

Lecture 5: Data Streaming

We study some interesting sublinear space algorithms, including detecting heavy hitters, estimating number of distinct elements, and estimating frequency moments.

We will see that ideas from hashing and k -wise independence are very useful in this setting.

Sublinear algorithms In some applications where we have a massive data set, the data is too much that we can't afford to store them all or read them once. Yet we can design sublinear space or sublinear time algorithms to do something nontrivial. Randomness is crucial in these settings, as it can be shown that deterministic algorithms cannot guarantee anything nontrivial for most problems.

For sublinear space algorithms, we focus on the data streaming model, in which the algorithm can take one pass on the data but has very limited space for computation. A good motivating example is in computer networks, where routers would like to get good statistics of the traffic but could not afford to store all the data. We will study some basic problems such as estimating the number of distinct elements, the second moment of the data, etc.

Since we use so little space, we could not hope to solve the problem exactly, and we work with a weaker requirement that the algorithm can distinguish between YES-instances and instances that are "far from YES" (e.g. $\# \text{ distinct elements} \geq D$, or $\# \text{ distinct elements} \leq D$).

Of course, we also allow some error probability δ ; we will make them precise in the following

Heavy Hitters

Given a data stream $X_1, X_2, X_3, \dots, X_T$, each $X_t = (i_t, c_t)$ where i_t is the ID of the t -th data and c_t is the weight associated with it, e.g. i_t is the IP address and c_t is the number of bytes of the data transmitted.

Our goal is to find all the heavy hitters of the data stream. Let $Q = \sum_t c_t$ be the total count.

Given a number q , we say an ID i is a heavy hitter if $\sum_{t: i_t = i} c_t \geq q$. (e.g. $q \geq 0.01 Q$)

Let $\text{Count}(i, T) = \sum_{t=1, i_t=i}^T c_t$ be the total count for ID i .

We will show a sublinear space algorithm with the following properties:

- ① All heavy hitters are reported, i.e. those i with $\text{Count}(i, T) \geq q$,
- ② If $\text{Count}(i, T) \leq q - \epsilon Q$, then i is reported with probability at most δ .

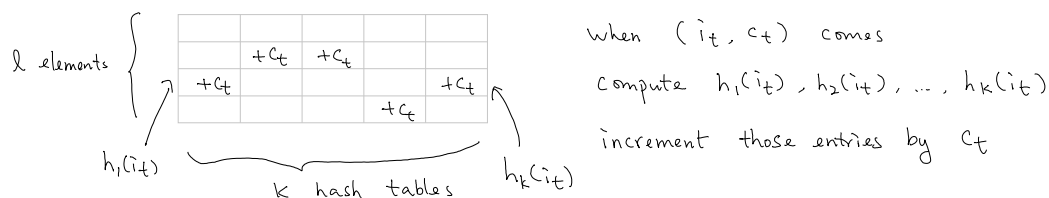
Note that the performance guarantee is nontrivial only when $q = \Omega(Q)$. In that range, the algorithm will report all heavy hitters, only having a small probability reporting an ID that is far from a heavy hitter.

The idea is to use 2-universal hash functions, similar to the construction of a "bloom filter".

We will use k hash functions, h_1, h_2, \dots, h_k , each maps the universe into $\{1, 2, \dots, \ell\}$.

We maintain $k \cdot \ell$ counters, each counter $C_{a,j}$ adds the weight of the items that are mapped to the j -th entry by the a -th hash function. Initially $C_{a,j} = 0$ for $1 \leq a \leq k$, $1 \leq j \leq \ell$.

The algorithm is simple. Given $x_t = (i_t, c_t)$, for each $1 \leq a \leq k$, increment $C_{a, h_a(i_t)}$ by c_t .



After we read all the data, we report all i with $\min_{j=h_a(i), 1 \leq a \leq k} C_{a,j} \geq q$. In words, we report an ID if the minimum counter associated with it is at least q .

Analysis

It should be clear that a heavy hitter is always reported, since by our algorithm all the counters associated with it must be at least q .

The key is to show that if an ID is not a heavy hitter, then it is reported with small probability.

Let us consider a specific ID i . Suppose $\text{Count}(i, T) \leq q - \epsilon Q$. What is the prob. that i is reported?

