

New Heavy Hitters Algorithms with Tail Bounds

Anonymous Author(s)

ABSTRACT

In a stream of m items belonging to $\{1, \dots, n\}$, the (ϵ, φ) -Heavy Hitters problem is to output a set S of items containing all items of frequency $\geq \varphi m$ and no item of frequency $< (\varphi - \epsilon)m$. It is one of the most heavily-studied problems in data streams, with a wide variety of applications.

In this work, we propose new randomized counter-based algorithms for the heavy-hitters problem. We show that not only do they have nearly optimal space complexity with respect to ϵ , φ , and n , they also exhibit strong error bounds in terms of the frequency tail of the input and perform favorably on realistic data sets.

CCS CONCEPTS

• **Information systems** \rightarrow **Data mining**; • **Theory of computation** \rightarrow *Design and analysis of algorithms*; • **Networks** \rightarrow *Network monitoring*.

KEYWORDS

data streaming, frequent items, randomized algorithms

ACM Reference Format:

Anonymous Author(s). 2019. New Heavy Hitters Algorithms with Tail Bounds. In *KDD '19: SIGKDD Conference on Knowledge Discovery and Data Mining*, Aug 04–06, 2019, Anchorage, AL. ACM, New York, NY, USA, 9 pages. <https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

1 INTRODUCTION

1.1 Background

The *data streaming* model is an important tool for analyzing algorithms on massive data sets [17]. Streaming algorithms take a long stream of items as input and produce a compact summary of the data as output. This summary can be used to answer queries about the properties of the input data stream. The algorithms are often constrained to take a single pass over the data.

A concrete application is the problem of monitoring IP network traffic. Here, data stream management systems monitor IP packets sent on communication links and perform detailed statistical analyses. These analyses are crucial for fault diagnoses and for verifying network performance and security. Examples of such systems are Gigascope at AT&T [10] and CMON at Sprint [19]. The main challenge in these systems is to quickly and accurately perform the needed analyses in the face of a very high rate of updates.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

KDD '19, August 04–06, 2019, Anchorage, AL

© 2019 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 978-x-xxxx-xxxx-x/YY/MM...\$15.00
<https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

In this work, we consider the *heavy-hitters* problem. Informally, the goal is to find the frequent elements in the stream; the same problem is also studied under the names of top-k, popular items, frequent items, elephants, or iceberg queries. It is very heavily studied and commonly used, e.g. in network analyses to find addresses that account for a large fraction of a network link utilization in a given time window. It is also often solved as a subroutine within more advanced data stream computations, e.g. approximating the entropy of a stream. Cormode and Hadjieleftheriou in [8] give an excellent survey of the work on this problem in a unified framework.

We formulate the heavy-hitters problem rigorously as follows:

Definition 1.1. In the (ϵ, φ) -Heavy Hitters Problem, we are given parameters $0 < \epsilon < \varphi \leq 1$ and a stream a_1, \dots, a_m of items a_j in $\{1, 2, \dots, n\}$. Let f_i denote the number of occurrences of item i , i.e., its frequency¹. An algorithm for the problem should make one pass over the stream and at the end of the stream output a set $S \subseteq \{1, 2, \dots, n\}$ which contains any i such that $f_i \geq \varphi m$ and does not contain any i such that $f_i < (\varphi - \epsilon)m$.

Furthermore, for each item $i \in S$, the algorithm should output an estimate \hat{f}_i of the frequency f_i which satisfies $|f_i - \hat{f}_i| \leq \epsilon m$.

A large number of algorithms have been proposed for finding the heavy hitters and their variants. Following [8], they can be broadly classified into *Counter-based* algorithms, *Quantile* Algorithms, and *Sketch* algorithms.

One of the first algorithms for the problem, as well as perhaps the most widely used, is the deterministic counter-based FREQUENT algorithm due to Misra and Gries [16] (independently re-discovered 20 years later with some refinements to the implementation by Karp et al. [13] and Demaine et al. [12]). There are several variants of the algorithm, notably LOSSYCOUNTING [14] and SPACESAVING [15], that sometimes perform better in terms of space usage or update time than the original version. All these algorithms use $O(\epsilon^{-1}(\log n + \log m))$ bits of space; note the bound does *not* depend on φ . Sketching algorithms, such as the CountSketch [7] or the Count-Min Sketch [9], require more storage (but they can handle more general streams with both insertions and deletions that we do not study here).

Recently, in [5], Bhattacharyya et al. introduced two new algorithms for the (ϵ, φ) -heavy hitters problem that improve on the space usage of FREQUENT, particularly when $\varphi \gg \epsilon$. One of their algorithms uses $O(\epsilon^{-1} \log \varphi^{-1} + \varphi^{-1} \log n + \log \log m)$ bits of storage, which is optimal upto constant factors. Moreover, the algorithms take a constant amount of time per update. These algorithms are fundamentally counter-based and use the Misra-Gries algorithm as a subroutine, but they also use random sampling and hashing in crucial ways.

¹We assume throughout this paper that there are only insertions (no deletion) in the stream and that each insertion of an item increases its frequency by 1.

1.2 Our Work

Given the strong theoretical guarantees of the algorithms proposed in [5], it is natural to wonder how well they fare in practice.

Berinde et al. [4] observe that one of the main reasons for the real-world success of FREQUENT is that it displays “better-than-advertised” performance on streams which have a *skewed* frequency distribution. In particular, they showed that the error made by FREQUENT in frequency estimation satisfies a “tail bound”, so that it is very small when heavy hitters constitute most of the stream.

The starting point of this work is the observation that the algorithms proposed in [5] do *not* exhibit a tail bound, so that the estimation error can be large even when the frequency distribution is very skewed. We design new randomized algorithms inspired by [5] that improve FREQUENT in terms of bits of storage while still achieving a tail bound. We show two different algorithms with this feature. Although they are novel to the best of our knowledge, they are both quite simple to implement and analyze. The first algorithm combines a FREQUENT data structure with a Count-Min sketch [9]. The second algorithm maintains two FREQUENT data structures, one with a smaller number of counters containing items from the original universe and another with a larger number of counters recording hashes of the original item id’s. Armed with the tail guarantee and following the analysis of Berinde et al. [4], we then show that our proposed algorithms enjoy several other appealing properties, such as guarantees for skewed Zipfian streams, for the sparse recovery problem and for merging multiple streams.

