

Lecture 10: Count-Min Sketch and HyperLogLog Counting

Sam McCauley

October 21, 2021

1 Goal for Today

Computers are increasingly often tasked with dealing with *extremely* large amounts of data. This is particularly true on the internet: every user, every action, every message is a piece of data that needs to be considered and processed very quickly.

In the last class, we learned about Bloom filters, which allow us to compress data. But you may have noticed that there were serious limits to this compression: generally, each data item is compressed to the size of a byte. At absolute minimum, a Bloom filter requires a bit or two of space for every stored item. That isn't going to work when the number of items we want to deal with number in the trillions or even quadrillions.¹

What we want is a much more aggressive form of compression. In fact, we want it to be so aggressive that we can have a guaranteed size of our data structure that's only *logarithmic* in the size of our data—or even better than that.

How could such a thing be possible? And, if our data structure is that small, won't we be losing almost all of the information we're looking for? Surprisingly, a simple and carefully built data structure can give us extremely valuable information about the data, even under such draconian space requirements.

2 What is a Stream?

A *stream* is a very long sequence of data which is given to the data structure one item at a time. Once an item is given to the data structure, the data structure updates its representation and the item is never seen again. You can imagine this as sitting next to a large stream. You are trying to sample small pieces of the stream for analysis. But once something in the stream goes past you, that's it—if you didn't sample it, you won't see it again.

Formalism A stream is a list of items x_1, \dots, x_N (a stream may contain duplicates, where $x_i = x_j$ but $i \neq j$). As in filters, we assume that items are of the same type for the sake of hashing: for some universe U , $x_i \in U$ for all i , and all queries q must also satisfy $q \in U$. Each item in the stream is presented to the data structure one at a time, in order. We will assume that an **insert** function is called on x_1 , then x_2 , and so on until **insert** is called on x_N .

The data structure is required to have small space, usually $O(\log N)$ or smaller (sometimes $O(1)$). Today we will be looking at an example of each: the Count-Min Sketch requires $O(\log N)$ bits of space for constant error, whereas the HyperLogLog data structure requires $O(\log \log N)$ bits.

After all items have been inserted a single time into the data structure, we can make queries to learn information about the stream. We will focus on two queries today:

- For a query q , how many times did q appear in the stream? In set notation, this asks us to determine $|\{i \in \{1, \dots, N\} \mid x_i = q\}|$.

¹The Brazil Internet Exchange alone processes over 7 trillion bits, on average, every second.

- How many distinct items appeared in the stream overall? In set notation, this asks us to determine $|\{x_i \mid i \in \{1, \dots, N\}\}|$.

Discussion Streaming algorithms have been a major area of research throughout the 2000s and 2010s. There exist entire courses on the topic.² Streaming algorithms are an excellent example of how randomized algorithms can make otherwise-impossible problems tractable.

Why are streaming algorithms so useful? One reason is that data streams are, in fact, fairly common nowadays: packets being transmitted through a node in a network (say, a router on the internet) essentially form a data stream. Any entity interested in monitoring that traffic needs to solve a problem on a large data stream.

However, a further motivation comes (again!) from cache efficiency. As we discussed in class, the most cache-efficient way to go through data is to scan it.³ So if you have a massive amount of data stored somewhere that you want to analyze, the cheapest way to approach it is to scan through the data once, keeping track of statistics in a small-space data structure as you go. This is, again, a problem on a large data stream.⁴

3 Count Min Sketch

Our first data structure keeps an impressively accurate count of *every* item in a data stream while using extremely small space.

The Count-Min Sketch is one of the best-known streaming data structures. It looks somewhat like a Bloom filter, but with substantially different parameters. We will use two error terms ε and δ in describing the Count Min Sketch, as well as the length of the stream N . We assume every hash h_i outputs a number in $\{0, \dots, \lceil e/\varepsilon \rceil - 1\}$.

Algorithm 1 Insert x_i

```

for  $j = 0$  to  $\lceil 1/\delta \rceil - 1$  do
     $T[j][h_j(x_i)] = T[j][h_j(x_i)] + 1$ 
end for

```

Algorithm 2 Query q

```

min = 0
for  $j = 0$  to  $\lceil 1/\delta \rceil - 1$  do
    if min <  $T[j][h_j(q)]$  then
        min  $\leftarrow T[j][h_j(q)]$ 
    end if
end for
return min

```

Structure The Count Min Sketch (CMS) consists of a two-dimensional table T with $\lceil \ln 1/\delta \rceil$ rows. Each row consists of $\lceil e/\varepsilon \rceil$ entries. Each entry must be of size at least $\log N$ bits.⁵

A Count Min Sketch has $\lceil \ln 1/\delta \rceil$ hash functions, one for each row of the table. In this way, the CMS essentially consists of $\lceil \ln 1/\delta \rceil$ independent hash tables stored in one location.

To initialize the structure, we set all entries in T to 0.

Insert To insert an item x_i , we iterate through each row of the CMS. For each row, we hash x_i , and increment the counter stored at that hash location.

