# Finding Interesting Associations without Support Pruning

Edith Cohen[*]     Mayur Datar[†]     Shinji Fujiwara[‡]     Aristides Gionis[§]     Piotr Indyk[¶]

Rajeev Motwani[‖]     Jeffrey D. Ullman[**]     Cheng Yang[††]

## Abstract

*Association-rule mining has heretofore relied on the condition of high support to do its work efficiently. In particular, the well-known a-priori algorithm is only effective when the only rules of interest are relationships that occur very frequently. However, there are a number of applications, such as data mining, identification of similar web documents, clustering, and collaborative filtering, where the rules of interest have comparatively few instances in the data. In these cases, we must look for highly correlated items, or possibly even causal relationships between infrequent items. We develop a family of algorithms for solving this problem, employing a combination of random sampling and hashing techniques. We provide analysis of the algorithms developed, and conduct experiments on real and synthetic data to obtain a comparative performance analysis.*

## 1   Introduction

A prevalent problem in large-scale data mining is that of *association-rule mining*, first introduced by Agrawal, Imielinski, and Swami [1]. This challenge is sometimes referred to as the *market-basket* problem due to its origins in the study of consumer purchasing patterns in retail stores, although the applications extend far beyond this specific

setting. Suppose we have a relation $R$ containing $n$ tuples over a set of *boolean* attributes $A_1, A_2, \ldots, A_m$. Let $I = \{A_{i_1}, A_{i_2}, \ldots, A_{i_k}\}$ and $J = \{A_{j_1}, A_{j_2}, \ldots, A_{j_l}\}$ be two sets of attributes. We say that $I \Rightarrow J$ is an *association rule* if the following two conditions are satisfied: *support* — the set $I \cup J$ appears in at least an $s$-fraction of the tuples; and, *confidence* — amongst the tuples in which $I$ appears, at least a $c$-fraction also have $J$ appearing in them. The goal is to identify all valid association rules for a given relation.

To some extent, the relative popularity of this problem can be attributed to its paradigmatic nature, the simplicity of the problem statement, and its wide applicability in identifying hidden patterns in data from more general applications than the original market-basket motivation. Arguably though, this success has as much to do with the availability of a surprisingly efficient algorithm, the lack of which has stymied other models of pattern-discovery in data mining. The algorithmic efficiency derives from an idea due to Agrawal et al [1, 2], called *a-priori*, which exploits the support requirement for association rules. The key observation is that if a set of attributes $S$ appears in a fraction $s$ of the tuples, then any subset of $S$ also appears in a fraction $s$ of the tuples.

This principle enables the following approach based on pruning: to determine a list $L_k$ of all $k$-sets of attributes with high support, first compute a list $L_{k-1}$ of all $(k-1)$-sets of attributes of high support, and consider as *candidates* for $L_k$ only those $k$-sets that have all their $(k-1)$-subsets in $L_{k-1}$. Variants and enhancements of this approach underlie essentially all known efficient algorithms for computing association rules or their variants. Note that in the worst case, the problem of computing association rules requires time exponential in $m$, but the a-priori algorithm avoids this pathology on real data sets. Observe also that the confidence requirement plays no role in the algorithm, and indeed is completely ignored until the end-game, when the high-support sets are screened for high confidence.

Our work is motivated by the long-standing open question of devising an efficient algorithm for finding rules that have *extremely high confidence*, but for which there is *no (or extremely weak) support*. For example, in market-basket data the standard association-rule algorithms may be useful

for commonly-purchased (i.e., high-support) items such as "beer and diapers," but are essentially useless for discovering rules such as that "Beluga caviar and Ketel vodka" are always bought together, because there are only a few people who purchase either of the two items. We develop a body of techniques which rely on the confidence requirement alone to obtain efficient algorithms.

There are two possible objections to removing the support requirement. First, this may increase the number of rules that are produced and make it difficult for a user to pin-point the rules of interest. But note that in most of the applications described next, the output rules are not intended for human analysis but rather for an automated analysis. In any case, our intent is to seek low-support rules with confidence extremely close to 100%; the latter will substantially reduce the output size and yet leave in the rules that are of interest. Second, it may be argued that rules of low support are inherently uninteresting. While this may be true in the classical market-basket applications, there are many applications where it is essential to discover rules of extremely high confidence without regard for support. We discuss these applications briefly and give some supporting experimental evidence before turning to a detailed description of our results.

One motivation for seeking such associations, of high confidence but without any support requirement, is that most rules with high support are obvious and well-known, and it is the rules of low-support that provide interesting *new* insights. Not only are the support-free associations a natural class of patterns for data mining in their own right, they also arise in a variety of applications such as: *copy detection* — identifying identical or similar documents and web pages [4, 13]; *clustering* — identifying similar vectors in high-dimensional spaces for the purposes of clustering data [6, 9]; and, *collaborative filtering* — tracking user behavior and making recommendations to individuals based on similarity of their preferences to those of other users [8, 16]. Note that each of these applications can be formulated in terms of a table whose columns tend to be sparse, and the goal is to identify column pairs that appear to be similar, without any support requirement. There are also other forms of data mining, e.g., detecting causality [15], where it is important to discover associated columns, but there is no natural notion of support.

We describe some experimental results for one such application: mining for pairs of words that occur together in news articles obtained from Reuters. The goal was to check whether low-support and high-confidence pairs provided any interesting information. Indeed the similar pairs proved to be extremely interesting as illustrated by the representative samples provided in Figure 1. A large majority of the output pairs were names of famous international personalities, cities, terms from medicine and other fields, phrases

| Category | Examples |
|---|---|
| Names | (Dalai, Lama) |
| | (Meryl, Streep) |
| | (Bertolt, Brecht) |
| | (Buenos, Aires) |
| | (Darth, Vader) |
| Terminology | (pneumocystis, carinii) |
| | (meseo, oceania) |
| | (fibrosis, cystic) |
| Phrases | (avant, garde) |
| | (mache, papier) |
| | (cosa, nostra) |
| | (hors, oeuvres) |
| | (presse, agence) |
| Misc Relations | (encyclopedia, Britannica) |
| | (Salman, Satanic) |
| | (Mardi, Gras) |
| | (emperor, Hirohito) |

**Figure 1. Examples of different types of similar pairs found in the news articles**

from foreign languages and other miscellaneous items like author-book pairs and organization names. We also obtained clusters of words, i.e., groups of words in which most pairs have high similarity. An example is the cluster (CHESS, TIMMAN, KARPOV, SOVIET, IVANCHUK, POLGER) which represents a chess event. It should be noted that the pairs discovered have very low support and would not be discovered under the standard definition of association rules. Of course, we can run the a-priori algorithm with very low support definition, but this would be really slow as indicated in the running time comparison provided in Section 5.

## 2 Summary of Results

The notion of confidence is asymmetric or unidirectional, and it will be convenient for our purpose to work with a symmetric or bi-directional measure of interest. At a conceptual level, we view the data as a 0/1 matrix $M$ with $n$ rows and $m$ columns. Typically, the matrix is fairly sparse and we assume that the *average* number of 1s per row is $r$ and that $r << m$. (For the applications we have in mind, $n$ could be as much as $10^9$, $m$ could be as large as $10^6$, and $r$ could be as small as $10^2$). Define $C_i$ as the set of rows that have a 1 in column $c_i$; also, define the *density* of column $c_i$ as $d_i = |C_i|/n$. We define the *similarity* of two columns $c_i$ and $c_j$ as

$$S(c_i, c_j) = \frac{|C_i \cap C_j|}{|C_i \cup C_j|}.$$

That is, the similarity of $c_i$ and $c_j$ is the fraction of rows, amongst those containing a 1 in either $c_i$ or $c_j$, that contain a 1 in both $c_i$ and $c_j$. Observe that the definition of similarity is symmetric with respect to $c_i$ and $c_j$; in contrast, the confidence of the rule $\{c_i\} \Rightarrow \{c_j\}$ is given by

$$\text{Conf}(c_i, c_j) = \frac{|C_i \cap C_j|}{|C_i|}.$$

To identify all pairs of columns with similarity exceeding a prespecified threshold is easy when the matrix $M$ is small and fits in main memory, since a brute-force enumeration algorithm requires $O(m^2 n)$ time. We are more interested in the case where $M$ is large and the data is disk-resident.