We also experimentally evaluate our proposed algorithms and compare them to existing ones. We consider two skewed datasets. One is generated by a ‘planted’ model where 5 items are randomly chosen to be the heavy hitters and they constitute 10% of the stream. The other is generated by a Zipfian distribution which is commonly used to model real-world frequencies. We show that indeed, we can improve the space usage in the regime where $m \ll n$ and $\varepsilon \ll \varphi$ while having high recall and precision rates. We note that existing algorithms can have better update times, leading to a direction for future work.

2 PRELIMINARIES

2.1 General

For a stream consisting of m insertions from the universe $[n] := \{1, \dots, n\}$, we let the n -dimensional vector $f = (f_1, \dots, f_n)$ denote the frequencies of the n items. So, $\sum_{i=1}^n f_i = m$. Without loss of generality, we will assume that $f_1 \geq f_2 \geq \dots \geq f_n$.

For any integer $k \geq 0$, the *residual k -tail* of the stream is defined as²:

$$F_1^{\text{res}(k)} = \sum_{i=k+1}^n f_i.$$

That is, the residual k -tail is the sum of the frequencies except the k largest ones.

Recall from Definition 1.1 that an algorithm for the (ε, φ) -heavy hitters problem is supposed to output estimates \tilde{f}_i for each i in the output set such that:

$$|\tilde{f}_i - f_i| \leq \varepsilon \cdot m = \varepsilon \cdot F_1^{\text{res}(0)}.$$

²The 1 in the subscript refers to the fact that this is the 1st moment of the residual stream.

Definition 2.1 (Tail Bound). An algorithm for the (ε, φ) -heavy hitters problem is said to satisfy the *K -tail bound* if for every $0 \leq k \leq K$:

$$|\tilde{f}_i - f_i| \leq \varepsilon \cdot F_1^{\text{res}(k)}$$

for all i in the output set of the algorithm.

The larger the K we can get for the tail bound, the stronger is the obtained guarantee.

Our proposed algorithms use the notion of a *universal hash family*, which we define next.

Definition 2.2. A family \mathcal{H} consisting of functions $h : U \rightarrow [m]$ is said to be a *universal hash family* if for all distinct $x, y \in U$:

$$\Pr_{h \leftarrow \mathcal{H}} [h(x) = h(y)] \leq \frac{1}{m}$$

where the probability is over h drawn uniformly at random from \mathcal{H} .

It is folklore that if $m \leq |U|$, there exists a universal hash family whose members can be represented uniquely using $O(\log |U|)$ bits.

2.2 The FREQUENT algorithm

Algorithm 1: FREQUENT, [16]

Data: Number of counters t

Initialization: **begin**

 Empty set T

 Empty array c of length n

end

Function *Insert*(i)

if $i \in T$ **then**

$c_i \leftarrow c_i + 1$;

end

else if $|T| < t$ **then**

$T \leftarrow T \cup \{i\}$;

$c_i \leftarrow 1$;

end

else

foreach $j \in T$ **do**

$c_j \leftarrow c_j - 1$;

if $c_j = 0$ **then** $T \leftarrow T \setminus \{j\}$;

end

end

Function *Report*()

foreach $i \in T$ **do** Output (i, c_i) ;

FREQUENT, shown in Algorithm 1, maintains t (item, counter) pairs. Its reported frequencies are accurate to within $m/(t+1)$ where m is the length of the stream. In particular, it reports all items with frequency at least $m/(t+1)$.

A simple grouping argument can be used to verify the correctness of the algorithm. With $t = \lceil 1/\varepsilon \rceil$ counters, FREQUENT solves the (ε, φ) -Heavy hitters problem using $O(\varepsilon^{-1}(\log n + \log m))$ bits of space. Berinde et al. [4] showed a stronger tail bound:

THEOREM 2.3. ([4]) For any $t \geq \frac{1}{\varepsilon}$, the FREQUENT algorithm with t counters satisfies the $(t - \frac{1}{\varepsilon})$ -tail bound and uses $O(t(\log n + \log m))$ bits of storage.

2.3 The work of [5]

In the recent work [5], Bhattacharyya, Dey and Woodruff introduce a couple of randomized counter-based algorithms that improve upon FREQUENT in terms of the number of bits needed for storage. Their first algorithm solves the (ε, φ) -Heavy hitters problem with space $O(1/\varepsilon \log 1/\varepsilon + 1/\varphi \log n + \log \log m)$ bits, while their second algorithm uses $O(1/\varepsilon \log 1/\varphi + 1/\varphi \log n + \log \log m)$ bits. When $\varepsilon \ll \varphi$, the space requirement for these algorithms is asymptotically smaller than the worst-case space requirement for FREQUENT. In fact, [5] show that the second algorithm achieves the optimal (up to constant factors) number of bits for solving the (ε, φ) -Heavy hitters problem.

The first algorithm of [5], which we dub COMPRESSEDFREQUENT, is shown below as Algorithm 2. COMPRESSEDFREQUENT runs Misra and Gries' FREQUENT algorithm on a sampled stream of length $O(1/\varepsilon^2)$ on a hashed universe of size $O(1/\varepsilon^4)$. By the choice of parameters, there are no hash collisions among the elements sampled from the stream with high probability. Also, by standard concentration results, the relative frequency of any element in the original stream and the sampled stream differ by at most $\varepsilon/2$. Some additional bookkeeping needs to be done to ensure that we have the original (unhashed) id's of the heavy hitters, but otherwise, the analysis is complete.

The second algorithm of [5] is more intricate. It also operates on $O(1/\varepsilon^2)$ samples from the original stream and then maintains a compressed histogram for hashes of the original items. It is somewhat reminiscent of the Count-Min sketch algorithm of Cormode and Muthukrishnan [9], but it also subsamples items from the sample with subsampling rate tuned to a multiplicative approximation of the frequency. We omit a more detailed description of this algorithm because we do not study it in this work.

