

Streaming Algorithms

We now consider a slightly different computational model called the *data streaming* model. In this model we see elements going past in a “stream”, and we have very little space to store things. For example, we might be running a program on an Internet router with limited space, and the elements might be IP Addresses. We certainly don’t have space to store all the elements in the stream. The question is: which functions of the input stream can we compute with what amount of time and space? While we focus on space, similar questions can be asked for update times.

We denote the stream elements by

$$a_1, a_2, a_3, \dots, a_t, \dots$$

We assume each stream element is from alphabet U , and takes $b = \lceil \log_2 U \rceil$ bits to represent. For example, the elements might be 32-bit integers IP addresses. We imagine we are given some function, and we want to compute it continually, on every prefix of the stream. Let us denote $a_{[1:t]} = \langle a_1, a_2, \dots, a_t \rangle$. For example, if we have seen the integers:

$$3, 1, 17, 4, -9, 32, 101, 3, -722, 3, 900, 4, 32, \dots \quad (12.1)$$

1. Can we compute the sum of all the integers seen so far? I.e., $F(a_{[1:t]}) = \sum_{i=1}^t a_i$. We want the outputs to be

$$3, 4, 21, 25, 16, 48, 149, 152, -570, -567, 333, 337, 369, \dots$$

If we have seen T numbers so far, the sum is at most $T2^b$ and hence needs at most $O(b + \log T)$ space. So we can just keep a counter, and when a new element comes in, we add it to the counter.

2. How about the maximum of the elements so far? $F(a_{[1:t]}) = \max_{i=1}^t a_i$. Even easier. The outputs are:

$$3, 1, 17, 17, 17, 32, 101, 101, 101, 101, 900, 900, 900$$

We just need to store b bits.

3. The median? The outputs on the various prefixes of (12.1) now are

$$3, 1, 3, 3, 3, 3, 4, 3, \dots$$

And doing this with small space is a lot more tricky.

4. (“distinct elements”) Or the number of distinct numbers seen so far? We’d want to output:

$$1, 2, 3, 4, 5, 6, 7, 7, 8, 8, 9, 9, 9, \dots$$

5. (“heavy hitters”) Or the elements that have appeared most often so far? Hmm...

We can imagine the applications of the data-stream model. An Internet router might see a lot of packets whiz by, and may want to figure out which data connections are using the most space? Or how many different connections have been initiated since midnight? Or the median (or the 90th percentile) of the file sizes that have been transferred. Which IP connections are “elephants” (say the ones that have used more than 0.01% of our bandwidth)? Even if we are not working at “line speed”, but just looking over the server logs, we may not want to spend too much time to find out the answers, we may just want to read over the file in one quick pass and come up with an answer. Such an algorithm might also be cache-friendly. But how to do this?

Such a router might see tens of millions of packets per second.

Two of the recurring themes will be:

1. Approximate solutions: in several cases, it will be impossible to compute the function exactly using small space. Hence we’ll explore the trade-offs between approximation and space.
2. Hashing: this will be a very powerful technique.

12.1 Streams as Vectors, and Additions/Deletions

An important abstraction will be to view the stream as a vector (in high dimensional space). Since each element in the stream is an element of the universe U , we can imagine the stream at time t as a vector $\mathbf{x}^t \in \mathbb{Z}^{|U|}$. Here

$$\mathbf{x}^t = (x_1^t, x_2^t, \dots, x_{|U|}^t)$$

and x_i^t is the number of times the i^{th} element in U has been seen until time t . (Hence, $x_i^0 = 0$ for all $i \in U$.) When the next element comes in and it is element j , we increment x_j by 1.

This brings us an extension of the model: we could have another model where each element of the stream is either a new element, or an old element departing. Formally, each time we get an *update* a_i , it looks like (add, e) or (del, e) . We usually assume that for each element, the number of deletes we see for it is at most the number of adds we see — the running counts of each element is non-negative. As an example, suppose the stream looks like:

$$(\text{add}, A), (\text{add}, B), (\text{add}, A), (\text{del}, B), (\text{del}, A), (\text{add}, C), \dots$$

and if A is the first element of U , then the first coordinate x_1 of the vector \mathbf{x} would be $1, 1, 2, 2, 1, 1, \dots$. This vector notation allows us to formulate some of the problems more easily:

1. The total number of elements currently in the system is just $\|\mathbf{x}\|_1 := \sum_{i=1}^{|U|} x_i$. (This is easy.)
2. We might want to estimate the norms $\|\mathbf{x}\|_2, \|\mathbf{x}\|_p$ of the vector \mathbf{x} .
3. The number of distinct elements is the number of non-zero entries in \mathbf{x} is denoted by $\|\mathbf{x}\|_0$.