In this paper, our primary focus is on the problem of identifying all *pairs of columns* with similarity exceeding a pre-specified threshold $s^*$. Restricting ourselves to this most basic version of the problem will enable us clearly to showcase our techniques for dealing with the main issue of achieving algorithmic efficiency in the absence of the support requirement. It is possible to generalize our techniques to more complex settings, and we discuss this briefly before moving on to the techniques themselves. It will be easy to verify that our basic approach generalizes to the problem of identifying high-confidence association rules on pairs of columns, as discussed in Section 6. We omit the analysis and experimental results from this version of the paper, as the results are all practically the same as for high-similarity pairs. It should be noted that several recent papers [3, 14, 15] have expressed dissatisfaction with the use of confidence as a measure of interest for association rules and have suggested various alternate measures. Our ideas are applicable to these new measures of interest as well. A major restriction in our work is that we only deal with pairs of columns. However, we believe that it should be possible to apply our techniques to the identification of more complex rules; this matter is discussed in more detail in Section 6.

All our algorithms for identifying pairs of similar columns follow a very natural three-phase approach: *compute signatures*, *generate candidates*, and *prune candidates*. In the first phase, we make a pass over the table generating a small hash-signature for each column. Our goal is to deal with large-scale tables sitting in secondary memory, and this phase produces a "summary" of the table that will fit into main memory. In the second phase, we operate in main memory, generating candidate pairs from the column signatures. Finally, in the third phase, we make another pass over the original table, determining for each candidate pair whether it indeed has high similarity. The last phase is identical in all our algorithms: while scanning the table data, maintain for each candidate column-pair $(c_i, c_j)$ the counts of the number of rows having a 1 in at least one of the two columns and also the number of rows having a 1 in both columns. Consequently, we limit the ensuing discussion to

the proper implementation of only the first two phases.

The key ingredient, of course, is the hashing scheme for computing signatures. On the one hand, it needs to be extremely fast, produce small signatures, and be able to do so in a single pass over the data. Competing with this goal is the requirement that there are not too many *false-positives*, i.e., candidate pairs that are not really highly-similar, since the time required for the third phase depends on the number of candidates to be screened. A related requirement is that there are extremely few (ideally, none) *false-negatives*, i.e., highly-similar pairs that do not make it to the list of candidates.

In Section 3 we present a family of schemes based on a technique called *Min-Hashing (MH)* which is inspired by an idea used by Cohen [5] to estimate the size of transitive closure and reachability sets (see also Broder [4]). The idea is to implicitly define a random order on the rows, selecting for each column a signature that consists of the first row index (under the ordering) in which the column has a 1. We will show that the probability that two columns have the same signature is proportional to their similarity. To reduce the probability of false-positives and false-negatives, we can collect $k$ signatures by independently repeating the basic process or by picking the first $k$ rows in which the column has 1's. The main feature of the Min-Hashing scheme is that, for a suitably large choice of $k$, the number of false-positives is fairly small and the number of false-negatives is essentially zero. A disadvantage is that as $k$ rises, the space and time required for the second phase (candidate generation) increases.

Our second family of schemes, called *Locality-Sensitive Hashing (LSH)*, is presented in Section 4 and is inspired by the ideas used by Gionis, Indyk, and Motwani [7] for high-dimensional nearest neighbors (see also Indyk and Motwani [11]). The basic idea here is to implicitly partition the set of rows, computing a signature based on the pattern of 1's of a column in each subtable; for example, we could just compute a bit for each column in a subtable, denoting whether the number of 1's in the column is greater than zero or not. This family of schemes suffers from the disadvantage that reducing the number of false-positives increases the number of false-negatives, and vice versa, unlike in the previous scheme. While it tends to produce more false-positives or false-negatives, it has the advantage of having much lower space and time requirements than Min-Hashing.

We have conducted extensive experiments on both real and synthetic data, and the results are presented in Section 5. As expected, the experiments indicate that our schemes outperform the a-priori algorithm by nearly an order of magnitude. They also illustrate the point made above about the trade-off between accuracy and speed in our two algorithms. If it is important to avoid any false-negatives, than we recommend the use of the Min-Hashing schemes

which tend to be slower. However, if speed is more important than complete accuracy in generating rules, than the Locality-Sensitive Hashing schemes are to be preferred. We conclude in Section 6 by discussing the extensions of our work alluded to earlier, and by providing some interesting directions for future work.

## 3   Min-Hashing Schemes

The Min-Hashing scheme used an idea due to Cohen [5], in the context of estimating transitive closure and reachability sets. The basic idea in the Min-Hashing scheme is to randomly permute the rows and for each column $c_i$ compute its hash value $h(c_i)$ as the index of the first row under the permutation that has a 1 in that column. For reasons of efficiency, we do not wish to explicitly permute the rows, and indeed would like to compute the hash value for each column in a single pass over the table. To this end, while scanning the rows, we will simply associate with each row a hash value that is a number chosen independently and uniformly at random from a range $R$. Assuming that the number of rows is no more than $2^{16}$, it will suffice to choose the hash value as a random 32-bit integer, avoiding the "birthday paradox" [12] of having two rows get identical hash value. Furthermore, while scanning the table and assigning random hash values to the rows, for each column $c_i$ we keep track of the *minimum* hash value of the rows which contain a 1 in that column. Thus, we obtain the Min-Hash value $h(c_i)$ for each column $c_i$ in a single pass over the table, using $O(m)$ memory.

**Proposition 1** *For any column pair* $(c_i, c_j)$, $\Pr[h(c_i) = h(c_j)] = S(c_i, c_j) = \dfrac{|C_i \cap C_j|}{|C_i \cup C_j|}$.

This is easy to see since two columns will have the same Min-Hash value if and only if, in the random permutation of rows defined by their hash values, the first row with a 1 in column $c_i$ is also the first row with a 1 in column $c_j$. In other words, $h(c_i) = h(c_j)$ if and only if in restriction of the permutation to the rows in $C_i \cup C_j$, the first row belongs to $C_i \cap C_j$.

In order to be able to determine the degree of similarity between column-pairs, it will be necessary to determine multiple (say $k$) independent Min-Hash values for each column. To this end, in a single pass over the input table we select (in parallel) $k$ independent hash values for each row, defining $k$ distinct permutations over the rows. Using $O(mk)$ memory, during the single pass we can also determine the corresponding $k$ Min-Hash values, say $h_1(c_j), \ldots, h_k(c_j)$, for each column $c_j$ under each of $k$ row permutations. In effect, we obtain a matrix $\widehat{M}$ with $k$ rows, $m$ columns, and $\widehat{M}_{ij} = h_i(c_j)$, where the $k$ entries in a column are the

Min-Hash values for it. The matrix $\widehat{M}$ can be viewed as a compact representation of the matrix $M$. We will show in Theorem 1 below that the similarity of column-pairs in $M$ is captured by their similarity in $\widehat{M}$.

