms the proposed algebraic cost model. Finally, we showed that the algorithmic ideas developed in the TAPER algorithm can be extended to efficiently compute uncentered Pearson's correlation coefficients as well as negative correlation.

There are several potential directions for future research. First, we are interested in studying how to specify statistically meaningful correlation thresholds. Second, we plan to efficiently compute the top-k correlated item pairs. Finally, we will investigate an efficient solution for correlation detection using Bayesian and Markov models.

## ACKNOWLEDGMENTS

This work was partially supported by US National Science Foundation (NSF) grant IIS-0308264, NSF grant ITR ACI-0325949, and by the US Army High Performance Computing Research Center under the auspices of the Department of the Army, Army Research Laboratory cooperative agreement number DAAD19-01-2-0014. The content of this work does not necessarily reflect the position or policy of the government and no official endorsement should be inferred. Access to computing facilities was provided by the AHPCRC and the Minnesota Supercomputing Institute. In addition, the authors thank Professor Joydeep Ghosh (University of Texas, Austin) and Professor Ke Wang (Simon Fraser University, Canada) for valuable comments. Finally, they are grateful to Kim Koffolt for her timely and detailed feedback to help improve the readability of this paper.

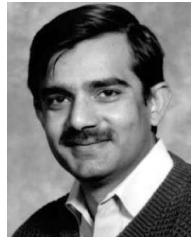
## REFERENCES

- [1] R. Agrawal, T. Imielinski, and A. Swami, "Mining Association Rules Between Sets of Items in Large Databases," *Proc. 1993 ACM SIGMOD Int'l Conf. Management of Data*, pp. 207-216, 1993.
- [2] C. Alexander, *Market Models: A Guide to Financial Data Analysis*. John Wiley & Sons, 2001.
- [3] R. Bayardo, R. Agrawal, and D. Gunopulos, "Constraint-Based Rule Mining in Large, Dense Databases," *Data Mining and Knowledge Discovery J.*, pp. 217-240, 2000.
- [4] S. Brin, R. Motwani, and C. Silverstein, "Beyond Market Baskets: Generalizing Association Rules to Correlations," *Proc. ACM SIGMOD Int'l Conf. Management of Data*, pp. 265-276, 1997.

- [5] C. Bucila, J. Gehrke, D. Kifer, and W.M. White, "Dualminer: A Dual-Pruning Algorithm for Itemsets with Constraints," *Data Mining and Knowledge Discovery J.*, pp. 241-272, 2003.
- [6] D. Burdick, M. Calimlim, and J. Gehrke, "Mafia: A Maximal Frequent Itemset Algorithm for Transactional Databases," *Proc. 17th Int'l Conf. Data Eng. (ICDE)*, pp. 443-452, 2001.
- [7] E. Cohen, M. Datar, S. Fujiwara, A. Gionis, P. Indyk, R. Motwani, J.D. Ullman, and C. Yang, "Finding Interesting Associations without Support Pruning," *Proc. 16th Int'l Conf. Data Eng. (ICDE)*, pp. 489-499, 2000.
- [8] P. Cohen, J. Cohen, S.G. West, and L. Aiken, *Applied Multiple Regression/Correlation Analysis for the Behavioral Science*, third ed. Lawrence Erlbaum Assoc., 2002.
- [9] W. Dumouchel and D. Pregibon, "Empirical Bayes Screening for Multi-Item Associations," *Proc. Seventh ACM SIGKDD Int'l Conf. Knowledge Discovery and Data Mining*, pp. 67-76, 2001.
- [10] G. Grahne, L.S. Lakshmanan, and X. Wang, "Efficient Mining of Constrained Correlated Sets," *Proc. 16th Int'l Conf. Data Eng.*, pp. 512-521, 2000.
- [11] J. Han, J. Pei, and Y. Yin, "Mining Frequent Patterns without Candidate Generation," *Proc. 2000 ACM SIGMOD Int'l Conf. Management of Data*, 2000.
- [12] I.F. Ilyas, V. Markl, P.J. Haas, P. Brown, and A. Aboulmaga, "Cords: Automatic Discovery of Correlations and Soft Functional Dependencies," *Proc. 2004 ACM SIGMOD Int'l Conf. Management of Data*, 2004.
- [13] C. Jermaine, "The Computational Complexity of High-Dimensional Correlation Search," *Proc. 2001 IEEE Int'l Conf. Data Mining (ICDM)*, pp. 249-256, 2001.
- [14] C. Jermaine, "Playing Hide-and-Seek with Correlations," *Proc. Ninth ACM SIGKDD Int'l Conf. Knowledge Discovery and Data Mining*, 2003.
- [15] S.K. Kachigan, *Multivariate Statistical Analysis: A Conceptual Introduction*, second ed., Radius Press, 1991.
- [16] M. Kendall and J.D. Gibbons, *Rank Correlation Methods*, fifth ed. Oxford Univ., 1990.
- [17] W.P. Kuo, T. Jenssen, A.J. Butte, L. Ohno-Machado, and I.S. Kohane, "Analysis of Matched mRNA Measurements from Two Different Microarray Technologies," *Bioinformatics*, vol. 18, no. 3, 2002.
- [18] E.L. Lehmann and H.J.M. D'Abbrera, *Nonparametrics: Statistical Methods Based on Ranks*. Prentice Hall, 1998.
- [19] B. Liu, W. Hsu, and Y. Ma, "Mining Association Rules with Multiple Minimum Supports," *Proc. ACM SIGKDD*, 1999.
- [20] R. Ng, L.V.S. Lakshmanan, J. Han, and A. Pang, "Exploratory Mining via Constrained Frequent Set Queries," *Proc. ACM SIGMOD Int'l Conf. Management Data*, 1999.
- [21] R. Rastogi and K. Shim, "Mining Optimized Association Rules with Categorical and Numeric Attributes," *IEEE Trans. Knowledge and Data Eng.*, vol. 14, no. 1, Jan. 2002.
- [22] H.T. Reynolds, *The Analysis of Cross-Classifications*. New York: The Free Press, 1977.
- [23] C.J. Van Rijsbergen, *Information Retrieval*, second ed. London: Butterworths, 1979.
- [24] R. Rymon, "Search through Systematic Set Enumeration," *Proc. Third Int'l Conf. Principles of Knowledge Representation and Reasoning*, pp. 539-550, 1992.
- [25] H.V. Storch and F.W. Zwiers, *Statistical Analysis in Climate Research*, reprint ed. Cambridge Univ. Press, Feb. 2002.
- [26] K. Wang, Y. He, D. Cheung, and Y.L. Chin, "Mining Confident Rules without Support Requirement," *Proc. 2001 ACM Int'l Conf. Information and Knowledge Management (CIKM)*, 2001.
- [27] X. Wu, C. Zhang, and S. Zhang, "Mining Both Positive and Negative Association Rules," *Proc. 19th Int'l Conf. Machine Learning (ICML)*, pp. 658-665, 2002.
- [28] H. Xiong, S. Shekhar, P. Tan, and V. Kumar, "Exploiting a Support-Based Upper Bound of Pearson's Correlation Coefficient for Efficiently Identifying Strongly Correlated Pairs," *Proc. 10th ACM SIGKDD Int'l Conf. Knowledge Discovery and Data Mining*, pp. 334-343, Aug. 2004.
- [29] G.K. Zipf, *Human Behavior and Principle of Least Effort: An Introduction to Human Ecology*, Cambridge, Mass: Addison Wesley Press, 1949.



