Bottom-k and Priority Sampling, Set Similarity and Subset Sums with Minimal Independence*

Mikkel Thorup
AT&T Labs—Research and University of Copenhagen
mikkel2thorup@gmail.com

June 12, 2013

Abstract

We consider bottom-k sampling for a set X, picking a sample $S_k(X)$ consisting of the k elements that are smallest according to a given hash function k. With this sample we can estimate the frequency f = |Y|/|X| of any subset Y as $|S_k(X) \cap Y|/k$. A standard application is the estimation of the Jaccard similarity $f = |A \cap B|/|A \cup B|$ between sets A and B. Given the bottom-k samples from A and B, we construct the bottom-k sample of their union as $S_k(A \cup B) = S_k(S_k(A) \cup S_k(B))$, and then the similarity is estimated as $|S_k(A \cup B) \cap S_k(A) \cap S_k(B)|/k$.

We show here that even if the hash function is only 2-independent, the expected relative error is $O(1/\sqrt{fk})$. For $fk = \Omega(1)$ this is within a constant factor of the expected relative error with truly random hashing.

For comparison, consider the classic approach of repeated min-wise hashing, where we use k independent hash functions $h_1, ..., h_k$, storing the smallest element with each hash function. For min-wise hashing, there can be a constant bias with constant independence, and this is not reduced with more repetitions k. Recently Feigenblat et al. showed that bottom-k circumvents the bias if the hash function is 8-independent and k is sufficiently large. We get down to 2-independence for any k. Our result is based on a simple union bound, transferring generic concentration bounds for the hashing scheme to the bottom-k sample, e.g., getting stronger probability error bounds with higher independence.

For weighted sets, we consider priority sampling which adapts efficiently to the concrete input weights, e.g., benefiting strongly from heavy-tailed input. This time, the analysis is much more involved, but again we show that generic concentration bounds can be applied.

1 Introduction

The concept of min-wise hashing (or the "MinHash algorithm" according to ¹) is a basic algorithmic tool suggested by Broder et al. [6, 8] for problems related to set similarity and containment. After the initial application of this algorithm in the early Altavista search engine to detecting and clustering similar documents, the scheme has reappeared in numerous other applications¹ and is now a standard tool in data mining where it is used for estimating similarity [6, 8, 9], rarity [13], document duplicate detection [7, 21, 23, 38], etc [2, 4, 10, 31].

^{*}A short preliminary version of this paper was presented at STOC'13 [35].

¹See http://en.wikipedia.org/wiki/MinHash

In an abstract mathematical view, we have two sets, A and B, and we are interested in understanding their overlap in the sense of the Jaccard similarity $f=\frac{|A\cap B|}{|A\cup B|}$. In order to do this by sampling, we need sampling correlated between the two sets, so we sample by hashing. Consider a hash function $h:A\cup B\to [0,1]$. For simplicity we assume that h is fully random, and has enough precision that no collisions are expected. The main mathematical observation is that $\Pr[\operatorname{argmin} h(A) = \operatorname{argmin} h(B)]$ is precisely $f = |A\cap B| / |A\cup B|$. Thus, we may sample the element with the minimal hash from each set, and use them in $[\operatorname{argmin} h(A) = \operatorname{argmin} h(B)]$ for an unbiased estimate of f. Here, for a logical statement S, [S] = 1 if S is true; otherwise [S] = 0.

For more concentrated estimators, we use repetition with k independent hash functions, $h_1, ..., h_k$. For each set A, we store $M^k(A) = (\operatorname{argmin} h_1(A), ..., \operatorname{argmin} h_k(A))$, which is a sample with replacement from A. The Jaccard similarity between sets A and B is now estimated as $|M^k(A) \cap M^k(B)|/k$ where $|M^k(A) \cap M^k(B)|$ denotes the number of agreeing coordinates between $M^k(A)$ and $M^k(B)$. We shall refer to this approach as repeated min-wise or $k \times min$.

For our discussion, we consider the very related application where we wish to store a sample of a set X that we can use to estimate the frequency $f = \frac{|Y|}{|X|}$ of any subset $Y \subseteq X$. The idea is that the subset Y is not known when the sample from X is made. The subset Y is revealed later in the form of a characteristic function that can tell if (sampled) elements belong to Y. Using the $k \times \min$ sample $M^k(X)$, we estimate the frequency as $|M^k(X) \cap Y|/k$ where $|M^k(X) \cap Y|$ denotes the number of samples from $M^k(X)$ in Y.

Another classic approach for frequency estimation is to use just one hash function h and use the k elements from X with the smallest hashes as a sample $S_k(X)$. This is a sample without replacement from X. As in [12], we refer to this as a bottom-k sample. The method goes back at least to [20]. The frequency of Y in X is estimated as $|Y \cap S_k(X)|/k$. Even though surprisingly fast methods have been proposed to compute $k \times \min[3]$, the bottom-k signature is much simpler and faster to compute. In a single pass through a set, we only apply a single hash function k to each element, and use a max-priority queue to maintain the k smallest elements with respect to k.

It is standard¹ to use bottom-k samples to estimate the Jaccard similarity between sets A and B, for this is exactly the frequency of the intersection in the union. First we construct the bottom-k sample $S_k(A \cup B) = S_k(S_k(A) \cup S_k(B))$ of the union by picking the k elements from $S_k(A) \cup S_k(B)$ with the smallest hashes. Next we return $|S_k(A) \cap S_k(B) \cap S_k(A \cup B)|/k$.

Stepping back, for subset frequency, we generally assume that we can identify samples from the subset. In the application to set similarity, it important that the samples are coordinated via hash functions, for this is what allows us to identify samples from the intersection as being sampled in both sets. In our mathematical analysis we will focus on the simpler case of subset frequency estimation, but it the application to set similarity that motivates our special interest in sampling via hash functions.

Limited independence The two approaches $k \times \min$ and bottom-k are similar in spirit, starting from the same base $1 \times \min$ bottom-1. With truly random hash functions, they have essentially the same relative standard deviation (standard deviation divided by expectation) bounded by $1/\sqrt{fk}$ where f is the set similarity or subset frequency. The two approaches are, however, very different from the perspective of pseudo-random hash functions of limited independence [37]: a random hash function h is d-independent if the hash values of any d given elements are totally random.

With min-wise hashing, we have a problem with bias in the sense of sets in which some elements have a better than average chance of getting the smallest hash value. It is known that 1 + o(1) bias requires $\omega(1)$ -independence [28]. This bias is not reduced by repetitions as in $k \times \min$. However, recently Porat et al. [19] proved that the bias for bottom-k vanishes for large enough $k \gg 1$ if we use 8-independent hashing.

Essentially they get an expected relative error of $O(1/\sqrt{fk})$, and error includes bias. For $fk = \Omega(1)$, this is only a constant factor worse than with truly random hashing. Their results are cast in a new framework of "d-k-min-wise hashing", and the translation to our context is not immediate.

Results In this paper, we prove that bottom-k sampling preserves the expected relative error of $O(1/\sqrt{fk})$ with 2-independent hashing, and this holds for any k including k=1. We note that when fk=o(1), then $1/\sqrt{fk}=\omega(1)$, so our result does not contradict a possible large bias for k=1.

We remark that we also get an $O(1/\sqrt{(1-f)k})$ bound on the expected relative error. This is important if we estimate the dissimilarity 1-f of sets with large similarity f=1-o(1).

For the more general case of weighted sets, we consider priority sampling [18] which adapts near-optimally to the concrete input weights [33], e.g., benefiting strongly from heavy-tailed input. We show that 2-independent hashing suffices for good concentration.

Our positive finding with 2-independence contrasts recent negative results on the insufficiency of low independence, e.g., that linear probing needs the 5-independence [28] that was proved sufficient by Pagh et al. [27].

Implementation For 2-independent hashing we can use the fast multiplication-shift scheme from [14], e.g., if the elements are 32-bit keys, we pick two random 64-bit numbers a and b. The hash of key x is computed with the C-code (a*x+b) >> 32, where * is 64-bit multiplication which as usual discards overflow, and >> is a right shift. This is 10-20 times faster than the fastest known 8-independent hashing based on a degree 7 polynomial tuned for a Mersenne prime field [36]².

Practical relevance We note that Mitzenmacher and Vadhan [24] have proved that 2-independence generally works if the input has enough entropy. However, the real world has lots of low entropy data. In [36] it was noted how consecutive numbers with zero entropy made linear probing with 2-independent hashing extremely unreliable. This was a problem in connection with denial-of-service attacks using consecutive IP-addresses. For our set similarity, we would have similar issues in scenarios where small numbers are more common, hence where set intersections are likely to be fairly dense intervals of small numbers whereas the difference is more likely to consists of large random outliers. Figure 1 presents an experiment showing what happens if we try to estimating such dissimilarity with 2-independent hashing.

Stepping back, the result Mitzenmacher and Vadhan is that 2-independence works for sufficiently random input. In particular, we do not expect problems to show up in random tests. However, this does not imply that 2-independent hashing can be trusted on real data unless we have specific reasons to believe that the input has high entropy. In Figure 1, bottom-k performs beautifully with 2-independent hashing, but no amount of experiments can demonstrate general reliability. However, the mathematical result of this paper is that bottom-k can indeed be trusted with 2-independent hashing: the expected relative error is $O(1/\sqrt{fk})$ no matter the structure of the input.

Techniques To appreciate our analysis, let us first consider the trivial case where we are given a non-random threshold probability p and sample all elements that hash below p. As in [17] we refer to this as threshold sampling. Since the hash of a element x is uniform in [0,1], this samples x with probability p. The sampling of x depends only on the hash value of x, so if, say, the hash function is d-independent, then the

²See Table 2 in [36] for comparisons with different key lengths and computers between multiplication-shift (TwoIndep), and tuned polynomial hashing (CWtrick). The table considers polynomials of degree 3 and 4, but the cost is linear in the degree, so the cost for degree 7 is easily extrapolated.

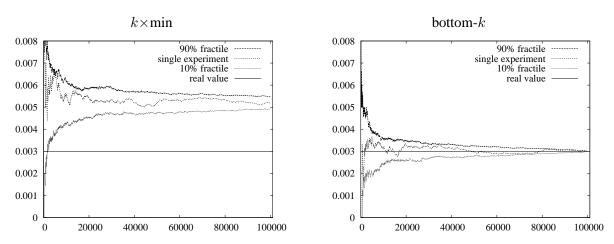


Figure 1: Experiment with set consisting of 100300 32-bit keys. It has a "core" consisting of the consecutive numbers 1, ... 100000. In addition it has 300 random "outliers". Using k samples from the whole set, we want to estimate the frequency of the outliers. The true frequency is $\frac{300}{100300} \approx 0.003$. We used k=1, ..., 100000 in $k \times \min$ and bottom-k and made one hundred experiments. For each k, we sorted the estimates, plotting the 10th and 90th value, labeled as 10% and 90% fractile in the figures. We also plotted the results from a single experiment. For readability, only one in every 100 values of k is plotted. Both schemes converge, but due to bias, $k \times \min$ converges to a value that is 70% too large. Since bottom-k does sampling without replacement, it becomes exact when the number of samples is the size of the whole set. The bias is a function of the structure of the subset within the whole set, e.g., the core set must have a negative bias complimenting the positive bias of the outliers. It is therefore not possible to correct for the bias if one only has the sample available.

number of samples is the sum of d-independent 0-1 variables. This scenario is very well understood (see, e.g., [15, 32]).

We could set p = k/n, and get an expected number of k samples. Morally, this should be similar to a bottom-k sample, which is what we get if we end up with exactly k samples, that is, if we end up with $h_{(k)} where <math>h_{(i)}$ denotes the ith smallest hash value. What complicates the situation is that $h_{(k)}$ and $h_{(k+1)}$ are random variables depending on all the random hash values.

An issue with threshold sampling is that the number of samples is variable. This is an issue if we have bounded capacity to store the samples. With k expected samples, we could put some limit $K \gg k$ on the number of samples, but any such limit introduces dependencies that have to be understood. Also, if we have room for K samples, then it would seem wasteful not fill it with a full bottom-K sample.

Our analysis of bottom-k samples is much simpler than the one in [19] for 8-independent hashing with $k \gg 1$. With a union bound we reduce the analysis of bottom-k samples to the trivial case of threshold sampling. Essentially we only get a constant loss in the error probabilities. With 2-independent hashing, we then apply Chebyshev's inequality to show that the expected relative error is $O(1/\sqrt{fk})$. The error probability bounds are immediately improved if we use hash functions with higher independence.

It is already known from [5] that we can use a 2-independent bottom-k sample of a set to estimate its size n with an expected error of $O(\sqrt{n})$. The estimate is simply the inverse of the kth smallest sample. Applying this to two $\Theta(n)$ -sized sets and their union, we can estimate $|A|, |B|, |A \cup B|$ and $|A \cap B| = |A| + |B| - |A \cup B|$ each with an expected error of $O(\sqrt{n})$. However, $|A \cap B|$ may be much smaller than $O(\sqrt{n})$. If we instead multiply our estimate of the similarity $f = |A \cap B|/|A \cup B|$ with the estimate of $|A \cup B|$, the resulting estimate of $|A \cap B|$ is

$$(1 \pm O(1/\sqrt{fk})f(|A \cup B| \pm O(\sqrt{n})) = |A \cap B| \pm O(\sqrt{|A \cap B|}).$$

The analysis of priority sampling for weighted sets is much more delicate, but again, using union bounds, we show that generic concentration bounds apply.