**Definition 1** *Let* $\widehat{S}(c_i, c_j)$ *be the fraction of Min-Hash values that are identical for* $c_i$ *and* $c_j$, *i.e.,*

$$
\begin{aligned}
\widehat{S}(c_i, c_j) &= \frac{|\{l \mid 1 \le l \le k \text{ and } \widehat{M}_{li} = \widehat{M}_{lj}\}|}{k} \\
&= \frac{|\{l \mid 1 \le l \le k \text{ and } h_l(c_i) = h_l(c_j)\}|}{k}.
\end{aligned}
$$

We have defined $\widehat{S}(c_i, c_j)$ as the fraction of rows of $\widehat{S}$ in which the Min-Hash entries for columns $c_i$ and $c_j$ are identical. We now show that $\widehat{S}(c_i, c_j)$ is a good estimator of $S(c_i, c_j)$. Recall that we set a threshold $s^*$ such that two columns are said to be *highly-similar* if $S(c_i, c_j) \ge s^*$. Assume that $s^*$ is lower bounded by some constant $c$. The following theorem shows that we are unlikely to get too many false-positives and false-negatives by using $\widehat{S}$ to determine similarity of column-pairs in the original matrix $M$.

**Theorem 1** *Let* $0 < \delta < 1$, $\epsilon > 0$, *and* $k \ge 2\delta^{-2}c^{-1}\log\epsilon^{-1}$. *Then, for all pairs of columns* $c_i$ *and* $c_j$, *we have the following two properties.*

a) *If* $S(c_j, c_j) \ge s^* \ge c$, *then* $\widehat{S}(c_i, c_j) \ge (1 - \delta)s^*$ *with probability at least* $1 - \epsilon$.

b) *If* $S(c_j, c_j) \le c$, *then* $\widehat{S}(c_i, c_j) \le (1 + \delta)c$ *with probability at least* $1 - \epsilon$.

We sketch the proof of the first part of the theorem; the proof of the second part is quite similar and is omitted. Fix any two columns $c_i$ and $c_j$ having similarity $S(c_j, c_j) \ge s^*$. Let $X_l$ be a random variable that takes on value 1 if $h_l(c_i) = h_l(c_j)$, and value 0 otherwise; define $X = X_1 + \ldots + X_k$. By Proposition 1, $E[X_l] = S(c_i, c_j) \ge s^*$; therefore, $E[X] \ge ks^*$. Applying the Chernoff bound [12] with the random variable $X$, we obtain that

$$
\begin{aligned}
\Pr[X < (1-\delta)ks^*] &\le \Pr[X < (1-\delta)E[X]] \\
&\le e^{-\frac{\delta^2 E[X]}{2}} \le e^{-\frac{\delta^2 ks^*}{2}} \le e^{-\frac{\delta^2 kc}{2}} < \epsilon.
\end{aligned}
$$

To establish the first part of the theorem, simply notice that $\widehat{S}(c_i, c_j) = X/k$.

Theorem 1 establishes that for sufficiently large $k$, if two columns have high similarity (at least $s^*$) in $M$ then they agree on a correspondingly large fraction of the Min-Hash values in $\widehat{M}$; conversely, if their similarity is low (at most $c$) in $M$ then they agree on a correspondingly small fraction of the Min-Hash values in $\widehat{M}$. Since $\widehat{M}$ can be computed

in a single pass over the data using $O(km)$ space, we obtain the desired implementation of the first phase (signature computation). We now turn to the task of devising a suitable implementation of the second phase (candidate generation).

## 3.1 Candidate Generation from Min-Hash Values

Having computed the signatures in the first phase as discussed in the previous section, we now wish to generate the candidate column-pairs in the second phase. At this point, we have a $k \times m$ matrix $\widehat{M}$ containing $k$ Min-Hash values for each column. Since $k << n$, we assume the $\widehat{M}$ is much smaller than the original data and fits in main memory. The goal is to identify all column-pairs which agree in a large enough fraction (at least to $(1-\delta)s^*$) of their Min-Hash values in $\widehat{M}$. A brute-force enumeration will require $O(k)$ time for each column-pair, for a total of $O(km^2)$. We present two techniques that avoid the quadratic dependence on $m$ and are considerably faster when (as is typically the case) the *average similarity* $\overline{S} = \sum_{1 \le i,j \le m} S(c_i, c_j)/m^2$ is low.

Row-Sorting: For this algorithm, view the rows of $\widehat{M}$ as a list of tuples containing a Min-Hash value and the corresponding column number. We sort each row on the basis of the Min-Hash values. This groups together identical Min-Hash values into a sequence of "runs." We maintain for each column an index into the position of its Min-Hash value in each sorted row. To estimate the similarity of column $c_i$ with all other columns, we use the following algorithm: use $m$ counters for column $c_i$ where the $j$th counter stores the number of rows in which the Min-Hash values of columns $c_i$ and $c_j$ are identical; for each row $1, \ldots, k$, index into the run containing the Min-Hash value for $c_i$, and for each other column represented in this run, increment the corresponding counter. To avoid $O(m^2)$ counter initializations, we re-use the same $O(m)$ counters when processing different columns, and remember and re-initialize only counters that were incremented at least once. We estimate the running time of this algorithm as follows. Sorting the rows requires total time $O(km \log m)$; thereafter, indexes on the columns can be built in time $O(km)$. The remaining time amounts to the total number of counter increments. When processing a row with column $c_i$, the number of counter increments is in fact the length of a run. The *expected* length of a run equals the sum of similarities $\sum_{j=1}^m S(c_i, c_j)$. Hence, the expected counter-increment cost when processing $c_i$ is $O(k \sum_{j=1}^m S(c_i, c_j))$, and the expected combined increments cost is $O(k \sum_{1 \le i,j \le m} S(c_i, c_j)) = O(k\overline{S}m^2)$. Thus, the expected total time required for this algorithm is is $O(km \log m + km^2\overline{S})$. Note that the average similarity $\overline{S}$ is typically a small fraction, and so the latter term in the running time is not really quadratic in $m$ as it appears to be.

Hash-Count: The next section introduces the K-Min-Hashing algorithm where the signatures for each column $c_i$ is a set $\text{SIG}_i$ of at most, but not exactly, $k$ Min-Hash values. The similarity of a column-pair $(c_i, c_j)$ is then estimated by computing the size of $\text{SIG}_i \cap \text{SIG}_j$; clearly, it suffices to consider ordered pairs $(c_j, c_i)$ such that $j < i$. This task can be accomplished via the following hash-count algorithm. We associate a bucket with each Min-Hash value. Buckets are indexed using a hash function defined over the Min-Hash values, and store column-indexes for all columns $c_i$ with some element of $\text{SIG}_i$ hashing into that bucket. We consider the columns $c_1, c_2, \ldots, c_m$ in order, and for column $c_i$ we use $i-1$ counters, of which the $j$th counter stores $\text{SIG}_j \cap \text{SIG}_i$. For each Min-Hash value $v \in \text{SIG}_i$, we access its hash-bucket and find the indexes of all columns $c_j$ $(j < i)$ which have $v \in \text{SIG}_j$. For each column $c_j$ in the bucket, we increment the counter for $(c_j, c_i)$. Finally, we add $c_i$ itself to the bucket.

Hash-Count can be easily adapted for use with the original Min-Hash scheme where we instead want to compute for each pair of columns the number of $\widehat{M}$ rows in which the two columns agree. To this end, we use a different hash table (and set of buckets) for each row of the matrix $\widehat{M}$, and execute the same process as for K-Min-Hash. The argument used for the row-sorting algorithm shows that hash-count for Min-Hashing takes $O(k\overline{S}m^2)$ time. The running time of Hash-Count for K-Min-Hash amounts to the number of counter increments. The number of increments made to a counter $(c_i, c_j)$) is exactly the size of $|\text{SIG}_i \cap \text{SIG}_j|$. A simple argument (see Lemma 1) shows that the expected size $E\{|\text{SIG}_i \cap \text{SIG}_j|\}$ is between $\min\{k, |C_i \cup C_j|\}S(c_i, c_j)$ and $\min\{2k, |C_i \cup C_j|\}S(c_i, c_j)$. Thus, the expected total running time of the hash-table scheme is $O(k\overline{S}m^2)$ in both cases.