3 NEW ALGORITHMS

We start with the observation that the algorithms proposed in [5] cannot satisfy a non-trivial tail guarantee. Consider the COMPRESSEDFREQUENT algorithm shown in Figure 2, and suppose the input stream only consists of two items with each inserted $\Omega(m)$ times. Although $F_1^{\text{res}(2)} = 0$, the frequency estimates output by the algorithm will still have nonzero error because due to random sampling, the relative frequency of the two items in the sampled stream will (with high probability) differ from that in the original stream. The same argument holds for the second algorithm of [5].

3.1 SKETCHFREQUENT

Our first algorithm, shown above, runs the Count-Min sketch on the output of a FREQUENT counter solving the $(\varphi, \varphi/2)$ -Heavy hitters problem.

THEOREM 3.1. Let $\ell = \log(4/\delta\varphi)$ as in the algorithm above. Suppose that $\varepsilon < \frac{\varphi^2\delta}{2\ell}$.

Algorithm 2: COMPRESSEDFREQUENT, [5]

Data: Parameters ε and φ with $0 < \varepsilon \leq \varphi < 1$.

Initialization: begin

$\ell \leftarrow 6 \log(6/\delta)/\varepsilon^2$.

Hash function h drawn uniformly from a universal family $H \subseteq \{[n] \rightarrow [4\ell^2/\delta]\}$.

An empty FREQUENT data structure \mathcal{T}_1 with ε^{-1} counters.

An empty set \mathcal{T}_2 .

end

Function Insert(x)

With probability $6\ell/m$, **continue**. Else, **return**;

Call Insert($h(x)$) on \mathcal{T}_1 ;

if $h(x)$ is among the top $1/\varphi$ valued items in \mathcal{T}_1 **then**

if $x \notin \mathcal{T}_2$ **then**

if $|\mathcal{T}_2| = 1/\varphi$ **then**

 For some y in \mathcal{T}_2 such that $h(y)$ is not among the top $1/\varphi$ valued items in \mathcal{T}_1 , replace y with x

else

 Insert x in \mathcal{T}_2 ;

end

end

end

Function Report()

foreach x in \mathcal{T}_2 **do**

if $\mathcal{T}_1[h(x)] \geq (\varphi - \varepsilon/2)\ell$ **then** Output $(x, \mathcal{T}_1[h(x)])$;

end

Algorithm 3: SKETCHFREQUENT

Data: Width w , Input parameters $\varphi, \delta \in (0, 1)$.

Initialization: begin

$\ell \leftarrow \log(4/\delta\varphi)$

Hash functions h_1, \dots, h_ℓ drawn uniformly from a universal family $H \subseteq \{[n] \rightarrow [w]\}$.

Empty FREQUENT data structure \mathcal{T} with $2/\varphi$ counters.

Empty matrix C with ℓ rows and w columns.

end

Function Insert(x)

 Insert x into \mathcal{T} ;

for $i = 1$ to ℓ **do**

 Increment $C[i][h_i(x)]$;

end

Function Report()

foreach $x \in \mathcal{T}$ **do**

$\hat{f}_x \leftarrow \min_{i \in [\ell]} C[i][h_i(x)]$;

if $\hat{f}_x \geq \varphi m$ **then**

 Output (x, \hat{f}_x) ;

end

end

For any $w \geq \frac{2}{\varepsilon}$ and $\delta \in [0, 1]$, *SKETCHFREQUENT* satisfies the $\sqrt{\frac{\delta w}{\ell}}$ -tail bound with probability $1 - \delta$ and uses $O(\varphi^{-1} \log n + w \log(\delta \varphi)^{-1} \log m)$ bits of storage.

In particular, for $w = \frac{2}{\varepsilon}$ and $\delta = \frac{1}{5}$, *SKETCHFREQUENT* solves the (ε, φ) -heavy hitters problem with probability at least $\frac{4}{5}$ using $O(\varphi^{-1} \log n + \varepsilon^{-1} \log \varphi^{-1} \log m)$ bits³.

PROOF. The space complexity is immediate. We will focus on the tail bound. We argue first that on the top $k = \sqrt{\delta w / \ell}$ elements of the stream, h_1, \dots, h_ℓ do not collide with probability at least $1 - \delta/2$.

LEMMA 3.2. Let $k = \sqrt{\delta w / \ell}$. For any fixed $T \subseteq [n]$ of size k ,

$$\Pr[\exists i \in [\ell], \exists x \neq y \in T, h_i(x) = h_i(y)] < \delta/2.$$

PROOF. By definition of a universal hash family, for any fixed $i \in [\ell]$ and $x \neq y \in T$, $\Pr[h_i(x) = h_i(y)] < 1/w$. Using the union bound:

$$\Pr[\exists i \in [\ell], \exists x \neq y \in T, h_i(x) = h_i(y)] \leq \ell \cdot \binom{t}{2} \cdot \frac{1}{w} < \frac{\delta}{2}$$

□

Condition on the above event. Note that by our assumption, $k > 2/\varphi$. We can now run the usual analysis of the count-min sketch, which we reproduce here for completeness.

Consider an $x \in [n]$ that is stored in one of \mathcal{T} 's counters. Because of the above, and using again universality of the hash family, for any $i \in [\ell]$:

$$\mathbb{E}[C[i][h_i(x)]] = f_x + \frac{1}{w} \sum_{y \neq x} f_y = f_x + \frac{F_1^{\text{res}(k)}}{w} \leq f_x + \frac{\varepsilon}{2} F_1^{\text{res}(k)}.$$

Then, by Markov's inequality: $\Pr[C[i][h_i(x)] - f_x > \varepsilon F_1^{\text{res}(k)}] \leq \frac{1}{2}$. Therefore:

$$\Pr[\hat{f}_x - f_x > \varepsilon F_1^{\text{res}(k)}] \leq \frac{1}{2\ell} = \frac{\varphi\delta}{4}.$$

Doing a union bound over $\frac{2}{\varphi}$ elements shows that with probability at least $1 - \delta/2$, for each $x \in \mathcal{T}$, $\hat{f}_x - f_x \leq \varepsilon F_1^{\text{res}(k)}$. Moreover, observe that every x with $f_x \geq \varphi m$ is output by the algorithm. □

3.2 DOUBLEFREQUENT

Our second algorithm runs two *FREQUENT* algorithms in parallel. The first is a gross estimation of the frequencies with relative error up to $\varphi/2$. This is done to filter out a list of $O(1/\varphi)$ elements on which we would like the estimate to be more accurate. For this, we use more counters but we also hash the universe size down to $\text{poly}(1/\varepsilon)$ to save space. The resulting algorithm inherits the desirable tail behavior of *FREQUENT*. The analysis which follows is an interesting mix of the analyses for *FREQUENT* and the Count-Min sketch.