Consider $C_{a, h_a(i)}$ for a specific a . Since i is reported, we must have $C_{a, h_a(i)} \geq q$.

Let Z_a be the value of this counter that is incremented by other IDs.

Since h_a is chosen from a 2-universal family, another item is mapped to $h_a(i)$ with probability $\leq \frac{1}{\ell}$.

Therefore $E[Z_a] \leq Q/\ell$.

By Markov's inequality, $\Pr(Z_a \geq \varepsilon Q) \leq \frac{1}{\varepsilon l}$.

So $\Pr(\min Z_a \geq \varepsilon Q) \leq (\frac{1}{\varepsilon l})^k$

Set $l = \frac{e}{\varepsilon}$ and $k = \ln(\frac{1}{\delta})$, this probability is at most δ .

So, we just need to maintain $O(\frac{1}{\varepsilon} \ln \frac{1}{\delta})$ counters for this task, which is just a constant for constant ε and δ .

The hash functions can be stored efficiently, using $O(k) = O(\ln \frac{1}{\delta})$ cells.

And the evaluation time is $O(k) = O(\ln \frac{1}{\delta})$ word operations.

Distinct Elements

The input is a stream of a sequence a_1, a_2, \dots, a_n , where each a_i is chosen from $\{1, 2, \dots, m\}$.

Our algorithm is allowed to take one pass of the stream, using very limited space (e.g. $O(\log(m+n))$), and return a good estimate on the number of distinct elements in the data stream.

We will give bounds in terms of $\log(m)$ instead of $\log(n)$. It is good in the case when $n \gg m$.

In the case when $m > n$, we can first hash each item to a table of size n^2 or larger, then with high probability there will be no collisions, and thus we can assume $\log m = O(\log n)$.

We show how to use a 2-universal hash family to solve the problem.

The idea is to map the universe $\{1, 2, \dots, m\}$ to $\{1, 2, \dots, m^3\}$

As we have seen last time, there will be no collisions with high probability.

Ideally, we want the hash values $h(1), h(2), \dots, h(m)$ to be evenly distributed in $\{1, 2, \dots, m^3\}$.

Then, if there are D distinct elements, the t -th smallest hash value should be close to $t \cdot \frac{m^3}{D}$.

Therefore, if we know the t -th smallest non-empty bin is T , then $T \approx t \cdot \frac{m^3}{D}$ and thus $D \approx \frac{t \cdot m^3}{T}$.

We must be careful about the space requirement. First of all, we are not going to store the whole hash table, as it is really unnecessary.

We only need to store t numbers, corresponding to the t smallest non-empty bins.

Storing more numbers will help get a better estimate, but of course it requires more space.

We will find a good t to balance both the error probability and the space requirement.

Let's state the algorithm again.

- Pick a random hash function h from a 2-universal family.
- For each item a_i in the data stream, compute $h(a_i)$. Maintain the list of the t smallest hash values seen so far.
- Let T be the t -th smallest hash value after all the data read. Return $Y = \frac{tm^3}{T}$

Theorem: By setting $t = O(\frac{1}{\epsilon^2})$, we have $(1-\epsilon)D \leq Y \leq (1+\epsilon)D$ with constant probability.

Case 1: We first bound the probability that $Y > (1+\epsilon)D$.

This implies that $T < \frac{tm^3}{(1+\epsilon)D}$, which is at most $\frac{tm^3(1-\frac{\epsilon}{2})}{D}$ for small enough ϵ .

This means that there are at least t hash values that are smaller than $\frac{tm^3(1-\frac{\epsilon}{2})}{D}$.

On the other hand, let X_i be the indicator variable that $h(a_i) \leq \frac{tm^3(1-\frac{\epsilon}{2})}{D}$.

Then $E[X_i] = \Pr(h(a_i) \leq \frac{tm^3(1-\frac{\epsilon}{2})}{D}) \leq \frac{t(1-\frac{\epsilon}{2})}{D}$ since $h(a_i)$ is equally likely to be in $\{1, \dots, m^3\}$.

Therefore, if there are D distinct elements, $E[\# \text{ of elements with hash values } \leq \frac{tm^3(1-\frac{\epsilon}{2})}{D}] \leq t(1-\frac{\epsilon}{2})$.

