

Lecture 9: The L_2 Heavy Hitters Problem

Source: Lecture notes by
Aaron Roth and Adam Smith

Lecturer: Uri Stemmer

Recall the formulation of the heavy-hitters streaming problem we studied in Lecture 7:

Given a (prefix of a) stream $S = (s_1, s_2, \dots, s_i)$, where each $s_\ell \in [n]$, return a list of size k containing all elements that appear more than i/k times in the stream, i.e., all elements $x \in [n]$ such that

$$\text{weight}_i(x) \triangleq |\{\ell \in [i] : s_\ell = x\}| > \frac{i}{k}$$

We saw a deterministic algorithm with space $\approx k$ for this problem, which is optimal in the sense that every algorithm for this problem must use $\Omega(k)$ space in the worst-case (as there could be k “heavy” elements that the algorithm “must remember”).

However, there is still room for improvement, in the sense that if no element has weight i/k then the algorithm from Lecture 7 guarantees nothing. That is, even if we are OK with paying space $\approx k$, we could still hope to “catch” elements with smaller weights under some conditions.

Let’s begin with the following more general formulation of the problem that allows for *deletions* in the stream.

Definition [The L_1 Heavy-Hitters Problem]:

- Every update in the stream is a pair (\mathbf{u}_i, Δ_i) where $\mathbf{u}_i \in \mathbb{R}^n$ is a standard basis vector and $\Delta_i \in \mathbb{R}$ is its weight
- After t updates, denote the frequency vector of the stream as

$$\mathbf{f}^{(t)} = \Delta_1 \cdot \mathbf{u}_1 + \dots + \Delta_t \cdot \mathbf{u}_t$$
- At query time, say after t updates, the goal is to return a set of $\mathcal{O}(k)$ coordinates (called keys) that include all the L_1 -heavy hitters of $\mathbf{f}^{(t)}$, that is, all keys $i \in [n]$ such that

$$|\mathbf{f}^{(t)}[i]| > \frac{1}{k} \|\mathbf{f}^{(t)}\|_1 = \frac{1}{k} (|\mathbf{f}^{(t)}[1]| + \dots + |\mathbf{f}^{(t)}[n]|)$$
- The output is correct when all heavy hitters are reported

The L_1 heavy hitters problem can be solved deterministically with small space (even though in Lecture 7 we did not address deletions). Now consider the following variant, where we aim to identify elements which are “heavy” w.r.t. the L_2 norm of the frequency vector rather than w.r.t. the L_1 norm:

Definition [The L_2 Heavy-Hitters Problem]:

- Every update in the stream is a pair (\mathbf{u}_i, Δ_i) where $\mathbf{u}_i \in \mathbb{R}^n$ is a standard basis vector and $\Delta_i \in \mathbb{R}$ is its weight
- After t updates, denote the frequency vector of the stream as

$$\mathbf{f}^{(t)} = \Delta_1 \cdot \mathbf{u}_1 + \dots + \Delta_t \cdot \mathbf{u}_t$$
- At query time, say after t updates, the goal is to return a set of $\mathcal{O}(k)$ coordinates (called keys) that include all the L_2 -heavy hitters of $\mathbf{f}^{(t)}$, that is, all keys $i \in [n]$ such that

$$|\mathbf{f}^{(t)}[i]|^2 > \frac{1}{k} \|\mathbf{f}^{(t)}\|_2^2 = \frac{1}{k} (|\mathbf{f}^{(t)}[1]|^2 + \dots + |\mathbf{f}^{(t)}[n]|^2)$$
- The output is correct when all heavy hitters are reported

Note:

(1) There could be at most k L2-heavy keys, as otherwise $\|f^{(t)}\|_2^2 \geq \left(|f^{(t)}[i_1]|^2 + \dots + |f^{(t)}[i_{k+1}]|^2 \right) \geq \left(\frac{1}{k} \|f^{(t)}\|_2^2 + \dots + \frac{1}{k} \|f^{(t)}\|_2^2 \right) = \frac{k+1}{k} \|f^{(t)}\|_2^2$

(2) If a key i is L1-heavy with parameter \sqrt{k} then it must also be L2-heavy with parameter k because

$$|f^{(t)}[i]| > \frac{1}{\sqrt{k}} \|f^{(t)}\|_1 \implies |f^{(t)}[i]|^2 > \frac{1}{k} \|f^{(t)}\|_1^2 \geq \frac{1}{k} \|f^{(t)}\|_2^2$$

where the last inequality follows since the L1 norm is always bigger or equal to the L2 norm.