THEOREM 3.3. Assume $\varepsilon < \varphi/2$. For any $t \geq \frac{1}{\varepsilon}$ and $\delta \in [0, 1]$, *DOUBLEFREQUENT* satisfies the $(t - \frac{1}{\varepsilon})$ -tail bound with probability $1 - \delta$ and uses $O(\varphi^{-1} \log n + t \log m + t \log(t/\delta))$ bits of storage.

³The assumption that $\varepsilon < \varphi^2/2\ell$ is not needed for the last part of the theorem.

Algorithm 4: DOUBLEFREQUENT

Data: Number of counters t , Input parameters $\varepsilon, \varphi, \delta \in (0, 1)$ with $\varepsilon \leq \varphi/2$.

Initialization: begin

Hash function h drawn uniformly from a universal family $H \subseteq \{[n] \rightarrow [2t^2/\delta]\}$.

An empty *FREQUENT* data structure \mathcal{T}_1 with $2/\varphi$ counters.

An empty *FREQUENT* data structure \mathcal{T}_2 with t counters.

end

Function Insert(x)

Insert x into \mathcal{T}_1 ;

Insert $h(x)$ into \mathcal{T}_2 ;

Function Report()

foreach $x \in \mathcal{T}_1$ **do**

if $\mathcal{T}_1(x) \geq \frac{\varphi m}{2} \wedge \mathcal{T}_2[h(x)] \geq (\varphi - \varepsilon)m$ **then**

 Output $(x, \mathcal{T}_2[h(x)])$;

end

end

In particular, setting $t = 1/\varepsilon$ and $\delta = 1/5$, *DOUBLEFREQUENT* solves the (ε, φ) -heavy hitters problem with probability at least $\frac{4}{5}$ using space $O(\varphi^{-1} \log n + \varepsilon^{-1} \log m + \varepsilon^{-1} \log \varepsilon^{-1})$ bits of storage.

PROOF. First, observe that there are likely going to be no hash collisions among the top $1/\varepsilon$ elements in the stream.

LEMMA 3.4. For any fixed $S \subseteq [n]$ of size t :

$$\Pr[\exists x \neq y \in S, h(x) = h(y)] \leq \frac{\delta}{2}.$$

PROOF. By the union bound and definition of universal hash family,

$$\Pr[\exists x \neq y \in S, h(x) = h(y)] \leq \binom{t}{2} \cdot \frac{1}{2t^2/\delta} \leq \frac{\delta}{2}$$

□

Henceforth, condition on Lemma 3.4 holding true for the top t elements. Define:

$$\bar{f}_x = \sum_{y: h(y)=h(x)} f_y.$$

Clearly, $\bar{f}_x \geq f_x$ for any x . Also, without loss of generality, suppose that $f_1 \geq f_2 \geq \dots \geq f_n$ so that:

$$F_1^{\text{res}(k)} = \sum_{i=k+1}^n f_i.$$

LEMMA 3.5. For any $0 \leq k \leq t$:

$$\Pr[\exists x, \mathcal{T}_1[x] \geq \varphi m/2 \wedge \bar{f}_x > f_x + \varepsilon \cdot F_1^{\text{res}(k)}] < \frac{\delta}{2}$$

PROOF. We observe that for x is among the top $2/\varphi \leq 1/\varepsilon \leq t$ elements. For any such x :

$$\mathbb{E}\bar{f}_x = f_x + \mathbb{E} \sum_{y \neq x: h(y)=h(x)} f_y = f_x + \mathbb{E} \sum_{y > t: h(y)=h(x)} f_y$$

where the last equality is from conditioning on no collisions in the top t elements. Again, using the definition of universal hash family, we get that:

$$\mathbb{E}\bar{f}_x \leq f_x + \frac{\delta}{2t^2} F_1^{\text{res}(k)}$$

for any $k \leq t$. By Markov's inequality and the union bound:

$$\Pr[\exists x, \mathcal{T}_1[x] \geq \varphi m/2 \wedge \bar{f}_x > f_x + \varepsilon \cdot F_1^{\text{res}(k)}] < \frac{2}{\varphi} \frac{1}{\varepsilon} \frac{\delta}{2t^2} \leq \frac{\delta}{2}$$

□

We now bound the error from the use of the FREQUENT algorithm in \mathcal{T}_2 .

LEMMA 3.6. For all x and all $0 \leq k \leq t - \frac{1}{\varepsilon}$:

$$\bar{f}_x - \mathcal{T}_2[h(x)] \leq \varepsilon F_1^{\text{res}(k)}.$$

PROOF. This directly follows from Theorem 2.3 shown in [4] and the fact that the residual k -tail of \bar{f} is at most the residual k -tail of f . □

Putting the lemmas together, we get that with probability at least $1 - \delta$, for all $0 \leq k \leq t - \frac{1}{\varepsilon}$, if x is such that $\mathcal{T}_1[x] \geq \varphi m/2$:

$$f_x \leq \bar{f}_x \leq \mathcal{T}_2[h(x)] + \varepsilon F_1^{\text{res}(k)}$$

from Lemma 3.6 and

$$f_x \geq \bar{f}_x - \varepsilon F_1^{\text{res}(k)} \geq \mathcal{T}_2[h(x)] - \varepsilon F_1^{\text{res}(k)}$$

from Lemma 3.5.

