

Today we'll talk about a topic that is both very old (as far as computer science goes), and very current. It's a model of computing where the amount of space we have is much less than the amount of data we examine, and hence we can store only a "summary" or "sketch" of the data. Back at the dawn of CS, data was stored on tapes and the amount of available RAM was very small, so people considered this model. And now, even when RAM is cheap and our machines have gigabytes of RAM and terabytes of disk space, it may be listening to an Ethernet cable carrying many gigabytes of data *per second*. *What can we compute in this model?*

- An introduction to the data streaming model, and its concerns.
- An algorithm for heavy hitters in the arrivals-only model.
- An algorithm for heavy hitters with both arrivals and departures.

1 Introduction

Today's lecture will be about a slightly different computational model called the *data streaming* model. In this model you see elements going past in a "stream", and you have very little space to store things. For example, you might be running a program on an Internet router, the elements might be IP Addresses, and you have limited space. You certainly don't have space to store all the elements in the stream. The question is: which functions of the input stream can you compute with what amount of time and space? (For this lecture, we will focus on space, but similar questions can be asked for update times.)

We will denote the stream elements by

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

We assume each stream element is from alphabet Σ and takes b 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 $a_{[1:t]} = \langle a_1, a_2, \dots, a_t \rangle$.

Let us consider some examples. Suppose we have seen the integers

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

- Computing the sum of all the integers seen so far? $F(a_{[1:t]}) = \sum_{i=1}^t a_i$. 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 keep a counter, and when a new element comes in, we add it to the counter.

- 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.

- The median? The outputs on the various prefixes of (\diamond) now are

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

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

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

1, 2, 3, 4, 5, 6, 7, 7, 8, 8, 9, 9, 9 . . .

- (“heavy hitters”) Or the elements that have appeared at least ε fraction of the times so far? Hmm...

You can imagine the applications of this 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 your bandwidth)? Even if you are not working at “line speed”,¹ but just looking over the server logs, you may not want to spend too much time to find out the answers, you 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? Two of the recurring themes will be:

- 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.
- Hashing: this will be a very powerful technique.

2 Finding ε -Heavy Hitters

Let’s formalize things a bit. We have a data stream with elements a_1, a_2, \dots, a_t seen by time t . Think of these elements as “arriving” and being added to a multiset S_t , with $S_0 = \emptyset$, $S_1 = \{a_1\}$, $\dots, S_i = \{a_1, a_2, \dots, a_i\}$, etc. Let

$$\text{count}_t(e) = \{i \in \{1, 2, \dots, t\} \mid a_i = e\}$$

be the number of times element e has been seen in the stream so far. The *multiplicity* of e in S_t .

Call element $e \in \Sigma$ is called a ε -heavy hitter at time t if $\text{count}_t(e) > \varepsilon t$. That is, e constitutes strictly more than ε fraction of the elements that arrive by time t .

The goal is simple — given an threshold $\varepsilon \in [0, 1]$, maintain a data structure that is capable of outputting ε -heavy-hitters. At any point in time, we can query the data structure, and it outputs a set of at most $1/\varepsilon$ elements that contains all the ε -heavy hitters.²

Crucial note: *it’s OK to output “false positives” but we are not allowed “false negatives”. I.e., we’re not allowed to miss any heavy-hitters, but we could output non-heavy-hitters. (Since we only output $1/\varepsilon$ elements, there can be $\leq 1/\varepsilon$ -false positives.)*

For example, if we’re looking for $\frac{1}{3}$ -heavy-hitters, and the stream is

$$E, D, B, D, D_5, D, B, A, C, B_{10}, B, E, E, E, E_{15}, E \dots \quad (\star)$$

(the subscripts are not part of the stream, just to help you count) then

- at time 5, the element D is the only $\frac{1}{3}$ -heavy-hitter,
- at time 11 both B and D are $\frac{1}{3}$ -heavy-hitters, and

¹Such a router might see tens of millions of packets per second.

²Clearly, at any moment in time, there are at most $1/\varepsilon$ elements that are ε -heavy-hitters, so this request is not unreasonable.