2 Bottom-k samples

We are given a set of X of n elements. A hash function maps the elements uniformly and collision free into [0,1]. Our bottom-k sample S consists of the k elements with the lowest hash values. The sample is used to estimate the frequency f=|Y|/|X| of any subset Y of X as $|Y\cap S|/k$. With 2-independent hashing, we will prove the following error probability bound for any $r \leq \bar{r} = \sqrt{k}/3$:

$$\Pr\left[||Y \cap S| - fk| > r\sqrt{fk}\right] \le 4/r^2. \tag{1}$$

The result is obtained via a simple union bound where stronger hash functions yield better error probabilities. With d-independence with d an even constant, the probability bound is $O(1/r^d)$.

It is instructive to compare d-independence with the idea of storing d independent bottom-k samples, each based on 2-independence, and use the median estimate. Generally, if the probability of a certain deviation is p, the deviation probability for the median is bounded by $(2ep)^{d/2}$, so the $4/r^2$ from (1) becomes $(2e4/r^2)^{d/2} < (5/r)^d$, which is the same type of probability that we get with a single d-independent hash function. The big advantage of a single d-independent hash function is that we only have to store a single bottom-k sample.

If we are willing to use much more space for the hash function, then we can use twisted tabulation hashing [29] which is very fast, and then we get exponential decay in r though only down to an arbitrary polynomial of the space used.

In order to show that the expected relative error is $O(1/\sqrt{fk})$, we also prove the following bound for $fk \le 1/4$:

$$\Pr[|Y \cap S| \ge \ell] = O(fk/\ell^2 + \sqrt{f}/\ell). \tag{2}$$

From (1) and (2) we get

Proposition 1 For bottom-k samples based on 2-independent hashing, a fraction f subset is estimated with an expected relative error of $O(1/\sqrt{fk})$.

Proof The proof assumes (1) and (2). For the case fk > 1/4, we will apply (1). The statement is equivalent to saying that the sample error $||Y \cap S| - fk|$ in expectation is bounded by $O(\sqrt{fk})$. This follows immediately from (1) for errors below $\bar{r}\sqrt{fk} = k\sqrt{|Y|/n}$. However, by (1), the probability of a larger error is bounded by $4/\bar{r}^2 = O(1/k)$. The maximal error is k, so the contribution of larger errors to the expected error is O(1). This is $O(\sqrt{fk})$ since fk > 1/4.

We will now handle the case $fk \le 1/4$ using (2). We want to show that the expected absolute error is $O(\sqrt{fk})$. We note that only positive errors can be bigger than fk, so if the expected error is above 2fk, the expected number of samples from Y is proportional to the expected error. We have $\sqrt{fk} \ge 2fk$, so for the expected error bound, it suffices to prove that the expected number of samples is $|Y \cap S| = O(\sqrt{fk})$. Using (2) for the probabilities, we now sum the contributions over exponentially increasing sample sizes.

$$\begin{split} \mathsf{E}[|Y\cap S|] &\leq \sum_{i=0}^{\lfloor \lg k \rfloor} \left(2^{i+1}\Pr[|Y\cap S| \geq 2^i]\right) \\ &= \sum_{i=0}^{\lfloor \lg k \rfloor} O\left(2^i (fk/2^{2i} + \sqrt{f}/2^i)\right) \\ &= O\left(fk + \sqrt{f}(1 + \lg k)\right) = O\left(\sqrt{fk}\right). \end{split}$$

2.1 A union upper bound

First we consider overestimates. For positive parameters a and b to be chosen, we will bound the probability of the overestimate

$$|Y \cap S| > \frac{1+b}{1-a}fk. \tag{3}$$

Define the threshold probability

$$p = \frac{k}{n(1-a)}.$$

Note that p is defined deterministically, independent of any samples. It is easy to see that the overestimate (3) implies one of the following two threshold sampling events:

- (A) The number of elements from X that hash below p is less than k. We expected pn = k/(1-a) elements, so k is a factor (1-a) below the expectation.
- **(B)** Y gets more than (1+b)p|Y| hashes below p, that is, a factor (1+b) above the expectation.

To see this, assume that both (A) and (B) are false. When (A) is false, we have k hashes from X below p, so the largest hash in S is below p. Now if (B) is also false, we have at most $(1+b)p|Y| = (1+b)/(1-a) \cdot fk$ elements from Y hashing below p, and only these elements from Y could be in S. This contradicts (3). By the union bound, we have proved

Proposition 2 The probability of the overestimate (3) is bounded by $P_A + P_B$ where P_A and P_B are the probabilities of the events (A) and (B), respectively.

Upper bound with 2-independence Addressing events like (A) and (B), let m be the number of elements in the set Z considered, e.g., Z=X or Z=Y. We study the number of elements hashing below a given threshold $p \in [0,1]$. Assuming that the hash values are uniform in [0,1], the mean is $\mu=mp$. Assuming 2-independence of the hash values, the variance is $mp(1-p)=(1-p)\mu$ and the standard deviation is $\sigma=\sqrt{(1-p)\mu}$. By Chebyshev's inequality, we know that the probability of a deviation by $r\sigma$ is bounded by $1/r^2$. Below we will only use that the relative standard deviation σ bounded by $1/\sqrt{\mu}$.

For any given $r \leq \sqrt{k}/3$, we will fix a and b to give a combined error probability of $2/r^2$. More precisely, we will fix $a = r/\sqrt{k}$ and $b = r/\sqrt{fk}$. This also fixes p = k/(n(1-a)). We note for later that $a \leq 1/3$ and $a \leq b$. This implies

$$(1+b)/(1-a) \le (1+3b) = 1 + 3r/\sqrt{fk}.$$
(4)

In connection with (A) we study the number of elements from X hashing below p. The mean is $pn \ge k$ so the relative standard deviation is less than $1/\sqrt{k}$. It follows that a relative error of $a = r/\sqrt{k}$ corresponds to at least r standard deviations, so

$$P_A = \Pr\left[\#\{x \in X | h(x) < p\} < (1-a)np\right] < 1/r^2.$$

In connection with (B) we study the number of elements from Y hashing below p. Let m=|Y|. The mean is pm=km/(n(1-a)) and the relative standard deviation less than $1/\sqrt{pm}<1/\sqrt{km/n}$. It follows than a relative error of $b=r/\sqrt{km/n}$ is more than r standard deviations, so

$$P_B = \Pr\left[\#\{y \in Y | h(y) < p\} > (1+b)mp\right] < 1/r^2.$$

By Proposition 2 we conclude that the probability of (3) is bounded by $2/r^2$. Rewriting (3) with (4), we conclude that

$$\Pr\left[|Y \cap S| > fk + 3r\sqrt{fk}\right] \le 2/r^2. \tag{5}$$

This bounds the probability of the positive error in (1). The above constants 3 and 2 are moderate, and they can easily be improved if we look at asymptotics. Suppose we want good estimates for subsets Y of frequency at least f_{\min} , that is, $|Y| \ge f_{\min}|X|$. This time, we set $a = r/\sqrt{fk}$, and then we get $P_A \le f/r^2$. We also set $b = r/\sqrt{fk}$ preserving $P_B \le 1/r^2$. Now for any $Y \subseteq X$ with |Y| > fn, we have

$$\Pr\left[|Y \cap S| > (1+\varepsilon)fk\right] = (1+f)/r^2$$
where $\varepsilon = \frac{1+r/\sqrt{fk}}{1-r/\sqrt{k}} - 1 = \frac{r/\sqrt{k} + r/\sqrt{fk}}{1-r/\sqrt{k}}.$
(6)

With f=o(1) and $k=\omega(1)$, the error is $\varepsilon=(1+o(1))r/\sqrt{fk}$, and the error probability is $P_{\varepsilon}=(1+f)/r^2=(1+o(1))/r^2$. Conversely, this means that if we for subsets of frequency f and a relative positive error ε want an error probability around P_{ε} , then we set $r=\sqrt{1/P_{\varepsilon}}$ and $k=r^2/(f\varepsilon^2)=1/(fP_{\varepsilon}\varepsilon^2)$.

2.2 A union lower bound

We have symmetric bounds for underestimates:

$$|Y \cap S| < \frac{1 - b'}{1 + a'} fk. \tag{7}$$

This time we define the threshold probability $p' = \frac{k}{n(1+a')}$. It is easy to see that the overestimate (3) implies one of the following two events:

- (A') The number of elements from X below p' is at least k. We expected p'n = k/(1+a') elements, so k is a factor (1+a') above the expectation.
- (B') Y gets less than (1-b')p|Y| hashes below p', that is, a factor (1-b') below the expectation.

To see this, assume that both (A') and (B') are false. When (A') is false, we have less than k hashes from X below p', so S must contain all hashes below p'. Now if (B) is also false, we have at least $(1-b)p'|Y| = (1-b)/(1+a) \cdot fk$ elements from $Y \subseteq X$ hashing below p', hence which must be in S. This contradicts (7). By the union bound, we have proved

Proposition 3 The probability of the underestimate (7) is bounded by $P_{A'} + P_{B'}$ where $P_{A'}$ and $P_{B'}$ are the probabilities of the events (A') and (B'), respectively.

Lower bound with 2-independence Using Proposition 3 we will bound the probability of underestimates, complementing our previous probability bounds for overestimates from Section 2.1. We will provide bounds for the same overall relative error as we did for the overestimates; namely

$$\varepsilon = \frac{1+b}{1-a} - 1 = (a+b)/(1-a)$$

However, for the events (A') and (B') we are going to scale up the relative errors by a factor (1 + a), that is, we will use a' = a(1 + a) and b' = b(1 + a). The overall relative negative error from (7) is then

$$\varepsilon' = 1 - \frac{1 - b'}{1 + a'} = (a' + b')/(1 + a')$$

$$< (1 + a)(a + b)/(1 + a') < (a + b) < \varepsilon.$$

Even with this smaller error, we will get better probability bounds than those we obtained for the overestimates. For (A) we used $1/\sqrt{k}$ as an upper bound on the relative standard deviation, so a relative error of a was counted as $s_A = a\sqrt{k}$ standard deviations. In (A') we have mean $\mu' = np' = k/(1+a')$, so the relative standard deviation is bounded by $1/\sqrt{k/(1+a')} = \sqrt{1+a+a^2}/\sqrt{k}$. This means that for (A'), we can count a relative error of a' = a(1+a) as

$$s'_A = a(1+a)\sqrt{k}/\sqrt{1+a+a^2}$$

= $s_A(1+a)/\sqrt{1+a+a^2} > s_A$

standard deviations. In Section 2.1 we bounded P_A by $1/s_A^2$, and now we can bound $P_{A'}$ by $1/s_A^2 \le 1/s_A^2$. The scaling has the same positive effect on our probability bounds for (B'). That is, in Section 2.1, a relative

error of b was counted as $s_B = b\sqrt{fk}$ standard deviations. With (B') our relative error of b' = b(1+a) is counted as

$$s'_B = b(1+a)\sqrt{fk}/\sqrt{1+a+a^2}$$

= $s_B(1+a)/\sqrt{1+a+a^2} > s_B$

standard deviations, and then we can bound $P_{B'}$ by $1/{s'_B}^2 \le 1/{s_B}^2$. Summing up, our negative relative error ε' is smaller than our previous positive error ε , and our overall negative error probability bound $1/{s'_A}^2 + 1/{s'_B}^2$ is smaller than our previous positive error probability bound $1/{s_A}^2 + 1/{s_B}^2$. We therefore translate (5) to

$$\Pr\left[|Y \cap S| < fk - 3r\sqrt{fk}\right] \le 2/r^2. \tag{8}$$

which together with (5) establishes (1). Likewise (6) translates to

$$\Pr\left[||Y \cap S| - fk| > \varepsilon fk\right] \le 2(1+f)/r^2$$
where $\varepsilon = \frac{1 + r/\sqrt{fk}}{1 - r/\sqrt{k}} - 1$. (9)

As for the positive error bounds we note that with f = o(1) and $k = \omega(1)$, the error is $\varepsilon = (1 + o(1))r/\sqrt{fk}$ and the error probability is $P_{\varepsilon} = (2 + o(1))/r^2$. Conversely, this means that if we for a target relative error ε want an error probability around P_{ε} , then we set $r = \sqrt{2/P_{\varepsilon}}$ and $k = r^2/(f\varepsilon^2) = 2/(fP_{\varepsilon}\varepsilon^2)$.

2.3 Rare subsets

We now consider the case where the expected number fk of samples from Y is less than 1/4. We wish to prove (2)

$$\Pr[|Y \cap S| \ge \ell] = O(fk/\ell^2 + \sqrt{f}/\ell).$$

For some balancing parameter $c \geq 2$, we use the threshold probability p = ck/n. The error event (A) is that less than k elements from X sample below p. The error event (B) is that at least ℓ elements hash below p. As in Proposition 2, we observe that ℓ bottom-k samples from Y implies (A) or (B), hence that $\Pr[|Y \cap S| \geq \ell] \leq P_A + P_B$.

The expected number of elements from X that hash below p is ck. The error event (A) is that we get less than k, which is less than half the expectation. This amounts to at least $\sqrt{ck}/2$ standard deviations, so by Chebyshev's inequality, the probability of (A) is $P_A \leq 1/(\sqrt{ck}/2)^2 = 4/(ck)$.