Let's consider the (non-trivial) problems one by one.

12.2 Computing Moments

Recall that \mathbf{x}^t was the vector of frequencies of elements seen so far. Several interesting problems can be posed as computing various norms of \mathbf{x}^t : in particular the Euclidean or 2-norm

$$\|\mathbf{x}^t\|_2 = \sqrt{\sum_{i=1}^{|U|} (x_i^t)^2},$$

and the 0-norm (which is not really a norm)

$$\|\mathbf{x}^t\|_0 := \text{number of non-zeroes in } \mathbf{x}^t.$$

Henceforth, we use the notation that $F_0 := \|\mathbf{x}^t\|_0$ is the number of non-zero entries in \mathbf{x} . For $p \geq 1$, we consider the *p-moment*, that is, the p^{th} -power of the p -norm:

$$F_p := \sum_{i=1}^{|U|} (x_i^t)^p. \tag{12.2}$$

We'll develop an algorithm to compute F_2 , and to compute F_0 ; we may see extensions from F_2 to F_p in the homeworks.

In data stream jargon, the addition-only model is called the *cash-register* model, whereas the model with both additions and deletions is called the *turnstile model*. I will not use this jargon.

12.2.1 Computing the Second Moment F_2

The “second moment” F_2 of the stream is often called the “surprise number” (since it captures how uneven the data is). This is also the *size of the self-join*. Clearly we can store the entire vector \mathbf{x} and compute F_2 , but that requires storing $|U|$ counts. Here’s an algorithm that uses much less space:

Pick a random hash function $h : U \rightarrow \{-1, +1\}$ from family H .
 Maintain counter C , which starts off at zero.
 On update $(add, i) \in U$, increment the counter $C \rightarrow C + h(i)$.
 On update $(delete, i) \in U$, decrement the counter $C \rightarrow C - h(i)$.
 On query about the value of F_2 , reply with C^2 .

This estimator was given by Noga Alon, Yossi Matias, and Mario Szegedy, in their Gödel-award winning paper on streaming computation.

This estimator is often called the “tug-of-war” estimator: the hash function randomly partitions the elements into two parties (those mapping to 1, and those to -1), and the counter keeps the difference between the sizes of the two parties.

Alon, Matias, Szegedy (2000)

12.2.2 Properties of the Hash Family

The choice of the hash family will be crucial: we want a small family so that we require only a small amount of space to store the hash function, but we want it to be rich enough for the subsequent analysis to go through.

Definition 12.1 (*k*-universal hash family). H is *k*-universal (also called *uniform* and *k*-wise independent) mapping universe U to some range R if all *distinct* elements $i_1, \dots, i_k \in U$ and for values $\alpha_1, \dots, \alpha_k \in R$,

$$\Pr_{h \leftarrow H} \left[\bigwedge_{j=1..k} (h(i_j) = \alpha_j) \right] = \frac{1}{|R|^k}. \quad (12.3)$$

In our application, we want the hash family to be 4-universal from U to the two-element range $R = \{-1, 1\}$. This means that for any element i ,

$$\Pr_{h \leftarrow H} [h(i) = 1] = \Pr_{h \leftarrow H} [h(i) = -1] = \frac{1}{2}.$$

Moreover, for four distinct elements i, j, k, l , their maps behave independently of each other, and hence

$$\begin{aligned} \mathbb{E}[h(i) \cdot h(j) \cdot h(k) \cdot h(l)] &= \mathbb{E}[h(i)] \cdot \mathbb{E}[h(j)] \cdot \mathbb{E}[h(k)] \cdot \mathbb{E}[h(l)]. \\ \mathbb{E}[h(i) \cdot h(j)] &= \mathbb{E}[h(i)] \cdot \mathbb{E}[h(j)]. \end{aligned}$$

We will discuss constructions of such hash families soon, but let us use them to analyze the tug-of-war estimator.