- at time 15, there is no $\frac{1}{3}$ -heavy-hitter, and
- at time 16, only E is a $\frac{1}{3}$ -heavy-hitter.

Note that as time passes, the set of frequent elements may change completely, so an algorithm would have to be adaptive. We cannot keep track of the counts for all the elements we've seen so far, there may be a lot of different elements that appear over time, and we have limited space.

Any ideas for how to find some “small” set containing all the ε -heavy-hitters?

Hmm, one trick that is useful in algorithm design, as in problem solving, is to try simple cases first. Can you find a 1-heavy-hitter? (Easy: there is no such element.) How about a 0.99-heavy-hitter? Random sampling will also work (maybe you'll see this in a homework or recitation), but let's focus on a deterministic algorithm for right now.

2.1 Finding a Majority Element

Let's first try to solve the problem of finding a 0.5-heavy-hitter, a.k.a. a majority element, an element that occurs (strictly) more than half the time. Keeping counts of all the elements is very wasteful! How can we find a majority element while using only a little space?

Here's an algorithm (due to R. Boyer and J.S. Moore) that keeps very little information. (One element from Σ in a variable called `memory`, and one counter.)

Initially, `memory = empty` and `counter = 0`.

```
When element a_t arrives
  if (counter == 0)
    set memory = a_t and counter = 1
  else
    if a_t == memory
      counter++
    else
      counter--
      discard a_t
```

At the end, return the element in `memory`.

- Suppose there is no majority element, then we'll output a false positive. As we said earlier, that's OK: we want to output a set of size 1 that contains the majority element, *if any*.
- If there is a majority element, we *will* indeed output it. Why? Observe: when we discard an element a_t , we also throw away another (different) element as well. (Decrementing the counter is like throwing away one copy of the element in `memory`.) So every time we throw away a copy of the majority element, we throw away another element too. Since there are fewer than half non-majority elements, we cannot throw away all the majority elements.

A different (albeit somewhat gruesome) way of looking at this is as a “gang war”. Every element is a gang member. When we have two members from different gangs, they shoot each other in the standoff. If there is a gang with more than $n/2$ members (a majority gang), it will be the only one whose members survive.

2.2 Finding an ε -heavy-hitter

It's possible to extend this idea to finding ε -heavy-hitters.³ Set $k = \lceil 1/\varepsilon \rceil - 1$; for $\varepsilon = 1/2$, we'd get $k = 1$.

*Keep an array $T[1 \dots k]$, each location can hold one element from Σ .
And an array $C[1 \dots k]$ each location can hold a non-negative integer.*

Initialize $C[i] = 0$ and $T[i] = \perp$ for all i .

When element a_t arrives.

If there exists $j \in \{1, 2, \dots, k\}$ such that $a_t == T[j]$, then $C[j] ++$.

Else if some counter $C[j] == 0$ then set $T[j] \leftarrow a_t$ and $C[j] \leftarrow 1$.

Else decrement all the counters by 1 (and discard element a_t).

Exercise: Check that this algorithm is identical to the one above for $\varepsilon = 1/2$.

At any point t , we will maintain estimates of the counts. Our estimates are:

$$\mathbf{est}_t(e) = \begin{cases} C[j] & \text{if } e == T[j] \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

The main claim naturally extends the proof for the case of majority.

Lemma 1 *The estimated counts satisfy:*

$$0 \leq \mathbf{count}_t(e) - \mathbf{est}_t(e) \leq \frac{t}{k+1} \leq \varepsilon t$$

Proof: The estimated count $\mathbf{est}_t(e)$ is at most the actual count, since we never increase a counter for e unless we see e . So $\mathbf{count}_t(e) - \mathbf{est}_t(e) \geq 0$. (In other words, $\mathbf{est}_t(e) \leq \mathbf{count}_t(e)$, it is always an underestimate.)