## 3.2 The K-Min-Hashing Algorithm

One disadvantage of the Min-Hashing scheme outlined above is that choosing $k$ independent Min-Hash values for each column entailed choosing $k$ independent hash values for each row. This has a negative effect on the efficiency of the signature-computation phase. On the other hand, using $k$ Min-Hash values per column is essential for reducing the number of false-positives and false-negatives. We now present a modification called K-Min-Hashing (K-MH) in which we use only a single hash value for each row, setting the $k$ Min-Hash values for each column to be the hash values of the first $k$ rows (under the induced row permutation) containing a 1 in that column. (A similar approach was also mentioned in [5] but without an analysis.) In other words, for each column we pick the $k$ smallest hash values for the rows containing a one in that column. If a column $c_i$ has fewer 1s than $k$, we assign as Min-Hash values all

hash values corresponding to rows with 1s in that column. The resulting set of (at most) $k$ Min-Hash values forms the signature of the column $c_i$ and is denoted by $\text{SIG}_i$.

**Proposition 2** *In the K-Min-Hashing scheme, for any column $c_i$, the signature $\text{SIG}_i$ consists of the hash values for a uniform random sample of distinct rows from $C_i$.*

We remark that if the number of 1s in each column is significantly larger than $k$, then the hash values may be considered independent and the analysis from Min-Hashing applies. The situation is slightly more complex when the columns are sparse, which is the case of interest to us.

Let $\text{SIG}_{i \cup j}$ denote the $k$ smallest elements of $C_i \cup C_j$; if $|C_i \cup C_j| < k$ then $\text{SIG}_{i \cup j} = C_i \cup C_j$. We can view $\text{SIG}_{i \cup j}$ as the signature of the "column" that would correspond to $C_i \cup C_j$. Observe that $\text{SIG}_{i \cup j}$ can be obtained (in $O(k)$ time) from $\text{SIG}_i$ and $\text{SIG}_j$ since it is in fact the set of the smallest $k$ elements from $\text{SIG}_i \cup \text{SIG}_j$. Since $\text{SIG}_{i \cup j}$ corresponds to a set of rows selected uniformly at random from all elements of $C_i \cup C_j$, the expected number of elements of $\text{SIG}_{i \cup j}$ that belong to the subset $C_i \cap C_j$ is exactly $|\text{SIG}_{i \cup j}| \times |C_i \cap C_j| / |C_i \cup C_j| = |\text{SIG}_{i \cup j}| \times S(c_i, c_j)$. Also, $\text{SIG}_{i \cup j} \cap C_i \cap C_j = \text{SIG}_{i \cup j} \cap \text{SIG}_i \cap \text{SIG}_j$, since the signatures are just the smallest $k$ elements. Hence, we obtain the following theorem.

**Theorem 2** *An unbiased estimator of the similarity $S(c_i, c_j)$ is given by the expression*

$$\frac{|\text{SIG}_{i \cup j} \cap \text{SIG}_i \cap \text{SIG}_j|}{|\text{SIG}_{i \cup j}|}.$$

Consider the computational cost of this algorithm. While scanning the data, we generate one hash value per row, and for each column we maintain the minimum $k$ hash values from those corresponding to rows that contain 1 in that column. We maintain the $k$ minimum hash values for each column in a simple data structure that allows us to insert a new value (smaller than the current maximum) and delete the current maximum in $O(\log k)$ time. The data structure also makes the maximum element amongst the $k$ current Min-Hash values of each column readily available. Hence, the computation for each row is constant time for each 1 entry and additional $\log k$ time for each column with 1 entry where the hash value of the row was amongst the $k$ smallest seen so far. A simple probabilistic argument shows that the expected number of rows on which the $k$-Min-Hash list of a column $c_i$ gets updated is $O(k \log |C_i|) = O(k \log n)$. It follows that the total computation cost is a single scan of the data and $O(|M| + mk \log n \log k)$, where $|M|$ is the number of 1s in the matrix $M$.

In the second phase, while generating candidates, we need to compute the sets $\text{SIG}_{i \cup j}$ for each column-pair using merge join ($O(k)$ operations) and while we are merging we

can also find the elements that belong to $\text{SIG}_i \cap \text{SIG}_j$. Hence, the total time for this phase is $O(km^2)$. The quadratic dependence on the number of columns is prohibitive and is caused by the need to compute $\text{SIG}_{i \cup j}$ for each column-pair. Instead, we first apply a considerably more efficient *biased approximate* estimator for the similarity. The biased estimator is computed for all pairs of columns using Hash-Count in $O(k \overline{S} m^2)$ time. Next we perform a *main-memory candidate pruning* phase, where the unbiased estimator of Theorem 2 is explicitly computed for all pairs of columns where the approximate biased estimator exceeds a threshold.

The choice of threshold for the biased estimator is guided by the following lemma.

**Lemma 1**

$$E\{|\text{SIG}_i \cap \text{SIG}_j|\} / \min\{2k, |C_i \cup C_j|\} \leq S(c_i, c_j) \leq$$
$$\leq E\{|\text{SIG}_i \cap \text{SIG}_j|\} / \min\{k, |C_i \cup C_j|\}.$$

Alternatively, the biased estimator and choice of threshold can be derived from the following analysis. Let $(c_i, c_j)$ be a column-pair with $|C_i| \geq |C_j|$; define $C_{ij} = C_i \cap C_j$. As before, for each column $c_i$, we choose a set $\text{SIG}_i$ of $k$ Min-Hash values. Let $\text{SIG}_{ij} = \text{SIG}_i \cap C_{ij}$ and $\text{SIG}_{ji} = \text{SIG}_j \cap C_{ij}$. Then, the expected sizes of $\text{SIG}_{ij}$ and $\text{SIG}_{ji}$ are given by $k|C_{ij}|/|C_i|$ and $k|C_{ij}|/|C_j|$. Also, $|\text{SIG}_i \cap \text{SIG}_j| = min(|\text{SIG}_{ij}|, |\text{SIG}_{ji}|)$. Hence can compute the expected value as

$$E[|\text{SIG}_i \cap \text{SIG}_j|] =$$

$$\sum_{x=0}^{k} \sum_{y=0}^{k} \Pr[|\text{SIG}_{ij}| = x] \Pr[|\text{SIG}_{ji}| = y \mid |\text{SIG}_{ij}| = x] \min(x, y).$$

Since $|C_i| > |C_j|$, we have $E(|\text{SIG}_{ij}|) \leq E(|\text{SIG}_{ji}|)$. We assume that $\Pr[|\text{SIG}_{ij}| > |\text{SIG}_{ji}|] \approx 0$ or $\sum_{y=x}^{k} P[|\text{SIG}_{ji}| = y \mid |\text{SIG}_{ij}| = x] \approx 1$. Then, the above equation becomes

$$E[|\text{SIG}_i \cap \text{SIG}_j|]$$

$$= \sum_{x=0}^{k} \sum_{y=x}^{k} \Pr[|\text{SIG}_{ij}| = x] \Pr[|\text{SIG}_{ji}| = y \mid |\text{SIG}_{ij}| = x] x$$

$$= \sum_{x=0}^{k} \Pr[|\text{SIG}_{ij}| = x] x \sum_{y=x}^{k} \Pr[|\text{SIG}_{ji}| = y \mid |\text{SIG}_{ij}| = x]$$

$$= E[\text{SIG}_{ij}].$$

Thus, we obtain the estimator $E[|\text{SIG}_i \cap \text{SIG}_j|] \approx k|C_{ij}|/|C_i|$. We use this estimate to calculate $|C_{ij}|$ and use that to estimate the similarity since we know $|C_i|$ and $|C_j|$. We compute $|\text{SIG}_i \cap \text{SIG}_j|$ using the hash table technique that we have described earlier in Section 3.1. The time required to compute the hash values is $O(|M| + mk \log n \log k)$ as described earlier, and the time for computing $|\text{SIG}_i \cap \text{SIG}_j|$ is $O(k \overline{S} m^2)$.