12.2.3 A Direct Analysis

Hence, having seen the stream that results in the frequency vector $\mathbf{x} \in \mathbb{Z}_{\geq 0}^{|U|}$, the counter will have the value

$$C := \sum_{i \in U} x_i h(i).$$

Remember, the resulting estimate is C^2 : so we need to show that $\mathbb{E}[C^2] = F_2$, and variance that is small enough that Chebyshev's inequality ensures we are correct with reasonable probability.

$$\begin{aligned} \mathbb{E}[C^2] &= \mathbb{E}\left[\sum_{i,j} (h(i)x_i \cdot h(j)x_j)\right] = \sum_{i,j} x_i x_j \mathbb{E}[(h(i) \cdot h(j))] \\ &= \sum_i x_i^2 \mathbb{E}[h(i) \cdot h(i)] + \sum_{i \neq j} \sum_{i,j} x_i x_j \mathbb{E}[h(i)] \cdot \mathbb{E}[h(j)] \\ &= \sum_i x_i^2 = F_2. \end{aligned}$$

So in expectation we are correct! Next, recall that the variance is defined as $\text{Var}(C^2) = \mathbb{E}[(C^2)^2] - \mathbb{E}[C^2]^2$:

$$\begin{aligned} \mathbb{E}[(C^2)^2] &= \mathbb{E}\left[\sum_{p,q,r,s} h(p)h(q)h(r)h(s)x_p x_q x_r x_s\right] = \\ &= \sum_p x_p^4 \mathbb{E}[h(p)^4] + 6 \sum_{p < q} x_p^2 x_q^2 \mathbb{E}[h(p)^2 h(q)^2] + \text{other terms} \\ &= \sum_p x_p^4 + 6 \sum_{p < q} x_p^2 x_q^2. \end{aligned}$$

This is because all the other terms have expectation zero. Why? The terms like $\mathbb{E}[h(p)h(q)h(r)h(s)]$ where p, q, r, s are all distinct, all become zero because of 4-universality. Terms like $\mathbb{E}[h(p)^2 h(r)h(s)]$ become zero for the same reason. It is only terms like $\mathbb{E}[h(p)^2 h(q)^2]$ and $\mathbb{E}[h(p)^4]$ that survive, and since $h(p) \in \{-1, 1\}$, they have expectation 1. So

$$\text{Var}(C^2) = \sum_p x_p^4 + 6 \sum_{p < q} x_p^2 x_q^2 - \left(\sum_p x_p^2\right)^2 = 4 \sum_{p < q} x_p^2 x_q^2 \leq 2\mathbb{E}[C^2]^2.$$

What does Chebyshev say then?

$$\Pr[|C^2 - \mathbb{E}[C^2]| > \varepsilon \mathbb{E}[C^2]] \leq \frac{\text{Var}(C^2)}{(\varepsilon \mathbb{E}[C^2])^2} \leq \frac{2}{\varepsilon^2}.$$

This is pretty pathetic: since ε is usually less than 1, the RHS usually more than 1.

12.2.4 Reduce the Variance by Repetition

The idea is the simplest one: if we have an estimator with mean μ and variance σ^2 , then taking the average of k independent copies of

this estimator has mean μ and variance σ^2/k . (Why? Summing the independent copies sums the variances and so increases it by k , but dividing by k reduces it by k^2 .)

So if we k such independent counters C_1, C_2, \dots, C_k , and return their average $\bar{C} = \frac{1}{k} \sum_i C_i$, we get

$$\Pr[|\bar{C}^2 - \mathbb{E}[\bar{C}^2]| > \epsilon \mathbb{E}[\bar{C}^2]] \leq \frac{\text{Var}(\bar{C}^2)}{(\epsilon \mathbb{E}[\bar{C}^2])^2} \leq \frac{2}{k\epsilon^2}.$$

Taking $k = \frac{2}{\epsilon^2\delta}$ independent counters gives a probability δ of error on any query. Each counter uses a 4-universal hash function, which requires $O(\log U)$ random bits to store.

12.2.5 Estimating the p -Moments

To fix, please skip. A bunch of students (Jason, Anshu, Aram) proposed that for the p^{th} -moment calculation we should use $2p$ -wise independent hash functions from U to R , where $R = \{1, \omega, \omega^2, \dots, \omega^{p-1}\}$, the p primitive roots of unity. Again, we set $C := \sum_{i \in U} x_i h(i)$, and return the real part of C^p as our estimate. This approach has been explored by Ganguly in [this paper](#). Some calculations (and elbow-grease) show that $\mathbb{E}[C^p] = F_p$, but it seems that naively $\text{Var}(C^p)$ tends to grow like F_2^p instead of F_k^p ; this leads to pretty bad bounds. Ganguly's paper gives some ways of controlling the variance.