The event (B) is that at least ℓ elements from Y hash below p, while the expectation is only fck. Assuming that $\ell \geq 2fck$, the error is by at least $(\ell/2)/\sqrt{fck}$ standard deviations. By Chebyshev's inequality, the probability of (B) is $P_B \leq 1/((\ell/2)/\sqrt{fck})^2 = 4fck/\ell^2$. Thus

$$P_A + P_B \le 4/(ck) + 4fck/\ell^2.$$

We wish to pick c for balance, that is,

$$4/ck = 4fck/\ell^2 \iff c = \ell/(\sqrt{f}k)$$

However, we have assumed that $c \geq 2$ and that $\ell \geq 2fck$. The latter is satisfied because $2fck = 2fk\ell/(\sqrt{f}k) = 2\sqrt{f}\ell$ and $f \leq 1/4$. Assuming that $c = \ell/(\sqrt{f}k) \geq 2$, we get

$$P_A + P_B \le 8/(k(\ell/(\sqrt{f}k))) = 8\sqrt{f}/\ell.$$

When $\ell/(\sqrt{f}k) < 2$, we set c = 2. Then

$$P_A + P_B \le 2/k + 8fk/\ell^2 \le 16fk/\ell^2$$
.

Again we need to verify that $\ell \ge fck = 2fk$, but that follows because $\ell \ge 1$ and $fk \le 1/4$. We know that at one of the above two cases applies, so we conclude that

$$P[|Y \cap S| \ge \ell] \le P_A + P_B = O(fk/\ell^2 + \sqrt{f}/\ell),$$

completing the proof of (2).

3 Priority sampling

We now consider the more general situation where we are dealing with a set I of weighted items with w_i denoting the weight of item $i \in I$. Let $\sum I = \sum_{i \in I} w_i$ denote the total weight of set I.

Now that we are dealing with weighted items, we will use *priority sampling* [18] which generalizes the bottom-k samples we used for unweighted elements. Each item or element i is identified by a unique key which is hashed uniformly to a random number $h_i \in (0,1)$. The item is assigned a *priority* $q_i = w_i/h_i > w_i$. In practice, hash values may have some limited precision b, but we assume that b is large enough that the resulting rounding can be ignored. We assume that all priorities end up distinct and different from the weights. If not, we could break ties based on an ordering of the items. The priority sample S of size k contains the k samples of highest priority, but it also stores a threshold τ which is the (k+1)th highest priority. Based on this we assign a weight estimate \widehat{w}_i to each item i. If i is not sampled, $\widehat{w}_i = 0$; otherwise $\widehat{w}_i = \max\{w_i, \tau\}$. A basic result from [18] is that $E[\widehat{w}_i] = w_i$ if the hash function is truly random (in [18], the h_i were described as random numbers, but here they are hashes of the keys).

We note that priority sampling generalize the bottom-k sample we used for unweighted items, for if all weights are unit, then the k highest priorities correspond to the k smallest hash values. In fact, priority sampling predates [12], and [12] describes bottom-k samples for weighted items as a generalization of priority sampling, picking the first k items according to an arbitrary randomized function of the weights.

The original objective of priority sampling [18] was subset sum estimation. A subset $J \subseteq I$ of the items is selected, and we estimate the total weight in the subset as $\widehat{w}_J = \sum \{\widehat{w}_i | i \in J \cap S\}$. By linearity of expectation, this is an unbiased estimator. A cool application from [18] was that as soon as the signature of the Slammer worm [25] was identified, we could inspect the priority samples from the past to track its history and identify infected hosts. An important point is that the Slammer worm was not known when the samples were made. Samples are made with no knowledge about which subsets will later turn out to be of interest.

Trivially, if we want to estimate the relative subset weight $\sum J/\sum I$ and we do not know the exact total, we can divide \widehat{w}_J with the estimate \widehat{w}_I of the total. As with the bottom-k sampling for unweighted items, we can easily use priority sampling to estimate the similarity of sets of weighted items: given the priority sample from two sets, we construct the priority sample of their union, and estimate the intersection as a subset. This is where it is important that we use a hash function so that the sampling from different sets is coordinated, e.g., we could not use iterative sampling procedures like the one in [11]. In the case of histogram similarity, it is natural to allow the same item to have different weights in different sets. More specifically, allowing zero weights, every possible item has a weight in each set. For the similarity we take the sum of the minimum weight for each item, and divide it by the sum of the maximum weight for each item. This requires a special sampling that we shall return to in Section 3.9.

Priority sampling is not only extremely easy to implement on-line with a standard min-priority queue; it also has some powerful universal properties in its adaption to the concrete input weights. As conjectured in [18] and proved in [33], given one extra sample, priority sampling has smaller variance sum $\sum_i \text{Var}[\widehat{w}_i]$ than any off-line unbiased sampling scheme tailored for the concrete input weights. In particular, priority sampling benefits strongly if there are dominant weights w_i in the input, estimated precisely as $\widehat{w}_i = \max\{w_i, \tau\} = w_i$. In the important case of heavy tailed input distributions [1], we thus expect most of the total weight to be estimated without any error. The quality of a priority sample is therefore often much better than what can be described in terms of simple parameters such as total weight, number of items, etc. The experiments in [18] on real and synthetic data show how priority sampling often gains orders of magnitude in estimate quality over competing methods.

The quality of a priority estimate depends completely on the distribution of weights in the input, and often we would like to know how much we can trust a given estimate. What we really want from a sample is not just an estimate of a subset sum, but a confidence interval [34]: from the information in the sample, we want to derive lower and upper bounds that capture the true value with some desired probability. Some applications of such concervative bounds are given in [16].

What makes priority sampling tricky to analyze is that the priority threshold τ is a random variable depending on all the random priorities. It may be very likely that the threshold τ ends up smaller than some dominant weight w_i , but it could also be bigger, so we do have variance on all weight estimates \widehat{w}_i .

All current analysis of priority sampling [18, 33, 34] is heavily based on true randomness, assuming that the priorities are independent random variables, e.g., the unbiasedness proof from [18] that $E[\widehat{w}_i] = w_i$ starts by fixing the priorities q_j of all the other items $j \neq i$. However, in this paper, we want to use hash functions with independence as low as 2, and then any such analysis breaks down. In fact, bias may now be introduced. To see this, consider the following extreme case of 2-independent hashing of n keys: divide (0,1] into n subintervals $I_i = (i/n, (i+1)/n]$. With probability 1-1/n, the keys are all mapped to different random subintervals, and with probability 1/n, all keys are mapped to the same random subinterval. Within the subintervals, the hashing is totally random. This scheme is clearly 2-independent, but highly restricted for n > 2. As a simple example of bias, consider a priority sample of k = 2 out of n = 3 unit weight keys. A messy computer calculation shows that the expected weight estimates are 1.084.

Relation to threshold sampling Generalizing the pattern for unweighted sets, our basic goal is to relate the error probabilities with priority sampling to the much simpler case of threshold sampling for weighted items. In threshold sampling, we are not given a predefined sample size. Instead we are given a fixed threshold t. We use exactly the same random priorities as in priority sampling, but now an item is sampled if and only if $q_i > t$. The weight estimate is

$$\widehat{w}_i^t = \begin{cases} 0 & \text{if } q_i \le t \\ \max\{w_i, t\} & \text{if } q_i > t \end{cases}$$
 (10)

In statistics, threshold sampling is known as Poisson sampling with probability proportional to size [30]. The name threshold sampling is taken from [17].

The \widehat{w}_i^t notation from (10) is well-defined also when t is a variable, and if priority sampling leads to threshold τ , then the priority estimate for item i is $\widehat{w}_i = \widehat{w}_i^{\tau}$.

With a fixed threshold t, it is trivial to see that the estimates are unbiased, that is, $\mathsf{E}[\widehat{w}_i^t] = w_i$; for if $w_i \geq t$, we always have $\widehat{w}_i^t = w_i$, and if $w_i < t$, then

$$\mathsf{E}[\widehat{w}_i^t] = t \Pr[q_i > t] = t \Pr[h_i < w_i/t] = w_i.$$

The unbiasedness with fixed threshold t only requires that each h_i is uniform in (0,1). No independence is required. This contrasts the bias we may get with the variable priority threshold τ with limited dependence.

With threshold sampling, concentration bounds for subset sum estimates are easily derived. For a subset $J \subseteq I$, the threshold estimate $\sum_{i \in J} \widehat{w}_i$ is naturally divided in an exact part for large weights and a variable part for small weights:

$$\sum_{i \in J, w_i \ge t} \widehat{w}_i^t = \sum_{i \in J, w_i \ge t} w_i$$

$$\sum_{i \in J, w_i < t} \widehat{w}_i^t = t \sum_{i \in J, w_i < t} X_i, \text{ where } X_i = [h_i < w_i/t] \in \{0, 1\}.$$

$$(11)$$

Each X_i depends on h_i only, so if the h_i are d-independent, then so are the h_i . Let $X = \sum_{i \in J, w_i < t} X_i$ and $\mu = \mathsf{E}[X] = \sum_{i \in J, w_i < t} w_i/t$. As in the unweighted case, if the hash function is 2-independent, then $\mathsf{Var}[X] < \mu$, and by Chebyshev's inequality $\Pr[|X - \mu| \ge r\sqrt{\mu}] \le 1/r^2$.

Informally speaking, for bounded errors and modulo constant factors, our main result is that concentration bounds for threshold sampling apply to priority sampling as if the variable priority threshold was fixed. As in the unweighted case, the result is obtained by a union bound over threshold sampling events. In the unweighted case, we only needed to consider the four threshold sampling error events (A), (B), (A'), and (B'). However, now with weighted items, we are going to reduce a priority sampling error event to the union of an unbounded number of threshold sampling error events that happen with geometrically decreasing probabilities.

3.1 Notation and definitions

Before formally presenting our priority sampling results, we introduce some notation and definitions.

Fractional subsets and inner products It is both convenient and natural to generalize our estimates from regular subsets to *fractional subsets*, where for each $i \in I$, there is a fraction $f_i \in [0, 1]$ specifying that item i contributes $f_i w_i$ to the weight of fractional subset. A regular subset corresponds to the special case where $f_i \in \{0, 1\}$.

We are now interested in inner products between the fraction vector $f=(f_i)_{i\in I}$ and the vectors of weights or weight estimates. Our goal is to estimate $fw=\sum_{i\in I}f_iw_i$. With threshold t, we estimate fw as $f\widehat{w}^t=\sum_{i\in I}f_i\widehat{w}_i^t=\sum_{i\in S}f_i\widehat{w}_i^t$. With fixed threshold t, we have $\mathsf{E}[\widehat{w}_i^t]=w_i$, so $\mathsf{E}[f_i\widehat{w}_i^t]=f_iw_i$ and $\mathsf{E}[f\widehat{w}^t]=fw$.

As an example, suppose we sampled grocery bills. For each bill sampled, we could check the fraction spent on candy, and based on that estimate the total amount spent on candy.

To emulate a standard subset J, we let f be the characteristic function of J, that is, $f_i = 1$ if $i \in J$; otherwise $f_i = 0$. In fact, we will often identify a set with its characteristic vector, so the weight of J can be written as Jw and estimated as $J\hat{w}^t$.

Using inner products will simplify a lot of notation in our analysis. The generalization to fractional subsets comes for free in our analysis which is all based on concentration bounds for sums of random variables $X_i \in [0, 1]$.

Notation for small and sampled weights With threshold t, we know that variability in the estimates is from items i with weight below t. We will generally use a subscript < t to denote the restriction to items i

with weights $w_i < t$, e.g., $I_{< t} = \{i \in I | w_i < t\}$, $w_{< t} = (w_i)_{i \in I_{< t}}$, and $fw_{< t} = \sum_{i \in I_{< t}} f_i w_i$. Notice that $fw_{< t}$ does not include i with $w_i \ge t$ even if $f_i w_i < t$.

Above we defined $w_{< t}$ to denote the vector $(w_i)_{i \in I_{< t}}$ of weights below t, and used it for the inner product $fw_{< t} = \sum_{i \in I_{< t}} f_i w_i$. When it is clear from the context that we need a number, not a vector, we will use $w_{< t}$ to denote the sum of these weights, that is, $w_{< t} = \sum_{i \in I_{< t}} w_i = \underline{1} w_{< t}$ where $\underline{1}$ is the all 1s vector. Since $f_i \leq 1$ for all i, we always have $fw_{< t} \leq w_{< t}$.

We shall use subscript $\leq t$, $\geq t$, and > t to denote the corresponding restriction to items with weight $\leq t$, $\geq t$, and > t, respectively.

We also introduce a superscript t notation to denote the restriction to items sampled with threshold t, that is, items i with $q_i > t$, so $I_{\leq t}^t$ denotes the set of items with weights below t that ended up sampled. Identifying this set with its characteristic vector, we can write our estimate with threshold t as

$$f\widehat{w}^t = fw_{\geq t} + t(fI_{\leq t}^t). \tag{12}$$

Error probability functions As mentioned earlier, we will reduce the priority sampling error event to the union of an unbounded number of threshold sampling error events that happen with geometrically decreasing probabilities. Our reduction will hold for most hash functions, including 2-independent hash functions, but to make such a claim clear, we have to carefully describe what properties of the hash functions we rely on.