(3) The converse is not necessarily true. That is, a key could be very L2-heavy while being L1-heavy. For example, for the vector $f = (\sqrt{n}, 1, \dots, 1)$ we have $f[1]^2 = n$ and $\|f\|_2^2 \approx 2n$. So key 1 is very L2-heavy. But it is not L1 heavy (\sqrt{n} compared to $\approx n$).

(4) This means that, in spirit, aiming to identify all L2-heavy-hitters is a harder task than identifying all L1-heavy-hitters, since the L2-heavy-hitters are “less heavy” (using the same space).¹

What is known about the L2-HH problem in the classical/oblivious setting? Can be solved with space $O(k \cdot \log n)$ using the CountSketch algorithm (to be surveyed next)

What will we get from the generic construction for the adaptive setting based on DP-stability?

Recall our generic construction:

Input: Collection of P random strings $R = (r_1, \dots, r_P) \in (\{0, 1\}^*)^P$ for some parameter P

1. Initiate P copies on a non-adaptive algorithm \mathcal{A} with random strings r_1, \dots, r_P
2. Every time an update comes, feed it to all of the copies of \mathcal{A}
3. On query time:
 - a) Query each of the copies of \mathcal{A} to obtain P intermediate outputs y_1, \dots, y_P
 - b) Return a DP-stable aggregation of $\{y_1, \dots, y_P\}$ as the current output

- Informally, composition arguments for DP-stability allow us to release t aggregations using a dataset of size $P \approx \sqrt{t}$
- In our context, we have T rounds, during each of which we need to release $O(k)$ heavy elements
- This amounts to a total of $t = Tk$ aggregations, for which we need $P \approx \sqrt{Tk}$ intermediate outputs
- This gives a total space of $O(\sqrt{T} \cdot k^{1.5})$. So the “price for robustness” is $\approx \sqrt{T} \cdot k$

Can we do better?

¹ This is only “in spirit” as there could be keys which are L1 heavy but not L2 heavy (up to the difference between \sqrt{k} and k mentioned above). For example, for $f = (10, 90)$ we have that $f[1] \geq \frac{\|f\|_1}{10} = 10$, but $(f[1])^2 < \frac{\|f\|_2^2}{10} = \frac{100+8100}{10} = 820$. So key 1 is L1 heavy but not L2 heavy (both with parameter $k = 10$).

Algorithm CountSketch

- Parameters: n =dimension of input vector, d =size of sketch, $b=O(k)$ =the “width” (where b/k controls “accuracy” as will become clear later)
- Initialize $\ell = \frac{d}{b}$ pairs of hash functions $\rho = ((h_1, s_1), \dots, (h_\ell, s_\ell))$ where $h_j: [n] \rightarrow [b]$ and $s_j: [n] \rightarrow \{\pm 1\}$
- Think of ρ as defining $\frac{d}{b} \times b$ “buckets” $C[\cdot, \cdot]$ (initially empty)
- Given an update (i, Δ) : For every $j \in [d/b]$ add $s_j(i) \cdot \Delta$ to the bucket $C[j, h_j(i)]$

Update	$h_1(i)=5, s_1(i)=+1$				+ Δ				
(i, Δ)	$h_2(i)=2, s_2(i)=-1$	- Δ							
	$h_3(i)=8, s_3(i)=-1$						- Δ		
	$h_4(i)=3, s_4(i)=+1$		+ Δ						

Utility analysis:

- For every $j \in [d/b]$ and $i \in [n]$ we have that $s_j(i) \cdot C[j, h_j(i)]$ is an unbiased estimator for $f[i]$
- Fix a bucket. Every key i is mapped to that bucket with value $f[i]$ or $-f[i]$ with probabilities $\frac{1}{2b}$
- Hence, the noise in the bucket is the sum of n “Bernoulli RVs” with variance $\approx \frac{1}{b} \sum_i (f[i])^2 = \frac{1}{b} \|f\|_2^2$
- If key i is an L2-heavy hitter then it overtakes the noise (since $|f[i]|$ is bigger than the standard deviation of the noise). Formally, fix a key i and consider its j th bucket $C[j, h_j(i)]$