# 4 Locality-Sensitive Hashing Schemes

In this section we show how to obtain a significant improvement in the running time with respect to the previous algorithms by resorting to *Locality Sensitive Hashing (LSH)* technique introduced by Indyk and Motwani [11] in designing main-memory algorithms for nearest neighbor search in high-dimensional Euclidean spaces; it has been subsequently improved and tested in [7]. We apply the LSH framework to the Min-Hash functions described in earlier section, obtaining an algorithm for similar column-pairs. This problem differs from nearest neighbor search in that the data is known in advance. We exploit this property by showing how to optimize the running time of the algorithm given constraints on the quality of the output. Our optimization is *input-sensitive*, i.e., takes into account the characteristics of the input data set.

The key idea in LSH is to hash columns so as to ensure that for each hash function, the probability of collision is much higher for similar columns than for dissimilar ones. Subsequently, the hash table is scanned and column-pairs hashed to the same bucket are reported as similar. Since the process is probabilistic, both false positives and false negatives can occur. In order to reduce the former, LSH amplifies the difference in collision probabilities for similar and dissimilar pairs. In order to reduce false negatives, the process is repeated a few times, and the union of pairs found during all iterations are reported. The fraction of false positives and false negatives can be analytically controlled using the parameters of the algorithm.

Although not the main focus of this paper, we mention that the LSH algorithm can be adapted to the *on-line* framework of [10]. In particular, it follows from our analysis that each iteration of our algorithm reduces the number of false negatives by a fixed factor; it can also add new false positives, but they can be removed at a small additional cost. Thus, the user can monitor the progress of the algorithm and interrupt the process at any time if satisfied with the results produce so far. Moreover, the higher the similarity, the earlier the pair is likely to be discovered. Therefore, the user can terminate the process when the output produced appears to be less and less interesting.

## 4.1 The Min-LSH Scheme

We present now the Min-LSH (M-LSH) scheme for finding similar column-pairs from the matrix $\widehat{M}$ of Min-Hash values. The M-LSH algorithm splits the matrix $\widehat{M}$ into $l$ sub-matrices of dimension $r \times m$. Recall that $\widehat{M}$ has dimension $k \times m$, and here we assume that $k = lr$. Then, for each of the $l$ sub-matrices, we repeat the following. Each column, represented by the $r$ Min-Hash values in the current sub-matrix, is hashed into a table using as hashing key the concatenation of all $r$ values. If two columns are similar, there is a high probability that they agree in all $r$ Min-Hash values and so they hash into the same bucket. At the end of the phase we scan the hash table and produce pairs of columns that have been hashed to the same bucket. To amplify the probability that similar columns will hash to the same bucket, we repeat the process $l$ times. Let $P_{r,l}(c_i, c_j)$ be the probability that columns $c_i$ and $c_j$ will hash to the same bucket at least once; since the value of $P$ depends only upon $s = S(c_i, c_j)$, we simplify notation by writing $P(s)$.

**Lemma 2** *Assume that columns $c_i$ and $c_j$ have similarity $s$, and also let $s^*$ be the similarity threshold. For any $0 < \delta, \epsilon < 1$, we can choose the parameters $r$ and $l$ such that:*

- *For any $s \geq (1 + \delta)s^*$, $P_{r,l}(c_i, c_j) \geq 1 - \epsilon$*

- *For any $s \leq (1 - \delta)s^*$, $P_{r,l}(c_i, c_j) \leq \epsilon$*

Proof: By Proposition 1, the probability that columns $c_i$, $c_j$ agree on one Min-Hash value is exactly $s$ and the probability that they agree in a group of $r$ values is $s^r$. If we repeat the hashing process $l$ times, the probability that they will hash at least once to the same bucket would be $P_{r,l}(c_i, c_j) = 1 - (1 - s^r)^l$. The lemma follows from the properties of the function $P$. ∎

Lemma 2 states that for large values of $r$ and $l$, the function $P$ approximates the unit step function translated to the point $C = s^*$, which can be used to filter out *all and only* the pairs with similarity at most $s^*$. On the other hand, the time/space requirements of the algorithm are proportional to $k = lr$, so the increase in the values of $r$ and $l$ is subject to a quality-efficiency trade-off. In practice, if we are willing to allow a number of false negatives ($n_-$) and false positives ($n_+$), we can determine optimal values for $r$ and $l$ that achieve this quality.

Specifically, assume that we are given (an estimate of) the similarity distribution of the data, defined as $d(s_i)$ to be the number of pairs having similarity $s_i$. This is not an unreasonable assumption, since we can approximate this distribution by sampling a small fraction of columns and estimating all pairwise similarity. The expected number of false negatives would be $\sum_{s_i \geq s_0} d(s_i)(1 - P(s_i))$, and the expected number of false positives would be $\sum_{s_i < s_0} d(s_i)P(s_i)$. Therefore, the problem of estimating optimal parameters turns into the following minimization problem:

$$\text{minimize} \qquad l \cdot r$$
$$\text{subject to} \quad \begin{cases} \sum_{s_i \geq s_0} d(s_i)(1 - P(s_i)) \leq n_- \\ \sum_{s_i < s_0} d(s_i)P(s_i) \leq n_+ \end{cases}$$

This is an easy problem since we have only two parameters to optimize, and their feasible values are small integers. Also, the histogram $d(\cdot)$ is typically quantified in 10-20 bins.

One approach is to solve the minimization problem by iterating on small values of $r$, finding a lower bound on the value of $l$ by solving the first inequality, and then performing binary search until the second inequality is satisfied. In most experiments, the optimal value of $r$ was between 5 and 20.

## 4.2 The Hamming-LSH Scheme

We now propose another scheme, Hamming-LSH (H-LSH), for finding highly-similar column-pairs. The idea is to reduce the problem to searching for column-pairs having small *Hamming distance*. In order to solve the latter problem we employ the techniques similar to those used in [7] to solve the nearest neighbor problem. We start by establishing the correspondence between the similarity and Hamming distance (the proof is easy).

Lemma 3  $S(C_i, C_j) = \frac{|C_i| + |C_j| - d_H(c_i, c_j)}{|C_i| + |C_j| + d_H(c_i, c_j)}$.

It follows that when we consider pairs $(c_i, c_j)$ such that the sum $\rho = |C_i| + |C_j|$ is fixed, then the high value of $S(c_i, c_j)$ corresponds to small values of $d_H(c_i, c_j)$ and vice versa. Hence, we partition columns into groups of similar density and for each group we find pairs of columns that have small Hamming distance. First, we briefly describe how to search for pairs of columns with small Hamming distance. This scheme is similar to to the technique from [7] and can be analyzed using the tools developed in there. This scheme finds highly-similar columns assuming that the density of all columns is roughly the same. This is done by partitioning the rows of database into $p$ subsets. For each partition, process as in the previous algorithm. We declare a pair of columns as a candidate if they agree on any subset. Thus this scheme is exactly similar to the earlier scheme, except that we are dealing with the actual data instead of Min-Hash values.

However there are two problems with this scheme. One problem is that if the matrix is sparse, most of the subsets just contain zeros and also the columns do not have similar densities as assumed. The following algorithm (which we call *H-LSH*) improves on the above basic algorithm.