Query To query q , we iterate through each row of the CMS. For each row, we hash q , and keep track of the value stored at that location. We return the minimum such value found. This

value is an estimate of how many times q occurred during the stream: that is to say, it estimates $|\{i \mid x_i = q\}|$.

²For example, this 2007 course at MIT: <http://stellar.mit.edu/S/course/6/fa07/6.895/materials.html>.

³Scans also tend to avoid branch mispredictions and other practical inefficiencies.

⁴You may notice that if our goal cache-efficiency, there's an opportunity for a tradeoff: I can scan the data multiple times, resulting in an increase in cache misses, but with a second chance to see important data items. Can seeing a stream several times lead to more efficient methods? In short, the answer is frequently yes. This is called the "multi pass" streaming model. While this model is fairly popular we will not be looking at it in this class.

⁵In practice, one would just choose 16, 32, or 64 bit entries depending on the size of the stream.

3.1 Bounds

A Count Min Sketch gives the following guarantee. Let $q \in U$ be any query, let o_q be the value given by the query algorithm on q , and let \hat{o}_q be the true number of occurrences of q in the stream (i.e. $\hat{o}_q = |\{i \mid x_i = q\}|$). Then we have the following two guarantees:

1. $\hat{o}_q \leq o_q$, and
2. with probability at least $1 - \delta$, $o_q \leq \hat{o}_q + \varepsilon N$.

Multiplying the number of rows, number of columns, and size of each entry, we can see that the Count Min Sketch requires $\lceil e/\varepsilon \rceil \lceil \ln 1/\delta \rceil \lceil \log N \rceil$ bits.

3.2 Analysis

The insertion and query algorithms of the count-min sketch immediately give us the lower bound $\hat{o}_q \leq o_q$. How can we upper bound o_q ?

You'll prove the following useful lemma on your assignment.

Lemma 1. *Let X be a positive random variable (i.e. a random variable where all possible values are positive). Then*

$$\Pr[X \geq e \cdot E[X]] \leq \frac{1}{e}$$

To begin, let's focus on a single row j of the table. Using linearity of expectation we can calculate the expected size of the entry of j hashed to by a given query q .

Lemma 2. *For any query q , for any row j of the hash table,*

$$\Pr[T[j][h_j(q)] \geq \hat{o}_q + \varepsilon N] \leq 1/e$$

Proof. Let's split all elements x_i with $h_j(x_i) = h_j(q)$ into two categories: elements where $x_i = q$; and elements where $x_i \neq q$. From the insertion algorithm above, we have

$$T[j][h_j(q)] = \hat{o}_q + (\# i \text{ such that } x_i \neq q \text{ but } h_j(x_i) = h_j(q)).$$

Let's examine the last term in the above. For all $j \in \{0, \dots, \lceil 1/\delta \rceil - 1\}$ and all $i \in \{1, \dots, N\}$, define a random variable $X_{j,i,q} = 1$ if $x_i \neq q$ and $h_j(q) = h_j(x_i)$; $X_{j,i,q} = 0$ otherwise. Therefore, we can rewrite

$$T[j][h_j(q)] = \hat{o}_q + \sum_{i=1}^N X_{j,i,q} \tag{1}$$

Let's look at the last term. By linearity of expectation,

$$\mathbb{E} \left[\sum_{i=1}^N X_{j,i,q} \right] \leq \sum_{i=1}^N \mathbb{E}[X_{j,i,q}]$$

To calculate $\mathbb{E}[X_{j,i,q}]$, we split into two cases. If $q = x_i$, then $X_{j,i,q} = 0$; therefore, $\mathbb{E}[X_{j,i,q}] = 0$. If $q \neq x_i$, then $\Pr(h_j(q) = h_j(x_i)) = \frac{1}{\lceil e/\varepsilon \rceil} \leq \varepsilon/e$; therefore, $\mathbb{E}[X_{j,i,q}] \leq \varepsilon/e$. Combining these two cases, we can conclude that $\mathbb{E}[X_{j,i,q}] \leq \varepsilon/e$.

Therefore,

$$\mathbb{E} \left[\sum_{i=1}^N X_{j,i,q} \right] = \sum_{i=1}^N \mathbb{E}(X_{j,i,q}) \leq \frac{N\varepsilon}{e}.$$

Applying Lemma 1,

$$\Pr \left[\sum_{i=1}^N X_{j,i,q} \geq N\varepsilon \right] \leq 1/e$$

Combining with Equation 1,

$$\Pr [T[j][h_j(q)] \geq \hat{o}_q + N\varepsilon] \leq 1/e$$

□

With Lemma 2 we're almost done. Since each row T is generated by a different hash function, the rows are independent. Therefore, we can multiply them to get the final answer.