Assume that the threshold t is fixed. With reference to (12), the variability in our estimate is all from

$$fI_{\le t}^t = \sum_{i \in I_{\le t}} X_i$$
, where $X_i = f_i[h_i < w_i/t] \in [0, 1]$.

As for the regular subsets in (11), let $X = \sum_{i \in I_{<t}} X_i$ and $\mu = \mathsf{E}[X]$. We are interested in an error probability function \wp such that for $\mu > 0$, $\delta > 0$, if $\mu = \mathsf{E}[X]$, then

$$\Pr[|X - \mu| > \delta \mu] \le \wp(\mu, \delta). \tag{13}$$

The error probability function \wp that we can use depends on the quality of the hash function. For example, if the hash function is 2-independent, then $Var[X] \le \mu$, and then by Chebyshev's inequality, we can use

$$\wp^{\text{Chebyshev}}(\mu, \delta) = 1/(\delta^2 \mu). \tag{14}$$

In the case of full randomness, for $\delta \leq 1$, we could use a standard 2-sided Chernoff bound (see, e.g., [26])

$$\wp^{\text{Chernoff}_{\delta \le 1}}(\mu, \delta) = 2\exp(-\delta^2 \mu/3). \tag{15}$$

For most of our results, it is more natural to think of δ as a function of μ and some target error probability $P \in (0,1)$, defining $\delta(\mu, P)$ such that

$$\mu(\mu, \delta(\mu, P)) = P. \tag{16}$$

Returning to threshold sampling with threshold t, by (12) the error is $f\widehat{w}^t - fw = t(fI_{< t}^t) - fw_{< t}$. Hence

$$\Pr[|f\widehat{w}^t - fw| > \delta(fw_{< t}/t, P)fw_{< t}] \le P.$$
(17)

When we start analyzing priority sampling, we will need to relate the probabilities of different threshold sampling events. This places some constraints on the error probability function \wp . Mathematically, it is convenient to allow \wp to attain values above 1, but only values below 1 are probabilistically interesting.

Definition 4 An error probability function $\wp : \mathbb{R}_{\geq 0} \times R_{\geq 0} \to R_{\geq 0}$ is well-behaved if

- (a) \wp is continuous and strictly decreasing in both arguments.
- (b) If with the same absolute error we decrease the expectancy, then the probability goes down. Formally if $\mu' < \mu$ and $\mu'\delta' \ge \mu\delta$, then $\wp(\mu', \delta') < \wp(\mu, \delta)$.

We also have an optional condition for cases where we only care for $\delta \leq 1$ as in (15)

(c) If $\delta \leq 1$ and $\wp(\mu, \delta) < P_{\wp}^{(c)}$ for some constant $P_{\wp}^{(c)}$ depending on \wp , then $\wp(\mu, \delta)$ falls at least proportionally to $\mu\delta^2$. Formally, if $\delta_0, \delta_1 \leq 1$, $\wp(\mu_0, \delta_0) < 1$, and $\mu_0\delta_0^2 < \mu_1\delta_1^2$, then

$$\wp(\mu_0, \delta_0) \ge \frac{\mu_0 \delta_0^2}{\mu_1 \delta_1^2} \wp(\mu_1, \delta_1). \tag{18}$$

We will use condition (c) to argue that probabilities of different events fall geometrically. The condition is trivially satisfied with our 2-independent Chebyshev bound (14), so we can just set $P_{\wp^{\text{Chebyshev}}}^{(c)}=1$. The restrictions in (c) are necessary for the Chernoff bound (15), $\wp^{\text{Chernoff}_{\delta \leq 1}}=2\exp(-\delta^2\mu/3)$, which only falls fast enough for $\delta^2\mu/3 \geq 1$, hence with $\wp^{\text{Chernoff}_{\delta \leq 1}}(\mu,\delta) \leq P_{\wp^{\text{Chernoff}_{\delta \leq 1}}}^{(c)}=2/e$. As a further illustration, with 4-independence, we have the 4th moment bound $\frac{\mu+3\mu^2}{(\delta\mu)^4}$ (see, e.g., [22, Lemma 4.19]). For $\mu \geq 1$, this is upper bounded by $\wp^{\text{4th-moment}_{\mu \geq 1}}(\mu,\delta) = \left(\frac{2}{\mu\delta^2}\right)^2$. For $\delta \leq 1$, the condition $\mu \geq 1$ is satisfied if $\wp^{\text{4th moment}_{\mu \geq 1}}(\mu,\delta) \leq 4$, so we can just use $P_{\wp^{\text{4th moment}}}^{(c)}=1$.

Threshold confidence intervals In the case of threshold sampling with a fixed threshold t, we get some trivial confidence intervals for the true value fw. The sample gives us the exact value $fw_{\geq t}$ for weights at least as big as t, and an estimate $f\widehat{w}_{\leq t}^t$ for those below. Setting

$$\begin{split} f\widehat{w}_{< t}^- &= \min\{x \,|\, (1 + \delta(x/t, P))x \geq f\widehat{w}_{\geq t}^t\} \\ f\widehat{w}_{< t}^+ &= \max\{x \,|\, (1 - \delta(x/t, P))x \leq f\widehat{w}_{\geq t}^t\} \end{split}$$

we get

$$\Pr\left[fw_{>t} + f\widehat{w}_{< t}^{-} \le fw \le fw_{>t} + f\widehat{w}_{< t}^{+}\right] \ge 1 - 2P.$$

We are going to show that similar bounds can be obtained for priority sampling.

3.2 Priority sampling: the main result

We are now ready to present our main technical result. We are considering a random priority sample of size k, and let τ denote the resulting priority threshold. The sample size k and the target error probability P are both fixed in advance of the random sampling.

Theorem 5 Let the error probability function \wp satisfy Definition 4 including (c). With target error probability $P \leq P_{\wp}^{(c)}$, let

$$\delta = 6 \, \delta(f w_{<\tau}/(3\tau), P).$$

If $\delta \leq 2$, then

$$\Pr[|f\widehat{w}^{\tau} - fw| > \delta fw_{<\tau}] \le 6P.$$

The above constants are not optimized, but with O-notation, some of our statements would be less clear. Ignoring the constants and the restriction $\delta \leq 2$, we see that our error bound for priority sampling with threshold τ is of the same type as the one in (17) for threshold sampling with fixed threshold $t = \tau$.

The proof of Theorem 5 is rather convoluted. We consider a single priority sampling event with k samples and priority threshold τ . It assigns a random priority q_i to each item, and this defines a sample for any given threshold t. In particular, $\tau \geq t$ if and only if we get at least k+1 samples with threshold t. Note that $\Pr[\tau = t] = 0$ for any given t. We define

$$t_{\max} = \min\{t \mid \Pr[\tau \ge t] \le P\}$$
 (19)

$$t_{\min} = \max\{t \mid \Pr[\tau < t] \le P\}$$
 (20)

By definition, $\Pr[\tau \notin [t_{\min}, t_{\max})] \leq 2P$. By union, to prove Theorem 5, it suffices to show that the following good event errors with probability at most 4P:

$$\forall t \in [t_{\min}, t_{\max}) : \delta = \delta(fw_{< t}/(3t), P) \le 1/3$$

$$\implies |f\widehat{w}^t - fw| > 6\delta fw_{< t}$$
(21)

Note that our variable δ is 6 times smaller than the one in the statement of Theorem 5. This parameter change will be more convenient for the analysis. What makes (21) very tricky to prove is that $\delta(fw_{< t}/(3t), P)$ can vary a lot for different $t \in [t_{\min}, t_{\max})$.

Priority confidence intervals The format of Theorem 5 makes it easy to derive confidence intervals like those for threshold sampling. A priority sample with priority threshold τ gives us the exact value $fw_{\geq \tau}$ for weights at least as big as τ , and an estimate $f\widehat{w}_{\leq \tau}^{\tau}$ for those below. For an upper bound on $fw_{\leq \tau}$, we compute

$$f\widehat{w}_{<\tau}^+ = \max\{x \mid \delta = 6\,\delta(x/(3\tau), P)) \wedge (1 - \delta)x \le f\widehat{w}_{\ge \tau}^\tau\}.$$

Note that here in the upper bound, we only consider $\delta \leq 1$, so we do not need to worry about the restriction $\delta < 2$ in Theorem 5. For a lower bound on $fw_{\leq \tau}$, we use

$$f\widehat{w}_{<\tau}^- = \min\{x \mid \delta = 6\,\delta(x/(3\tau), P)) \le 2 \wedge (1+\delta)x \ge f\widehat{w}_{>\tau}^\tau\}.$$

Here in the lower bound, the restriction $\delta = 6 \, \delta(x/(3\tau), P)) \leq 2$ prevents us from deriving a lower bound $x = f \widehat{w}_{\leq \tau}^- \leq f \widehat{w}_{\geq \tau}^\tau/3$. In such cases, we use the trivial lower bound $x = f \widehat{w}_{\leq \tau}^- = 0$ which in distance from $f \widehat{w}_{\leq \tau}^\tau$ is at most 1.5 times bigger. Now, by Theorem 5,

$$\Pr\left[fw_{\geq \tau} + f\widehat{w}_{<\tau}^{-} \le fw \le fw_{\geq \tau} + f\widehat{w}_{<\tau}^{+}\right] \ge 1 - 12P.$$

In cases where the exact part $fw_{\geq \tau}$ of an estimate is small compared with the variable part $f\widehat{w}_{\geq \tau}^{\tau}$, we may be interested in a non-zero lower bound $f\widehat{w}_{\leq \tau}^{-}$ even if it is smaller than $f\widehat{w}_{\geq \tau}^{\tau}/3$. To do this, we need bounds for larger δ .

Large errors We are now going to present bounds that works for arbitrarily large relative errors δ . We assume a basic error probability function \wp satisfying Definition 4 (a) and (b) while (c) may not be satisfied. The bounds we get are not as clean as those from Theorem 5. In particular, they involve t_{\max} from (19). Since we are only worried about errors $\delta > 1$, we only have to worry about positive errors.

Theorem 6 Set

$$\ell = \lg\lceil (t_{\text{max}}/\tau) \rceil$$

$$\delta = \delta(fw_{<\tau}/\tau, P/\ell^2).$$

Then

$$\Pr[f\widehat{w}_{<\tau}^{\tau} > (2+2\delta)fw_{<\tau}] < 3P.$$

Complementing Theorem 5, we only intend to use Theorem 6 for large errors where $(2+2\delta)=O(\delta)$. We wish to provide a probabilistic lower bound for $fw_{<\tau}$. Unfortunately, we do not know t_{\max} which depends on the whole weight vector $(w_i)_{i\in I}$. However, based our priority sample, it is not hard to generate a probabilistic upper bound \overline{t}_{\max} on t_{\max} such that $\Pr[\overline{t}_{\max} < t_{\max}] \leq P$. We set

$$\overline{\ell} = \lceil \lg(\overline{t}_{\text{max}}/\tau) \rceil \tag{22}$$

$$f\widehat{w}_{<\tau}^{-} = \min\{x \mid \overline{\delta} = \delta(x/\tau, P/\overline{\ell}^{2}) \land 2(1+\overline{\delta})x \ge f\widehat{w}_{<\tau}^{\tau}\}.$$
 (23)

Then by Theorem 6,

$$\Pr\left[fw_{<\tau} \geq f\widehat{w}_{<\tau}^{-}\right] \geq 1 - 4P.$$

To see this, let $f\widehat{w}_{<\tau}^*$ be the value we would have obtained if we had computed $f\widehat{w}_{<\tau}^-$ using the real t_{\max} . Our error event is that $\overline{t}_{\max} < t_{\max}$ or $fw_{<\tau} < f\widehat{w}_{<\tau}^*$. The former happens with probability at most P, and Theorem 6 states that the latter happens with probability at most 3P. Hence none of these error events happen with probability at least 1-4P, but then $\overline{t}_{\max} \geq t_{\max}$, implying $f\widehat{w}_{<\tau}^- \leq f\widehat{w}_{<\tau}^* \leq fw_{<\tau}$.

Theorem 5, our main result, is proved in Sections 3.3–3.5. Theorem 6, which is much easier, is proved in Section 3.6. In Section 3.7 we show how to compute the \overline{t}_{max} used for confidence lower bound with Theorem 6. Finally, in Section 3.8 we will argue that we for typical weight distributions expect to get $\overline{\ell}=1$.

3.3 The priority threshold

To prove Theorem 5 we need a handle on the variable priority threshold. With priority sampling we specify the number k of samples, and use as threshold τ the (k+1)th priority. Recall for any threshold t that the subscript $\leq t$ indicates restriction to items with weight below t. To relate this notation to a priority sample of some specified size k, we let $k \leq t$ denote k minus the number of items with weight bigger than t. We define k < t accordingly. With threshold t, the expected number of samples is $E[|I^t|] = k - k \leq t + w \leq t/t = k - k < t + w < t/t$. The last equality is because weights $w_i = t$ cancel out.

The ideal threshold We define the *ideal threshold* t^* to be the one leading to an expected number of exactly k samples, that is $w_{< t^*} = t^* k_{< t^*}$.

Lemma 7 $tk_{\leq t} - w_{\leq t}$ is strictly increasing in t, so (a) t^* is uniquely defined, (b) $w_{\leq t} > k_{\leq t}t$ for any $t < t^*$, and (c) $w_{\leq t} < k_{\leq t}t$ for any $t > t^*$.