- By Chebyshev’s inequality: $\Pr \left[|s_j(i)C[j, h_j(i)] - f[i]| > \sqrt{\frac{3}{b}} \cdot \|f\|_2 \right] \leq \frac{\frac{1}{b} \|f\|_2^2}{\left(\sqrt{\frac{3}{b}} \cdot \|f\|_2\right)^2} \leq \frac{1}{3}$

- That is, for every $j \in [d/b]$, with probability at least $2/3$ we have that $s_j(i)C[j, h_j(i)]$ is an accurate estimate for $f[i]$ up to additive error of $\sqrt{\frac{3}{b}} \cdot \|f\|_2$.
- Therefore, by the Chernoff bound (and a union bound over the n keys), if the number of rows $[d/b]$ is at least $\approx \log n$ then w.h.p. for every key i we have that

Median $_{j \in [d/b]} \{ s_j(i) \cdot C[j, h_j(i)] \}$ is within $\sqrt{\frac{3}{b}} \cdot \|f\|_2$ of $f[i]$.

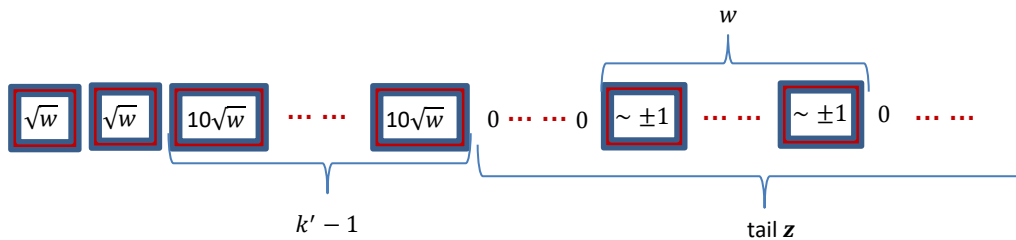
- That is, from the sketch we can estimate the value of key i as **Median** $_{j \in [d/b]} \{ s_j(i) \cdot C[j, h_j(i)] \}$
- Note that for an L2-heavy key i it holds that $\sqrt{\frac{3}{b}} \cdot \|f\|_2 < \sqrt{\frac{3k}{b}} \cdot |f[i]|$, which means multiplicative error at most $\left(1 + \sqrt{\frac{3k}{b}}\right)$ for heavy keys.
- To identify the list of heavy elements we report $O(k)$ keys with maximal estimated values:
 - To see that this list indeed includes all heavy elements, note that if key i is heavy with parameter k (i.e., $|f[i]| \geq \frac{1}{\sqrt{k}} \cdot \|f\|_2$) and key j is *not* heavy with parameter $2k$ (i.e., $|f[j]| <$

$\frac{1}{\sqrt{2k}} \cdot \|f\|_2$), then our estimation for $f[i]$ is bigger than our estimation for $f[j]$, because for all $b \geq 5k$ we have $\frac{1}{\sqrt{k}} \cdot \|f\|_2 - \sqrt{\frac{3}{b}} \cdot \|f\|_2 > \frac{1}{\sqrt{2k}} \cdot \|f\|_2 + \sqrt{\frac{3}{b}} \cdot \|f\|_2$

- That is, only keys that are at least “ $2k$ heavy” might outshine keys that are “ k heavy”. As there could be at most $2k$ keys that are “ $2k$ heavy”, we get that all “ k heavy” keys must be included in the list of $2k$ keys with top estimations.

