

Comparing Data Streams Using Hamming Norms (How to Zero In)

Graham Cormode, Mayur Datar, Piotr Indyk, and S. Muthukrishnan

Abstract—Massive data streams are now fundamental to many data processing applications. For example, Internet routers produce large scale diagnostic data streams. Such streams are rarely stored in traditional databases and instead must be processed “on the fly” as they are produced. Similarly, sensor networks produce multiple data streams of observations from their sensors. There is growing focus on manipulating data streams and, hence, there is a need to identify basic operations of interest in managing data streams, and to support them efficiently. We propose computation of the Hamming norm as a basic operation of interest. The Hamming norm formalizes ideas that are used throughout data processing. When applied to a single stream, the Hamming norm gives the number of distinct items that are present in that data stream, which is a statistic of great interest in databases. When applied to a pair of streams, the Hamming norm gives an important measure of (dis)similarity: the number of unequal item counts in the two streams. Hamming norms have many uses in comparing data streams. We present a novel approximation technique for estimating the Hamming norm for massive data streams; this relies on what we call the “ l_0 sketch” and we prove its accuracy. We test our approximation method on a large quantity of synthetic and real stream data, and show that the estimation is accurate to within a few percentage points.

Index Terms—Data stream analysis, approximate query processing, data structures and algorithms, data reduction.

1 INTRODUCTION

DATA streams are now fundamental to many data processing applications. For example, telecommunication network elements such as switches and routers periodically generate records of their traffic—telephone calls, Internet packet, and flow traces—which are data streams [3], [26], [37]. Atmospheric observations—weather measurements, lightning stroke data, and satellite imagery—also produce multiple data streams [38], [44]. Emerging sensor networks produce many streams of observations, for example, highway traffic conditions [35], [42]. Sources of data streams—large scale transactions, Web clicks, ticker tape updates of stock quotes, toll booth observations—are ubiquitous in daily life.

For many applications that produce data streams, it is useful to visualize the underlying data seen so far, or underlying state of the system as a very high-dimensional vector (note that in most cases the vector is *not* explicitly represented or materialized). For example, if we consider a data stream created by network flows that originate from a source IP address, the state at any time t can be visualized as a vector $\mathbf{a} = a_1 \dots a_n$ that is indexed by the destination IP

address. The entry along each dimension (a_i) represents the total number of flows that were observed between the given source IP address and the destination IP address a_i . The input is usually presented in the order it arrives (rather than sorted on any attribute) and consists of updates only. For instance, in the example above, we may not receive data sorted on the destination IP address and each data element would represent an additional flow along an arbitrary destination IP address. Formally, each update to \mathbf{a} is represented by a pair (i, d_k) , which is interpreted as “add the value d_k to the i th coordinate.” Thus, at any time t , the value of a_i is the sum of d_{ks} that were added to the i th coordinate. The sequence of updates we see on the stream therefore implicitly represents \mathbf{a} . Thus, we call \mathbf{a} the (implicit) state vector of the data stream.

Data stream processing entails a special constraint. Despite the exponential growth in the capacity of storage devices, it is not common for such streams to be stored. Nor is it desirable or helpful to store them, since the cost of any simple processing—even just sorting the data—would be too great. The main challenge in handling data streams is to perform necessary computations “on the fly” using a small amount of storage, while maintaining low total running time.

Since there is growing focus on manipulating data streams, the database and data processing infrastructure needed to handle stream data is now being investigated. Also, there is a need to identify basic operations of interest in managing data streams and to support them efficiently. The database community has just begun to investigate the challenges involved [2], [4], [27] complementing the efforts emerging in the other communities—algorithms [16], [18], [29], [30], [32], networking [3], physical sciences [44], and elsewhere.

In this paper, we propose *Hamming norm computation* as a basic operation of interest for data stream processing.

- G. Cormode is with the Center for Discrete Mathematics and Computer Science, Rutgers University, 96 Frelinghuysen Road, Piscataway NJ 08854, E-mail: graham@dimacs.rutgers.edu.
- M. Datar is with the Computer Science Department, 482 Gates Computer Science, Stanford University, Stanford CA 94305. E-mail: datar@stanford.edu.
- P. Indyk is with MIT Laboratory for Computer Science, 545 Technology Square, NE43-373, Cambridge MA 02139. E-mail: indyk@theory.lcs.mit.edu.
- S. Muthukrishnan is with the Division of Computer and Information Sciences, Rutgers University, 110 Frelinghuysen Road, Piscataway NJ 08854, and AT&T Research, 180 Park Avenue, Florham Park, NJ 07932. E-mail: muthu@research.att.com.

Manuscript received 2 Aug. 2002; revised 29 Nov. 2002; accepted 4 Dec. 2002. For information on obtaining reprints of this article, please send e-mail to: tkde@computer.org, and reference IEEECS Log Number 117895.

Consider a data stream with state vector \mathbf{a} . The *Hamming norm* of vector \mathbf{a} , written $|\mathbf{a}|_H$, is the number of values i for which $a_i \neq 0$. That is, $|\mathbf{a}|_H = |\{i | a_i \neq 0\}|$.

There are two compelling reasons for proposing Hamming norms. First, in the special case where \mathbf{a} represents a vector of counts similar to our earlier example, $|\mathbf{a}|_H$ is the number of distinct items in the data stream (the number of distinct values i seen); computing $|\mathbf{a}|_H$ on a stream is equivalent to maintaining the number of distinct values in a database relation in the presence of inserts and deletes to it. As such, this is an important problem in traditional databases. Second, when we apply Hamming norm to two (or more) data streams, we get very interesting measures. For two streams that represent state vectors \mathbf{a} and \mathbf{b} , respectively, we may consider the Hamming norm of the sum of the vectors or of their difference. The Hamming norm of the sum $|\mathbf{a} + \mathbf{b}|_H = |\{i | (a_i + b_i) \neq 0\}|$ represents the union of the two streams. The Hamming norm of the difference $|\mathbf{a} - \mathbf{b}|_H = |\{i | (a_i - b_i) \neq 0\}| = |\{i | a_i \neq b_i\}|$ is the number of dimensions in which they differ.¹ Both of these values are of fundamental interest. We will see a few of the large number of uses for the Hamming norm in more detail in Section 2.

The main features of this work are as follows:

1. The initiation of the study of Hamming norm computations for data streams.
2. A novel algorithm for calculating a very small summary for any data stream (what we call the l_0 sketch) such that the Hamming norm of that stream can be found up to a user-specified approximation factor (with high probability) using only the l_0 sketches. This is the first known algorithm for the problem of Hamming norm computation. For the special case of estimating the number of distinct items, algorithms are known to exist and we demonstrate that our algorithm can outperform them. Our algorithm has following properties.
 - The approximation factor is a priori guaranteed to be $1 \pm \epsilon$ with constant probability, and the sketch requires only space $O(1/\epsilon^2)$. Note that this is of constant size for a fixed fraction ϵ and is independent of the size of the data stream of the signal.
 - The l_0 sketches can be maintained efficiently in presence of a stream consisting of dynamic updates. We can estimate Hamming norms without requiring to rescan the entire relation even under an arbitrary mixture of additions and deletions.
 - The l_0 sketches can be computed separately and then combined in a number of ways: They can be added or subtracted to find the union or difference of the corresponding streams. Hamming norm information can be computed in time proportional to the size of the l_0 sketch, which is effectively constant. This procedure is therefore fast, much faster than sampling.
 - The l_0 sketch is an embedding of the Hamming norm into a very small number of dimensions.

1. Notice that the generated "difference stream," $\mathbf{a} - \mathbf{b}$, will usually contain negative values corresponding to points where $b_i > a_i$.

As such l_0 sketches can be used to answer nearest neighbors and other proximity and similarity queries amongst data streams for Hamming norms.

3. We perform a thorough set of experiments with both synthetic data and real NetFlow data drawn from a large ISP, demonstrating the power of the l_0 sketches on computing various Hamming norms. We show experiments where we estimate Hamming norms accurately, correct to within a few percentage points of the correct answer. For this we use a working space of only 8Kb and handle datasets of tens of megabytes, and we can fix this working space while scaling up to gigabytes of data and larger size with the same accuracy guarantees. For finding the Hamming norm of a single stream, which corresponds to the maintenance of the number of distinct elements under insertions and deletions to the databases, our solution is more accurate than existing methods. For multiple general data streams where both insertions and deletions are allowed to occur arbitrarily within both streams, we present the first known solutions for union and difference problems.