To prove the other inequality, think of the arrays C and T saying that we have $C[1]$ copies of $T[1]$, $C[2]$ copies of $T[2]$, etc., in our hand. Look at some time when the difference between $\mathbf{count}_t(e)$ and $\mathbf{est}_t(e)$ increases by 1. Each time this happens, we decrement k different counters by 1 and discard an element (which is not currently present in the array T). This is like dropping $k+1$ distinct elements (one of which is e). We can drop at most t elements until time t . So the gap can be at most $t/(k+1)$. And since $k = \lceil 1/\varepsilon \rceil - 1 \geq 1/\varepsilon - 1$, we get that $t/(k+1) \leq \varepsilon t$. ■

Corollary 2 *For every n , the set of items in the array T contains all the ε -heavy-hitters.*

Proof: After we've seen n elements, the ε -heavy-hitters have occurred at least $\mathbf{count}_t(e) > \varepsilon t$ times. If e is a ε -heavy-hitter, by Lemma 1, the estimate

$$\mathbf{est}_t(e) \geq \mathbf{count}_t(e) - \varepsilon t > 0.$$

So element e must be in the array T . ■

To summarize: if we set $k = \lceil 1/\varepsilon \rceil - 1$, we get an algorithm that gives us element counts on a stream of length t to within an additive error of at most $(\varepsilon \cdot t)$ with space $O(1/\varepsilon) \cdot (\log \Sigma + \log t)$ bits. As a corollary we get an algorithm to find a set containing the ε -heavy-hitters.

³The extension to finding ε -heavy-hitters was given by J. Misra and D. Gries, and independently by R.M. Karp, C.H. Papadimitriou and S. Shenker.

3 Heavy Hitters with Deletions

In the above problem, we assumed that the elements were only being added to the current set at each time. We maintained an approximate count for the elements (up to an error of εt). Now suppose we have a stream where elements can be both added and removed from the current set. How can we give estimates for the counts?

Formally, each time we get an *update*, it looks like (add, e) or (del, e) . We will 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 looked like:

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

then the sets at various times are

$$S_0 = \emptyset, S_1 = \{A\}, S_2 = \{A, B\}, S_3 = \{A, A, B\}, S_4 = \{A, A\}, S_5 = \{A\}, S_6 = \{A, C\}, \dots$$

The counts of the element are defined in the natural way: $\text{count}_3(A) = 2$ and $\text{count}_5(A) = 1$. Observe that the “active” set S_t has size $|S_t| = \sum_{e \in \Sigma} \text{count}_t(e)$, and its size can grow and shrink.

What do we want now? We want a data structure to answer count queries approximately. Specifically, at any point in time t , for any element, we should be able to ask “What is $\text{count}_t(e)$?” The data structure should respond with an estimate $\text{est}_t(e)$ such that

$$\Pr \left[\underbrace{|\text{est}_t(e) - \text{count}_t(e)|}_{\text{error}} \leq \underbrace{\varepsilon |S_t|}_{\text{is small}} \right] \geq \underbrace{1 - \delta}_{\text{with high probability}}.$$

Again, we would again like to use small space — perhaps even close to the

$$O(1/\varepsilon) \cdot (\log \Sigma + \log |S_t|) \text{ bits of space}$$

we used in the above algorithm. (As you’ll see, we’ll get close.)

3.1 A Hashing-Based Solution: First Cut

We’re going to be using hashing for this approach, simple and effective. We’ll worry about what properties we need for our hash functions later, for now assume we have a hash function $h : \Sigma \rightarrow \{0, 1, \dots, k - 1\}$ for some suitably large integer k . Maintain an array $A[1 \dots k]$ capable of storing non-negative integers.

```
When update a_t arrives
  If (a_t == (add, e))
    then A[h(e)]++;
  else // a_t == (del, e)
    A[h(e)]--;
```