Proof If t increases without passing any weight value, then $k_{\leq t}$ and $w_{\leq t}$ are not changed, and the statement is trivial. When t reaches the value of some weight w_i , then both $tk_{\leq t}$ and $w_{\leq t}$ are increased by the same value w_i (if there are j weights with value w_i , the increase is by jw_i).

We would like to claim that the priority sampling threshold τ is concentrated around t^* , but this may be far from true. To illustrate what makes things tricky to analyze, consider the case where, say, we have k-1 weights of size t^* , and then a lot of small weights that sum to t^* . In this case we get an upper bound on τ which is close to t^* , but we do not get any good lower bound on τ even if we have full randomness. On the other hand, in this case, it is only little weight that is affected by the downwards variance in τ .

3.4 Tightening the gap

The following lemmas give us a much tighter understanding of $t \in [t_{\min}, t_{\max})$.

Lemma 8 If $t \ge t^*$ then $t(1 + \delta(w_{\le t}/t_{\max}, P)) \ge t_{\max}$.

Proof Let $\delta = \delta(w_{\leq t}/t_{\max}, P)$ and $T = (1 + \delta)t$. Assume for a contradiction that $T < t_{\max}$. By Lemma 7 (c), since $t \geq t^*$, we have $w_{\leq t} \leq t k_{\leq t}$, so

$$k_{\le t} > w_{\le t}/t = (1 + \delta)w_{\le t}/T$$

For the priority threshold to be as big as T we need $|I_{\leq t}^T| \geq k_{\leq t} + 1$ and $\mathsf{E}[|I_{\leq t}^T|] = w_{\leq t}/T$, so

$$\Pr[\tau \ge T] \le \wp(w_{< t}/T, \delta) \le \wp(w_{< t}/t_{\max}, \delta) = P.$$

But this contradicts the minimality of t_{max} from (19).

Lemma 9 If $t \in [t_{\min}, t^*)$ then $t \geq (1 - \delta(w_{\leq t}/t, P))t^*$.

Proof Let $\delta = \delta(w_{\leq t}/t, P)$. The proof is by contradiction against our maximal lower bound t_{\min} from (20). The priority threshold is smaller than t if $|I_{\leq t}^t| \leq k_{\leq t}$. The expectancy is $\mathsf{E}[|I_{\leq t}^t|] = w_{\leq t}/t$. Suppose for a contradiction that $(1-\delta)w_{\leq t}/t > k_{\leq t}$. Then

$$\Pr[\tau < t] < \wp(w_{\le t}/t, \delta) = P,$$

implying that $t \geq t_{\min}$ is a lower bound. If $t > t_{\min}$ this contradicts the maximality of t_{\min} . Otherwise, we pick an infinitesimally larger t' > t with no weights in (t,t'], that is, $w_{\leq t'} = w_{\leq t}$ and $k_{\leq t'} = k_{\leq t}$. By Definition 4 (a), we get a corresponding infinitesimally larger $\delta' = \delta(w_{\leq t'}/t', P) > \delta$, and then we still have $(1-\delta')w_{\leq t'}/t' > k_{\leq t'}$, implying that $t' > t = t_{\min}$ is a lower bound contradicting the maximality of t_{\min} . Thus we conclude that $(1-\delta)w_{\leq t}/t \leq k_{\leq t}$, or equivalently, $t \geq (1-\delta)w_{\leq t}/k_{\leq t}$. Finally by Lemma 7 (b), since $t < t^*$, we have $w_{\leq t}/k_{\leq t} \geq t^*$. This completes the proof that $t \geq (1-\delta_\ell)t^*$.

In order to give a joint analysis for t bigger and smaller than t^* , we make a conservative combination of Lemma 8 and 9.

Lemma 10 Suppose $[t^-, t^+) = [t_{\min}, t^*)$ or $[t^-, t^+) = [t^*, t_{\max})$. If $t \in [t^-, t^+)$ then $t \geq (1 - \delta(w_{\leq t}/t^+, P))t^+$.

Proof Let $\delta = \delta(w_{\leq t}/t^+, P)$. If $[t^-, t^+) = [t^*, t_{\max})$, by Lemma 8, $t \geq t^+/(1+\delta) \geq (1-\delta)t^+$. If $[t^-, t^+) = [t_{\min}, t^*)$, by Lemma 9, $t \geq (1-\delta(w_{\leq t}/t, P))t^+ > (1-\delta)t^+$. The last inequality uses that $\delta(w_{\leq t}/t, P) \leq \delta$. This is because $w_{\leq t}/t \geq w_{\leq t}/t^+$, so by Definition 4 (a), $\wp(w_{\leq t}/t, \delta) \geq \wp(w_{\leq t}/t^+, \delta) = P$.

Loosing a factor 2 in the error probability to cover t bigger or smaller than t^* , our good event (21) reduces to

$$\forall t \in [t^-, t^+) : \delta = \delta(fw_{< t}/(3t), P) \le 1/3$$

$$\implies |f\widehat{w}^t - fw| > 6\delta fw_{< t}$$
(24)

By Lemma 10, the bound on δ implies that $t \geq (2/3)t^+$. Therefore (24) is implied by

$$\forall t \in [t^-, t^+) : \delta = \delta(fw_{< t}/(2t^+), P) \le 1/3$$

$$\implies |f\widehat{w}^t - fw| \le 6\delta fw_{< t}. \tag{25}$$

One advantage of dealing with $f\widehat{w}_{< t}/t^+$ instead of $f\widehat{w}_{< t}/t$ is that $f\widehat{w}_{< t}/t^+$ is proportional to $f\widehat{w}_{< t}$ hence increasing in t whereas $f\widehat{w}_{< t}/t$ may not be monotone in t. Then $\delta(fw_{< t'}^-/(2t^+), P)$ is decreasing in t. If $\delta(fw_{< t'}^-/(2t^+), P) > 1/3$, we let t'^- be the smallest value such that $\delta(fw_{< t'}^-/(2t^+), P) \le 1/3$; otherwise we set $t'^- = t^-$. Then (25) is equivalent to

$$\forall t \in [t'^-, t^+) : \delta = \delta(fw_{< t}/(2t^+), P)$$

$$\implies |f\widehat{w}^t - fw| \le 6\delta fw_{< t}. \tag{26}$$

3.5 Dividing into layers

We now define a sequence $t_0 > t_1 > \cdots > t_{L+1}$ of decreasing thresholds with $t_0 = t^+$ and $t_{L+1} = t'^-$. For $\ell = 0, ..., L$ we require

$$fw_{< t_{\ell+1}} \ge fw_{< t_{\ell}}/2.$$
 (27)

For $\ell < L$, we pick $t_{\ell+1}$ smallest possible satisfying (27). Then

$$fw_{< t_{\ell+1}} < fw_{< t_{\ell}}/2.$$
 (28)

We arrive $\ell = L$ when $fw_{\leq t'^-} \geq fw_{\leq t_\ell}/2$, and then we set $t_{L+1} = t'^-$.

For each "layer" $\ell \leq \overline{L}$, we define

$$\delta_{\ell} = \delta(fw_{< t_{\ell}}/(2t^{+}), P), \tag{29}$$

noting that this is the same value as we would use for $t=t_\ell$ in (26). By definition, for all $\ell \leq L$, we have $\wp(fw_{< t_\ell}/(2t^+), \delta_\ell) = P$. Since $t_\ell > t'^-$, we have $\delta_\ell \leq 1/3 \leq 1$, so it follows from Definition 4 (c) that there is a constant C such that $fw_{< t_\ell} \delta_\ell^2 = C$ for all $\ell \leq L$. For $\ell = 1, ..., L$, by (28), we have $fw_{< t_\ell} < fw_{< t_{\ell-1}}/2$. Therefore

$$\delta_{\ell} > \sqrt{2} \, \delta_{\ell-1} \tag{30}$$

$$\delta_{\ell} f w_{\langle t_{\ell}} \langle \delta_{\ell-1} f w_{\langle t_{\ell-1}} / \sqrt{2} \tag{31}$$

This will correspond to an effect where the relative errors are geometrically increasing while the absolute errors are geometrically decreasing. Another important thing to notice is that by (27), $fw_{\leq t_{\ell+1}} \geq fw_{< t_{\ell}}/2$, so $\delta(fw_{< t_{\ell+1}}/t^+, P) \leq \delta(fw_{< t_{\ell}}/(2t^+), P) = \delta_{\ell}$. Therefore, by Lemma 10,

$$t_{\ell+1} \ge (1 - \delta_{\ell})t^{+}. \tag{32}$$

Good layers For each layer $\ell < L$, our good event will be that for weights in $[t_{\ell+1}, t_{\ell})$, the relative estimate error is bounded by $2\delta_{\ell}$. Formally

$$\forall t \in [t_{\ell+1}, t^+) : \left| f \widehat{w}_{[t_{\ell+1}, t_{\ell})}^t - f w_{[t_{\ell+1}, t_{\ell})} \right|$$

$$\leq 2\delta_{\ell} f w_{[t_{\ell+1}, t_{\ell})}$$
(33)

Above, the subscript $[t_{\ell+1},t_{\ell}]$ denotes the restriction to items i with weights $w_i \in [t_{\ell+1},t_{\ell}]$. The last layer L is special in that we want to consider all weights below t_L . Here the good event is that

$$\forall t \in [t_{L+1}, t^+) : \left| f \widehat{w}_{\leq t_L}^t - f w_{\leq t_L} \right|$$

$$\leq 3\delta_L f w_{\leq t_L}$$
(34)

To prove Theorem 5, we are going to prove two statements.

- Assume that all layers are good satisfying (33) and (34). If for any $t \in (t'^-, t^+]$ we add up the errors from all relevant layers, then the total error is bounded by $6\delta(fw_{< t}/(2t^+), P)fw_{< t}$, so (26) satisfied.
- If P_{ℓ} is the probability that a layer ℓ fails, then the P_{ℓ} are geometrically increasing and $P_L = O(P)$, so by union, the probability that any layer fails is O(P).

Adding layer errors Assuming that all layers are good satisfying (33) and (34), we pick an arbitrary threshold $t \in [t'^-, t^+)$. We which to bound the estimate error $|f\widehat{w}^t - fw|$.

Let h be the layer such that $t \in [t_{h-1}, t_h)$. We can only have estimate errors from weights below $t < t_h$, so

$$|f\widehat{w}^{t} - fw| \leq \sum_{\ell=h}^{L-1} \left| f\widehat{w}_{[t_{\ell+1}, t_{\ell})}^{t} - fw_{[t_{\ell+1}, t_{\ell})} \right| + \left| f\widehat{w}_{< t_{L}}^{t} - fw_{< t_{L}} \right|$$

$$\leq \sum_{\ell=h}^{L-1} 2\delta_{\ell} w_{[t_{\ell+1}, t_{\ell})} + 3\delta_{L} fw_{< t_{L}}$$

$$= \sum_{\ell=h}^{L-1} 2\delta_{\ell} (fw_{< t_{\ell}} - fw_{< t_{\ell+1}}) + 3\delta_{L} fw_{< t_{L}}$$

By (30), the δ_ℓ are increasing, so in the above sum, every $fw_{< t_\ell}$ appears with a positive coefficient. It follows that we could only get a larger sum if (28) was more than tight with $fw_{< t_\ell} = fw_{< t_{\ell-1}}/2$ for $\ell \le L$. Then we would have $fw_{< t_\ell} - fw_{< t_{\ell+1}} = fw_{< t_\ell}/2$ and corresponding to (31), $\delta_\ell fw_{< t_\ell} = \delta_{\ell-1} fw_{< t_{\ell-1}}/\sqrt{2}$. Thus we get

$$|f\widehat{w}^{t} - fw| \leq \sum_{\ell=h}^{L-1} 2\delta_{\ell} (fw_{< t_{\ell}} - fw_{< t_{\ell+1}}) + 3\delta_{L} fw_{< t_{L}}$$

$$\leq \sum_{\ell=h}^{L-1} \delta_{h} fw_{< t_{h}} / \sqrt{2}^{\ell-h} + 3\delta_{h} fw_{< t_{h}} / \sqrt{2}^{L-h}$$

$$< 3\delta_{h} fw_{< t_{h}}$$

The last inequality exploits that $\sum_{i=1}^{\infty} 1/\sqrt{2} = 1/(1-1/\sqrt{2}) < 3$. Since $t \leq t_h$, we have $\delta(fw_{< t}/(2t^+), P) \geq \delta(fw_{< t_h}/(2t^+), P) = \delta_h$. Also, $t > t_{h+1}$, so by (27), $fw_{< t} \geq fw_{\leq t_{h+1}} \geq fw_{< t_h}/2$, so

$$\delta(fw_{< t}/(2t^+), P) fw_{< t} \ge \delta_h fw_{< t_h}/2.$$

Therefore

$$|f\widehat{w}^t - fw| \le 3\delta_h f w_{< t_h} \le 6 \, \delta(f w_{< t}/(2t^+), P) \, f w_{< t}.$$

Thus we conclude that (26) follows from (33) and (34).