Attack on CountSketch (median estimator)

- CountSketch reports exactly k' elements with highest estimated values
- Attack vectors are a concatenation of a fixed prefix of length $k' + 1$ of “heavy” keys and a “tail”
- Tails have disjoint supports
- The reported candidate HH are $H = \{1, 3, \dots, k' + 1\}$ OR $\{2, 3, \dots, k' + 1\}$



$a \leftarrow 0$

Repeat r times:

Query with a random tail z

If $1 \in H$ $a \leftarrow a + z$

else $a \leftarrow a - z$

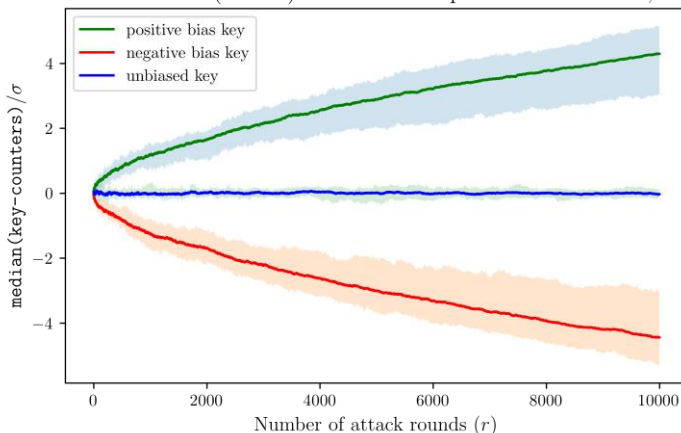
Return a

- a biases “towards” 1 and “away from” 2
- Adversarial vector: “0” on keys 1,2 and heavy on k' keys and “tail” a
- Keys 1,2 get large estimates and are reported in H

Attack Simulations

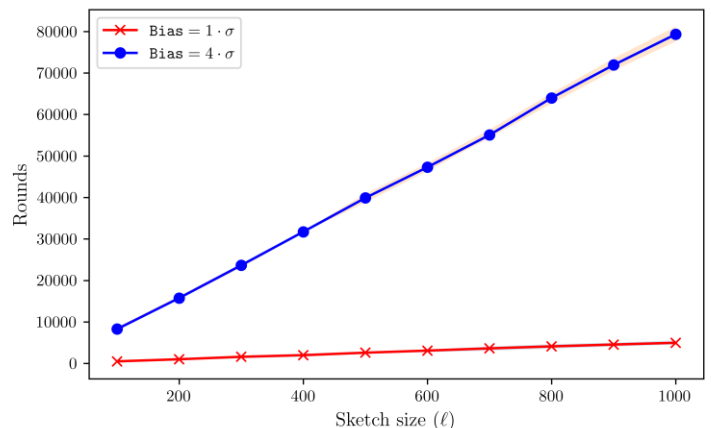
Simulations show that $r = 5 \cdot \tau^2 \cdot \frac{d}{b}$ attack queries suffice for generating an adversarial vector f with keys biased by $\frac{\tau}{\sqrt{k}} \cdot \|f\|_2$

Median estimator ($k = 10$) attack: sketch parameters: $l = 100, b = 30$



Bias-to-noise ratio (τ) for number of rounds, average of 10 simulations with different initializations (shaded region between the minimum and maximum)

Attack rounds vs. sketch size



Attack rounds versus sketch size $\ell = d/b$ to obtain bias-to-noise ratio $\tau \in \{1, 4\}$, averaged over 20 and 5 simulations respectively (shaded region between minimum and maximum)

Attack Analysis (simplified)

- Suppose for simplicity that we can get estimates for weights, and consider the following attack:

```

a ← 0
Repeat r times:
    Query value of 1 on a random tail z (w disjoint keys)
    If Val(1) < 0 then a ← a + z
Return a
    
```

- Estimates are computed as the median of $\ell = \frac{d}{b}$ buckets in which 1 participates. The value of every such bucket is a “binomial-like” RV consisting of the sum of w “bits” which are 1 or -1 each wp $\frac{1}{2b}$, and 0 otherwise. The standard deviation of the value of every bucket is $\approx \sqrt{w/b}$.
- In expectation, 50% the ℓ buckets in which 1 participates are positive and 50% are negative. The standard deviation of the difference between number of positives and number of negatives is $\approx \sqrt{\ell} = \sqrt{d/b}$
- When we collect a tail, there are more negatives than positives buckets
- Whp, we collected $R = \Omega(r)$ “good tails” in which there are $\sqrt{\ell}$ more negatives than positives
- Fix a bucket in which 1 participates, and sample a “good tail”. The bucket is negative wp $\geq \frac{\ell + \sqrt{\ell}}{\ell} = \frac{1}{2} + \frac{1}{\sqrt{\ell}}$
- After R such samples, the bucket is negative in $\geq \frac{R}{2} + \frac{R}{\sqrt{\ell}} \pm \sqrt{R}$ and is positive in $\leq \frac{R}{2} - \frac{R}{\sqrt{\ell}} \pm \sqrt{R}$
- Assuming $R \gg \ell$ we get about $R/\sqrt{\ell}$ extra negative times, which gives bias