The approach of l_0 sketches to approximate Hamming norms allows us to zero-in to estimate distinct values and differences in massive data streams using a very small summary structure.

We motivate Hamming norms in more detail in Section 2, and discuss previous work in Section 3. We present preliminaries in Section 4 and our solution in Section 5. The results of experimental evaluations are shown in Section 6 and conclusions given in Section 7.

2 MOTIVATING HAMMING NORMS OF DATA STREAMS

We will motivate Hamming norms in more detail by considering two applications.

2.1 Maintaining Distinct Values in Traditional Databases

The Hamming norm of a stream² is of large interest in itself. It follows from the definition above that this quantity is precisely the number of distinct items in the stream. For example, let \mathbf{a} be a stream representing any attribute of a given database, so a_i is the number of tuples in the database with value i in the attribute of interest. Computing the Hamming norm of \mathbf{a} provides the number of distinct values of that attribute taken by the tuples. This is a foundational problem. We consider a traditional database table which is subject to a sequence of insertions and deletions of rows. It is of importance to query optimization and otherwise to know the number of distinct values that each attribute of the table assumes. The importance of this problem is highlighted in [9]: "A principled choice of an execution plan by an optimizer heavily depends on the availability of

2. To simplify the exposition, here and in the remainder of this paper we will write "a norm of a stream" instead of "a norm of the implicit state vector of a stream."

statistical summaries like histograms and the number of distinct values in a column for the tables referenced in the query.” Distinct values are also of importance in statistics and scientific computing (see [23], [24], [31]). Unfortunately, it is provably impossible to approximate this statistic without looking at a large fraction of the database (e.g., via sampling) [9]. Our algorithm avoids this problem by *maintaining* the desired statistics under database updates, so that we never have to *compute* it from scratch.

2.2 Monitoring and Auditing Network Databases

Network managers view information from multiple data stream sources. Routers periodically send traffic information: traces of IP packets and IP flows (which are aggregated IP packet flows) [37]; there are management routine updates: SNMP traps, card/interface/link status updates, route reachability via pings, and other alarms [28]; configuration information: topology and various routing tables [3]. Network managers need ways to take this continuous stream of diagnostic information and extract meaningful information. The infrastructure for collecting this information is often error-prone because of unreliable transfer (typically UDP and not TCP is used for data collection), network elements fail (links go down), configuration tables have errors, and data is incomplete (not all network elements might be configured to provide diagnostic data).

Continuous monitoring tools are needed to audit different data sources to ensure their integrity. This calls for “slicing and dicing” different data streams and corroborating them with alternate data sources. We give three examples of how this can be accomplished by computing the Hamming distance between data streams.

1. Let a_i be the number of *transit* IP packets sent from IP address i that enter a part of the network, and b_i be the number of IP packets that exit that part from i . We would like to determine the Hamming Distance between these counts to find out how many transit flows are losing packets within this part of the network.
2. There are published methods for constructing flows from IP packet traces. We can take traces of IP packets, aggregate them, and generate a flow log from this. This can then be compared with the flows generated by routers [14], [37] to spot discrepancies in the network.
3. Denial of Service attacks involve flooding a network with a large number of requests from spoofed IP addresses. Since these addresses are faked, responses are not acknowledged. So, the Hamming difference between a vector of addresses which issued requests and those which sent acknowledgments will be high in this situation [36]. The Hamming norm of the difference between these two vectors provides a quick check for the presence of sustained Denial of Service attacks and other network abnormality, and could be incorporated into network monitoring toolkits.

There are many other potential applications of Hamming norm computation such as in database auditing and data cleaning. Data cleaning requires finding columns that are

mostly similar [15]; the Hamming norm of columns in a table can quickly identify such candidates, even if the rows are arranged in different orders. It is beyond the scope of this paper to go into detail on all these applications, so we do not elaborate further on them.

3 PRIOR WORK

3.1 Work on Data Streams and Sketches

Previous work on data streams has addressed problems of finding approximately the l_2 (Euclidean) norm and the l_1 norm of massive vectors whose entries are listed in an arbitrary order [1], [18], [32]. We study a related problem, of finding the Hamming norm of a vector, and the Hamming distance between pairs of vectors. No previous results were known for these problems in their general form as stated here. As mentioned in the introduction, our algorithms are based on the technique called *sketching*. The basic idea is to represent the whole dataset using only very small amount of space, while preserving important information. When combined with data streams, these sketches must be produced online as the data arrives.

The sketching technique has its roots in the field of mathematics called functional analysis [34]. We consider in particular the application of sketching to approximate l_p norms of the data. The l_p norm of a vector a is equal to

$$\|a\|_p = \left(\sum_i |a_i|^p \right)^{\frac{1}{p}}.$$

Sketching was first applied to tracking approximate size of a self-join of a relation in [1]; in our language, this corresponds to maintaining the l_2 norm of the vector represented by a stream. The techniques of Alon et al. [1] and Johnson and Lindenstrauss [34] were later generalized by Indyk [32] to maintain the l_p norm of the stream vector for any $p \in (0, 2]$. In the context of databases, the group of techniques covered by the umbrella term “sketching” have been applied to finding representative trends in massive one and multidimensional time series data [12], [33]. They have also been applied to multidimensional histograms [43] and data cleaning [15]. But, our concept of l_0 sketch is a novel approach to estimating the Hamming norm and l_0 sketches have not been used previously in the literature.

Besides comparing multiple streams, there has been work on the problem of one pass clustering of data streams [30]. There has been a lot of work in computing over data streams for purposes such as set resemblance, data mining, creating histograms, and so on [11], [17], [29]. Particularly relevant is some recent work [23], [25] which studies the problem of finding the size of the union of two streams. Here, the streams define multisets of elements, and it is the size of the union of the supporting sets that is of interest. Their method is only applicable to streams which consist solely of inserts—they fail when deletions are allowed. The method we present solves this problem when both insertions and deletions are allowed to occur arbitrarily within both streams.

```

initialize  $c[1, 1] \dots c[m, \log n] = 0$ 
for all tuples  $(i, dk)$  do
  for  $j = 1$  to  $m$  do
     $c[j, \text{hash}_j(i)] = c[j, \text{hash}_j(i)] + dk$ 
for  $j = 1$  to  $m$  do
  for  $k = \log n$  downto  $1$  do
    if  $c[j, k] = 0$  then
       $\text{minzero} = k$ 
     $\text{total} = \text{total} + \text{minzero}$ 
return  $(1.2928 \times 2^{\text{total}/m})$ 

```

Fig. 1. The Flajolet-Martin algorithm for computing the Hamming norm of a stream.

3.2 Maintaining Distinct Elements Estimates

There have been two main styles of approach to counting distinct elements: these are sampling-based and synopsis-based. Sampling methods attempt to do a small amount of probing of a complete list of the items in an attempt to find how many distinct values there are [7], [9], [10], [31]. However, sampling-based approaches are known to be inaccurate and substantial lower bounds on the sample size required have been shown [9], proving that a large fraction of the data must be sampled. The alternative paradigm consists of synopsis-based approaches which keep a small summary of the stream, and update this every time an item is added or removed. We focus on these synopsis methods, since they can work in our data streams model, whereas sampling is not suited to dynamic modification of the data. The most widely applicable synopsis method is that of Flajolet and Martin [21], [22], which we describe in outline to enable comparison with our algorithm.

The algorithm is shown in Fig. 1. The crucial part is the set of m hash functions hash_j , which map item values onto the range $[1 \dots \log n]$. hash_j is designed so that the probability $\Pr[\text{hash}_j(i) = \ell] = 2^{-\ell}$, over all choices of hash functions from the family from which hash_j is drawn. Intuitively, this procedure works because if the probability that any item is mapped onto counter ℓ is $2^{-\ell}$, then if there are d nonzero entries in the vector, then we expect $d/2$ to be mapped to the first entry, $d/4$ to be mapped to the second, and so on until there are none expected in the $(\log_2 d)$ th. Several repetitions of this procedure are done independently and the result scaled by an appropriate factor (established in [22] as 1.2928).