The basic idea is as follows. We perform computation on a *sequence* of matrices with increasing densities; we denote them by $M_0, M_1, M_2, \ldots$. The matrix $M_{i+1}$ is obtained from the matrix $M_i$ by randomly pairing all rows of $M_i$, and placing in $M_{i+1}$ the "OR" of each pair. [1] One can see that for each $i$, $M_{i+1}$ contains half the rows of $M_i$ (for illustration purposes we assume that the initial number of rows is a power of 2). The algorithm is applied to all matrices

---

[1]Notice, that the "OR operation" gives similar results to hashing each columns to a set of increasingly smaller hash table; this provides an alternative view of our algorithm.

in the set. A pair of columns can become a candidate only on a matrix $M_i$ in which they are both sufficiently dense and both their densities belong to a certain range. False negatives are controlled by repeating each sample $l$ times, and taking the *union* of the candidate sets across all $l$ runs. Hence, $kr$ rows are extracted from each compressed matrix. Note that this operation may increase false positives.

We now present the algorithm that was implemented. Experiments show that this scheme is better than the Min-Hashing algorithms in terms of running time, but the number of false positives is much larger. Moreover, the number of false positives is increases rapidly if we try to reduce the number of false negatives. In the case of Min-Hashing algorithms, if we decreased the number of false negatives by increasing $k$, the number of false positives would also decrease.

The Algorithm:

1. Set $M_0 = M$ and generate $M_1, M_2, \ldots$ as described above.

2. For each $i \geq 0$, select $k$ sets of $r$ sample rows from $M_i$.

3. A column pair is a candidate if there exists an $i$, such that (i) the column pair has density in $(1/t, (t-1)/t)$ in $M_i$, and (ii) has identical hash values (essentially, identical $r$-bit representations) in at least one of the $k$ runs.

Note that $t$ is a parameter that indicates the range of density for candidate pairs, and we use $t = 4$ in our experiments.

## 5 Experiments

We have conducted experiments to evaluate the performance of the different algorithms. In this section we report the results for the different experiments. We use two sets of data namely synthetic data and real data.

Synthetic Data: The data contains $10^4$ columns and the number of rows vary from $10^4$ to $10^6$. The column densities vary from $1\%$ to $5\%$ and for every 100 columns we have a pair of similar columns. We have 20 pairs similar columns whose similarity fall in the ranges $(85, 95)$, $(75, 85)$, $(65, 75)$, $(55, 65)$ and $(45, 55)$.

Real Data: The real data set consists of the log of HTTP requests made over a period of 9 days to the Sun Microsystems Web server (www.sun.com). The columns in this case are the URL's and the rows represent distinct client IP addresses that have recently accessed the server. An entry is set to 1 if there has been at least one hit for that URL from that particular client IP. The data set has about thirteen thousand columns and more than 0.2 million rows. Most of the
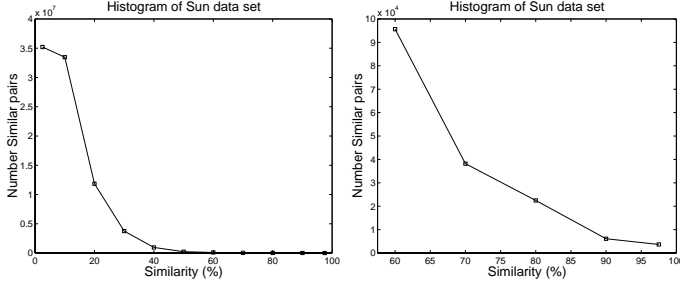
**Figure 2. The first figure shows the similarity distribution of the Sun data. The second shows again the same distribution but it focuses on the region of similarities that we are interested in.**

| Support threshold | | 0.1% | 0.15% | 0.2% |
|---|---|---|---|---|
| Number of columns after pruning | | 15559 | 11568 | 9518 |
| A-priori | (sec) | - | 96.05 | 79.94 |
| MH | (sec) | 71.4 | 44.8 | 25.8 |
| K-MH | (sec) | 87.6 | 52.0 | 36.0 |
| H-LSH | (sec) | 15.6 | 6.7 | 6.0 |
| M-LSH | (sec) | 10.7 | 9.7 | 5.1 |

**Figure 3. Running times for the news articles data set**

columns are sparse and have density less than 0.01%. The histogram in Figure 2 shows the number of column pairs for different values of similarity. Typical examples of similar columns that we extracted from this data were URLs corresponding to `gif` images or `Java` applets which are loaded automatically when a client IP accesses a parent URL.

To compare our algorithms with existing techniques, we implemented and executed the a-priori algorithm [1, 2]. Of course, the a-priori is not designed for this setting of low support, but it is the only existing technique and gives us a benchmark against which we can compare the improvements afforded by our algorithms. The comparison was done for the news articles data that we have mentioned in Section 1.

We conducted experiments on the news article data and our results are summarized in Figure 3. The a-priori algorithm cannot be run on the original data since it runs out of memory. Therefore we performed support pruning to remove columns that have very few ones in them. It is evident that our techniques give nearly an order of magnitude improvement in running time; for support threshold below 0.1%, a-priori runs out of memory on our systems and does a lot of thrashing. Note that although our algorithms are probabilistic they report the same set of pairs as reported by a-priori.

### 5.1 Results

We implemented the four algorithms described in the previous section, namely MH, K-MH, H-LSH, and M-LSH. All algorithms were compared in terms of the running time and the quality of the output. Due to the lack of space we report experiments and give graphs for the Sun data, which in any case are more interesting, but we have also performed tests for the synthetic data, and all algorithms behave similarly.

The quality of the output is measured in terms of false positives and false negatives generated by each algorithm. To do that, we plot a curve that shows the ratio of the number of pairs found by the algorithm over the real number of pairs (computed once off-line) for a given similarity range (e.g. Figure 7). The result is typically an "S"-shaped curve, that gives a good visual picture for the false positives and negatives of the algorithm. Intuitively, the area below the curve and left to a given similarity cutoff corresponds to the number of false positives, while the area above the curve and right to the cutoff corresponds to the number of false negatives.

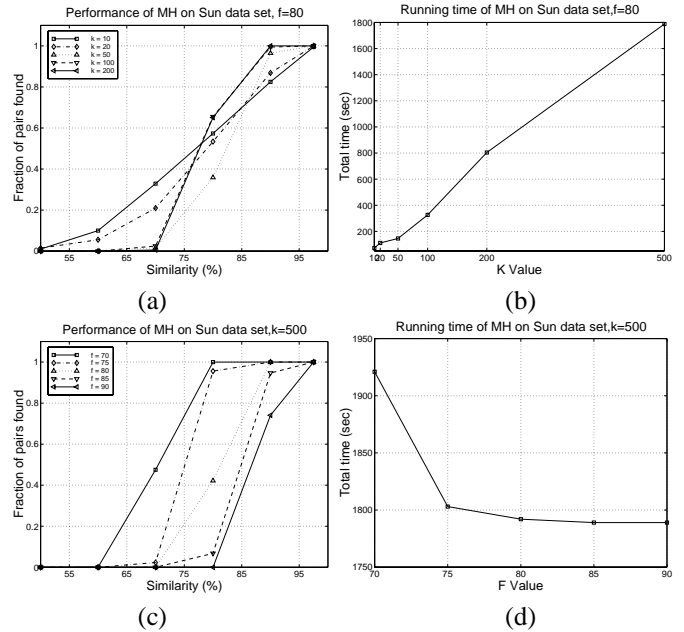We now describe the behavior of each algorithm as their parameters are varied .