The space bound for DOUBLEFREQUENT is immediate from the space bound for FREQUENT. □

3.3 Consequences

Berinde et al. [4] show that the tail guarantee implies desirable behavior for a few other settings. We state the results below for the DOUBLEFREQUENT algorithm, though similar results should also hold for the SKETCHFREQUENT algorithm.

3.3.1 *Zipfian distribution.* A stream of length m on a universe of size n is said be *Zipfian* if there exists $\alpha \geq 1$ such that for every $i \in [n]$,

$$f_i = \frac{m}{i^\alpha \zeta(\alpha)} \quad \text{where} \quad \zeta(\alpha) = \sum_{i=1}^n \frac{1}{i^\alpha}.$$

In fact, for the results here, it is only required that the tail of the distribution be upper-bounded by the above Zipfian distribution. Of course, the order of arrivals in the stream is arbitrary.

THEOREM 3.7. Given a Zipfian stream with parameter $\alpha \geq 1$, the error in the frequency for any item reported by the DOUBLEFREQUENT algorithm with parameters δ, ε , and φ is at most $O(\varepsilon^\alpha m)$ with probability at least $1 - \delta$.

PROOF. Corollary of Theorem 8 in [4]. □

3.3.2 *k-sparse recovery.* In the k -sparse recovery problem, given a frequency vector f , we want to find a vector f' such that f' is k -sparse (has only k nonzero entries) and the p -norm error $\|f - f'\|_p = (\sum_i |f_i - f'_i|^p)^{1/p}$ is minimized. The next theorem shows that if $k \leq 1/\varphi$, taking the top k of the output of DOUBLEFREQUENT is a good approximation.

THEOREM 3.8. If $k \leq \frac{1}{\varphi}$ and f' is the k -sparse vector obtained by taking the top k elements of the output of DOUBLEFREQUENT (run with parameters δ, ε , and φ), with probability at least $1 - \delta$, for any $p \geq 1$:

$$\|f - f'\|_p \leq (F_p^{\text{res}(k)})^{1/p} + O\left(\frac{\varepsilon F_1^{\text{res}(k)}}{k^{2-1/p}}\right)$$

where $(F_p^{\text{res}(k)})^{1/p}$ is the smallest L_p error of any k -sparse recovery of f .

PROOF. Corollary of Theorem 5 in [4]. □

The following result also gives a guarantee if we would like to output exactly the top k elements.

THEOREM 3.9. Given a Zipfian stream with parameter $\alpha > 1$, if $k = O((\alpha/\varphi^\alpha)^{1/(\alpha+1)})$, then with probability at least $1 - \delta$, DOUBLEFREQUENT can retrieve the top- k elements of the stream in correct order.

PROOF. Corollary of Theorem 9 of [4]. □

4 EXPERIMENTAL EVALUATIONS

4.1 Setup

Existing algorithms for the heavy hitters problem have been carefully implemented in a few different places: the MassDAL Public Code Bank [1], its extension by Cormode [2], and the Yahoo Sketches Library [3]. We chose to build on the second option, Cormode's C++ extension of the MassDAL codebase, to include our implementations. We deployed our experiments on a standard laptop with an Intel Core i7-6700HQ processor at 260 GHz and 16 GB of RAM, running Ubuntu 4.8.4. The C++ code was compiled using gcc. (Due to the double-blind reviewing policy, we are unable to give a pointer to the code currently; however, it will be made public in the future.)

Among previous algorithms, we considered FREQUENT and LCDelta. LCDelta is a version of the Lossy Counting algorithm of Manku and Motwani [14] that is described in Cormode and Hadjieleftheriou's survey [8] and that outperforms FREQUENT in many of their experiments. For our new algorithms, in addition to the above described SKETCHFREQUENT and DOUBLEFREQUENT, we also considered SKETCHLC and DOUBLELC which invoke the LCDelta algorithm in place of FREQUENT. The analyses of SKETCHLC and DOUBLELC mimic that of SKETCHFREQUENT and DOUBLEFREQUENT respectively, and we do not reproduce them.

We ran experiments on two sets of synthetically generated data. The first set is generated from a planted distribution where 5 randomly chosen elements each constitute 2% of the stream, while the other $n - 5$ elements are inserted with uniform probability for the rest 90% of the stream. Here, we set the universe size $n = 10^5$ and the length of each stream $m = 100$. The second set is generated

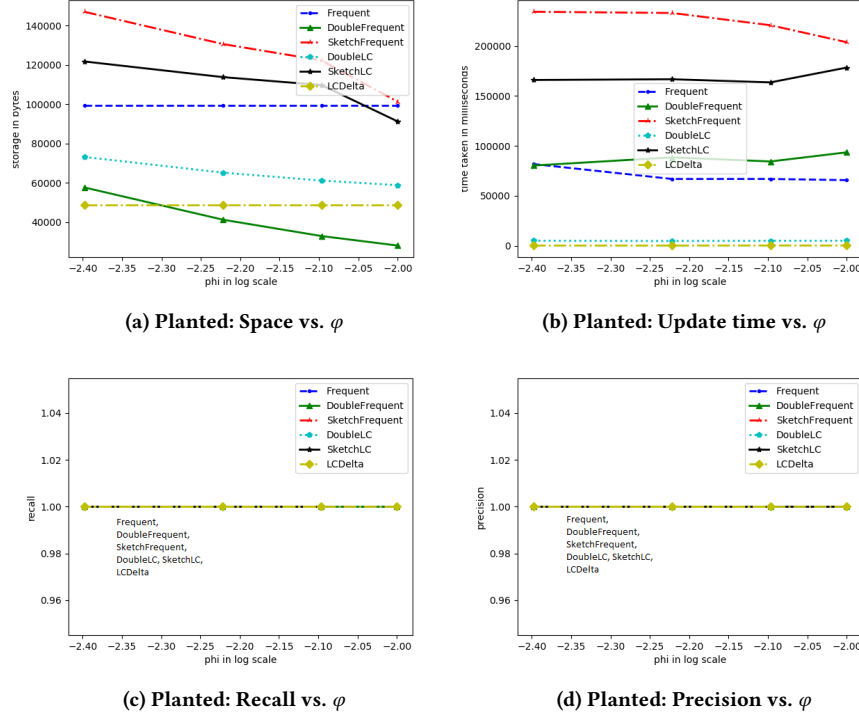


Figure 1: Comparison between the algorithms when run on data generated from the planted model. Here, ϵ is fixed to 10^{-3} .

from skewed Zipfian distributions, where the skew ranges from 1.2 to 2. In this case, the streams are drawn from a universe of size 10^5 , while the length of each stream is 16. For both datasets, we varied ϵ from 10^{-5} to 10^{-2} and varied ϕ from 0.004 to 0.01 (while keeping $\phi \geq \epsilon$).