Theorem 1 Due to Flajolet and Martin, (Theorem 2 of [21]).

This procedure gives an unbiased estimate of the number of distinct values that are seen in a stream.

This solves the problem of finding the number of distinct values in a stream of values, which as we know is the Hamming norm of that stream. However, this method fails to find the Hamming norm of general vectors since it relies on the input conforming to certain conditions: The result is not defined if the implicit vector \mathbf{a} has any entries that are negative. This can lead to decreasing a counter below zero,

or to producing an estimate of the number of distinct elements that is highly inaccurate. In particular then, this method cannot be used to find the Hamming norm of the difference of two streams, $|\mathbf{a} - \mathbf{b}|_H$, since this leads to negative values in the “difference stream,” $\mathbf{a} - \mathbf{b}$.

There has been much recent work extending these results. Gibbons and Tirthapura [25], [23] propose a method which uses a similar kind of hashing to keep a small sample of the stream, from which the number of distinct values can be approximated. Bar-Yossef et al. [6], [5] describe and analyze a number of algorithms for this problem. In particular, they prove (Theorem 1 of [5]) that the probabilistic counting algorithm of Flajolet and Martin gives a $1 \pm \epsilon$ approximation with probability $1 - \delta$ using space proportional to $1/\epsilon^3 \log 1/\delta \log n$. They also give new algorithms which use $\tilde{O}((1/\epsilon^3 + \log n) \log 1/\delta)$ space,³ a significant improvement in the dependency on n and ϵ . However, one disadvantage is that these new algorithms do not allow values to be deleted (corresponding to a negative value for a d_k). Our algorithm is not limited in this way.

3.3 Database Work on Network Monitoring

Database issues in network monitoring are beginning to get explored. Specific data processing problems in networking databases have been studied such as constructing traffic matrices [20]. The database architecture necessary for processing multiple configuration and data files arising in networks has been discussed in [3], [19]. In the specific case of sensor network data streams, a system architecture has been presented in [35]. At least two different philosophies seem to exist for dealing with network traffic data streams: One is to appropriately sample to decrease them to manageable size and to collate them in massive data warehouses [19], and the other is to deploy a database infrastructure where querying and summarization can be pushed to network elements and processing is distributed (for example, [14], [26]). However, we are not aware of any prior work using the Hamming norm in this context.

4 PRELIMINARIES

4.1 Data Stream Model

We assume a very general, abstracted model of data streams where our input arrives as a stream of data values, (i, d_k) . This indicates that we should add the integer d_k to the count for item i . Clearly, we can accommodate subtractions by allowing d_k to be negative. Update operations have the effect that each tuple (i, d_k) causes $a_i \leftarrow a_i + d_k$. The accumulation of all these pieces of information defines the implicit vector, \mathbf{a} such that $a_i = l$ means that over all tuples for item i the total of the d_k s is l .

An important factor is any restrictions on how the information in the stream arrives. For the most part, we expect the data to arrive in no particular order, since it is unrealistic to expect it to be sorted. Another question is whether every attribute value will be seen at most once, or whether there can be multiple such data items spread out arbitrarily within the stream. Here, we assume the

3. The \tilde{O} notation hides factors in $\log 1/\epsilon$ and $\log \log n$.

most general case, that the data arrives unordered and the same value can appear multiple times within the stream. This is termed the *unordered, unaggregated model* (cash register) in [27].

The processing of massive data streams requires the use of a more restricted model of computation. In this model, we must process a stream of data, with the demand that each item in the stream must be processed completely and then discarded before the next is received. It is not possible to backtrack on the stream, so once a data item has been seen it cannot be retrieved unless it is explicitly stored in the working space. For a method in this model to be useful in practice, it must therefore use an amount of working space much smaller than the total size of the data, and also process each item rapidly. There has been a great deal of interest in processing data streams recently, see, for example, [4], [35].

Example. We consider the following stream of IP flows as $\langle \text{source}, \text{dest} \rangle$ pairs:

$\langle 10.0.0.59, 10.0.0.21 \rangle, \langle 10.0.0.105, 10.0.0.109 \rangle,$
 $\langle 10.0.0.59, 10.0.0.21 \rangle, \langle 10.0.0.105, 10.0.0.17 \rangle,$
 $\langle 10.0.0.59, 10.0.0.105 \rangle, \langle 10.0.0.252, 10.0.0.253 \rangle.$

In this example, we see that the same source address appears many times throughout the stream and the same pair can appear more than once. In a realistic setting, IP address pairs can come in arbitrary order (and are drawn from a space of $(2^{32})^2$ possible pairs).

4.2 Stable Distributions

A vital part of our solution is the use of what are known as stable distributions. We consider statistical distributions that are (strictly) stable, with a stability parameter p . Strictly stable distributions with stability parameter p have the following property: If random variables X_1, X_2, \dots, X_l have stable distributions with stability parameter p , then $a_1 X_1 + a_2 X_2 + \dots + a_l X_l$ is distributed as $(\sum_i |a_i|^p)^{1/p} X_0$, where X_0 is also a random variable with p stable distribution. This property will let us use the stable distributions to compute l_p norms. For example, the Gaussian distribution is stable with stability parameter 2, and the Cauchy distribution is stable with parameter 1. See, for example, [39] for more details of stable distributions.

5 OUR HAMMING NORM COMPUTATION

We first show the algorithm for computing a sketch to approximate the Hamming norm of a single stream. This takes a number of steps: how the Hamming norm can be found via l_p norms, how to create the small sketch that will be used in all our computations, how this can be found in the streaming model, and how values can be drawn from the necessary statistical distributions. We then show how these results can be easily extended to compute norms of combinations (union and differences) of streams.

5.1 The l_0 Norm

Theorem 2. *The Hamming norm $|a|_H$ can be approximated by finding the l_p norm of the vector \mathbf{a} for sufficiently small p ($0 < p \leq \epsilon / \log U$) provided we have an upper bound (U) on the size of each entry in the vector, so $\forall i : |a_i| < U$.*

Proof. We provide another mathematical definition for the Hamming norm which is crucial for our algorithms. We want to find $|\{i | a_i \neq 0\}|$. Observe that $|a_i|^0 = 1$ if $a_i \neq 0$; we can define $|a_i|^0 = 0$ for $a_i = 0$. Thus, the Hamming norm of a vector \mathbf{a} is given by $\sum_i |a_i|^0$. This is similar to the definition of the l_p norm of a vector given in Section 3.1. We must define $l_0 = \sum_i |a_i|^0 = |a|_H$. Hence, the l_0 norm as defined here is our Hamming norm, and we refer to our sketches as l_0 sketches.

We show that the l_0 norm of a vector can be well-approximated by $(l_p)^p$ if we take p small enough. We consider $\sum_i |a_i|^p = (l_p)^p$ for a small value of p ($p > 0$). If, for all i we have that $|a_i| \leq U$ for some upper bound U , then

$$\begin{aligned} |a|_H &= \sum_i |a_i|^0 \leq \sum_i |a_i|^p \leq \sum_i U^p |a_i|^0 \\ &\leq U^p \sum_i |a_i|^0 \leq (1 + \epsilon) \sum_i |a_i|^0 = (1 + \epsilon) |a|_H. \end{aligned}$$

We use the fact that a_i is an integer and $\forall i : |a_i| \leq U$. The last inequality uses $U^p \leq (1 + \epsilon)$ which follows if we set $p \leq \log(1 + \epsilon) / \log U \approx \epsilon / \log(U)$. \square

5.2 Creating the l_0 Sketch

Definition 1. *We define an l_0 sketch vector $sk(\mathbf{a})$ with dimension m as follows: $sk(\mathbf{a})$ is the dot product $\mathbf{x} \cdot \mathbf{a}^T$ (where \mathbf{x} is a matrix of values $x_{i,j}$), so*

$$sk(\mathbf{a})_j = \sum_{i=1}^n x_{i,j} a_i.$$

Each $x_{i,j}$ is drawn independently from a random stable distribution with parameter p , with p as small as possible. Here, $1 \leq i \leq n$ and $1 \leq j \leq m$; n is the dimension of the underlying vector \mathbf{a} , and m is the dimension of the sketch vector.