Intermediate layer error probabilities We now consider the intermediate layers $\ell = 0, ..., L - 1$. We want to show that the probability P_{ℓ} of violating (33) increases geometrically with ℓ , yet remains bounded by P. First we consider the upper bound part of (33)

$$\forall t \in [t_{\ell+1}, t^+) : f\widehat{w}_{[t_{\ell+1}, t_{\ell})}^t \le (1 + 2\delta_{\ell}) f w_{[t_{\ell+1}, t_{\ell})}. \tag{35}$$

We claim that it can never be violated. The worst that can happen is that every item i in the layer gets sampled, and the estimate is at most f_it^+ . However, the items all have weight at least $t_{\ell+1}$ and by (32), $t_{\ell+1} \geq (1-\delta_\ell)t^+$. The increase is thus by at most a factor $t^+/t_{\ell+1} \leq 1/(1-\delta_\ell)$, which for $\delta_\ell \leq 1/3$ is at most $(1+2\delta_\ell)$. Thus (35) is satisfied regardless of the random choices.

We now consider the lower bound part of (33)

$$\forall t \in [t_{\ell+1}, t^+) : f\widehat{w}_{[t_{\ell+1}, t_{\ell})}^t \ge (1 - 2\delta_{\ell}) f w_{[t_{\ell+1}, t_{\ell})}. \tag{36}$$

This event could happen. To bound the probability, we will focus on the loss $fw_{[t_{\ell+1},t_{\ell})}-f\widehat{w}_{[t_{\ell+1},t_{\ell})}^t$. When bounding the loss, we do not consider the gain from possible overestimates of sampled items. We only consider the actual losses f_iw_i from unsampled items i. Conservatively, we consider an item i lost if $q_i \leq t^+$. This includes any item unsampled with some threshold $t \leq t^+$. The loss for every threshold $t \leq t^+$ is thus bounded as

$$fw_{[t_{\ell+1},t_{\ell})} - f\widehat{w}_{[t_{\ell+1},t_{\ell})}^t \le \sum_{i:w_i \in [t_{\ell+1},t_{\ell})} [q_i \le t^+] f_i w_i.$$

We know that $w_i \geq t_{\ell+1} \geq (1 - \delta_{\ell})t^+$. Therefore

$$\Pr[q_i \le t^+] = \Pr[h_i \ge w_i/t^+] \le \Pr[h_i \ge t_{\ell+1}/t^+]$$

= 1 - t_{\ell+1}/t^+ \le \delta_\ell.

The expected loss from layer ℓ is thus bounded by

$$\sum_{i:w_i \in [t_{\ell+1},t_\ell)} \delta_\ell f_i w_i = \delta_\ell f w_{[t_{\ell+1},t_\ell)}.$$

For (36) to fail, we need a loss that is twice this big, that is, $2\delta_{\ell} f w_{[t_{\ell+1},t_{\ell})}$. We know that items *i*'s loss contribution $[q_i \leq t^+] f_i w_i$ depends only on h_i and that it is at most t^+ . The probability of violating (36) is therefore bounded by

$$\wp(\delta_{\ell} f w_{[t_{\ell+1},t_{\ell})}/t^+, 1) \le \wp(\delta_{\ell} f w_{< t_{\ell}}/(2t^+), 1).$$

Let $\mu_{\ell} = \delta_{\ell} f w_{< t_{\ell}}/(2t^{+})$. Then $P_{\ell} = \wp(\mu_{\ell}, 1)$ is our bound on the error probability. From (31), we know that $\delta_{\ell} f w_{< t_{\ell}} < \delta_{\ell-1} f w_{< t_{\ell-1}}/\sqrt{2}$, so $\mu_{\ell} < \mu_{\ell-1}/\sqrt{2}$. It follows from Definition 4 (c) that $\wp(\mu_{\ell}, 1) > \sqrt{2} \wp(\mu_{\ell-1}, 1)$, so the P_{ℓ} are geometrically increasing, and their sum is bounded by $3P_{L-1}$.

Finally, by definition (29), $\wp(fw_{< t_{L-1}}/(2t^+), \delta_{L-1}) = P$. To compare P with P_{L-1} , by Definition 4 (c), we compare $fw_{< t_{L-1}}/(2t^+)\delta_{L-1}^2$ with $\mu_{L-1}1^2 = \delta_{L-1}fw_{< t_{L-1}}/(2t^+)$, and conclude that $P_{L-1} \leq \delta_{L-1}P \leq P/3$. The probability that any intermediate layer $\ell < L$ fails (36) is thus at most P. Since (36) was always satisfied, we conclude that (34) is satisfied for all layers $\ell < L$ with probability P.

The last layer We now consider items i with weights below t_L . On the upper bound side, the good event (34) states that

$$\forall t \in [t_{L+1}, t^+) : f \widehat{w}_{\leq t_L}^t \leq (1 + 3\delta_L) f w_{\leq t_L}. \tag{37}$$

We will show that (37) fails with probability less than P. For an upper bound on the estimate with any threshold $t \in [t_{L+1}, t^+)$, we include item i if $q_i > t_{L+1}$, and if so, we give it at an estimate of $f_i t^+$ which is bigger than the sampled estimate with threshold $t \le t^+$. The result is at most a factor t^+/t_{L+1} bigger than in the sampled estimate with threshold t_{L+1} , and by (32), $t_{L+1} \ge (1 - \delta_L)t^+$. Thus, regardless of the random choices made, we conclude that

$$\forall t \in [t_{L+1}, t^+) : f \widehat{w}_{\leq t_L}^t \leq f \widehat{w}_{\leq t_L}^{t_{L+1}} / (1 - \delta_L).$$

Consider the following error event:

$$f\widehat{w}_{\leq t_L}^{t_{L+1}} > (1 + \delta_L) f w_{\leq t_L} \tag{38}$$

The maximal item contribution to $f\widehat{w}_{\leq t_L}^{t_{L+1}}$ is bounded by $t_L \leq t^+$, so the probability of (38) is bounded by

$$\wp(fw_{\leq t_L}/t^+, \delta_L) \leq \wp(fw_{\leq t_L}/(2t^+), \delta_L)/2 = P/2.$$

If (38) does not happen, then since $\delta_L \leq 1/3$,

$$\forall t \in [t_{L+1}, t^+) : f \widehat{w}_{\leq t_L}^t \leq (1 + \delta_L) f w_{\leq t_L} / (1 - \delta_L)$$

$$\leq (1 + 3\delta_L) f w_{\leq t_L},$$

which is the statement of (37). We conclude that (37) fails with probability at most P/2.

We now consider the lower bound side of (34) which states that

$$\forall t \in [t_{L+1}, t^+) : f \widehat{w}_{< t_L}^t \ge (1 - 3\delta_L) f w_{< t_L}$$
(39)

The analysis is very symmetric to the upper bound case. For a lower bound for weights $w_i < t_L$ with any threshold $t \in [t_{L+1}, t^+)$, we only include i if $q_i > t^+$, and if so, we only give i the estimate $f_i t_{L+1}$ which is smaller than the sampled estimate with threshold $t > t_{L+1}$. Our samples are exactly the same as those we would get with threshold t^+ , and our estimates are smaller by a factor $t_{L+1}/t^+ \ge (1 - \delta_L)$, so we conclude that regardless of the random choices,

$$\forall t \in [t_{L+1}, t^+) : f \widehat{w}_{\leq t_L}^t \ge f \widehat{w}_{\leq t_L}^{t^+} (1 - \delta_L)$$

We now consider the error event:

$$f\widehat{w}_{\leq t_L}^{t^+} < (1 - \delta_L) f w_{\leq t_L} \tag{40}$$

The maximal item contribution to $f\widehat{w}_{\leq t_L}^{t^+}$ is t^+ , so as for (38), we get that the probability of (40) is bounded by $\wp(fw_{\leq t_L}/t^+, \delta_L) \leq P/2$.

If (40) does not happen, then since $\delta_L \leq 1/3$,

$$\forall t \in [t_{L+1}, t^+) : f \widehat{w}_{\leq t_L}^t \ge (1 - \delta_L) f w_{\leq t_L} (1 - \delta_L)$$

> $(1 - 2\delta_L) f w_{\leq t_L}$,

which implies (39). We conclude that (37) fails with probability at most P/2. Including the probability of an upper bound error (38), we get that (34) fails with probability at most P.

Summing up Above we proved that the probability that (33) failed for any layer $\ell < L$ was at most P. We also saw that (34) failed with probability P. If none of them fail, we proved that (26) and hence (24) was satisfied, so (24) fails with probability at most 2P. However, for (21) we need (24) both for $[t^-,t^+)=[t_{\min},t^*)$ and for $[t^-,t^+)=[t^*,t_{\max})$, so (21) fails with probability at most 4P. Finally, we need to consider both the case that $\tau>t_{\max}$ and $\tau\leq t_{\min}$. Either of these events happens with probability at most P, so we conclude that the overall error probability is at most 6P. If no error happened and $\delta=\delta(fw_{< t}/(3t),P)\leq 1/3$, then $|f\hat{w}^t-fw|>6\delta fw_{< t}$. This completes the proof of Theorem 5.

3.6 Large errors

The limitation of Theorem 5 is that it can only be used to bound the probability that the estimate error $|f\widehat{w}^{\tau} - fw|$ is bigger than $2fw_{<\tau}$. Note that errors above $fw_{<\tau}$ can only be overestimates. We will now target larger errors and prove the statement of Theorem 6:

Set

$$\ell = \lceil \lg(t_{\text{max}}/\tau) \rceil$$
$$\delta = \delta(fw_{<\tau}/\tau, P/\ell^2).$$

Then

$$\Pr[f\widehat{w}_{<\tau}^{\tau} > (2+2\delta)fw_{<\tau}] < 3P$$

Proof of Theorem 6 Since we are targeting arbitrarily large relative errors δ , for the probability function \wp , we can only assume conditions (a) and (b) in Definition 4.

We will use some of the same ideas as we used for the last layer in Section 3.5, but tuned for our situation. We will study intervals based on $t_\ell = t_{\max}/2^\ell$ for $\ell = 1, 2, ...$ Interval ℓ is for thresholds $t \in [t_\ell, t_{\ell-1}) = [t_\ell, 2t_\ell)$, so $t < t_{\max}$ belongs to interval $\ell = \lceil \lg(t_{\max}/t) \rceil$. To define the error for interval ℓ , set

$$\delta_{\ell} = \wp(fw_{< t_{\ell}}/t_{\ell}, P/\ell^2).$$

The good non-error event for interval ℓ is that

$$f\widehat{w}_{< t_{\ell}}^{t_{\ell}} \le (1 + \delta_{\ell}) f w_{< t_{\ell}}. \tag{41}$$

By definition, the probability that (41) is violated is at most P/ℓ^2 , so the probability of failure for any ℓ is bounded by $\sum_{\ell=1}^{\infty} P/\ell^2 \leq P\pi^2/6 < 1.65P$. The probability that $\tau = t_{\text{max}}$ is zero, so by (19), the event

$$\tau < t_{\text{max}}.$$
 (42)

is violated with probability at most P. Our total error probability is thus bounded by 2.65P < 3P. Below we assume no errors, that is, (41) holds for all ℓ and so does (42).

Consider an arbitrary threshold $t \in (0, t_{\max})$ and let ℓ be such that $t \in [t_\ell, 2t_\ell)$. We can only have errors for weights $w_i < t$, so we want an upper bound on $f\widehat{w}_{< t}^t$. The basic idea for an upper bound is to say that we sample all items with priority above t_ℓ , just as in the estimate $f\widehat{w}_{< t_\ell}^{t_\ell}$, but instead of giving sampled item i estimate $\max\{w_i, t_\ell\}$, it gets value $\max\{w_i, t\}$ which is at most $t/t_\ell < 2$ times bigger. Thus, regardless of the random choices,

$$f\widehat{w}_{\leq t}^t < 2f\widehat{w}_{\leq t}^{t_\ell}.$$

Assuming no error as in (41), we get

$$f\widehat{w}_{< t}^{t_{\ell}} = fw_{[t_{\ell}, t]} + f\widehat{w}_{< t_{\ell}}^{t_{\ell}} \le fw_{< t} + \delta_{\ell} fw_{< t_{\ell}}.$$

For the further analysis, we need a general lemma.

Lemma 11 For thresholds t', t, and relative errors δ' , δ , if t' < t and $\wp(fw_{< t'}/t', \delta_{< t'}) = \wp(fw_{\le t}/t, \delta_{\le t})$, then

$$\delta' f w_{< t'} < \delta f w_{< t}$$
.

Hence, for any fixed target error probability Q in (17), the target error

$$\delta(fw_{< t}/t, Q)fw_{< t}$$

decreases together with the threshold t.

Proof We will divide the decrease from t to t' into a series of atomic decreases. The first atomic "decrease" is from $fw_{\leq t}$ to $fw_{< t}$. This makes no difference unless there are weights equal to t so that $fw_{< t} < fw_{\leq t}$. Assume this is the case and suppose $\wp(fw_{< t}/t, \delta_{< t}) = \wp(fw_{\leq t}/t, \delta_{\leq t})$. Since $fw_{\leq t}/t < fw_{< t}/t$, it follows directly from (b) that $\delta_{< t} fw_{< t}/t < \delta_{\leq t} w_{\leq t}/t$, hence that $\delta_{< t} fw_{< t} < \delta_{\leq t} w_{\leq t}$.