Note that although we do not test our algorithms on real-world datasets, the latter can often be approximated with Zipfian distributions. For example, it has been established that sizes of cities and word frequencies in text [20], citations of papers [18], web page accesses [6], and file transfer sizes [11] have distributions whose tails are bounded by Zipfian distributions. Thus, we can reasonably expect that the results of our experiments on the Zipfian dataset capture some behavior that would be observed in reality.

The algorithms are compared below according to the following attributes:

- **Space:** Storage space used, measured in bytes
- **Update time:** Time take to insert all the items in the stream, measured in milliseconds
- **Recall:** number of true heavy hitters reported, divided by the total number of true heavy hitters.
- **Precision:** number of true heavy hitters reported, divided by the total number of items reported.

For each of the above, we performed 5 runs per experiment.

4.2 Planted Data

First, we present results for the dataset generated using the planted model. We investigated how the different algorithms fare against

each other as we vary the parameters ϕ and ϵ . Figure 1 shows the data when $\epsilon = 0.001$ and ϕ is varied, while Figure 2 is for $\phi = 0.01$ and ϵ varying.

Consider the case when ϕ is varied. Figure 1a shows that as ϕ becomes larger, DOUBLEFREQUENT uses the least space. Second in the race is LCDelta which beats out DOUBLELC. FREQUENT and the algorithms using sketches need at least 50% more space. In terms of update time, Figure 1b shows that LCDelta and DOUBLELC are far more efficient than the rest, while those two are comparable. Figures 1c and 1d show that all the algorithms achieve perfect recall and precision for this setting of parameters.

Now, we turn to Figure 2 where ϵ is varied as ϕ is kept at 0.01. The range over which ϵ is varied is much bigger (10^{-5} to 10^{-2}) than the range over which ϕ was varied in Figure 1, so here, we see a more high-level view of the algorithms' behavior. As ϵ becomes small, the algorithm taking the least amount of space is DOUBLEFREQUENT followed by LCDelta and DOUBLELC (Figure 2a). Note that when ϵ is not much smaller than $\phi = 0.01$, all the algorithms have roughly the same space usage. In terms of update time (Figure 2b), again as above, LCDelta and DOUBLELC have by far the best performance, whether for small or large ϵ . For precision and recall (Figures 2c and 2d), all the algorithms have perfect precision and recall⁴.

Overall, we see that DOUBLEFREQUENT has the best performance in terms of space, although its updates are slow compared to LCDelta and DOUBLELC. Surprisingly, for our planted dataset, DOUBLELC

⁴Except for one datapoint at $\epsilon = 10^{-5}$ in Figure 2c which we believe is due to the probabilistic nature of DOUBLELC

does not offer any advantages over LCDELTA. The algorithms using sketching, SKETCHFREQUENT and SKETCHLC, are expensive in terms of both time and space.

4.3 Zipfian Data

Next, we present results for the Zipfian data. Here, other than ϵ and ϕ , we have a third parameter, the skew, that controls how heavy the tail of the distribution is. The higher the skew, the smaller the frequency tail.

Figure 3 compares the algorithms when the skew is varied while ϵ and ϕ are held fixed. In terms of space (Figure 3a), DOUBLEFREQUENT needs the minimum storage, while FREQUENT uses about 50% more and DOUBLELC and LCDELTA need about 5-6 times more. In terms of running time (Figure 3b), LCDELTA performs the best, followed (at some distance) by FREQUENT and DOUBLELC which have similar numbers. Figure 3c shows that recall is perfect for all the algorithms. As for precision, interestingly, we see in Figure 3d, LCDELTA does not do so well while DOUBLELC has perfect precision throughout.

In Figure 4, the skew is held fixed at 2, ϵ fixed at 0.001 while ϕ is varied. The story is similar to what we have seen so far in terms of space (winner: DOUBLEFREQUENT), time (winner: LCDELTA) and recall (winner: all). For precision though, we see that LCDELTA has low precision for small values of ϕ while SKETCHFREQUENT, SKETCHLC, and DOUBLELC has perfect precision throughout.

The picture is clarified in Figure 5 where we vary over a large range of ϵ while the skew and ϕ are held fixed. DOUBLEFREQUENT has the best space complexity while DOUBLELC and LCDELTA are similar in that regard (Figure 5a). With respect to time complexity, Figure 5b shows that LCDELTA is the most efficient while FREQUENT and DOUBLELC are similar. Recall is almost perfect for all the algorithms (Figure 5c). When it comes to precision (Figure 5d), DOUBLELC has precision 1 throughout while LCDELTA, DOUBLEFREQUENT and FREQUENT suffer from lower precision when ϵ is not too small.

Overall, our experiments suggest that on frequency distributions that may be bounded by a Zipfian and in our parameter regime where $n \gg m$ and $\phi \gg \epsilon$, DOUBLELC has good recall and precision, is similar to FREQUENT with respect to update time, and is similar to LCDELTA with respect to space usage. LCDELTA is very fast but it can have low precision, as also observed in [8].

5 CONCLUSION

We have proposed some simple algorithms for the heavy hitters problem that combine traditional counter-based approaches with hashing tricks used by sketching algorithms. Theoretically, they have excellent guarantees, in terms of space complexity and strong error bounds with respect to the residual frequency tail. Experimentally also, they perform favorably in terms of space complexity, precision and recall, though their update speed is worse than for existing solutions.