According to Section 4.2, each entry of $sk(\mathbf{a})$ is distributed as $(\sum_i |a_i|^p)^{1/p} X$, where X is a random variable chosen from a p -stable distribution. We will use $sk(\mathbf{a})$ to make our approximation of the l_0 norm. In particular, we use this sketch to find $\sum_i |a_i|^p$ for $0 < p \leq \epsilon / \log U$, from which we can approximate the Hamming norm up to a $(1 + \epsilon)$ factor. By construction, we can use any $sk(\mathbf{a})_j$ to estimate the $(l_p)^p = \sum_i |a_i|^p$. We combine these values to get a good estimator for the $(l_p)^p$ by taking the median of all entries $|sk(\mathbf{a})_j|^p$. This is a good approximation of the Hamming norm of the stream:

Fact: Let X be any real random variable with density function continuous over $(-\infty, \infty)$. For $m = C \log(1/\delta) / \epsilon^2$ (where C is a large constant), let $s_1 \dots s_m$ be independently chosen from the same distribution as X . Then, $M = \text{median}(s_1, \dots, s_m)$ satisfies

$$1/2 - \epsilon \leq \Pr[X < M] \leq 1/2 + \epsilon$$

with probability $\geq 1 - \delta$.

Proof. Denote $F(t) = \Pr_X[X < t]$. We need to show that if $m = C \log(1/\delta) / \epsilon^2$, then

$$\Pr_{s_1 \dots s_m} [1/2 - \epsilon \leq F(M) \leq 1/2 + \epsilon] \geq 1 - \delta.$$

We will show

$$\Pr_{s_1 \dots s_m} [F(M) < 1/2 - \epsilon] < \delta/2;$$

the probability of $F(M) > 1/2 + \epsilon$ can be bounded in the same way.

Let c^- be the smallest real number such that $F(c^-) = 1/2 - \epsilon$. The key observation is that, since $F(t)$ is nondecreasing in t , $F(M) < 1/2 - \epsilon$ implies $M < c^-$. Now, since M is the median of $s_1 \dots s_m$, $M < c^-$ implies that at least $m/2$ values of s_i s are smaller than c^- . And, this is sufficient since we can define indicator variables $Y_i \in \{0, 1\}$, such that $Y_i = 1$ if and only if $s_i < c^-$. $Y_1 \dots Y_m$ are independent, $\Pr[Y_i = 1] = 1/2 - \epsilon$, and we need to bound from the above probability that $\sum_i Y_i > m/2$, i.e., that the sum of Y_i s is larger from its expectation by a factor of $\frac{1/2}{1/2 - \epsilon} \approx 1 + 2\epsilon$. This is solved by applying the Chernoff bound. \square

We apply this fact to our situation with values drawn from stable distributions.

Theorem 3.

$$(1 - \epsilon)^p \text{median}_j |sk(\mathbf{a})_j|^p \leq \text{median} |X_0|^p \left(\sum_i |a_i|^p \right) \leq (1 + \epsilon)^p \text{median}_j |sk(\mathbf{a})_j|^p$$

with probability $1 - \delta$ if $m = O(1/\epsilon^2 \cdot \log 1/\delta)$, where X_0 is a random variable with symmetric, strictly p -stable distribution.

Proof. The proof of this theorem follows from the results of Indyk in [32] and the above fact. By the properties of stable distributions given in Section 4.2, if $X_0 \dots X_n$ are distributed identically and independently as stable distributions with parameter p , then $a_1 X_1 + \dots a_n X_n$ is distributed as $\|a\|_p X_0$. Then, $(|a_1 X_1 + \dots a_n X_n|)$ is distributed as $(\|a\|_p |X_0|)$. Hence, when we take the median of values $|sk(\mathbf{a})|$ we get $\|a\|_p \text{median}(|X_0|)$, where $\text{median}(|X_0|)$ is the median of absolute values from a stable distribution with parameter p . By the above fact and that the cumulative distribution function F has bounded derivative around the median, we perform this procedure independently $m = O(1/\epsilon^2 \log 1/\delta)$ times. This gives us with probability at least $1 - \delta$

$$(1 - \epsilon) \text{median}_j |sk(\mathbf{a})_j| \leq \text{median}_j |X_0| \left(\sum_i |a_i|^p \right)^{1/p} \leq (1 + \epsilon) \text{median}_j |sk(\mathbf{a})_j|.$$

Raising this to the power p gives the result as stated. \square

5.3 Computation in the Streaming Model

We now list the additional modifications to this procedure to produce a streaming algorithm with low space requirements.

Maintaining the Sketch Under Updates. The l_0 sketch is initially the zero vector, since this is the sketch of an empty stream. We can then build the sketch progressively as each

item in the data stream is received. Our update procedure on receiving tuple (i, d_k) is as follows: We add d_k times $x_{i,j}$ to each entry $sk(\mathbf{a})_j$ ($1 \leq j \leq m$) in the sketch. That is, given (i, d_k)

$$\forall 1 \leq j \leq m : sk(\mathbf{a})_j \leftarrow sk(\mathbf{a})_j + d_k x_{i,j}.$$

Clearly, this procedure ensures that at any point, the sketch is indeed the dot product of the vector \mathbf{a} with \mathbf{x} , as required.

Computation with Small Space Usage. We need to show that this technique can be implemented in small space. So, we do not wish to precompute and store all the values $x_{i,j}$. To do so would consume much more space than simply recording each a_i from which the number of distinct items could be easily found. Instead, we will generate $x_{i,j}$ when it is needed. Note that $x_{i,j}$ may be needed several times during the course of the algorithm, and must take the same value each time. We can achieve this by using pseudorandom generators for the generation of the values from the stable distributions. In other words, we will use standard techniques to generate a sequence of pseudorandom numbers from a seed value (see, for example, [41]). We will use i to seed a pseudorandom number generator $random()$. We will then use the stream of pseudorandom numbers generated by $random()$ to generate a sequence of p -stable distributed random variables $x_{i,1}, x_{i,2}, \dots, x_{i,m}$. This ensures that $x_{i,j}$ takes the same value each time it is used, since we use the same seed i each time, but that it still has the appearance of being drawn from a p -stable distribution. Results in [32] assure us that we can use random number generators in place of a true source of randomness with out any fear of loss of quality of the results.

5.4 Generating Values from Stable Distributions

We need to be able to generate values from a stable distribution with a very small stability parameter p . Standard methods such as those described by Chambers et al. [8], [39] can be used to draw values from a stable distribution with arbitrary parameters. These take uniform random numbers r_1, r_2 drawn from the range $[0 \dots 1]$ and output a value drawn from a stable distribution with parameter p . This is much like the Box-Muller procedure for drawing values from the Normal distribution. We denote this transform as a (deterministic) function $stable(r_1, r_2, p)$ computable in constant time. This function is defined as follows: First, we compute $\theta = \pi(r_1 - \frac{1}{2})$. Then,

$$stable(1/2 + \theta/\pi, r_2, p) = \frac{\sin p\theta}{\cos^{1/p} \theta} \left(\frac{\cos(\theta(1-p))}{-\ln r_2} \right)^{\frac{1-p}{p}}.$$

To use the result of Theorem 3, we need to find $\text{median} |X_0|^p$, the median of absolute values from a stable distribution with parameter p . We can do this in advance using numeric methods and then scale by this constant factor to find the desired result, denoted $scale_factor(p)$ in our algorithm.

However, under our initial implementation, the time cost of using stable distributions against using probabilistic counting was quite high. The time can be much reduced by considering other methods to draw values from stable distributions, since the majority of the processing time is in creating the values of the stable distribution using a transformation from uniform random variables. We now

TABLE 1
Timing Results for Different Methods

Time for each method / ms	Number of repetitions				
	64	128	256	512	1024
Stable Distributions	23.4	46.3	92.0	184	379
Uniform Simulation	0.97	1.6	2.9	5.1	10
Probabilistic Counting	0.40	0.7	1.6	3.2	6.5

describe some ways which allow these values to be drawn much more quickly, while retaining the ability to compute them in the stream with small space overheads. The following ideas are mostly due to John Nolan [39], [40].

