

Lecture 12: "Sampling"

Source: Lecture notes by
Aaron Roth and Adam Smith

Lecturer: Uri Stemmer

Recall the following algorithm and theorem (from Lecture 10), that operate in the oblivious streaming model.

Reservoir Sampling [Vitter, 1985]:

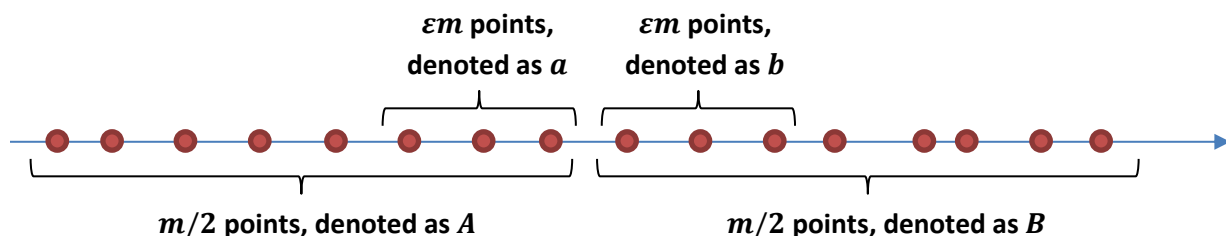
1. Init $s = \perp$
2. When the next item u_i is read, toss a biased coin and with probability $1/i$ let $s = u_i$ (note that we need to maintain both the element s and a counter for i)
3. Output s

Theorem: The output of the above algorithm is an element uniformly sampled from the stream. Specifically, if the stream is of length m , then for every $i \in [m]$ we have $\Pr[s = u_i] = \frac{1}{m}$.

Recall that we can maintain a sample of size k by simply running k copies of this algorithm in parallel.

Observation: In the oblivious model, if we maintain a sample of size $k = O\left(\frac{1}{\varepsilon^2}\right)$ then w.h.p. the median of the sample will be an ε -approximation to the median of the entire stream.¹

Proof sketch: Fix an input stream and consider the following plot where the red point correspond to the (sorted) elements in the stream:



Note that for a random sample from the stream, the probability of it falling in $A \cup b$ or in $a \cup B$ is exactly $\frac{1}{2} + \varepsilon$. Thus, by Chernoff/Hoeffding, w.h.p., a sample of size $k \gtrsim \frac{1}{\varepsilon^2}$ contains more than $\frac{k}{2}$ points from $A \cup b$ and more than $\frac{k}{2}$ points from $a \cup B$. Thus, the sample median cannot be outside of $a \cup b$. To see this, note that the sample median cannot be to the left of a because then there are more than $\frac{k}{2}$ sample points to the right of it.

¹ The term "approximate" here means that the median of the sampled set will be an element whose order among the elements of the full stream, when the elements are sorted by value from smallest to largest, is within the range $(1 \pm \varepsilon)m/2$ with high probability.

What happens in the adaptive/adversarial setting?

AdaptiveSamplingGame(A, M)

1. Instantiate the sampling algorithm M and let σ_0 denote its initial internal state
2. The adversary A gets σ_0
3. For $i = 1, 2, \dots, m$:
 - The adversary A chooses the next item in the stream u_i
 - The mechanism M gets u_i and updates its internal state to σ_i
 - The adversary A gets σ_i
4. At the end of the stream, the sampling algorithm M outputs a set of k elements

Will Reservoir Sampling “work” in the adaptive case in the sense that the sample median approximates the stream median?

Let k denote the sample size and let k' be the total number of elements that were sampled during the reservoir sampling process. That is, k' counts sampled elements that were possibly evicted at a future iteration. We bound k' as follows:

$$\mathbb{E}[k'] = \sum_{i=1}^m \frac{k}{i} = k \sum_{i=1}^m \frac{1}{i} \stackrel{\text{harmonic sum}}{\leq} k(\ln(m) + 1) \leq 2k \ln m$$

Thus, by Markov's inequality, with probability at least **0.9**, we will have $k' \leq 20k \ln m$. We next present an attack in which these k' elements are the smallest elements in the stream. The final sample set consists of some k elements among these k' elements. Hence, when the stream is of length $m \gg 2k'$, the sample median does not approximate the stream median.