So, the expected number of such hash values is $\leq t(1-\frac{\epsilon}{2})$, but we have $\geq t$ now.

To bound the deviation, we compute the variance. Let $X = \sum_{i=1}^D X_i$. Then $E[X] \leq t(1-\frac{\epsilon}{2})$.

Since X_i are pairwise independent, $\text{Var}[X] = \sum_{i=1}^D \text{Var}[X_i] \leq t(1-\frac{\epsilon}{2})$.

By Chebyshev's inequality, $\Pr(X > t) = \Pr(X > t(1-\frac{\epsilon}{2}) + \frac{\epsilon}{2} \cdot t) \leq \Pr(|X - E[X]| \geq \frac{\epsilon t}{2})$
 $\leq \frac{4\text{Var}[X]}{\epsilon^2 t^2} \leq \frac{4(1-\frac{\epsilon}{2})}{\epsilon^2 t} \leq \frac{4}{\epsilon^2 t} \leq \frac{1}{6}$ for $t = O(\frac{1}{\epsilon^2})$.

Case 2: We bound the probability that $Y < (1-\epsilon)D$. This is very similar to case 1.

The condition implies that there are less than t distinct elements with hash values $\leq \frac{(1+\epsilon)tm^3}{D}$.

By a very similar calculation, $E[X] \geq (1+\frac{\epsilon}{2})t$ and $\text{Var}[X] \leq (1+\frac{\epsilon}{2})t$.

By Chebyshev's inequality, $\Pr(X < t) \leq \Pr(|X - E[X]| \geq \frac{t\epsilon}{2}) \leq \frac{4\text{Var}[X]}{\epsilon^2 t^2} < \frac{1}{6}$ for $t = O(\frac{1}{\epsilon^2})$. \square

Therefore, with probability at most $1/3$, $(1-\epsilon)D \leq Y \leq (1+\epsilon)D$.

By running $O(\log \frac{1}{\delta})$ copies and taking the median, the error probability is at most δ .

The total space used is $O(\frac{1}{\epsilon^2} \log m \log(\frac{1}{\delta}))$, since each copy stores $O(\frac{1}{\epsilon^2})$ numbers each of $O(\log m)$ bits, and the hash function only requires $O(\log m)$ bits to store.

Note that the list can be implemented by balanced binary search trees or heaps to support fast update.

Since there are at most $O(\frac{1}{\epsilon^2})$ elements, each operation can be supported in $O(\log \frac{1}{\epsilon})$ steps.

Recently there is an optimal algorithm that uses only $O(\frac{1}{\epsilon^2} + \log m)$ space and $O(1)$ update time.

This is best possible; see the paper "An optimal algorithm for the distinct element problem", by Kane, Nelson, and Woodruff.

Frequency Moments

Again, let a_1, a_2, \dots, a_n be the items in the data stream, where each a_i is from $\{1, 2, \dots, m\}$.

For $1 \leq i \leq m$, let x_i be the number of items equal to i .

We would like to estimate $\sum_{i=1}^m x_i^p$ where p is a given constant.

Note that when $p=0$ this is just the distinct element problem, when $p=1$ it is trivial.

Here we focus on the problem when $p=2$, and later mention results for other p .

The algorithm is very simple and elegant. This is an important result in data streaming algorithms.

- Set r_1, r_2, \dots, r_m independently to be ± 1 with prob $1/2$ and -1 with prob $1/2$.
- Maintain $Z = \sum_{i=1}^m r_i x_i$ (this can be done by just storing the current value of Z)
- Return $Y = Z^2$

We will prove that Y is a good approximation to $\sum_{i=1}^m x_i^2$ with high probability.

First we show that $E[Y] = \sum_{i=1}^m x_i^2$.

Since $Y = Z^2 = (\sum_{i=1}^m r_i x_i)^2$, $E[Y] = E[Z^2] = \sum_{1 \leq i, j \leq m} E[r_i x_i r_j x_j] = \sum_{1 \leq i, j \leq m} x_i x_j E[r_i r_j]$

When $i=j$, $E[r_i r_j] = 1$; when $i \neq j$, $E[r_i r_j] = E[r_i] E[r_j] = 0$.

Therefore, $E[Y] = E[Z^2] = \sum_{i=1}^m x_i^2$.