Definition 2. A random variable X is said to be in the domain of attraction of Z if and only if there exist constants $c_n > 0$, $d_n \in \mathcal{R}$ with

$$c_n(X_1 + X_2 + \dots + X_n) - d_n \xrightarrow{d} Z,$$

where X_1, X_2, \dots, X_n are independent identically distributed copies of X and \xrightarrow{d} denotes convergence in distribution.

Define $X = \text{sign}(U) * |U|^{-1/p}$, where U is a *Uniform*($-1, 1$) random variable. It can be shown that X is in the “domain of attraction” of a p -stable distribution. This means that if X_1, X_2, \dots, X_n are i.i.d. copies of the random variable X , then $(X_1 + X_2 + \dots + X_n)/cn^{1/p}$ converges in distribution to a p -stable distribution. Note, according to the definition above $c_n = cn^{1/p}$ and $d_n = 0$, where $c = c(p) = \left(\frac{\pi}{2\Gamma(p)\sin(\pi p/2)}\right)^{1/p}$ is a constant independent of n . In general, $a_1X_1 + a_2X_2 + \dots + a_nX_n/c \|a\|_p$ converges to a p -stable distribution.

This indicates that, instead of using the i.i.d. random variables that are p -stable distributed, if we use i.i.d. random variables that are distributed as $X = \text{sign}(U) * |U|^{-1/p}$, where U is *Uniform*($-1, 1$), we get the same resulting distribution as n tends to infinity. Thus, for large values of n , the distribution $a_1X_1 + a_2X_2 + \dots + a_nX_n/c \|a\|_p$ is very close to the p -stable distribution and can be used to estimate the quantity of interest, $\|a\|_p$, using medians as above. A similar trick of using pseudorandom generators to reduce space requirement can be used, as before. Effectively, the method remains the same except for the step where we use the transform $\text{stable}(r_1, r_2, p)$: Instead, we will only generate one random variable r_1 which is uniform over $(-1, 1)$ and add the quantity $d_k * \text{sign}(r_1)|r_1|^{-1/p}$ to $sk[j]$. In order to extract the approximation of $\|a\|_p$, we need to do some additional scaling: Instead of multiplying by $\text{scale factor}(p)$, we multiply by $\text{scale factor}(p)/c(p)^p$. We note that, as p tends to zero, the quantity $c(p)^p$ tends to 1.

The advantage of using the above method is that the quantity $\text{sign}(r_1)|r_1|^{-1/p}$ is faster to compute. We incorporated this method into our implementation and compared it to generating sketches using Stable Distributions directly, as well as to the method of Probabilistic counting. Table 1 summarizes the per-item processing time in milliseconds

```

initialize  $sk[1..m] = 0.0$ 
for all tuples  $(i, dk)$  do
    initialize  $random$  with  $i$ 
    for  $j = 1$  to  $m$  do
         $r1 = random()$ 
         $r2 = random()$ 
         $sk[j] = sk[j] + dk * \text{stable}(r1, r2, p)$ 
    for  $j = 1$  to  $m$  do
         $sk[j] = \text{absolute}(sk[j])^p$ 
return median( $sk[j]$ ) *  $\text{scale factor}(p)$ 

```

Fig. 2. Algorithm to approximate the Hamming norm.

for each of the different methods. The time is independent of the data distribution, and varies linearly with the number of repetitions, that is, the size of the sketch. This is for illustrative purposes only, since we have not attempted to fully optimize the code, and is specific to the machine on which we ran the experiments. However, we observe that while the Stable Distributions method incurs significant per-item processing cost, approximating stable distributions is competitive with probabilistic counting in terms of time taken. In subsequent sections, we will see that it obtains similar levels of accuracy.

The full algorithm for processing a stream is shown in Fig. 2. We can now state our main theorem about computing this norm.

Theorem 4. We can compute a sketch, $sk(a)$ of a stream a using the algorithm in Fig. 2. The sketch has length $m = O(1/\epsilon^2 \cdot \log 1/\delta)$, and allows approximation of the Hamming norm within a factor of $1 \pm \epsilon$ of the true answer with probability $1 - \delta$. Processing each new item and computing the Hamming norm of a both take time linear in the size of the sketch, i.e., $O(1)$ for fixed ϵ and δ .

The proof of this claim relies on the results from the preceding sections and Theorems 2 and 3. The algorithm has two basic parts: The sketch is updated based on every item encountered in the stream; the approximation of the l_0 norm is found by returning the median value of the absolute values of the sketch vector, scaled appropriately.

5.5 Computing Norms of Multiple Streams

With relatively little modification, the above method can be used to find the union or difference of two or more streams.

Theorem 5. The Hamming difference, $|a - b|_H$ can be computed using sketches of size $O(1/\epsilon^2 \log 1/\delta)$. The difference is then approximated within a factor of $1 \pm \epsilon$ with probability $1 - \delta$.

Proof. The Hamming distance between two streams can be computed using just the Hamming norm, since it is equal to $|\{i|a_i \neq b_i\}| = |\{i|(a_i - b_i) \neq 0\}| = |a - b|_H$. Hence, we need to find $sk(a - b)$. As seen before, the sketch $sk(a)$ is the result of the dot product of the induced vector a with the matrix of values $x_{i,j}$. So $sk(a - b) = sk(a) - sk(b)$ follows immediately from the fact that the dot product is a linear function. Therefore, to find the approximate

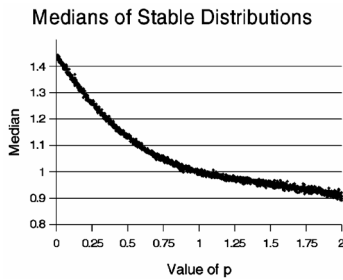


Fig. 3. The value of the median of stable distributions with parameter $p < 2$.

Hamming distance between two streams we have the following simple procedure: 1) Find sketches for each stream individually as before. 2) Compute a new sketch as $sk(a) - sk(b)$. 3) Find the Hamming norm of this compound sketch as described above, by finding the median value. \square

Theorem 6. *The Hamming norm of the union of multiple streams, $|a + b + \dots|_H$ can be computed using sketches of size $O(1/\epsilon^2 \log 1/\delta)$. The norm is then approximated within a factor of $1 \pm \epsilon$ with probability $1 - \delta$.*

This also results from the linearity of the dot product function, in the same way to the above proof. It follows that the union of multiple streams $a, b, \dots = a + b + \dots$ and so a sketch for this can be computed as

$$sk(a + b + \dots) = sk(a) + sk(b) + \dots$$

This, in particular, allows the computation of the number of distinct elements over several separate streams, without overcounting any elements common to two or more of the streams.

6 EXPERIMENTS

We implemented our method of using stable distributions to approximate the Hamming norm of a stream and Hamming norm of two or more streams. We implemented Probabilistic Counting as described in Section 3.2 for approximating the number of distinct elements, since this is the method that comes closest to being able to compute the Hamming norm of a sequence. The two methods of generating sketch values described in Section 5.4 were also tested. This allows us to test how well our methods perform for the two main motivating applications: counting the number of distinct items in a stream, and comparing streams of network data.

6.1 Implementation Issues

For computing with stable distributions, we implemented the method of Chambers et al. [8], [39] to generate stable distributions with arbitrary stability parameters, as well as the method using Uniform variables raised to the power $-1/p$, which converges in distribution to Stable.

Median of Stable Distributions. As noted above, the result we find from our comparison is the L_p norm of vectors, multiplied by the median of absolute values from a stable distribution with parameter p . So, to find the accurate answer, we need to scale this by an appropriate scaling factor. We know of no way to find the necessary scaling factor analytically, so we find it empirically instead. In Fig. 3, we show the results of using random simulation to find the median of stable distributions for different values of the stability parameter p . Each data point represents the

median of 80,000 values chosen at random from a distribution with that parameter. We then take this to the power $1/p$, giving the scaling factor necessary for finding the $(L_p)^p$ distance—this is what we will want when we are using very small values of p to approximate the Hamming norm of a sequence. We repeated the experiment nine times for each value of p as a multiple of 0.01 less than 2. For any given value of p , we can use the median of stable distributions for this parameter to scale our results accordingly.