(a)

(b)

(c)

(d)

**Figure 4. Quality of output and total running time for MH algorithm as $k$ and $s$ are varied**
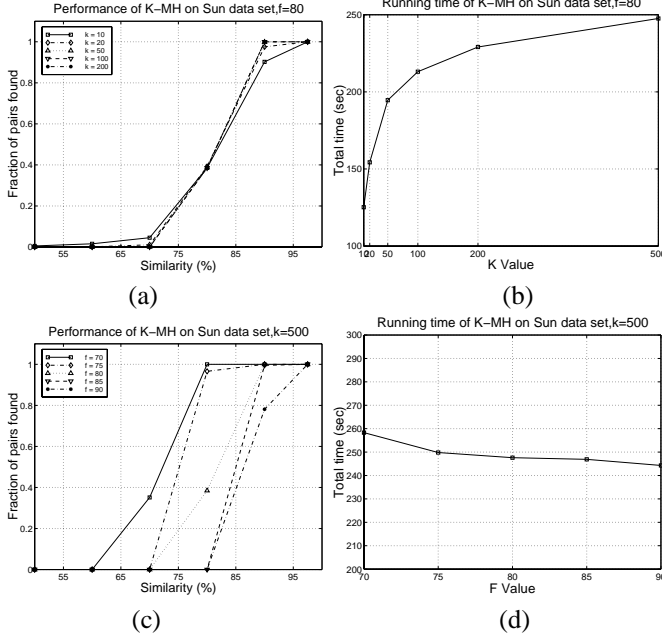
MH and K-MH algorithms have two parameters, $s^*$ the

**Figure 5. Quality of output and total running time for K-MH algorithm as $k$ and $s$ are varied**



**Figure 6. Quality of output and total running time for H-LSH algorithm as $r$ and $l$ are varied**

user specified similarity cutoff, and $k$, the number of Min-Hash values extracted to represent the signature of of each column. Figures 4(a) and 5(a) plot "S"-curves for different values of $k$ for the MH and K-MH algorithms. As the $k$ value increases the curve gets sharper indicating better quality. In Figures 4(c) and 5(c) we keep $k$ fixed and change the value $s^*$ of the similarity cutoff. As expected the curves shift to the right as the cutoff value increases. Figures 4(d) and 5(d) show that for a given value of $k$ the total running time decreases marginally since we generate fewer candidates. Figure 4(b) shows that the total running time for MH algorithm increases linearly with $k$. However this is not the case for K-MH algorithm as depicted by Figure 5(b). The sub-linear increase of the running time is due to the sparsity of the data. More specifically, the number of hash values extracted from each column is upper bounded by the number of ones of that column, and therefore, the hash values extracted do not increase linearly with $k$.

We do a similar exploration of the parameter space for the M-LSH and H-LSH algorithms. The parameters of this algorithm are $r$, and $l$. Figures 7(a) and 6(a) illustrate the fact that as $r$ increases the probability that columns mapped to the same bucked decreases, and therefore the number of false positives decreases but as a trade-off consequence the number of false negatives increases. On the other hand, Figure 7(c) and 6(c) shows that an increase in $l$, corresponds to an increase of the collision probability, and therefore the
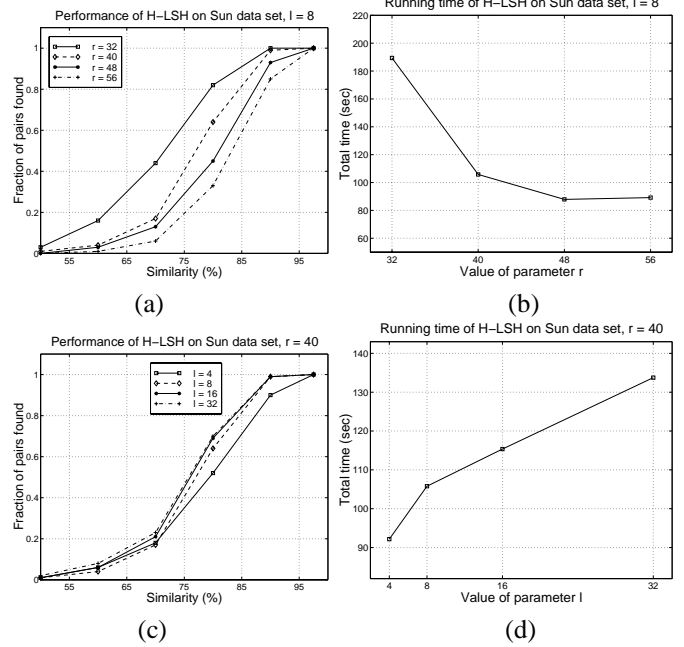
number of false negatives decrease but the number of false positives increases. Figures 7(d) and 6(d) show that the total running time increases with $l$ since we hash each column more times and this also results in an increase in the number of candidates. In our implementation of M-LSH, the extraction of min hash values dominates the total computation time, which increases linearly with the value of $r$. This is shown in Figure 7(b). On the other hand, in the implementation of H-LSH, checking for candidates dominates the running times, and as a result the total running time decreases as $r$ increases since less candidates are produced. This is shown in Figure 6(b).

We now compare the different algorithms that we have implemented. When comparing the time requirements of the algorithm we compare the CPU time for each algorithm since the time spent in I/O is same for all the algorithms. It is important to note that the for all the algorithms the number of false negatives is very important and this is the quantity that requires to be kept in control. As long as the number of false positives is not too large (i.e. all of candidates can fit in main memory) we can always eliminate them in the pruning phase. To compare the algorithms we fix the percentage of false negatives that can be tolerated. For each algorithm we pick the set of parameters for which the number of false negatives is within this threshold and the total running time is minimum. We then plot the total running time and the number of false positives against the false negative threshold.
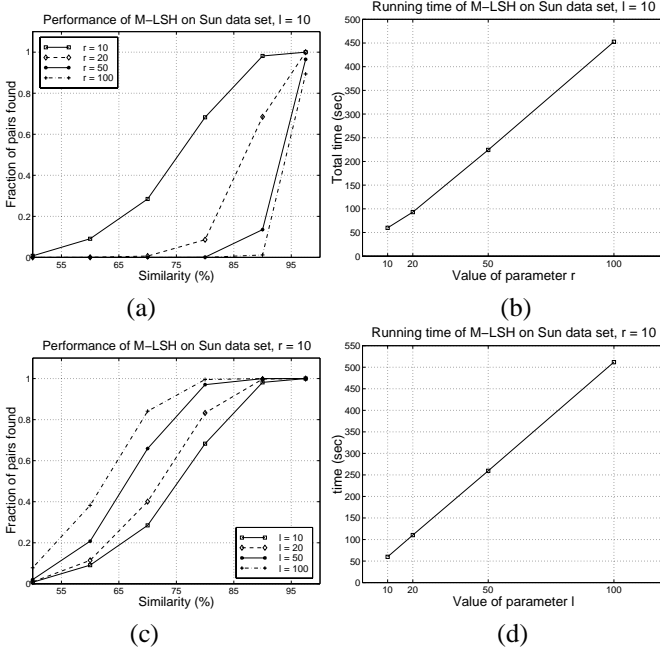
(a)  (b)

(c)  (d)

**Figure 7. Quality of output and total running time for M-LSH algorithm as $r$ and $l$ are varied**

(a)  (b)

(c)  (d)

**Figure 8. Comparison of different algorithms in terms of total running time and number of false positives for different negative thresholds.**

Consider Figures 8(a) and 8(c). The Figures show the total running time against the false negative threshold. We can see that the H-LSH algorithm requires a lot of time if the false negative threshold is less while it does better if the limit is high. In general the M-LSH and H-LSH algorithms do better than the MH and K-MH algorithms. However it should be noted that H-LSH algorithm cannot be used if we are interested in similarity cutoffs that are low. The graph shows that the best performance is shown by the M-LSH algorithm.