To show that it is a good approximation with high probability, we need to show that Z^2 is concentrated around its expected value, for which we'll use the Chebyshev's inequality.

So we compute $E[Y^2] = E[Z^4] = \sum_{1 \leq i, j, k, \ell \leq m} x_i x_j x_k x_\ell E[r_i r_j r_k r_\ell]$.

Note that $E[r_i r_j r_k r_\ell] = 0$ whenever one index appears only once.

Therefore there are only two situations when $E[r_i r_j r_k r_\ell] \neq 0$:

- when $i=j=k=l$, and this contributes $\sum_{i=1}^m x_i^4$ to $E[Y^2]$.
- when two indexes appear twice, and this contributes $6 \sum_{i < j} x_i^2 x_j^2$ to $E[Y^2]$.

$$\text{So } E[Y^2] = E[Z^4] = \sum_{i=1}^m x_i^4 + 6 \sum_{i < j} x_i^2 x_j^2.$$

$$\begin{aligned} \text{Thus } \text{Var}[Y] &= E[Y^2] - (E[Y])^2 = E[Z^4] - (E[Z^2])^2 \\ &= \sum_{i=1}^m x_i^4 + 6 \sum_{i < j} x_i^2 x_j^2 - \left(\sum_{i=1}^m x_i^2 \right)^2 \\ &= \sum_{i=1}^m x_i^4 + 6 \sum_{i < j} x_i^2 x_j^2 - \left(\sum_{i=1}^m x_i^4 + 2 \sum_{i < j} x_i^2 x_j^2 \right) \\ &= 4 \sum_{i < j} x_i^2 x_j^2 \leq 2 \left(\sum_{i=1}^m x_i^2 \right)^2 = 2(E[Y])^2 \end{aligned}$$

By Chebyshev's inequality, $\Pr(|Y - E[Y]| \geq c \sqrt{\text{Var}[Y]}) \leq \frac{1}{c^2}$.

This already implies that $\Pr(|Y - E[Y]| \geq c \sqrt{E[Y]}) \leq \frac{1}{c^2}$, getting a constant factor approximation.

To get a $(1 \pm \epsilon)$ -approximation, we need to find a variable Y' with $E[Y] = E[Y']$ but the variance of Y' is smaller. The standard way to do this is to compute Y_1, Y_2, \dots, Y_k and take the average.

Let $Y' = (\sum_{i=1}^k Y_i) / k$. Then $E[Y'] = E[Y]$ and $\text{Var}[Y'] = \sum_{i=1}^k \text{Var}[Y_i / k] = \frac{1}{k} \text{Var}[Y] \leq \frac{2}{k} (E[Y])^2$.

Therefore, $\Pr(|Y' - E[Y']| \geq c \sqrt{2/k} E[Y']) \leq 1/c^2$.

Setting c to be a constant and $k = O(1/\epsilon^2)$ will get an $(1 \pm \epsilon)$ -approximation with constant probability.

The space required to compute one Y is to store one number Z . So the total space required is $O(1/\epsilon^2)$ numbers.

Wait, how do we compute Z if we don't store r_1, r_2, \dots, r_m ? This requires m bits, although it is not as bad as storing all the data, it is still undesirable if the universe is large (e.g. IP address).

A key point about this approach is that we only need the variables to be 4-wise independent, so that the analysis about $E[r_i r_j r_k r_l]$ would hold.

m 4-wise independent bits can be generated from $O(\log m)$ random bits. So the extra storage required for the computation is only $O(\log m)$.

This approach is called sketching, which is very useful in data streaming algorithms. It is closely related to the idea of dimension reduction.

What about $\sum_{i=1}^m x_i^3$? It turns out that $O(n^{1-\frac{4}{p}} \text{polylog}(n))$ space is required to estimate $\sum_{i=1}^m x_i^p$ for any $p > 2$. On the other hand, for $0 \leq p \leq 2$ (p can be fractional), then $O(\text{polylog}(m))$ space is enough.

References Heavy hitter is from Chapter 13 of "Probability and Computing".

Distinct elements and frequency moments are from course notes of 6.895 in MIT by Piotr Indyk.

In the same course page, you can find other interesting sublinear algorithms such as Compressed sensing.