Lemma 3. *For any query q ,*

$$\Pr \left[\min_j T[j][h_j(q)] \geq \hat{o}_q + \varepsilon N \right] \leq \delta$$

Proof. If

$$\min_j T[j][h_j(q)] \geq \hat{o}_q + \varepsilon N,$$

then for all j , $T[j][h_j(q)] \geq \hat{o}_q + \varepsilon N$. By Lemma 2 this occurs for a given j with probability $1/e$. Since the value in each row is independent, the probability that this occurs for all $\lceil \ln 1/\delta \rceil$ rows is

$$\left(\frac{1}{e}\right)^{\lceil \ln 1/\delta \rceil} \leq \left(\frac{1}{e}\right)^{\ln 1/\delta} = \delta.$$

□

4 HyperLogLog

The HyperLogLog data structure is much smaller than Count Min Sketch. It does not attempt to retain the number of occurrences for each individual count. Rather, it attempts to estimate the number of *distinct* items in the stream. That is to say, it estimates $|\{x_i \mid i \in \{1, \dots, N\}\}|$.

Structure The HyperLogLog data structure consists of an array M of counters, where M is of length m (we assume m is a power of 2). Each counter should have at least $\log \log n$ bits; in practice 8-bit counters are sufficient for *any* application.⁶ The HyperLogLog data structure also requires a (single) hash function, which we will call h .

Insert Inserting x_i begins by hashing x_i . Processing this hash proceeds in two steps. First, we obtain an index into the table. Then, we perform a calculation on the hash.

The first step is to obtain an index j using the rightmost $\log_2 m$ bits of $h(x_i)$.⁷ These bits are then shifted off to obtain a remainder $r = h(x_i) \gg \log_2 m$.

Then, we count the number of trailing zeroes in the binary representation of r and store that number in an integer z . (Mathematically, we find the largest

Algorithm 3 Insert x_i

```

j ← h(xi) % m
r ← h(xi) >> log2 m
z ← # trailing 0s of r
if M[j] < z + 1 then
    M[j] ← z + 1
end if

```

z such that $r/2^z$ is an integer.)

We look up $M[j]$. If $M[j] < z + 1$, we set $M[j] = z + 1$. In other words, at every point in time $M[j]$ stores the largest number of zeroes (plus one!) of the remainder of any item that hashes to j .

⁶In particular, 8-bit counters are sufficient even if your stream consists of all particles in the observable universe—much larger than any stream you would see on Earth.

⁷Because m is a power of 2, $\log_2 m$ is always an integer.

Query The query attempts to estimate how many distinct items were seen in the stream based on M . At a high level, the following process constitutes taking a biased harmonic mean⁸ of the entries of M .

To begin, we set a double $Z = 0$. Then for each j from 0 to $m - 1$, we add $(1/2)^{M[j]}$ to Z .

Algorithm 4 Query

```

 $Z = 0$ 
for  $j = 0$  to  $m - 1$  do
     $Z \leftarrow Z + (1/2)^{M[j]}$ 
end for
return  $bm^2/Z$ 

```

We calculate a bias constant b , which depends on m . Good values of b can be found in Table 1.

Finally, we return bm^2/Z .

Parameters and caveats Using 8 hash bits per element is sufficient for a stream of any earthly size. In fact, as streams get large, the main parameter that needs to change is the size of the hash function output. The hash function output should be large enough that any two elements are unlikely to collide

(after all, if x and y have $h(x) = h(y)$, the HyperLogLog will treat them the same and will not count them as distinct elements). 32 bits was suggested in the original paper, and is fine for moderate-sized streams (in fact, for the assignment, I'd suggest you just use one 32-bit output from Murmurhash). However, streams with billions (or even hundreds of millions) of distinct elements—specifically, streams with close to 2^{32} elements or more—will incur a decrease in quality with 32 bit hashes, and a 64 bit hash should be used instead.

m	b
16	.673
32	.697
64	.709
≥ 128	$.7213/(1 + 1.079/m)$

Table 1

The value of m is relatively small in the assignment, but is often moderately large in practice. Counters are cheap, so one often sees $m = 1024$ or $m = 2048$ to give improved estimates.

For streams where the number of distinct elements is very small (less than $5m/2$) or very large (close to 2^{32} with 32-bit hashes), one can still use the same data structure, but should calculate the approximate set size using an alternate method. We won't be dealing with these methods in class, but they are important for implementing this structure in practice.

Guarantees? You are not expected to know the guarantees of this algorithm in this class, but I include them here for completeness. Let D be the number of distinct elements in the stream. The *expected* output of HyperLogLog is D itself. For concentration bounds, it's a little trickier. To get an idea of how close we usually are to D , the standard error⁹ of D is $1.04/\sqrt{m}$. It appears that the outputs given by HyperLogLog are close to a Gaussian distribution. Therefore, the estimate should be within a $1.04/\sqrt{m}$ factor of D 65% of the time, and should be within a $2.08/\sqrt{m}$ factor of D 95% of the time.

⁸That is to say, a special kind of average that happens to work particularly well in this setting.

⁹This is essentially the standard deviation multiplied by the correct estimate D .