The attack: Consider a setting where the stream consists of k points x_1, \dots, x_n in the one-dimensional range $[0, 1]$. The adversary keeps a “working range” at any point during the game, starting with the full range $[0, 1]$. In the first round, the adversary chooses the number $x_1 = 0.5$ as the first element in the stream. If x_1 is sampled, then the adversary moves to the range $[0.5, 1]$, and otherwise, to the range $[0, 0.5]$. Next, the adversary submits x_2 as the middle of the current range. This continues for m steps.

Note that at any point throughout the process, the adversary always submits an element that is larger than all elements in the current sampled set, and also smaller than all the non-sampled elements of the stream. Therefore, the end result is that after this process is over, the k' sampled elements are precisely the smallest k' elements in the stream.

An adaptive “sampling” algorithm for the median

Define a “merge” operation on two sets A, B of size k each: Sort $A \cup B$ and return the set C obtained by deleting every second item in the ordering. We think of $A \cup B$ as a *multiset* containing $2k$ items. Note that $|C| = k$.

Also note that for every number t it holds that

$$2 \sum_{x \in C} \mathbf{1}_{\{x \leq t\}} \in \sum_{x \in A \cup B} \mathbf{1}_{\{x \leq t\}} \pm \mathbf{1} \in \sum_{x \in A \cup B} \mathbf{1}_{\{x \leq t\}} \pm \varepsilon |A \cup B| \quad ((1))$$

Where the last transition follows by requiring $k \geq \frac{2}{\varepsilon}$, where ε is the approximation parameter.

Consider a complete binary tree whose nodes correspond to buckets of size k that are initially empty. Let \mathbf{p} be a pointer to the “current leaf” of the tree, initiated as the left-most leaf. Throughout the execution, the elements in the stream are added to \mathbf{p} 's bucket until it is full (containing k elements), at which point we move to the next leaf (meaning that \mathbf{p} now points to the next leaf). Whenever two sibling nodes in the tree are full, we “merge” them and place the resulting set of size k in their parent node.

Let S denote the input stream and let $m = |S|$ be the length of the stream. Note that in this case our tree has m/k leaves and its height is $\log \frac{m}{k}$.

Claim: Let $\vec{B} = (B_1, \dots, B_j)$ denote the buckets in level ℓ of the tree (where the leaves are at level 0). Then for every number t we have

$$2^\ell \cdot \sum_{x \in \vec{B}} \mathbf{1}_{\{x \leq t\}} \in \sum_{x \in S} \mathbf{1}_{\{x \leq t\}} \pm \varepsilon m \ell$$

Proof: The proof is by induction on ℓ . The base-case for $\ell = 0$ is trivial. Suppose that the claim holds for $\ell - 1$. For every bucket B_i in level ℓ , let B_i^0, B_i^1 denote its two children (which are buckets at level $\ell - 1$). Then,

$$\begin{aligned} 2^\ell \cdot \sum_{x \in \vec{B}} \mathbf{1}_{\{x \leq t\}} &= 2^{\ell-1} \left(2 \sum_{x \in B_1} \mathbf{1}_{\{x \leq t\}} + \dots + 2 \sum_{x \in B_j} \mathbf{1}_{\{x \leq t\}} \right) \\ &\stackrel{\text{by } ((1))}{\leq} 2^{\ell-1} \left(\left(\sum_{x \in B_1^0 \cup B_1^1} \mathbf{1}_{\{x \leq t\}} + \varepsilon |B_1^0 \cup B_1^1| \right) + \dots + \left(\sum_{x \in B_j^0 \cup B_j^1} \mathbf{1}_{\{x \leq t\}} + \varepsilon |B_j^0 \cup B_j^1| \right) \right) \\ &= 2^{\ell-1} \left(\sum_{x \in B_1^0 \cup B_1^1} \mathbf{1}_{\{x \leq t\}} + \dots + \sum_{x \in B_j^0 \cup B_j^1} \mathbf{1}_{\{x \leq t\}} \right) + 2^{\ell-1} \varepsilon \cdot (|B_1^0 \cup B_1^1| + \dots + |B_j^0 \cup B_j^1|) \end{aligned}$$