Ideally, we set the stability parameter p of the stable distribution to be as low as possible, to approximate as well as possible the actual Hamming norm. However, as p gets closer to zero, the values generated from stable distributions get significantly larger, gradually exceeding the range of floating point representation. Through experimentation, we found the smallest value of p that did not generate floating point overflow was 0.02. Hence, we set p to this value, where the median of the stable distribution as generated by this procedure is $1.425 = \text{scale factor}(0.02)$. Note that using $p = 0.02$ means that, even if every distinct element occurs a million times, then the contribution by every distinct element to the Hamming norm will be $(10^6)^{0.02} = 1.318$, so this gives a worst case overestimate of 32 percent. This could be a large source of error, although even this level of error is likely to be acceptable for many applications. In fact, we shall see that most of our experiments show an error of much less than 10 percent.

Experimental Environment. Experiments were run on a Sun Enterprise Server on one of its UltraSparc 400MHz processors. To test our methods, we used a mixture of synthetic data generated by random statistical distributions, and real data from network monitoring tools. For this, we obtained 26Mb of streaming NetFlow data [37], from an AT&T network. We performed a series of experiments, first to compare the accuracy of using sketches against existing methods for counting the number of distinct elements. We started by comparing our approach with the probabilistic counting algorithm for the insertions-only case (i.e., no deletions). We then investigated the problem for streams where both insertions and deletions were allowed. Next, we ran experiments on the more general situations presented by Network data with streams whose entries in the implicit vectors are allowed to be negative. As mentioned earlier, probabilistic counting techniques can fail dramatically when presented with this situation. Finally, we ran experiments for computing the Hamming distance between network data streams and on the union of multiple data streams. We were unable to use Probabilistic Counting in this case, since this method is not applicable in this case. Instead, we compared the two methods of drawing values to use in sketches.

For our experiments, the main measurement that we gathered is how close the approximation was to the correct value. This was done by using exact methods to find the correct answer (*exact*), and then comparing this to the approximation (*approx*). Then, a percentage error was calculated simply as

$$(\max(\text{exact}, \text{approx}) / \min(\text{exact}, \text{approx}) - 1) \times 100\%.$$

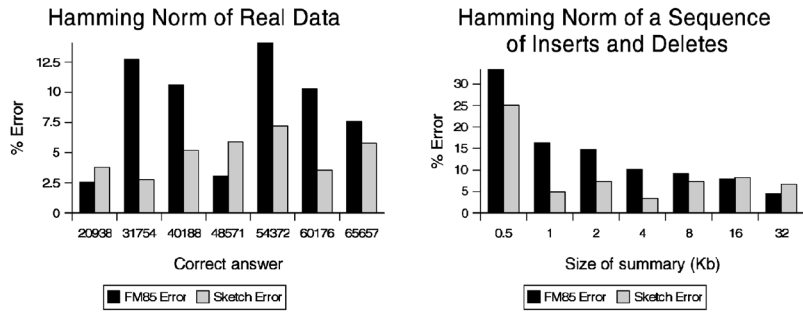


Fig. 4. Testing performance on finding the Hamming Norm.

6.2 Results of Hamming Norm Experiments

Hamming Norm of Network Data. We first examined finding the Hamming norm of the sequence of IP addresses, to find out how many distinct hosts were active. We used exact methods to be able to find the error ratio. We increased the number of items examined from the stream, and looked at between 100,000 and 700,000 items in 100,000 increments. The results presented on the left of Fig. 4 show that there were between 20,000 and 70,000 distinct IP addresses in the stream. Both probabilistic counting and sketches were used, given a workspace of 8Kb. Using sketches is highly competitive with probabilistic counting and is, on the whole, more reliable with an expected error of close to 5 percent against probabilistic counting, which is nearer to 10 percent.

This should also be compared against sampling-based methods as reported in [23], where the error ratio was frequently in excess of 200 percent. This shows that for comparable amounts of space usage, the two methods are competitive with each other for counting distinct items. The worst case for the l_0 sketch occurs when the number of distinct elements is very low (high skew), and here exact methods could be used with small additional space requirements.

Streams Based on Sequences of Inserts and Deletes. Our second experiment tested how the methods worked on more dynamic data, with a mixture of insertions and deletions. It also tested how much they depend on the amount of working space. We created a sequence of insertions and deletions of items, to simulate addition and removal of records from a database table. This was done by

inserting an element with one probability, p_1 , and removing an element with probability p_2 , while ensuring that for each element i , the number of such elements seen was never less than zero. Again, 100,000 transactions were carried out to test the implementation.

We ran a sequence of experiments, varying the amount of working space allocated to the counting programs, from 0.5Kb, up to 32Kb. The results are shown on the right of Fig. 4. The first observation is that the results outdo what we would expect from our theoretical limits from Theorem 4. Even with only 1Kb of working space, the sketching procedure using stable distributions was able to compute a very accurate approximation, correct to within a few percentage points. It is important to note that l_0 sketches were able to nearly equal or better the fidelity of probabilistic counting in every case, and also offer additional functionality. Although in this example the quality of the result does not improve as more working space is made available for it, we claim that this is due to the algorithm being fortunate in small space, since these procedures are in their nature strongly probabilistic. Certainly, with a workspace of only a few kilobytes, we can be sure of a result which is highly likely to be within a few percentage points of the correct answer. This is more than good enough for most of the applications we have already mentioned.

Insertions Only. We tested the algorithms on synthetic data generated from a Zipf distribution with varying levels of skewness. The results are shown in Fig. 5. We used sketches that were vectors with 512 entries, against Flajolet-

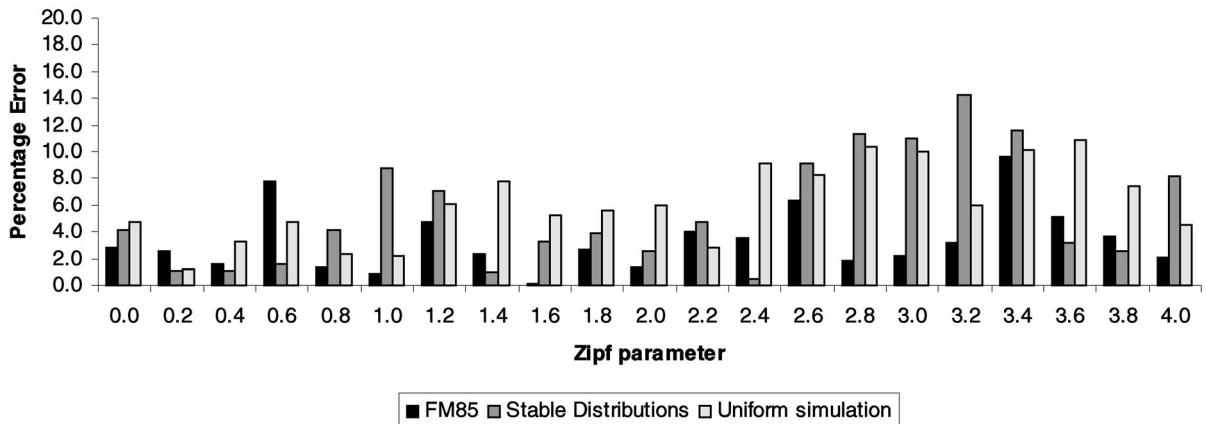


Fig. 5. Results for synthetic data.

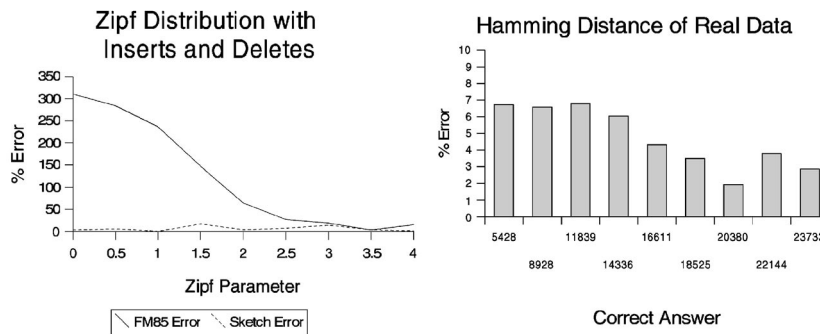


Fig. 6. Testing performance for Network Monitoring.