$$\approx \frac{R}{\sqrt{\ell}} \cdot \sqrt{\frac{w}{b}}$$

% This is similar to the expectation of the sum of $(\ell - \frac{R}{\sqrt{\ell}})$ unbiased buckets, which is zero, and the sum of $\frac{R}{\sqrt{\ell}}$ positive buckets, each of which is $\approx \sqrt{\frac{w}{b}}$

- We want this bias to overtake the random noise in the bucket, which is $\approx \sqrt{Rw/b}$, which happens when $R \gg \ell$
 % This is the standard deviation of the sum of Rw “bits” which are ± 1 wp $\frac{1}{2b}$ and 0 otherwise
- That is, we caused CountSketch to estimate $\text{Val}(1) \gtrsim \sqrt{Rw/b}$ even though $\text{Val}(1) = 0$. When $R \gg \ell$, then this is more than the allowed error level of $\frac{1}{\sqrt{k}} \sqrt{R \cdot w}$.

PS: We need $w = w(\ell)$. If w is constant then the green text is incorrect, because the probability of bucket being zero is $\frac{1}{\sqrt{w}}$ meaning that a constant fraction of the buckets are zero, and hence so is the median, and we get $\text{Val}(1) = 0$ (in which case no tail is collected).

Warmup towards a robust construction

- Consider a single copy of CountSketch
- Instead of estimating values using the actual median of the buckets, return a DP-stable estimate of it
- So every time we return an estimate we aggregate $\frac{d}{b}$ buckets
- Thus, for releasing $O(T \cdot k)$ estimates overall we need to set $\frac{d}{b} \approx \sqrt{T \cdot k}$, so sketch size $\approx \frac{d}{b} \times b \approx \sqrt{T} \cdot k^{1.5}$
- So we didn't gain much in terms of the sketch size
- **Conceptual improvement:** In the generic construction, stability was w.r.t. a complete sketch. Now stability is w.r.t. one row, i.e. one pair of hash functions. Maintaining stability w.r.t. "smaller" objects is potentially easier.
- **Can we take this idea further?** Would it suffice to maintain stability w.r.t. a single bucket?

BCountSketch

- Recall that the generalization properties of DP-stability (as we have seen them in Theorem 8 in Lecture 6) hold when the data is sampled in an iid manner. However, In CountSketch, different buckets in the same row are slightly correlated.
- This is a minor issue and there is in fact a variant of CountSketch, called **BCountSketch**, where the buckets are iid:
 - For every bucket $j \in [d]$ we have a hash function $g_j: [n] \rightarrow \{-1, 0, 1\}$ where $\forall i \in [n]$
$$\Pr[g_j(i) = 0] = 1 - \frac{1}{b} \quad \text{and} \quad \Pr[g_j(i) = 1] = \Pr[g_j(i) = -1] = \frac{1}{2b}$$
 - Very similar performances (basically the same analysis)

Fine-tuned stability analysis

- Fix a key i and let B be a bucket that i participates in. As we have seen in the analysis of CountSketch (see page 3), in the oblivious case, with constant probability (say 99%) B provides an accurate estimation for $f[i]$ (up to the allowed error).
- Assuming that our algorithm maintains DP-stability w.r.t. the buckets, then by the generalization properties of DP-stability it follows that (whp) even in the adaptive case, 98% of the buckets in which i participates provide good estimates for it. **Good!**
- **How should we design a DP-stable algorithm w.r.t. the buckets?**
- Note that by working with the buckets as our "stability unit", we now have a dataset of size d instead of d/b as we had when working with the rows. **Good!**
 - (But recall that every key participates in $\approx d/b$ buckets, so we need to use this big dataset carefully because when estimating its value we still only have d/b elements to aggregate...)