BTW, there is a lower bound saying that any algorithm that outputs a 2-approximation for F_k requires at least $|U|^{1-2/k}$ bits of storage. Hence, while we just saw that for $k = 2$, we can get away with just $O(\log |U|)$ bits to get a $O(1)$ -estimate, for $k > 2$ things are much worse.

12.3 A Matrix View of our Estimator

Here's a equivalent way of looking at this estimator, that also relates it to the previous chapter and the JL Theorem. Recall that the stream can be viewed as representing a vector \mathbf{x} of size $|U|$, and $F_2 = \|\mathbf{x}\|^2$.

Take a matrix M of dimensions $k \times D$, where $D = |U|$: again, M is a "fat and short" matrix, since $k = O(\epsilon^{-2}\delta^{-1})$ is small and $D = |U|$ is huge. Pick k independent hash functions h_1, h_2, \dots, h_k from the 4-universal hash family, and use each one to fill a row of M :

$$M_{ij} := h_i(j).$$

The k counters C_1, C_2, \dots, C_k are now nothing other than the entries of the matrix-vector product

$$M \mathbf{x}.$$

The estimate $\bar{C}^2 = \left(\frac{1}{k} \sum_{i=1}^k C_i\right)^2$ is nothing but

$$\frac{1}{k} \|M\mathbf{x}\|_2^2.$$

This is completely analogous to the construction for JL: we've got a slightly taller matrix with $k = O(\epsilon^{-2} \delta^{-1})$ rows instead of $k = O(\epsilon^{-2} \log \delta^{-1})$ rows. However, the matrix entries are not fully independent (as in JL), just 4-wise independent. I.e., we need to store only $O(k \log D)$ bits and can generate any entry of M quickly, whereas the construction for JL stored all kD bits.

Let us record two properties of this construction:

Theorem 12.2 (Tug-of-War Sketch). *Take a $k \times D$ matrix S whose columns are 4-wise independent $\{\frac{1}{\sqrt{k}}, \frac{-1}{\sqrt{k}}\}^k$ -valued r.v.s. Then for $\mathbf{x}, \mathbf{y} \in \mathbb{R}^D$,*

1. $\mathbb{E}[\langle S\mathbf{x}, S\mathbf{y} \rangle] = \langle \mathbf{x}, \mathbf{y} \rangle.$
2. $\text{Var}(\langle S\mathbf{x}, S\mathbf{y} \rangle) = \frac{2}{k} \cdot \|\mathbf{x}\|_2^2 \|\mathbf{y}\|_2^2.$

The proofs is similar to that in §12.2.3; using $\mathbf{y} = \mathbf{x}$ gives us exactly the results from that section. Moreover, an analogous theorem can also be given in the JL construction, with fewer rows but with completely independent entries.

12.4 Application: Approximate Matrix Multiplication

Suppose we want to multiply square matrices $A, B \in \mathbb{R}^{n \times n}$, but want to solve the problem faster, at the expense of getting only an approximate solution $C \approx AB$. How should we measure the error? Requiring that the answer be close entry-wise to the actual answer is a hard problem. Let's aim for something weaker: we want the "aggregate error" to be small.

Formally, the *Frobenius norm* of matrix M is

$$\|M\|_F := \sqrt{\sum_{i,j} M_{ij}^2}.$$

Our guarantee for approximate matrix multiplication will be

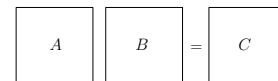
$$\|C - AB\|_F^2 \leq \text{small}.$$

Here's the idea: we want to do the matrix multiplication:

$$C = AB$$

Henceforth, we use $S = \frac{1}{\sqrt{k}}M$ to denote the "sketch" matrix.

It's as though we think of the matrix as just a vector and look at its Euclidean length.



This usually takes $O(n^3)$ time. Indeed, the ij^{th} entry of the product C is the dot-product of the i^{th} row $A_{i\star}$ of A with the j^{th} column $B_{\star j}$ of B , and the dot-product takes $O(n)$ time.

Suppose instead we use a “fat and short” $k \times n$ matrix S (for $k \ll n$), and calculate

$$\tilde{C} = AS^T SB.$$

By associativity of matrix multiplication, we could first compute (AS^T) and (SB) in times $O(n^2k)$, and then multiply the results in time $O(nk^2)$. Moreover, the matrix S from the previous section works pretty well, where we set $D = n$.