Martin probabilistic counting given the same amount of working space. We also examined the difference between the two methods of drawing values for the sketch vectors. 100,000 elements were generated from Zipf distributions with parameters ranging from 0 (uniform distribution) up to 4 (highly skewed). The variation in skew tests the ability of the methods to cope with differing numbers of distinct items, since low skew generates many distinct items, whereas high skew gives only a few.

Both algorithms produce answers that, in most cases, are within a few percentage points of the true answer. The worst case is still around 10 percent away from the correct answer. All three methods give similar results for the most part, with little to choose between them. The worst cases occur when the data is highly skewed (skew parameter of 2.5 and higher), when the true number of distinct values is one hundred or lower, so the absolute error is still quite small. The two methods of generating values for sketches appear quite similar in terms of accuracy.

Hamming Norm of Unrestricted Streams. We generated a set of synthetic data to test the method's performance on the more general problems presented by Network data. The main purpose of the next experiment was to highlight that existing methods are unable to cope with many data sets. Zipf distributions with a variety of skewness parameters were used. Additionally, when a value was generated, a coin was tossed: With probability $\frac{1}{2}$, the transaction is an insertion of that element, and with probability $\frac{1}{2}$ it is a deletion of that element. The results are presented on the left of Fig. 6. When we compare the results of probabilistic counting on this sequence to the Hamming norm of the induced vector, we see massive disparity. The error fraction ranges from 20 percent to 400 percent depending on the skew of the distribution, and it is on the uniform distribution on which this procedure performs the worst. On the other hand, using sketches gives a result which is consistently close to the correct answer, and in the same region of error as the previous experiments. Probabilistic counting is only competitive at computing the Hamming norm for distributions of high skew. This is when the number of nonzero entries is low (less than 100), and so it could be computed exactly without difficulty using a small amount of memory.

Hamming Distance Between Network Streams. Our second experiment on network data was to investigate finding the Hamming distance (dissimilarity) between two

streams. This we did by construction, to observe the effect as the Hamming distance increased. We fixed one stream, then constructed a new stream by merging this stream with a second. With probability p , we took item i from the first stream, and with probability $1 - p$ we took item i from the second, for 100,000 items. We then created sketches for both streams and found their difference using sketches of size 8Kb. The results are shown in the right hand chart of Fig. 6 as we varied the probability from 0.1 to 0.9. Here, it was not possible to compare to existing approximation methods, since no other method is able to find the Hamming distance between streams.

The performance of sketching mostly show high fidelity. For sketches generated using Stable Distributions, the answers are all within 7 percent of the true answer, and the trend is for the quality to improve as the size of the Hamming distance between the streams increases. This is because the worst performance of sketches for this problem is observed when the number of different items is low (when the norm is low). This is shown more clearly with the alternate method of generating values, which achieves quite poor accuracy when the difference is smallest, but has high accuracy when the difference increases. This suggests that this method of generating values is of most use when the difference is expected to be reasonably large. For consistent accuracy, the slower method of drawing values is needed.

Union of Multiple Data Streams. Finally, we tested the stated properties of summing sketches to merge data streams. We again used real network data, and for each experiment we split a stream of 100,000 values into a number of pieces. A sketch for each piece was computed separately, and then these sketches combined and compared against the result found when a single sketch was computed. The results found confirm the result of Theorem 6 completely: Perhaps remarkably, no matter how many pieces the stream is split into (2, 5, 10, 20, 50, or 100), the final result is exactly the same. In this case, the norm was approximated as 18745.07, an error of 3.96 percent from the real value. We might have been concerned that round off errors from floating point arithmetic could cause discrepancies between the answers, but this turns out not to be the case.

7 CONCLUDING REMARKS

We have proposed the computation of the Hamming norm of a stream as a basic operation of interest. Computing the Hamming norm of a single stream is a generalization of the problem of maintenance of the number of distinct values in a database under insertions and deletions, which is of fundamental interest. The Hamming norm of two (or more) streams gives the size of their union or the number of places where they differ, depending on whether the norm is computed on the sum of the streams or the difference, respectively. Both these estimations are of great interest as well. In spite of its importance, no algorithms were previously known for computing the Hamming Norm.

We presented a novel and efficient algorithm for computing the " l_0 sketch" for any data stream such that its Hamming norm can be estimated to an arbitrarily small factor using only the sketches. The sketches are very small in size and can be computed in a distributed manner. They can be added to obtain the sketch of the merged stream or subtracted to obtain the sketch of the "difference stream." Their size and accuracy does not depend upon the data distribution. Our experiments with real network flow data demonstrate the powerful accuracy of our algorithm, which outperformed existing methods (where applicable) significantly, and is comparable in terms of speed.

Our Hamming norm estimation technique is more general than is needed in the applications described. For example, we could allow entries in the implicit vector to become nonintegral. Also, our solution will work even if some entries become negative which can occur in certain situations, such as when finding the difference of two or more streams. It can be applied to generalizations of the Hamming norm problem, such as finding the Dominance Norm of multiple streams [13]. New applications may arise in the future where these features will find greater use.

The "sketch" we compute for estimating Hamming norms is suitable for indexing. That is, once we have computed the short "sketches" for data streams, we can cluster, perform similarity and other proximity searches on the streams using *only* the sketches. This is a powerful feature which will be desirable in emerging stream databases.

ACKNOWLEDGMENTS

The authors would like to thank John Nolan for his assistance regarding generating stable distributions.