Figure 8 gives the number of false positives generated by the algorithms against the tolerance limit. The false positives are plotted on a logarithmic scale. In case of H-LSH and M-LSH algorithms the number of false positives decreases if we are ready to tolerate more false negatives since in that case we hash every column fewer times. However the false positive graph for K-MH and MH is not monotonic. There exists a tradeoff in the time spent in the candidate generation stage and the pruning stage. To maintain the number of false negatives less than the given threshold we could either increase $k$ and spend more time in the candidate generation stage or else decrease the similarity cutoff $s$ and spend more time in the pruning stage as we get more false positives. Hence the points on the graph correspond to different values of similarity cutoff $s^*$ with which the algorithms are run to get candidates with similarity above a certain threshold. As a result we do not observe a monotonic behavior in case of
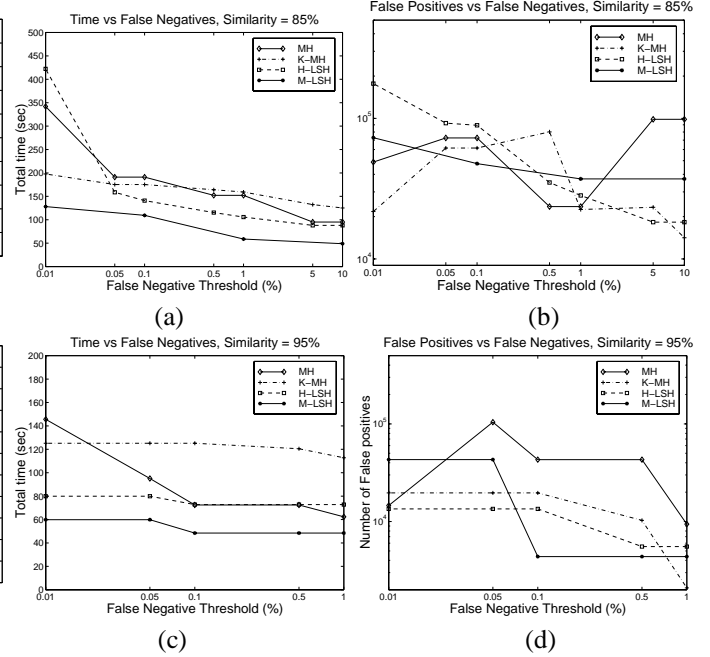
these algorithms.

We would like to comment that the results provided should be analyzed with caution. The reader should note that whenever we refer to time we refer to only the CPU time and we expect I/O time to dominate in the signature generation phase and pruning phase. If we are aware about the nature of the data then we can be smart in our choice of algorithms. For instance the K-MH algorithm should be used instead of MH for sparse data sets since it takes advantage of sparsity.

## 6 Extensions and Further Work

We briefly discuss some extensions of the results presented here as well as directions for future work. First, note that all the results presented here were for the discovery of bi-directional similarity measures. However, the Min-Hash technique can be extended to the discovery of column-pairs $(c_i, c_j)$ which form a high-confidence association rule of the type $\{c_i\} \Rightarrow \{c_j\}$ but without any support requirements. The basic idea is to generate a set of Min-Hash values for each column, and to determine whether the fraction of these values that are identical for $c_i$ and $c_j$ is proportional to the ratio of their densities, $d_i/d_j$. The analytical and the experimental results are qualitatively the same as for similar

column-pairs.

We can also use our Min-Hashing scheme to determine more complex relationships, e.g., $c_i$ is highly-similar to $c_j \vee c_{j'}$, since the hash values for the induced column $c_j \vee c_{j'}$ can be easily computed by taking the component-wise minimum of the hash value signature for $c_j$ and $c_{j'}$. Extending to $c_j \wedge c_{j'}$ is more difficult. It works as follows. First, observe that "$c_i$ implies $c_j \wedge c_{j'}$" means that "$c_i$ implies $c_j$" and "$c_i$ implies $c_{j'}$". The latter two implications can be generated as above. Now, we can conclude that "$c_i$ implies $c_j \wedge c_{j'}$" if (and only if) the cardinality of $c_i$ is roughly that of $c_j \wedge c_{j'}$. This presents problems when the cardinality of $c_i$ is really small, but is not so difficult otherwise. The case of small $c_i$ may not be very interesting anyway, since it is difficult to associate any statistical significance to the similarity in that case. It is also possible to define "anti-correlation," or mutual exclusion between a pair of columns. However, for statistical validity, this would require imposing a support requirement, since extremely sparse columns are likely to be mutually exclusive by sheer chance. It is interesting to note that our hashing techniques can be extended to deal with this situation, unlike a-priori which will not be effective even with support requirements. Extensions to more than three columns and complex boolean expressions are possible but will suffer from an exponential overhead in the number of columns.

## References

[1] R. Agrawal, T. Imielinski, and A. Swami. Mining Association Rules Between Sets of Items in Large Databases. In *Proceedings of the ACM SIGMOD Conference,* 1993, pp. 207–216.

[2] R. Agrawal and R. Srikant. Fast Algorithms for Mining Association Rules. In *Proceedings of the 20th International Conference on Very Large Databases,* 1994.

[3] S. Brin, R. Motwani, J.D. Ullman, and S. Tsur, Dynamic itemset counting and implication rules for market basket data. In *Proceedings of the ACM SIGMOD Conference on Management of Data,* 1997, pp. 255–264.

[4] A. Broder. On the resemblance and containment of documents. In *Compression and Complexity of Sequences (SEQUENCES'97)*, 1998, pp. 21–29.

[5] E. Cohen. Size-Estimation Framework with Applications to Transitive Closure and Reachability. *Journal of Computer and System Sciences* 55 (1997): 441–453.

[6] R.O. Duda and P.E. Hart. *Pattern Classification and Scene Analysis.* A Wiley-Interscience Publication, New York, 1973.

[7] A. Gionis, P. Indyk, and R. Motwani. Similarity Search in High Dimensions via Hashing. In *Proceedings of the 25th International Conference on Very Large Data Bases*, 1999, pp. 518–529.

[8] D. Goldberg, D. Nichols, B.M. Oki, and D. Terry. Using collaborative filtering to weave an information tapestry. *Communications of the ACM* 55 (1991): 1–19.

[9] S. Guha, R. Rastogi, and K. Shim. CURE - An Efficient Clustering Algorithm for Large Databases In *Proceedings of the ACM-SIGMOD International Conference on Management of Data,* 1998, pp. 73–84.

[10] J.M. Hellerstein, P.J. Haas, and H.J. Wang. Online Aggregation. In *Proceedings of the ACM-SIGMOD International Conference on Management of Data*, 1997.

[11] P. Indyk and R. Motwani. Approximate Nearest Neighbor: Towards Removing the Curse of Dimensionality. In *Proceedings of the 30th Annual ACM Symposium on Theory of Computing,* 1998, pp. 604–613.

[12] R. Motwani and P. Raghavan. *Randomized Algorithms.* Cambridge University Press, 1995.

[13] N. Shivakumar and H. Garcia-Molina. Building a Scalable and Accurate Copy Detection Mechanism. In *Proceedings of the 3rd International Conference on the Theory and Practice of Digital Libraries*, 1996.

[14] C. Silverstein, S. Brin, and R. Motwani. Beyond Market Baskets: Generalizing Association Rules to Dependence Rules. In *Proceedings of the ACM SIGMOD Conference on Management of Data,* 1997, pp. 265–276. In *Data Mining and Knowledge Discovery* 2 (1998): 69–96.

[15] C. Silverstein, S. Brin, R. Motwani, and J.D. Ullman. Scalable Techniques for Mining Causal Structures. In *Proceedings of the 24th International Conference on Very Large Data Bases*, 1998, pp. 594–605.

[16] H.R. Varian and P. Resnick, Eds. CACM Special Issue on Recommender Systems. *Communications of the ACM* 40 (1997).