Indeed, entries of the error matrix $Y = C - \tilde{C}$ satisfy

$$\mathbb{E}[Y_{ij}] = 0$$

and

$$\mathbb{E}[Y_{ij}^2] = \text{Var}(Y_{ij}) + \mathbb{E}[Y_{ij}]^2 = \text{Var}(Y_{ij}) \leq \frac{2}{k} \|A_{i\star}\|_2^2 \|B_{\star j}\|_2^2.$$

So

$$\begin{aligned} \mathbb{E}[\|AB - AS^T SB\|_F^2] &= \mathbb{E}[\sum_{ij} Y_{ij}^2] = \sum_{ij} \mathbb{E}[Y_{ij}^2] = \frac{2}{k} \sum_{ij} \|A_{i\star}\|_2^2 \|B_{\star j}\|_2^2 \\ &= \frac{2}{k} \|A\|_F^2 \|B\|_F^2. \end{aligned}$$

Finally, setting $k = \frac{2}{\varepsilon^2 \delta}$ and using Markov’s inequality, we can say that for any fixed $\varepsilon > 0$, we can compute an approximate matrix product $C := AS^T SB$ such that

$$\Pr \left[\|AB - C\|_F \leq \varepsilon \cdot \|A\|_F \|B\|_F \right] \geq 1 - \delta,$$

in time $O(\frac{n^2}{\varepsilon^2 \delta})$. (If we want to make δ very small, at the expense of picking more independent random bits in the sketching matrix S , we can use the JL matrices instead. Details will appear in a homework.) Finally, if the matrices A, B are sparse and contains only $\ll n^2$ entries, the time can be made to depend on $nnz(A, B)$.

The approximate matrix product question has been considered often, e.g., by Edith Cohen and David Lewis using a random-walks approach. The algorithm we present is due to Tamás Sarlós; his paper gives better results, as well as extensions to computing SVDs faster. Better bounds have subsequently been given by Clarkson and Woodruff. [More recent refs too.](#)

The intuition is that $S^T S$ is an almost-identity matrix, it has 1 on the diagonals and at most ε everywhere else. And hence it gives only a small error. Of course, we don’t multiply out $S^T S$, but instead compute AS^T and SB , and then multiply the smaller matrices.



The squared Frobenius norm of a matrix is the sum of squared Euclidean lengths of the columns, or of the rows.

Cohen and Lewis (1999)

12.5 Optional: Computing the Number of Distinct Elements

Our last example today will be to compute F_0 , the number of distinct elements seen in the data stream, but in the addition-only model, with no deletions. (We’ll see another approach in a HW.)

12.5.1 A Simple Lower Bound

Of course, if we store x explicitly (using $|U|$ space), we can trivially solve this problem exactly. Or we could store the (at most) t elements seen so far, again we could give an exact answer. And indeed, we cannot do much better if we want no errors. Here’s a proof sketch for deterministic algorithms (one can extend this to randomized algorithms with some more work).

Lemma 12.3 (A Lower Bound). *Suppose a deterministic algorithm correctly reports the number of distinct elements for each sequence of length at most N . Suppose $N \leq 2|U|$. Then it must use at least $\Omega(N)$ bits of space.*

Proof. Consider the situation where first we send in some subset S of $N - 1$ elements distinct elements of U . Look at the information stored by the algorithm. We claim that we should be able to use this information to identify exactly which of the $\binom{|U|}{N-1}$ subsets of U we have seen so far. This would require

$$\log_2 \binom{|U|}{N-1} \geq (N-1)(\log_2 |U| - \log_2(N-1)) = \Omega(N)$$

bits of memory.¹

OK, so why should we be able to uniquely identify the set of elements until time $N - 1$? For a contradiction, suppose we could not tell whether we’d seen S_1 or S_2 after $N - 1$ elements had come in. Pick any element $e \in S_1 \setminus S_2$. Now if we gave the algorithm e as the N^{th} element, the number of distinct elements seen would be N if we’d already seen S_2 , and $N - 1$ if we’d seen S_1 . But the algorithm could not distinguish between the two cases, and would return the same answer. It would be incorrect in one of the two cases. This contradicts the claim that the algorithm always correctly reports the number of distinct elements on streams of length N . \square

¹ We used the approximation that $\binom{m}{k} \geq (\frac{m}{k})^k$, and hence $\log_2 \binom{m}{k} \geq k(\log_2 m - \log_2 k)$.