Thus, contrary to popular belief, there is still room for new algorithmic approaches to this classical problem. Our future work will target the following questions:

- In [8], it was shown that the Space Saving algorithm of Metwally et al. [15] has the best overall performance among counter algorithms with respect to space usage, update time,

recall, precision and average error. We wonder whether a DOUBLESPPACING algorithm (similar to DOUBLEFREQUENT) can be better in practice than the algorithms proposed here.

- We plan to implement the algorithms proposed here in the Yahoo! Sketches library [3] API. The codebase has highly optimized routines for many of the counter-based algorithms. We also aim to do more extensive evaluations on large, realistic datasets arising from click streams and computer networks, to support the results we obtained on the synthetic datasets.

REFERENCES

- [1] 2005. MassDAL Public Code Bank. <https://www.cs.rutgers.edu/~muthu/massdal-code-index.html>. Accessed: 2019-02-03.
- [2] 2014. Finding Frequent Items. <http://hadjieleftheriou.com/frequent-items>. Accessed: 2019-02-03.
- [3] 2019. Sketches Library from YAHOO! <https://datasketches.github.io>. Accessed: 2019-02-03.
- [4] Radu Berinde, Piotr Indyk, Graham Cormode, and Martin J Strauss. 2010. Space-optimal heavy hitters with strong error bounds. *ACM Transactions on Database Systems (TODS)* 35, 4 (2010), 26.
- [5] Arnab Bhattacharyya, Palash Dey, and David P Woodruff. 2018. An Optimal Algorithm for ℓ_1 -Heavy Hitters in Insertion Streams and Related Problems. *ACM Transactions on Algorithms (TALG)* 15, 1 (2018), 2.
- [6] Lee Breslau, Pei Cao, Li Fan, Graham Phillips, and Scott Shenker. 1999. Web caching and Zipf-like distributions: Evidence and implications. In *INFOCOM'99. Eighteenth Annual Joint Conference of the IEEE Computer and Communications Societies. Proceedings. IEEE*, Vol. 1. IEEE, 126–134.
- [7] Moses Charikar, Kevin Chen, and Martin Farach-Colton. 2002. Finding frequent items in data streams. In *International Colloquium on Automata, Languages, and Programming*. Springer, 693–703.
- [8] Graham Cormode and Marios Hadjieleftheriou. 2009. Finding the frequent items in streams of data. *Commun. ACM* 52, 10 (2009), 97–105.
- [9] Graham Cormode and Shan Muthukrishnan. 2005. An improved data stream summary: the count-min sketch and its applications. *Journal of Algorithms* 55, 1 (2005), 58–75.
- [10] Chuck Cranor, Theodore Johnson, Oliver Spatschek, and Vladislav Shkapenyuk. 2003. Gigascope: a stream database for network applications. In *Proceedings of the 2003 ACM SIGMOD international conference on Management of data*. ACM, 647–651.
- [11] Mark E Crovella, Murad S Taqqu, and Azer Bestavros. 1999. Heavy-tailed probability distributions in the World Wide Web. In *A practical guide to heavy tails*. Vol. 1. Birkhauser, 3–26.
- [12] Erik D Demaine, Alejandro López-Ortiz, and J Ian Munro. 2002. Frequency estimation of internet packet streams with limited space. In *European Symposium on Algorithms*. Springer, 348–360.
- [13] Richard M Karp, Scott Shenker, and Christos H Papadimitriou. 2003. A simple algorithm for finding frequent elements in streams and bags. *ACM Transactions on Database Systems (TODS)* 28, 1 (2003), 51–55.
- [14] Gurmeet Singh Manku and Rajeev Motwani. 2002. Approximate frequency counts over data streams. In *Vldb'02: Proceedings of the 28th International Conference on Very Large Databases*. Elsevier, 346–357.
- [15] Ahmed Metwally, Divyakant Agrawal, and Amr El Abbadi. 2005. Efficient computation of frequent and top-k elements in data streams. In *International Conference on Database Theory*. Springer, 398–412.
- [16] J MISRA. 1982. FINDING REPEATED ELEMENTS. *Science of Computer Programming* 2 (1982), 143–152.
- [17] Shanmugavelayutham Muthukrishnan et al. 2005. Data streams: Algorithms and applications. *Foundations and Trends® in Theoretical Computer Science* 1, 2 (2005), 117–236.
- [18] Sidney Redner. 1998. How popular is your paper? An empirical study of the citation distribution. *The European Physical Journal B-Condensed Matter and Complex Systems* 4, 2 (1998), 131–134.
- [19] K To, Tao Ye, and S Bhattacharyya. 2004. CMON: A general purpose continuous ip backbone traffic analysis platform. *Research Report RR04-ATL-110309, Sprint ATL* (2004).
- [20] George Kingsley Zipf. 1949. *Human behavior and the principle of least effort: An introduction to human ecology*. Addison Wesley.

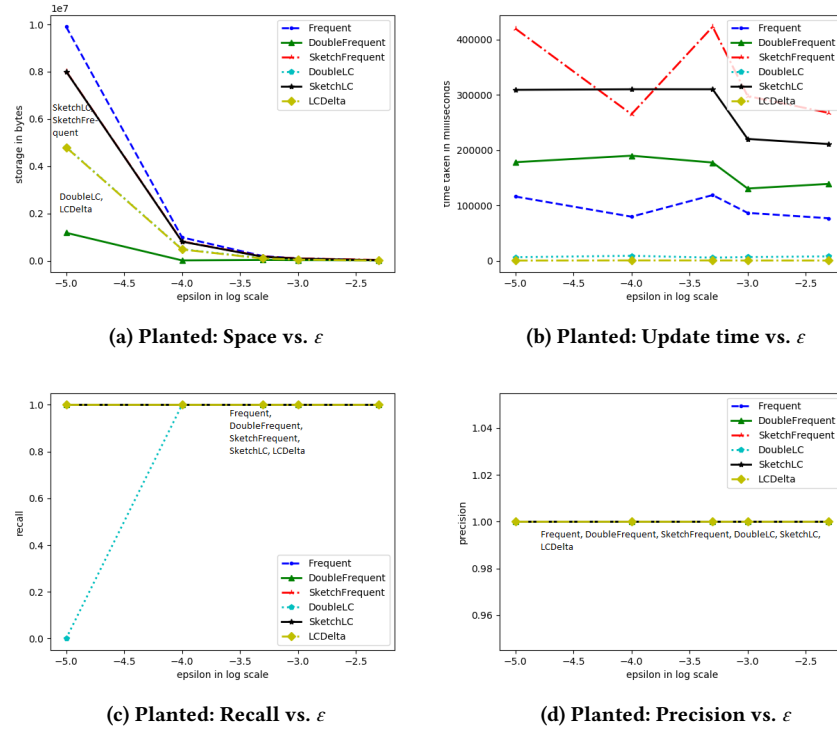


Figure 2: Comparison between the algorithms when run on data generated from the planted model. Here, φ is fixed to 10^{-2} .