$$\begin{aligned}
& \stackrel{\{\leq\}}{* \text{ see below}} 2^{\ell-1} \left(\sum_{x \in B_1^0 \cup B_1^1} \mathbf{1}_{\{x \leq t\}} + \dots + \sum_{x \in B_j^0 \cup B_j^1} \mathbf{1}_{\{x \leq t\}} \right) + 2^{\ell-1} \varepsilon \cdot \frac{m}{2^{\ell-1}} \\
&= 2^{\ell-1} \left(\sum_{x \in B_1^0 \cup B_1^1} \mathbf{1}_{\{x \leq t\}} + \dots + \sum_{x \in B_j^0 \cup B_j^1} \mathbf{1}_{\{x \leq t\}} \right) + \varepsilon \cdot m \\
&\stackrel{\text{induction}}{\{\leq\}} \sum_{x \in S} \mathbf{1}_{\{x \leq t\}} + \varepsilon m (\ell - 1) + \varepsilon \cdot m \\
&= \sum_{x \in S} \mathbf{1}_{\{x \leq t\}} + \varepsilon m \ell
\end{aligned}$$

The transition marked with * follows from the fact that $|B_1^0 \cup B_1^1| + \dots + |B_j^0 \cup B_j^1|$ is the width of level $(\ell - 1)$. There are $\frac{m/k}{2^{\ell-1}}$ buckets in this level, each of size k , so overall $|B_1^0 \cup B_1^1| + \dots + |B_j^0 \cup B_j^1| = \frac{m}{2^{\ell-1}}$.

This completes the proof of the claim.

Now let R denote the bucket in the root of the tree (the level of the root is $\ell = \log(m/k)$). By the above claim, for every number t it holds that

$$2^{\log(m/k)} \cdot \sum_{x \in R} \mathbf{1}_{\{x \leq t\}} \in \sum_{x \in S} \mathbf{1}_{\{x \leq t\}} \pm \varepsilon m \cdot \log(m/k)$$

That is,

$$\frac{1}{k} \sum_{x \in R} \mathbf{1}_{\{x \leq t\}} \in \frac{1}{m} \sum_{x \in S} \mathbf{1}_{\{x \leq t\}} \pm \varepsilon \cdot \log(m/k)$$

Therefore, by setting $k \gtrsim \frac{1}{\varepsilon}$ we ensure that the root provides a good approximation to the median of the stream (we need to set k to be slightly bigger than $\frac{1}{\varepsilon}$ in order to cancel out this $\log(m/k)$ factor).

Question: What is the space complexity of this algorithm?

Balls into bins

Suppose we throw n balls into n bins uniformly at random. The expected load of every bin is 1. A simple application of the Chernoff bound shows that w.h.p. the maximal load would be $O(\log n)$. In fact, classical results w.r.t. this problem show that w.h.p. the maximal load is bounded by $O\left(\frac{\log n}{\log \log n}\right)$.

Proof (even for $2n$ balls and n bins):

The probability that bin 1 receives at least M balls is at most

$$\binom{2n}{M} \left(\frac{1}{n}\right)^M < \left(\frac{e2n}{M}\right)^M \cdot \left(\frac{1}{n}\right)^M = \left(\frac{2e}{M}\right)^M$$

This follows from a union bound; there are $\binom{2n}{M}$ distinct sets of M balls, and for any set of M balls the probability that all land in bin 1 is $\left(\frac{1}{n}\right)^M$.