OK, so we need an approximation if we want to use little space. Let’s use some hashing magic.

12.5.2 The Intuition

Suppose there are $d = \|x\|_0$ distinct elements. If we randomly map d distinct elements onto the line $[0, 1]$, we expect to see the smallest mapped value at location $\approx \frac{1}{d}$. (I am assuming that we map these elements *consistently*, so that multiple copies of an element go to the same place.) So if the smallest value is δ , one estimator for the number of elements is $1/\delta$.

This is the essential idea. To make this work (and analyze it), we change it slightly: The variance of the above estimator is large. By the

same argument, for any integer s we expect the s^{th} smallest mapped value at $\frac{s}{d}$. We use a larger value of s to reduce the variance.

12.5.3 The Algorithm

Assume we have a hash family H with hash functions $h : U \rightarrow [M]$. (We'll soon figure out the precise properties we'll want from this hash family.) We will later fix the value of the parameter s to be some large constant. Here's the algorithm:

Pick a hash function h randomly from H .

If query comes in at time t

Consider the hash values $h(a_1), h(a_2), \dots, h(a_t)$ seen so far.

Let L_t be the s^{th} smallest distinct hash value $h(a_i)$ in this set.

Output the estimate $D_t = \frac{M \cdot s}{L_t}$.

The crucial observation is: it does not matter if we see an element e once or multiple times — the algorithm will behave the same, since the output depends on what *distinct* elements we've seen so far. Also, maintaining the s^{th} smallest element can be done by remembering at most s elements. (So we want to make s small.)

How does this help? As a thought experiment, if we had d distinct darts and threw them in the continuous interval $[0, M]$, we would expect the location of the s^{th} smallest dart to be about $\frac{s \cdot M}{d}$. So if the s^{th} smallest dart was at location ℓ in the interval $[0, M]$, we would be tempted to equate $\ell = \frac{s \cdot M}{d}$ and hence guessing $d = \frac{s \cdot M}{\ell}$ would be a good move. Which is precisely why we used the estimate

$$D_t = \frac{M \cdot s}{L_t}.$$

Of course, all this is in expectation—the following theorem argues that this estimate is good with reasonable probability.

Theorem 12.4. *Consider some time t . If H is a uniform 2-universal hash family mapping $U \rightarrow [M]$, and M is large enough, then both the following guarantees hold:*

$$\Pr[D_t > 2 \|\mathbf{x}^t\|_0] \leq \frac{3}{s}, \text{ and} \tag{12.4}$$

$$\Pr[D_t < \frac{\|\mathbf{x}^t\|_0}{2}] \leq \frac{3}{s}. \tag{12.5}$$

We will prove this in the next section. First, some observations. Firstly, we now use the stronger assumption that that the hash family *2-universal*; recall the definition from Section 12.2.2. Next, setting $s = 8$ means that the estimate D_t lies within $[\frac{\|\mathbf{x}^t\|_0}{2}, 2\|\mathbf{x}^t\|_0]$ with probability at least $1 - (1/4 + 1/4) = 1/2$. (And we can boost the

success probability by repetitions.) Secondly, we will see that the estimation error of a factor of 2 can be made $(1 + \varepsilon)$ by changing the parameters s and k .

12.5.4 Proof of Theorem 12.4

Now for the proof of the theorem. We'll prove bound (12.5), the other bound (12.4) is proved identically. Some shorter notation may help. Let $d := \|\mathbf{x}^t\|_0$. Let these d distinct elements be $T = \{e_1, e_2, \dots, e_d\} \subseteq U$.

The random variable L_t is the s^{th} smallest distinct hash value seen until time t . Our estimate is $\frac{sM}{L_t}$, and we want this to be at least $d/2$. So we want L_t to be at most $\frac{2sM}{d}$. In other words,

$$\Pr[\text{estimate too low}] = \Pr[D_t < d/2] = \Pr[L_t > \frac{2sM}{d}].$$

Recall T is the set of all d ($= \|\mathbf{x}^t\|_0$) distinct elements in U that have appeared so far. How many of these elements in T hashed to values greater than $2sM/d$? The event that $L_t > 2sM/d$ (which is what we want to bound the probability of) is the same as saying that fewer than s of the elements in T hashed to values smaller than $2sM/d$. For each $i = 1, 2, \dots, d$, define the indicator

$$X_i = \begin{cases} 1 & \text{if } h(e_i) \leq 2sM/d \\ 0 & \text{otherwise} \end{cases} \quad (12.6)$$

Then $X = \sum_{i=1}^d X_i$ is the number of elements seen that hash to values below $2sM/d$. By the discussion above, we get that

$$\Pr\left[L_t < \frac{2sM}{d}\right] \leq \Pr[X < s].$$

We will now estimate the RHS.