**Hui Xiong** received the BE degree in automation from the University of Science and Technology of China, the MS degree in computer science from the National University of Singapore, and the PhD degree in Computer Science from the University of Minnesota. He is currently an assistant professor in the management science and Information Systems Department at Rutgers, the State University of New Jersey. His research interests include data mining, statistical computing, geographic information systems, biomedical informatics, and information security. He has published more than 20 technical papers in peer-reviewed journals and conference proceedings and is the coeditor of the book *Clustering and Information Retrieval*. He has also served on the program committees for a number of conferences and workshops. He is a member of the IEEE, the IEEE Computer Society, and the ACM.



**Shashi Shekhar** received the BTech degree in computer science from the Indian Institute of Technology, Kanpur, India, in 1985, and the MS degree in Business Administration and the PhD degree in Computer Science from the University of California, Berkeley, in 1989. He is a McKnight Distinguished University Professor at the University of Minnesota, Minneapolis, MN, USA. His research interests include spatial databases, spatial data mining, geographic and information systems (GIS), and intelligent transportation systems. He is a coauthor of the textbook *Spatial Databases* (Prentice Hall, 2003) and has published more than 160 research papers in peer-reviewed journals, books, and conferences and workshops. He is a coeditor-in-chief of *Geo-Information: An International Journal on Advances in Computer Science for GIS*. He is serving as a member of the mapping science committee of the National Research Council National Academy of Sciences (2004-2006). Dr. Shekhar is a fellow of the IEEE and a member of the IEEE Computer Society and the ACM.



**Pang-Ning Tan** received his MS degree in Physics and PhD degree in Computer Science from the University of Minnesota. He is an assistant professor in the Department of Computer Science and Engineering at the Michigan State University. Prior to joining Michigan State University, he was a research associate at the Army High Performance Computing Center. His research interests include data mining, Web intelligence, and machine learning. He has published numerous technical papers in data-mining journals, conferences, and workshops. He coauthored the textbook *Introduction to Data Mining*, (Addison Wesley). He has also served on the program committees for many international conferences. He is a member of the IEEE, the IEEE Computer Society, and the ACM.



**Vipin Kumar** is currently William Norris Professor and Head of the Computer Science and Engineering Department at the University of Minnesota. His research interests include high-performance computing and data mining. He has authored more than 200 research articles and coedited or coauthored nine books, including the widely used textbook *Introduction to Parallel Computing and Introduction to Data Mining*, both published by Addison-Wesley. He has served as chair/co-chair for many conferences/workshops in the area of data mining and parallel computing, including the IEEE International Conference on Data Mining (2002) and the 15th International Parallel and Distributed Processing Symposium (2001). He serves as the chair of the steering committee of the SIAM International Conference on Data Mining, and is a member of the steering committee of the IEEE International Conference on Data Mining. Dr. Kumar serves or has served on the editorial boards of several journals including *Knowledge and Information Systems*, *Journal of Parallel and Distributed Computing*, and the *IEEE Transactions on Knowledge and Data Engineering* (1993-1997). He is a fellow of the ACM and the IEEE, and a member of SIAM.