- Fix a set H of k' keys (think of them as the keys reported as heavy at some time step). If we were to resample a bucket $j \in [d]$ (i.e., sample g_j), then in expectation this bucket contains $\frac{k'}{b} < 1$ keys from H
- Similarly, fix a sequence of sets H_1, \dots, H_T , each containing k' (possibly different) keys. If we were to resample a bucket $j \in [d]$ (i.e., sample g_j), then in expectation this bucket contains $\frac{k'}{b} < 1$ keys from each H_r and thus a total of $< T$ keys from H_1, \dots, H_T
- By Markov, the probability that a bucket contributes to more than $100T$ such keys is $< \frac{1}{100}$
- Suppose we delete a bucket once it contributes to $100T$ reported keys. Then in the oblivious setting the expected number of deleted buckets is at most $\frac{d}{100}$
- A similar argument shows that for every fixture of a key i , in the oblivious setting the expected number of buckets that are both deleted and contain i is at most $\approx \frac{d/b}{100}$
 - % To see this, note that the probability that a bucket B contains i and contains $100T$ keys from H_1, \dots, H_T is at most that probability that B contains i and contains $99T$ keys different than i from H_1, \dots, H_T , which happens w.p. at most $\frac{1}{b} \cdot \frac{1}{99} \approx \frac{1}{100b}$
- Assuming that our algorithm maintains DP-stability w.r.t. the buckets, then again by the generalization properties of DP-stability it follows that (whp) even in the adaptive case, 98% of the buckets in which i participates are not deleted. **Good!**
- **Midway summary:** As long as we maintain DP-stability w.r.t. the buckets, then (w.h.p.) for every key i throughout all of the execution:
 - (a) 98% of the buckets in which i participates provide good estimates for it
 - (b) 98% of the buckets in which i participates are not deleted
- We next leverage these properties to present an algorithm based on **ThresholdMonitor** from Lecture 8. Recall that in **ThresholdMonitor**, if we delete an element after its T 's contribution, then we need to add noise of magnitude $\approx \sqrt{T}$ to our counts.
- Our counts will be of the form: "How many buckets estimate the value of a key X to be more than some threshold Y ".
- As our noise is of the order of $\approx \sqrt{T}$, we need to ensure that the number of "alive and accurate" buckets containing the key X , which is roughly $0.98 \cdot 0.98 \cdot d/b$ is more than $\approx \sqrt{T}$
- This will mean that the size of our (robust) sketch needs to satisfy $d \gtrsim \sqrt{T} \cdot b \approx \sqrt{T} \cdot k$, instead of $\sqrt{T} \cdot k^{1.5}$ as we got from the generic construction

Algorithm Robust-BCountSketch

1. For every $j \in [d]$ sample a hash function $g_j: [n] \rightarrow \{-1, 0, 1\}$ and initiate an empty bucket $c_j = \mathbf{0}$
2. Initialize algorithm **ThresholdMonitor** on (g_1, \dots, g_d) with threshold $\approx \sqrt{T}$ and deletion parameter $\approx 100T$
3. When receiving an update (i, Δ) where $i \in [n]$ is the key and $\Delta \in \mathbb{R}$ is the weight:
 - a. For every $j \in [d]$ do $c_j \leftarrow c_j + g_j(i) \cdot \Delta$.
(Note that most buckets are unchanged as $g_j(i) = 0$)
4. When queried (at most T queries):
 - a. Set $\mathcal{T} = \emptyset$
 - b. For $v = m, (1 - \alpha)m, \dots, \mathbf{1}$ do: (here m bounds the length of the stream)
 - i. For $i \in [n]$ do:
 - Define the counting query $f_{v,i}: [d] \rightarrow \{0,1\}$ as

$$f_{v,i}(j) = \mathbf{1} \iff |g_j(i) \cdot c_j| \geq v$$
 (That is, $f(j) = \mathbf{1}$ if key i participates in the j th bucket and the value of this bucket is at least v in absolute value)
 - Submit the query $f_{v,i}$ to **ThresholdMonitor**. If the answer is T then add i to \mathcal{T} . If $|\mathcal{T}| = k' = O(k)$ then GOTO 4c
- c. Report \mathcal{H}