This was the update procedure. And what is our estimate for the number of copies of element e in our active set S_t ? It is

$$\text{est}_t(e) := A(h(e)).$$

In words, we look at the location $h(e)$ where e gets mapped using the hash function h , and look at the count $A[h(x)]$ stored at that location. What does it contain? It contains the current count for

element e for sure. But added to it is the current count for any other element that gets mapped to that location. In math:

$$A(h(e)) = \sum_{e' \in \Sigma} \text{count}_t(e') \cdot \mathbf{1}(h(e') = h(e)),$$

where $\mathbf{1}(\text{some condition})$ is a function that evaluates to 1 when the condition in the parentheses is true, and 0 if it is false. We can rewrite this as

$$A(h(e)) = \text{count}_t(e) + \sum_{e' \neq e} \text{count}_t(e') \cdot \mathbf{1}(h(e') = h(e)), \quad (2)$$

or using the definition of the estimate, as

$$\text{est}_t(e) - \text{count}_t(e) = \sum_{e' \neq e} \text{count}_t(e') \cdot \mathbf{1}(h(e') = h(e)). \quad (3)$$

A great situation will be if no other elements $e' \neq e$ hashed to location $h(e)$ and the error (i.e., the summation on the right) evaluates to zero. But that may be unlikely. On the other hand, we can show that the expected error is not too much.

What's the expected error? Now we need to assume something good about the hash functions. Assume that the hash function h is a random draw from a universal family. Recall the definition from Lecture #5.

Definition 3 A family H of hash functions from $\Sigma \rightarrow R$ is universal if for any pair of distinct keys $x_1, x_2 \in \Sigma$ with $x_1 \neq x_2$,

$$\Pr[h(x_1) = h(x_2)] \leq \frac{1}{|R|}.$$

In other words, if we just look at two keys, the probability that they collide is no more than if we chose to map them randomly into the range R . We gave a construction where each hash function in the family used $(\lg k) \cdot (\lg |\Sigma|)$ bits to specify.

Good. So we drew a hash function h from this universal hash family H , and we used it to map elements to locations $\{0, 1, \dots, k-1\}$. What is its expected error? From (3), it is

$$E \left[\sum_{e' \neq e} \text{count}_t(e') \cdot \mathbf{1}(h(e') = h(e)) \right] = \sum_{e' \neq e} \text{count}_t(e') \cdot E[\mathbf{1}(h(e') = h(e))] \quad (4)$$

$$\begin{aligned} &= \sum_{e' \neq e} \text{count}_t(e') \cdot \Pr[h(e') = h(e)] \\ &\leq \sum_{e' \neq e} \text{count}_t(e') \cdot (1/k) \\ &= \frac{|S_t| - \text{count}_t(e)}{k} \leq \frac{|S_t|}{k}. \end{aligned} \quad (5)$$

We used linearity of expectations in equality (4). To get (5) from the previous line, we used the definition of a universal hash family. Let's summarize:

Claim 4 The estimator $\text{est}_t(e) = A(h(e))$ ensures that

- (a) $\text{est}_t(e) \geq \text{count}_t(e)$, and
- (b) $E[\text{est}_t(e)] - \text{count}_t(e) \leq |S_t|/k$.

The space used is: k counters, and $O((\log k)(\log \Sigma))$ to store the hash function.

That’s pretty awesome. (But perhaps not so surprising, once you think about it.) Note that if we were only doing arrivals and no deletions, the size of S_t would be exactly t . So the expected error would be at most t/k , which is about the same as we were getting in Section 2 using an array of size $(k - 1)$. But now we can also handle deletions!

What’s the disadvantage? We only have a bound on the *expected error* $E[\mathbf{est}_t(e)] - \mathbf{count}_t(e)$. It is no longer deterministic. Also, the expected error being small is weaker than saying: we have small error with high probability. So let’s see how to improve things.

3.2 Amplification of the Success Probability

Any ideas how to *amplify* the probability that we are close to the expectation? The idea is simple: *independent repetitions*. Let us pick m hash functions h_1, h_2, \dots, h_m . Each $h_i : \Sigma \rightarrow \{0, 1, \dots, k-1\}$. How do we choose these hash functions? *Independently* from the universal hash family H .⁴ We have m arrays A_1, A_2, \dots, A_m , one for each hash function. The algorithm now just uses the i^{th} hash function to choose a location in the i^{th} array, and increments or decrements the same as before.

```

When update a_t arrives
  For each i from 1..m
    If (a_t == (add, e))
      then A_i[h_i(e)]++;
    else // a_t == (delete, e)
      A_i[h_i(e)]--;

```

And what is our new estimate for the number of copies of element e in our active set? It is

$$\mathbf{best}_t(e) := \min_{i=1}^m A_i(h_i(e)).$$

In other words, each (h_i, A_i) pair gives us an estimate, and we take the least of these. It makes perfect sense — the estimates here are all *overestimates*, so taking the least of these is the right thing to do.