Next, what is the chance that $X_i = 1$? The hash $h(e_i)$ takes on each of the M integer values with equal probability, so

$$\Pr[X_i = 1] = \frac{\lfloor sM/2d \rfloor}{M} \geq \frac{s}{2d} - \frac{1}{M}. \quad (12.7)$$

By linearity of expectations,

$$\mathbb{E}[X] = E\left[\sum_{i=1}^d X_i\right] = \sum_{i=1}^d E[X_i] = \sum_{i=1}^d \Pr[X_i = 1] \geq d \cdot \left(\frac{s}{2d} - \frac{1}{M}\right) = \left(\frac{s}{2} - \frac{d}{M}\right).$$

Let's imagine we set M large enough so that d/M is, say, at most $\frac{s}{100}$. Which means

$$\mathbb{E}[X] \geq \left(\frac{s}{2} - \frac{s}{100}\right) = \frac{49s}{100}.$$

So by Markov's inequality,

$$\Pr[X > s] = \Pr\left[X > \frac{100}{49}\mathbb{E}[X]\right] \leq \frac{49}{100}.$$

Good? Well, not so good. We wanted a probability of failure to be smaller than $2/s$, we got it to be slightly less than $1/2$. Good try, but no cigar.

12.5.5 Enter Chebyshev

Recall that $\text{Var}(\sum_i Z_i) = \sum_i \text{Var}(Z_i)$ for pairwise-independent random variables Z_i . (Why?) Also, if Z_i is a $\{0, 1\}$ random variable, $\text{Var}(Z_i) \leq \mathbb{E}[Z_i]$. (Why?) Applying these to our random variables $X = \sum_i X_i$, we get

$$\text{Var}(X) = \sum_i \text{Var}(X_i) \leq \sum_i \mathbb{E}[X_i] = E(X).$$

(The first inequality used that the X_i were pairwise independent, since the hash function was 2-universal.) Is this variance “low” enough? Plugging into Chebyshev's inequality, we get:

$$\Pr[X > s] = \Pr\left[X > \frac{100}{49}\mu_X\right] \leq \Pr\left[|X - \mu_X| > \frac{50}{49}\mu_X\right] \leq \frac{\sigma_X^2}{(50/49)^2\mu_X^2} \leq \frac{1}{(50/49)^2\mu_X} \leq \frac{3}{s}.$$

Which is precisely what we want for the bound (12.4). The proof for the bound (12.5) is similar and left as an exercise.

12.5.6 Final Bookkeeping

Excellent. We have a hashing-based data structure that answers “number of distinct elements seen so far” queries, such that each answer is within a multiplicative factor of 2 of the actual value $\|\mathbf{x}^t\|_0$, with small error probability.

Let's see how much space we actually used. Recall that for failure probability $1/2$, we could set $s = 12$, say. And the space to store the s smallest hash values seen so far is $O(s \lg M)$ bits. For the hash functions themselves, the standard constructions use $O((\lg M) + (\lg U))$ bits per hash function. So the total space used for the entire data structure is

$$O(\log M) + (\lg U) \text{ bits.}$$

What is M ? Recall we needed to M large enough so that $d/M \leq s/100$. Since $d \leq |U|$, the total number of elements in the universe, set $M = \Theta(U)$. Now the total number of bits stored is

$$O(\log U).$$

And the probability of our estimate D_t being within a factor of 2 of the correct answer $\|\mathbf{x}^t\|_0$ is at least $1/2$.

If we want the estimate to be at most $\frac{\|\mathbf{x}^t\|_0}{(1+\epsilon)}$, then we would want to bound $\Pr[X < \frac{\mathbb{E}[X]}{(1+\epsilon)}]$. Similar calculations should give this to be at most $\frac{3}{\epsilon^2 s}$, as long as M was large enough. In that case we would set $s = O(1/\epsilon^2)$ to get some non-trivial guarantees.