The other atomic decrease we consider is from $fw_{\leq t}$ to $fw_{\leq t'}$ where t' < t and with no weights in (t',t), hence with $fw_{\leq t'} = fw_{< t}$. Suppose $\wp(fw_{\leq t'}/t',\delta_{\leq t'}) = \wp(fw_{< t}/t,\delta_{< t})$. Since t' < t, $fw_{\leq t'}/t' > fw_{< t}/t$, so by (a), $\delta_{\leq t'} < \delta_{< t'}$. It follows that $\delta_{\leq t'}fw_{\leq t'} < \delta_{< t}fw_{< t}$. Alternating between these two atomic decreases, we can implement an arbitrary decrease in the threshold as required for the lemma.

Let

$$\delta = \delta(w_{< t}/t, P/\ell^2)$$

By Lemma 11, since $t > t_{\ell}$, we have $\delta_{\ell} f w_{< t_{\ell}} < \delta f w_{< t}$, so

$$f\widehat{w}_{\leq t}^{t_{\ell}} \leq fw_{\leq t} + \delta_{\ell} fw_{\leq t_{\ell}} < fw_{\leq t} + \delta fw_{\leq t_{\ell}}.$$

We thus conclude

$$\forall t \in (0, t_{\text{max}}) : f \widehat{w}_{< t}^t < 2f \widehat{w}_{< t}^{t_{\ell}} < 2(1 + \delta) f w_{< t}, \tag{43}$$

This completes the proof of Theorem 6.

3.7 Upper bounding the upper bound

Theorem 6 uses the threshold upper bound $t_{\rm max}$ which is a value that depends on all the input weights, and these are not known if we only have a sample. As described in Section 3.2, to get confidence bounds out of Theorem 6, it suffices if we based on our sample can compute a probabilistic upper bound $\overline{t}_{\rm max}^{\tau}$ on the upper bound $t_{\rm max}$ such that

$$\overline{t}_{\max}^{\tau} \ge t_{\max} \tag{44}$$

with probability at least 1-P. For better confidence lower bounds, we want \overline{t}_{max} to be small.

Theorem 12 Define $\delta_{\leq \tau}^{\downarrow}$ and $\delta_{\leq \tau}^{\uparrow}$ such that $\wp(k_{\leq \tau}/(1-\delta_{\leq \tau}^{\downarrow}),\ \delta_{\leq \tau}^{\downarrow})) = P$ and $\wp(k_{\leq \tau}/(1+\delta_{\leq \tau}^{\uparrow}),\ \delta_{\leq \tau}^{\uparrow})) = P$. Let

$$\overline{t}_{\max}^{\tau} = \frac{1 + \delta_{\leq \tau}^{\uparrow}}{1 - \delta_{\leq \tau}^{\downarrow}} \tau.$$

Then $\Pr[\overline{t}_{\max}^{\tau} < t_{\max}] \leq P$

Proof Our first step will be to compute a probabilistic upper bound $\overline{w}_{\leq \tau}^{\tau}$ on $w_{\leq \tau}$ such that

$$\overline{w}_{\leq \tau}^{\tau} \ge w_{\leq \tau}.\tag{45}$$

with probability at least 1-P. We are going to define $\overline{w}_{\leq t}^t$ for any possible threshold t as a function of only the values of t and $k_{\leq t}$. We define $\overline{\mu}_{\leq t}^t = k_{\leq t}/(1-\delta_{\leq t}^\downarrow)$, and

$$\overline{w}_{\leq t} = t \, \overline{\mu}_{\leq t}^t = t \, k_{\leq t} / (1 - \delta_{\leq t}^{\downarrow}).$$

The lemma below states that $\overline{w}_{\leq \tau}^{\tau}$ does give us the desired probabilistic upper bound on $w_{\leq \tau}$.

Lemma 13 For the random priority threshold τ , the probability that $w_{\leq \tau} > \overline{w}_{\leq \tau}^{\tau}$ is at most P, so (45) holds true with probability at least 1 - P.

Proof For any given set of input weights consider a threshold t such that $\overline{w}_{\leq t}^t \leq w_{\leq t}$. We claim that the random priority threshold τ is expected no smaller than t. Note that $\tau < t$ if and only if $|I_{\leq t}^t| \leq k_{\leq t}$. Since $\overline{w}_{\leq t}^t \leq w_{\leq t}$, we have $\mathsf{E}[|I_{\leq t}^t|] = w_{\leq t}/t \geq \overline{\mu}_{\leq t}^t$. Moreover, $k_{\leq t} = (1 - \delta_{\leq t}^\downarrow)\overline{\mu}_{\leq t}^t$, so

$$\Pr[\tau \le t] = \Pr[|I_{\le t}^t| \le k_{\le t}] = \Pr[|I_{\le t}^t| \le (1 - \delta_{\le t}^{\downarrow})\overline{\mu}_{\le t}^t]$$
$$\le \wp(\overline{\mu}_{\le t}, \delta_{\le t}^{\downarrow}) = \wp(k_{\le t}/(1 - \delta_{\le t}^{\downarrow}), \delta_{\le t}^{\downarrow}) \le P.$$

Let t^+ be the maximal value such that $\overline{w}_{\leq t^+}^{t^+} \leq w_{\leq t^+}$. The probability that $\tau \leq t^+$ is at most P, and if $\tau > t^+$, then $\overline{w}_{\leq \tau}^{\tau} > w_{\leq \tau}$. If the maximal value t^+ does not exist, we define instead t^+ as the limit where $\overline{w}_{\leq t}^{t} \leq w_{\leq t}$ for $t < t^+$ while $\overline{w}_{\leq t^+}^{t^+} > w_{\leq t^+}$. The probability that $\tau < t^+$ is at most P, and if $\tau \geq t^+$, we again have $\overline{w}_{\leq \tau}^{\tau} > w_{\leq \tau}$.

Lemma 14 For any threshold t, $t_{\text{max}} \leq (1 + \delta^{\uparrow}_{< t}) w_{\leq t} / k_{\leq t}$.

Proof Consider any two thresholds t and T. Then $|I_{\leq t}^T|$ is the number of weights below t with threshold above T, and $\tau \geq T$ implies $|I_{\leq t}^T| \geq k_{\leq t}$. Let $T = (1 + \delta_{\leq t}^{\uparrow}) w_{\leq t} / k_{\leq t}$. Then $\mathsf{E}[|I_{\leq t}^T|] \leq w_{\leq t} / T = k_{\leq t} / (1 + \delta_{\leq t}^{\uparrow})$ with equality if $T \geq t$. It follows that

$$\Pr[\tau \ge T] \le \Pr[|I_{\le t}^T| \ge k_{\le t}] \le P.$$

Hence $T \geq t_{\text{max}}$.

Assuming (45), with $t = \tau$, we get

$$t_{\max} \leq (1+\delta_{\leq \tau}^{\uparrow}) w_{\leq \tau}/k_{\leq \tau} \leq (1+\delta_{\leq \tau}^{\uparrow}) \overline{w}_{\leq \tau}/k_{\leq \tau} = (1+\delta_{\leq \tau}^{\uparrow}) \tau/(1-\delta_{\leq t}^{\downarrow}).$$

By Lemma 13, (45) holds true with probability 1 - P. This completes the proof of Theorem 12.

3.8 What to trust, and what to expect

All our theorems about confidence intervals are trustworthy in the sense that they hold true for any set of input weights. We will now discuss what to expect if the input follows a reasonable distribution. As we shall formalize below, we expect a typical priority sample to consist of only a few large weights above the priority threshold, and a majority of small weights that are significantly smaller than the priority threshold. This will imply that our estimated threshold upper bound $\overline{t}_{\max}^{\tau}$ is very close to the priority threshold τ .

This view has consequences for what we would consider worth optimizing for in our confidence intervals, e.g., one could try getting better confidence intervals for cases where the sampled items have weight below but close to the threshold (information that is currently ignored, and not even contained in the sample), but since we do not expect many such items, we do not optimize for this case.

As our formal model, we assume that each weight w_i is drawn independently from a Pareto distributions, that is, for a positive real parameter $\alpha = \Omega(1)$, and any real $x \ge 1$, we have the survival function

$$\overline{F}(x) = \Pr[w_i \ge x] = 1/x^{\alpha}. \tag{46}$$

Then all weights are at least 1. For $\alpha \to \infty$, all weights are 1. As α decreases, we get more heavy weights. The mean is infinite for $\alpha \le 1$, and the variance is infinite for $\alpha \le 2$. The probability density function f is the derivative of $1 - \overline{F}(x)$, so

$$f(x) = \alpha/x^{\alpha+1}. (47)$$

We are going to use n such input weights as input to a priority sample of size k, where $1 \ll k \ll n$. The priority sampling events assigns priorities $q_i = w_i/r_i$, $r_i \in U(0,1)$ to each item, and in the analysis, we will study these weights and priorities relative to any given threshold t. For any such given threshold t, we assume that our error probability function \wp from (13) holds for the number of samples i, $q_i > t$, for the combined event where we first assign the weights w_i independently and second assign the h_i and hence the q_i based on a hash of each i.

We assume that the number of priority samples k is so large that for some small error $\varepsilon = o(1)$ and target error probability P, we have

$$\wp(\Omega(k), O(\varepsilon)) = o(P). \tag{48}$$

The basic idea is that the number k of samples is so large, that we do not expect significant errors for the sample as a whole. However, we might still have significant errors in estimation of small subsets. More precisely, our analysis will imply the following result.

Theorem 15 Let t_k be the threshold leading to an expected number of k samples, and let τ be the actual priority threshold, that is, the (k+1)th largest priority when all random choices are made. Then, with probability 1 - o(P), we have $\tau = (1 \pm O(\varepsilon))t_k = (1 \pm o(1))t_k$. Moreover, we have $k \le \tau = \Omega(k)$ small weight samples, most of which are from weights below $\tau/2$.

Proof As a first simple observation, since all weights at least 1, the expected number of samples with threshold t is at least n/t. It follows that $k \ge n/t_k$, hence that

$$t_k \ge n/k = \omega(1). \tag{49}$$

The following analysis is for an arbitrary threshold t, not just $t = t_k$. We want to study the expected number of large weights $w_i > t$ that are sampled for sure, and the expected number of small weight samples

 $w_i \le t < q_i$. By linearity of expectation, this is n times the probability of these events for any given item i. By (46), $\Pr[w_i > t] = 1/t^{\alpha}$. Using the probability density function (47), we get

$$\Pr[w_i \le t < q_i] = \int_1^t f(x) \cdot x/t \, dx \tag{50}$$

$$= \alpha/t \cdot \int_{1}^{t} 1/x^{\alpha} \, dx \tag{51}$$

$$=\frac{\alpha/t}{1-\alpha} \left[x^{1-\alpha} \right]_1^t \tag{52}$$

$$=\frac{\alpha}{1-\alpha}(t^{1-\alpha}-1)/t\tag{53}$$

$$=\frac{\alpha}{1-\alpha}(1/t^{\alpha}-1/t)\tag{54}$$

This should be compared with the probability of a large weight sample $w_i \ge t$ which was $1/t^{\alpha}$. For $\alpha < 1$, the low weight sample probability is $\Omega(1/t^{\alpha})$, and for $\alpha > 1/2$ we start expecting more low weight samples than large weights. For $\alpha > 1$, the sampled weights dominate in that we expect more than 1/t if them.

Above we assumed $\alpha \neq 1$. For $\alpha = 1$, continuing from (51), we get

$$Pr[w_i \le t, q_i > t] = 1/t \cdot [\ln x]_1^t$$
$$= 1/t \cdot \ln t$$

which means that the small weight samples are dominant by a factor of $\ln t$ for $\alpha = 1$. For continuity, it is easily verified that (54) also converges to $(\ln t)/t$ for $\alpha \to 1$. For simplicity, we assume below that $\alpha \neq 1$.

The total expected number of samples is

$$s_{\alpha}(t) = \frac{1/t^{\alpha} - \alpha/t}{1 - \alpha}$$

By definition, $s_{\alpha}(t_k)=k$. Define $t_k^+>t_k$ such that $s_{\alpha}(t_k^+)=k/(1+\varepsilon)$ where $\varepsilon=o(1)$ is the error from (48). To get priority threshold $\tau\geq t_k^+$, we need at least k+1 samples t_k^+ . By (48), this happens with probability o(P). If $\alpha<1$, then $t_k^+<(1+\varepsilon)^{1/\alpha}\,t_k$, and if $\alpha>0$, $t_k^+<(1+\varepsilon)t_k$. Since $\alpha=\Omega(1)$, we conclude in both cases that

$$t_k^+ < (1 + O(\varepsilon))t_k$$
.

A symmetric argument shows that $\tau \geq (1 - O(\varepsilon))t_k$ with probability 1 - o(P). Thus $\tau = (1 \pm O(\varepsilon))t_k$ with probability 1 - o(P).

Next we need to argue that with probability 1-o(P), the number $k_{\leq \tau}$ of sampled small weights in the priority sample is $\Omega(k)$. We may assume that $\tau \leq t_k^+$, so $k_{\leq \tau}$ is at least as big as the number of sampled small weights with threshold t_k^+ . By definition of t_k^+ , the expected number of samples with threshold t_k^+ is $k/(1+\varepsilon)$, and we know for any given threshold, that the expected number of small weight samples is at least a constant fraction of the expected total, so we expect $\Omega(k)$ small weight samples with threshold t_k^+ . By (48), this implies that their actual number is $\Omega(k)$ with probability 1-o(P). Thus $k_{\leq \tau}=\Omega(k)$ with probability 1-o(P).