REFERENCES

- [1] N. Alon, Y. Matias, and M. Szegedy, "The Space Complexity of Approximating the Frequency Moments," *JCSS: J. Computer and System Sciences*, vol. 58, pp. 137-147, 1999.
- [2] Amazon Project, <http://www.cs.cornell.edu/database/amazon/amazon.htm>, 2002.
- [3] S. Babu, L. Subramanian, and J. Widom, "A Data Stream Management System for Network Traffic Management," *Proc. Workshop Network-Related Data Management*, 2001.
- [4] S. Babu and J. Widom, "Continuous Queries over Data Streams," *ACM Sigmod Record*, vol. 30, no. 3, pp. 109-120, 2001.
- [5] Z. Bar-Yossef, T.S. Jayram, R. Kumar, D. Sivakumar, and L. Trevisan, "Counting Distinct Elements in a Data Stream." *Proc. RANDOM 2002*, pp. 1-10, 2002.
- [6] Z. Bar-Yossef, R. Kumar, and D. Sivakumar, "Reductions in Streaming Algorithms, with an Application to Counting Triangles in Graphs," *Proc. ACM-SIAM Ann. Symp. Discrete Algorithms*, pp. 623-632, 2002.
- [7] J. Bunge and M. Fitzpatrick, "Estimating the Number of Species: A Review," *J. Am. Statistical Assoc.*, vol. 88, no. 421, p. 364, 1993.
- [8] J.M. Chambers, C.L. Mallows, and B.W. Stuck, "A Method for Simulating Stable Random Variables," *J. Am. Statistical Assoc.*, vol. 71, no. 354, pp. 340-344, 1976.
- [9] M. Charikar, S. Chaudhuri, R. Motwani, and V.R. Narasayya, "Towards Estimation Error Guarantees for Distinct Values," *Proc. 19th Symp. Principles of Database Systems*, pp. 268-279, 2000.
- [10] S. Chaudhuri, R. Motwani, and V. Narasayya, "Random Sampling for Histogram Construction: How Much is Enough?" *Proc. ACM SIGMOD Int'l Conf. Management of Data*, vol. 27, no. 2, pp. 436-447, 1998.
- [11] E. Cohen, M. Datar, S. Fujiwara, A. Gionis, P. Indyk, R. Motwani, J. Ullman, and C. Yang, "Finding Interesting Associations without Support Pruning," *Proc. 16th Int'l Conf. Data Eng.*, pp. 489-500, 2000.
- [12] G. Cormode, P. Indyk, N. Koudas, and S. Muthukrishnan, "Fast Mining of Tabular Data via Approximate Distance Computations," *Proc. Int'l Conf. Data Eng.*, 2002.
- [13] G. Cormode and S. Muthukrishnan, "Estimating Dominance Norms of Multiple Data Streams," Technical Report 2002-35, DIMACS, 2002.
- [14] C. Cranor, L. Gao, T. Johnson, and O. Spatscheck, "Gigascop: High Performance Network Monitoring with an SQL Interface," *Proc. ACM SIGMOD Int'l Conf. Management of Data*, p. 262, 2002.
- [15] P. Dasu, T. Johnson, S. Muthukrishnan, and V. Shkapenyuk, "Mining Database Structures or How to Build a Data Quality Browser," *Proc. ACM SIGMOD Int'l Conf. Management of Data*, pp. 240-251, 2002.
- [16] DIMACS Workshop Streaming Data Analysis and Mining; and DIMACS Workshop Sublinear Algorithms, Center for Discrete Mathematics and Computer Science, <http://dimacs.rutgers.edu/Workshops/Streaming/>, <http://dimacs.rutgers.edu/Workshops/Sublinear>, 2001.
- [17] P. Domingos, G. Hulten, and L. Spencer, "Mining Time-Changing Data Streams," *Proc. Seventh Int'l Conf. Knowledge Discovery and Data Mining*, pp. 97-106, 2001.
- [18] J. Feigenbaum, S. Kannan, M. Strauss, and M. Viswanathan, "An Approximate L_1 -Difference Algorithm for Massive Data Streams," *Proc. 40th Ann. Symp. Foundations of Computer Science*, pp. 501-511, 1999.
- [19] A. Feldmann, A. Greenberg, C. Lund, N. Reingold, and J. Rexford, "Netscope: Traffic Engineering for IP Networks," *IEEE Network Magazine*, pp. 11-19, 2000.
- [20] A. Feldmann, A. Greenberg, C. Lund, N. Reingold, J. Rexford, and F. True, "Deriving Traffic Demands for Operational IP Networks: Methodology and Experience," *IEEE/ACM Trans. Networking*, pp. 265-279, 2001.
- [21] P. Flajolet and G.N. Martin, "Probabilistic Counting," *Proc. 24th Ann. Symp. Foundations of Computer Science*, pp. 76-82, 1983.
- [22] P. Flajolet and G.N. Martin, "Probabilistic Counting Algorithms for Database Applications," *J. Computer and System Sciences*, vol. 31, pp. 182-209, 1985.
- [23] P. Gibbons, "Distinct Sampling for Highly-Accurate Answers to Distinct Values Queries and Event Reports," *Proc. 27th Int'l Conf. Very Large Databases*, pp. 541-550, 2001.
- [24] P. Gibbons and Y. Matias, "Synopsis Structures for Massive Data Sets," *DIMACS, Series in Discrete Math. and Theoretical Computer Science, A*, 1999.
- [25] P. Gibbons and S. Tirthapura, "Estimating Simple Functions on the Union of Data Streams," *Proc. 13th ACM Symp. Parallel Algorithms and Architectures*, pp. 281-290, 2001.
- [26] A. Gilbert, Y. Kotidis, S. Muthukrishnan, and M. Strauss, "QuickSAND: Quick Summary and Analysis of Network Data," Technical Report 2001-43, DIMACS, 2001.
- [27] A. Gilbert, Y. Kotidis, S. Muthukrishnan, and M. Strauss, "Surfing Wavelets on Streams: One-Pass Summaries for Approximate Aggregate Queries," *Proc. 27th Int'l Conf. Very Large Data Bases*, pp. 79-88, 2001.
- [28] M. Grossglauer, N. Koudas, Y. Park, and A. Varriot, "Falcon: Fault Management via Alarm Warehousing and Mining," *Proc. Workshop Network-Related Data Management*, 2001.

- [29] S. Guha, N. Koudas, and K. Shim, "Data Streams and Histograms," *Proc. Symp. Theory of Computing*, pp. 471-475, 2001.
- [30] S. Guha, N. Mishra, R. Motwani, and L. O'Callaghan, "Clustering Data Streams," *Proc. 41st Ann. Symp. Foundations of Computer Science*, 2000.
- [31] P.J. Haas, J.F. Naughton, S. Seshadri, and L. Stokes, "Sampling-Based Estimation of the Number of Distinct Values of an Attribute," *Proc. 21st Int'l Conf. Very Large Databases*, pp. 311-322, 1995.
- [32] P. Indyk, "Stable Distributions, Pseudorandom Generators, Embeddings and Data Stream Computation," *Proc. 40th Symp. Foundations of Computer Science*, pp. 189-197, 2000.
- [33] P. Indyk, N. Koudas, and S. Muthukrishnan, "Identifying Representative Trends in Massive Time Series Data Sets Using Sketches," *Proc. 26th Int'l Conf. Very Large Databases*, pp. 363-372, 2000.
- [34] W.B. Johnson and J. Lindenstrauss, "Extensions of Lipschitz Mapping into Hilbert Space," *Contemporary Math.*, vol. 26, pp. 189-206, 1984.
- [35] S. Madden and M.J. Franklin, "Fjording the Stream: An Architecture for Queries over Streaming Sensor Data," *Proc. 18th Int'l Conf. Data Eng.*, 2002.
- [36] D. Moore, G.M. Voelker, and S. Savage, "Inferring Internet Denial of Service Activity," *Proc. Usenix Security Symp.*, pp. 9-22, 2001.
- [37] Cisco NetFlow, more details at <http://www.cisco.com/warp/public/732/Tech/netflow/>, 2002.
- [38] NOAA, National Oceanic and Atmospheric Administration, US National Weather Service, <http://www.nws.noaa.gov/>, 2003.
- [39] J. Nolan Stable Distributions, available from <http://academici2.american.edu/~jpnolan/stable/chap1.pdf>, 2002.
- [40] J. Nolan, *Personal Comm.* 2002.
- [41] W.H. Press, S.A. Teukolsky, W.T. Vetterling, and B.P. Flannery, *Numerical Recipes in C: The Art of Scientific Computing*. second ed., Cambridge Univ. Press, 1992.
- [42] California PATH Smart-AHS, more details at <http://path.berkeley.edu/SMART-AHS/index.html>, 1997.
- [43] N. Thaper, S. Guha, P. Indyk, and N. Koudas, "Multidimensional Dynamic Histograms," *Proc. ACM SIGMOD Int'l Conf. Management of Data*, 2002.
- [44] Unidata, Atmospheric Data Repository, <http://www.unidata.ucar.edu/>, 2003.



Graham Cormode received the BA degree in computer science from the University of Cambridge in 1998. He completed the PhD, titled "Sequence Distance Embeddings," at the University of Warwick in 2002. He is currently a postdoctoral fellow at the Center for Discrete Mathematics and Computer Science (DIMACS), based at Rutgers University, Piscataway, New Jersey. He is working on small space algorithms for massive data sets.



Mayur Datar received the BTech degree in computer science from the Indian Institute of Technology (Bombay) in 1998. He was awarded the President of India, Gold Medal for being the most outstanding student of his graduating class. He joined Stanford University in 1998, where he is pursuing the PhD degree in computer science under the supervision of Professor Rajeev Motwani. His academic achievements and honors include: recipient of

the School of Engineering Fellowship for the year 1998, recipient of the Microsoft Graduate Fellowship for the years 2000-01, recipient of the Siebel Scholarship for the year 2002, and winner of the best student paper award at the European Symposium of Algorithms (ESA), 2002. His research interests include design and analysis of algorithms, databases, data mining, data stream algorithms, and combinatorial optimization.



Piotr Indyk received the Magister degree from Warsaw University in 1995, and the PhD degree in computer science from Stanford University in 2001. He has been an assistant professor at MIT since September 2000. His main research area is high-dimensional computational geometry, which involves design and analysis of algorithms for geometric problems (e.g., clustering, closest pair, etc.) in high-dimensional spaces. His other interests include computational geometry in low

dimensions, machine learning, explicit constructions of combinatorial objects (e.g., error-correcting codes), and design of algorithms in general.



S. Muthukrishnan received the PhD degree at the Courant Institute, New York University. Since then, he has been at DIMACS, University of Warwick, UK, Bell Labs, AT&T Labs-Research, and Rutgers University. He works on algorithms, databases, networking, computational biology, scheduling, pattern matching, and a few other assorted areas.

► **For more information on this or any other computing topic, please visit our Digital Library at <http://computer.org/publications/dlib>.**