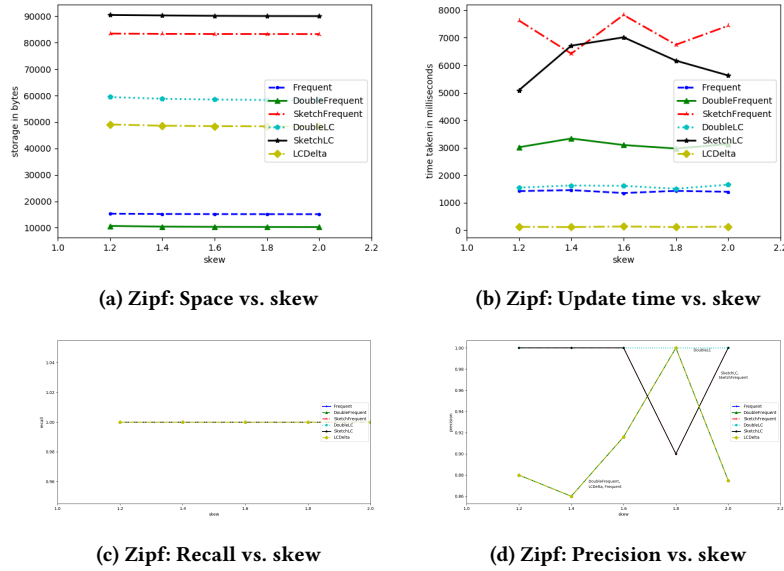


Figure 3: Comparison between the algorithms when run on Zipfian data. Skew varies as $\epsilon = 0.001$ and $\varphi = 0.01$.

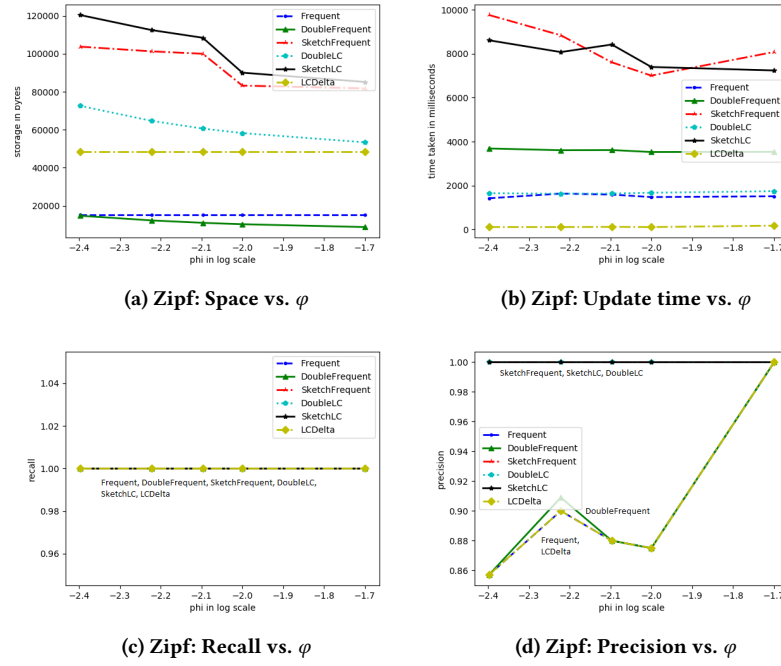


Figure 4: Comparison between the algorithms when run on Zipfian data. ϕ varies as $\epsilon = 0.001$ and skew = 2.

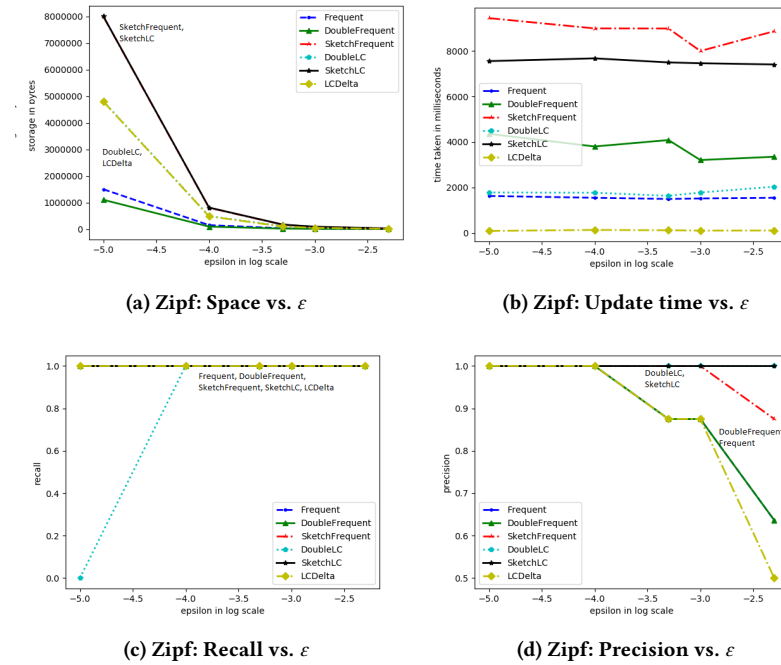


Figure 5: Comparison between the algorithms when run on Zipfian data. ϵ varies as $\phi = 0.01$ and skew = 2.