But how much better is this estimator? Let’s do the math.

What is the chance that one single estimator has error more than $2|S_t|/k$? Remember the expected error is at most $|S_t|/k$. So by Markov’s inequality

$$\Pr \left[\mathit{error} > 2 \cdot \frac{|S_t|}{k} \right] \leq \frac{1}{2}.$$

And what is the chance that all of the m copies have “large” (i.e., more than $2|S_t|/k$) error? The probability of m failures is

$$\begin{aligned} & \Pr[\text{each of } m \text{ copies have large error}] \\ &= \prod_{i=1}^m \Pr[i^{\text{th}} \text{ copy had large error}] \\ &\leq (1/2)^m. \end{aligned}$$

The first equality there used the independence of the hash function choices. (Only if events \mathcal{A}, \mathcal{B} are independent you can use $\Pr[\mathcal{A} \wedge \mathcal{B}] = \Pr[\mathcal{A}] \cdot \Pr[\mathcal{B}]$.) And so the minimum of the estimates will have “small” error (i.e., at most $2|S_t|/k$) with probability at least $1 - (1/2)^m$.

⁴If we use the hash function construction given in Lecture #5 (Section 5.4.1), this means the $(\lg k) \cdot (\lg |\Sigma|)$ -bit matrices for each hash function must be filled with independent random bits.

3.2.1 Final Bookkeeping

Let's set the parameters now. Set $k = 2/\varepsilon$, so that the error bound $2|S_t|/k = \varepsilon|S_t|$. And suppose we set $m = \lg 1/\delta$, then the failure probability is $(1/2)^m = \delta$, and our query will succeed with probability at least $1 - \delta$.⁵

Then on any particular estimate $\mathbf{best}_t(e)$ we ensure

$$\Pr \left[\left| \mathbf{best}_t(e) - \mathbf{count}_t(e) \right| \leq \varepsilon |S_t| \right] \geq 1 - \delta.$$

Just as we wanted. And the total space usage is

$$m \cdot k \text{ counters} = O(\log 1/\delta) \cdot O(1/\varepsilon) = O(1/\varepsilon \log 1/\delta) \text{ counters.}$$

Each counter has to store at most $\lg T$ -bit numbers after T time steps.⁶

Space for Hash Functions: We need to store the m hash functions as well. How much space does that use? The construction from Lecture #5 used $s := (\lg k) \cdot (\lg \Sigma)$ bits per hash function. Since $k = 1/\varepsilon$, the total space used for the functions is

$$m \cdot s = O(\log 1/\delta) \cdot (\lg 1/\varepsilon) \cdot (\lg \Sigma) \text{ bits.}$$

3.2.2 And in Summary...

Using about $1/\varepsilon \times$ poly-logarithmic factors space, and very simple hashing ideas, we could maintain the counts of elements in a data stream under both arrivals and departures (up to an error of $\varepsilon|S_t|$). As in the arrival-only case, these counts make sense only for very high-frequency elements.

⁵How small should you make δ ? Depends on how many queries you want to do. Suppose you want to make a query a million times a day, then you could make $\delta = 1/10^9 \approx 1/2^{30}$ to get a 1-in-1000 chance that even one of your answers has high error. Our space varies linearly as $\lg 1/\delta$, so setting $\delta = 1/10^{18}$ instead of $1/10^9$ doubles the space usage, but drops the error probability by a factor of billion.

⁶So a 32-counter can handle a data stream of length 4 billion. If that is not enough, there are ways to reduce this space usage as well, look online for "approximate counters".