Finally, among the sampled small weights $w_i \le \tau < q_i$, we want to see what fraction is below $\tau/2$. As usual, in our analysis, we first consider given thresholds rather than the variable priority threshold. Generally, for given thresholds $t_0 \le t_1$, and any given value of $w_i \le t_0$,

$$\Pr[q_i > t_1 | w_i] = \Pr[q_i > t_0 | w_i] \cdot t_0 / t_1.$$

Hence

$$\Pr[w_i \le t_0, q_i > t_1] = \Pr[w_i \le t_0 < q_0] \cdot t_0 / t_1 = \frac{\alpha}{1 - \alpha} (t_0^{1 - \alpha} - 1) / t_1.$$

For $\alpha = \Omega(1)$ and $t_1 \ge t_0 = \omega(1)$, we get

$$\Pr[w_i \le t_0, q_i > t_1] > (t_0/t_1)^{1-\Omega(1)} \cdot \Pr[w_i \le t_1, q_i > t_1].$$

We know that with probability 1-o(P), that $t_k^- \le \tau \le t_k^+$ where $t_k^+ = (1+o(1)t_k^-)$. A weight w_i is in the priority sample if $q_i > \tau$. We say weight w_i is over-sampled if $q_i \ge t_k^-$ and under-sampled if $q_i \ge t_k^+$. The over-sampled weights $w_i \le t_k^+$ include all weights $w_i \le \tau$ in the priority sample, and the under-sampled weights $w_i \le t_k^-/2$ are all included among the weights $w_i \le \tau/2$ in the priority sample.

The expected number of over-sampled weights $w_i \leq t_k^+$ is

$$n \cdot \Pr[w_i \le t_k^+, q_i > t_k^-] \le n \cdot \frac{\alpha}{1 - \alpha} ((t_k^+)^{1 - \alpha} - 1) / t_k^-.$$
 (55)

while the expected number of under-sampled weights $w_i \leq t_k^-/2$ is

$$n \cdot \Pr[w_i \le t_k^+, q_i > t_k^-] = n \cdot \frac{\alpha}{1 - \alpha} ((t_k^-/2)^{1 - \alpha} - 1) / t_k^+.$$
 (56)

With $\alpha = \Omega(1)$, (56) is within a factor $1/2 + \Omega(1)$ of (55). Moreover, from our analysis of $k_{\leq \tau}$, we know that (55) and hence (56) is $\Omega(k)$. It follows from (48) that with probability 1 - o(P), the expected bounds (55) and (56) end up both satisfied within a factor $1 \pm o(1)$. Then, among the priority sampled small weights $w_i \leq \tau$, at least half have weight below $\tau/2$.

We now return to our confidence bounds for large errors. By Theorem 15, with probability 1-o(P), we get $k_{\leq \tau} = \Omega(k)$. Hence by (48), we get $\delta^{\uparrow}_{\leq \tau}$, $\delta^{\downarrow}_{\leq \tau} = O(\varepsilon) = o(1)$ in Theorem 12, so

$$\overline{t}_{\max}^{\tau} = (1 + o(1))\tau$$

and then $\overline{\ell} = 1$ in (22).

3.9 Histogram similarity

We will now discuss estimators for the similarity of weighted sets. First consider the simple case where each key has a unique weight. The similarity is then just the total weight of the intersection divided by the weight of the union, and we estimate these two quantities independently.

As in the bottom-k sample for unweighted items, we note that given the size-k priority sample of two sets A and B, we can easily construct the size-k priority sample of their union, and identify which of these samples come from the intersection. Our analysis for subset sums now applies directly.

In the case of histogram similarity, it is natural to allow the same item to have different weights in different sets. More specifically, allowing zero weights, every possible item has a weight in each set. For the similarity we take the sum of the minimum weight for each item, and divide it by the sum of the maximum weight for each item. Formally, we are considering two sets A and B. Item i has weight w_i^A in A and weight w_i^B in B. Let $w_i^{\max} = \max\{w_i^A, w_i^B\}$ and $w_i^{\min} = \min\{w_i^A, w_i^B\}$. The histogram similarity is $w^{\min}/w^{\max} = (\sum_i w_i^{\min})/(\sum_i w_i^{\max})$.

This would seem a perfect application of our fractional subsets with $w_i = w_i^{\text{max}}$ and $f_i = w_i^{\text{min}}/w_i^{\text{max}}$. The issue is as follows. From our priority samples over the w_i^A and w_i^B , we can easily get the priority

sample for the $w_i = w_i^{\text{max}}$. However, for the items i sampled, we would typically not have a sample with w_i^{min} , and then we cannot compute f_i .

Our solution is to keep the instances of an item i in A and B separate as twins i^A and i^B with priorities $q_i^A = w_i^A/h_i$ and $q_i^B = w_i^B/h_i$. Note that it is the same hash value h_i we use to determine these two priorities. If $w_i^A = w_i^B$, we get $q_i^A = q_i^B$, and then we break the tie in favor of i^A . The priority sample for the union $A \cup B$ consists of the split items with the k highest priorities, and the priority threshold τ is the k+1 biggest among all priorities. Estimation is done as usual: for $C \in \{A, B\}$, if i^C is sampled, $\widehat{w}_i^C = \max\{w_i^C, \tau\}$. The important point here is the interpretation of the results. If $w_i^A \ge w_i^B$, then the priority of i^A is higher than that of i^B . Thus, in our sample, when we see an item i^C , $C \in \{A, B\}$, we count it for the union \widehat{w}^{\max} if it is not preceded by its twin; otherwise we count it for the intersection \widehat{w}^{\min} .

The resulting estimators \widehat{w}^{\min} and \widehat{w}^{\max} will no longer be unbiased even with truly random hashing. To see this, note that with sample size k=1, we always get $\widehat{w}^{\min}=0$. However, for our concentration bounds, we only lose a constant factor. The point is that the current analysis is using union bounds over threshold sampling events, using the fact that each hash value h_i contributes at most 1 to the number of items with priorities above a given threshold t. Now h_i affects at most 2 twins, but this is OK since all we really need is that the contribution of each random variable is bounded by a constant. The only effect on Theorem 5 is that we replace the relative error $6\delta(fw_{<\tau}/(3\tau), P)$ with $6\delta(fw_{<\tau}/(6\tau), P)$.

References

- [1] R. Adler, R. Feldman, and M. Taqqu. A Practical Guide to Heavy Tails. Birkhauser, 1998.
- [2] Y. Bachrach, R. Herbrich, and E. Porat. Sketching algorithms for approximating rank correlations in collaborative filtering systems. In *Proc. 16th SPIRE*, pages 344–352, 2009.
- [3] Y. Bachrach and E. Porat. Fast pseudo-random fingerprints. CoRR, abs/1009.5791, 2010.
- [4] Y. Bachrach, E. Porat, and J. S. Rosenschein. Sketching techniques for collaborative filtering. In *Proc.* 21st IJCAI, pages 2016–2021, 2009.
- [5] Z. Bar-Yossef, T. S. Jayram, R. Kumar, D. Sivakumar, and L. Trevisan. Counting distinct elements in a data stream. In *International Workshop on Randomization and Approximation Techniques in Computer Science (RANDOM)*, pages 1–10, 2002.
- [6] A. Z. Broder. On the resemblance and containment of documents. In Proc. Compression and Complexity of Sequences (SEQUENCES), pages 21–29, 1997.
- [7] A. Z. Broder. Identifying and filtering near-duplicate documents. In *Proc. 11th CPM*, pages 1–10, 2000.
- [8] A. Z. Broder, M. Charikar, A. M. Frieze, and M. Mitzenmacher. Min-wise independent permutations. *Journal of Computer and System Sciences*, 60(3):630–659, 2000. Announced at STOC'98.
- [9] A. Z. Broder, S. C. Glassman, M. S. Manasse, and G. Zweig. Syntactic clustering of the web. *Computer Networks*, 29:1157–1166, 1997.
- [10] E. Cohen, M. Datar, S. Fujiwara, A. Gionis, P. Indyk, R. Motwani, J. D. Ullman, and C. Yang. Finding interesting associations without support pruning. *IEEE Trans. Knowl. Data Eng.*, 13(1):64–78, 2001.

- [11] E. Cohen, N. Duffield, H. Kaplan, C. Lund, and M. Thorup. Efficient stream sampling for variance-optimal estimation of subset sums. *SIAM Journal on Computing*, 40(5):1402–1431, 2011. Announced at SODA'09.
- [12] E. Cohen and H. Kaplan. Summarizing data using bottom-k sketches. In *Proc. 26th PODC*, pages 225–234, 2007.
- [13] M. Datar and S. Muthukrishnan. Estimating rarity and similarity over data stream windows. In *Proc.* 10th ESA, pages 323–334, 2002.
- [14] M. Dietzfelbinger. Universal hashing and k-wise independent random variables via integer arithmetic without primes. In *Proc. 13th STACS*, pages 569–580, 1996.
- [15] M. Dietzfelbinger, J. Gil, Y. Matias, and N. Pippenger. Polynomial hash functions are reliable (extended abstract). In *Proc. 19th ICALP*, pages 235–246, 1992.
- [16] N. Duffield, C. Lund, and M. Thorup. Charging from sampled network usage. In *Proc. 1st ACM SIGCOMM Internet Measurement Workshop*, pages 245–256, 2001.
- [17] N. Duffield, C. Lund, and M. Thorup. Learn more, sample less: control of volume and variance in network measurements. *IEEE Transactions on Information Theory*, 51(5):1756–1775, 2005.
- [18] N. Duffield, C. Lund, and M. Thorup. Priority sampling for estimation of arbitrary subset sums. *J. ACM*, 54(6):Article 32, 2007. Announced at SIGMETRICS'04.
- [19] G. Feigenblat, E. Porat, and A. Shiftan. Exponential space improvement for minwise based algorithms. In *Proc. FSTTCS*, pages 70–85, 2012.
- [20] J. M. Hellerstein, P. J. Haas, and H. J. Wang. Online aggregation. In *Proc. SIGMOD*, pages 171–182, 1997.
- [21] M. R. Henzinger. Finding near-duplicate web pages: a large-scale evaluation of algorithms. In *Proc.* 29th SIGIR, pages 284–291, 2006.
- [22] C. P. Kruskal, L. Rudolph, and M. Snir. A complexity theory of efficient parallel algorithms. *Theor. Comput. Sci.*, 71(1):95–132, 1990.
- [23] G. S. Manku, A. Jain, and A. D. Sarma. Detecting near-duplicates for web crawling. In *Proc. 16th WWW*, pages 141–150, 2007.
- [24] M. Mitzenmacher and S. P. Vadhan. Why simple hash functions work: exploiting the entropy in a data stream. In *Proc. 19th ACM/SIAM Symposium on Discrete Algorithms (SODA)*, pages 746–755, 2008.
- [25] D. Moore, V. Paxson, S. Savage, C. Shannon, S. Staniford, and N. Weaver. Inside the slammer worm. *IEEE Security and Privacy Magazine*, 1(4):33–39, 2003.
- [26] R. Motwani and P. Raghavan. Randomized algorithms. Cambridge University Press, 1995.
- [27] A. Pagh, R. Pagh, and M. Ružić. Linear probing with constant independence. *SIAM Journal on Computing*, 39(3):1107–1120, 2009. Announced at STOC'07.

- [28] M. Pătraşcu and M. Thorup. On the *k*-independence required by linear probing and minwise independence. In *Proc. 36th ICALP, Part I, LNCS 6198*, pages 715–726, 2010.
- [29] M. Pătrașcu and M. Thorup. Twisted tabulation hashing. In *Proc. 23nd SODA*, pages 209–228, 2013.
- [30] C.-E. Särndal, B. Swensson, and J. Wretman. *Model Assisted Survey Sampling*. Springer, 1992.
- [31] S. Schleimer, D. S. Wilkerson, and A. Aiken. Winnowing: Local algorithms for document fingerprinting. In *Proc. SIGMOD*, pages 76–85, 2003.
- [32] J. P. Schmidt, A. Siegel, and A. Srinivasan. Chernoff-Hoeffding bounds for applications with limited independence. *SIAM Journal on Discrete Mathematics*, 8(2):223–250, 1995. Announced at SODA'93.
- [33] M. Szegedy. The DLT priority sampling is essentially optimal. In *Proc. 38th STOC*, pages 150–158, 2006.
- [34] M. Thorup. Confidence intervals for priority sampling. In *Proc. SIGMETRICS*, pages 252–263, 2006.
- [35] M. Thorup. Bottom-k and priority sampling, set similarity and subset sums with minimal independence. In *Proc. 45th STOC*, pages 371–378, 2013.
- [36] M. Thorup and Y. Zhang. Tabulation-based 5-independent hashing with applications to linear probing and second moment estimation. *SIAM Journal on Computing*, 41(2):293–331, 2012. Announced at SODA'04 and ALENEX'10.
- [37] M. N. Wegman and L. Carter. New classes and applications of hash functions. *Journal of Computer and System Sciences*, 22(3):265–279, 1981. Announced at FOCS'79.
- [38] H. Yang and J. P. Callan. Near-duplicate detection by instance-level constrained clustering. In *Proc.* 29th SIGIR, pages 421–428, 2006.