Applying a union bound again allows us to find that, for $M \geq \frac{6 \ln n}{\ln \ln n}$, the probability that any bin receives at least M balls is bounded above by

$$\begin{aligned} 2n \left(\frac{2e}{M}\right)^M &\leq 2n \left(\frac{2e \ln \ln n}{6 \ln n}\right)^{\frac{6 \ln n}{\ln \ln n}} \leq 2n \left(\frac{\ln \ln n}{\ln n}\right)^{\frac{6 \ln n}{\ln \ln n}} = e^{\ln(2n)} \cdot \left(e^{\ln \ln \ln n - \ln \ln n}\right)^{\frac{6 \ln n}{\ln \ln n}} \\ &= e^{\ln(2n)} \cdot e^{6 \ln n \cdot \ln \ln \ln n / \ln \ln n} \cdot e^{-6 \ln n} \leq \frac{1}{n} \end{aligned}$$

for n sufficiently large.

Now suppose we play the following adaptive game with an adversary:

AdaptiveBallsBins(n)

1. Randomly throw n balls into n bins
2. The adversary A gets to see the resulting allocation
3. For $i = 1, 2, \dots, n/2$:
 - The adversary A chooses a bin B_i , and gains the reward $r_i = \mathbf{Load}(B_i)$
 - The B_i is now excluded from the rest of the game and its balls are thrown randomly among the remaining bins
 - The adversary A gets to see the resulting allocation
4. The total reward of the adversary is $r_1 + r_2 + \dots + r_{n/2}$

In the oblivious setting, the sequence of bin deletions is fixed before the game begins. In this case, for every step $i \in [n]$ of the game, we have that the balls are distributed uniformly across the remaining bins, and thus the maximal load is bounded by $O\left(\frac{\log n}{\log \log n}\right)$ w.h.p., meaning that the total reward of the adversary is bounded by $O\left(n \cdot \frac{\log n}{\log \log n}\right)$. Actually, the total reward would be even smaller (since when the adversary eliminates the maximum-load bin then the second-highest bin would have significantly lower load).

What happens in the adaptive setting? Can the adversary somehow effect the distribution of the balls in a way that allows him to significantly increase its total reward?

Let us reconsider the oblivious case when the sequence of bin deletions is fixed before the game begins. To simplify the notations, let us assume that the order of bin deletions is simply Bin 1, Bin 2, Bin 3, ..., Bin $n/2$.

Let us focus on a single ball b and bound the reward it generates during its movements. As we have fixed the order of bin-deletions, the reward generated by ball b is exactly the number of times it is thrown into one of the bins $1, 2, \dots, n/2$. This is a Geometric RV with success probability 0.5.

Thus, the total reward across all n balls throughout the execution is distributed as the sum of n Geometric RV's with parameter 0.5 (this is known as the "Negative binomial" distribution). We will use the following concertation bound:

Theorem: Let $GS(n, p)$ be a negative binomially distributed random variable which is the sum of n i.i.d. geometrically distributed random variables, each with success probability p . Then, for $k > 1$

$$\Pr \left[GS(n, p) > \frac{k \cdot n}{p} \right] \leq \exp \left(- \frac{k(1 - 1/k)^2}{2} \cdot n \right)$$

Using this theorem, for our fixed ordering, we get that

$$\begin{aligned} \Pr \left[\text{total}_{\text{reward}} > 16n \log n \right] &= \Pr \left[GS(n, 0.5) > \frac{8 \cdot \ln(n) \cdot n}{0.5} \right] \leq \exp \left(- \frac{8 \cdot \ln(n) \left(1 - \frac{1}{4 \ln(n)}\right)^2}{2} \cdot n \right) \\ &\leq \exp(-2 \cdot n \cdot \ln(n)) \quad ((*) \end{aligned}$$

That is, for every fixed ordering of the deletions of the bins, with very high probability it holds that the total reward is bounded by $O(n \cdot \log n)$.

How many orders could there be? At most $n! \leq n^n$.

Note that that failure probability we derived in ((*)) is so low that it allows us to “union bound” over all possible orderings. But what does this mean exactly? What is the probability space that ties all of these orderings together?

We consider a simultaneous game against all possible choices of the adaptive adversary using a rooted game tree T . Each node in T represents an assignment of the balls into the bins which have not yet been deleted. The root represents the initial assignment of the balls into the n bins. A node v which has ℓ remaining bins (at depth $n - \ell$ of T) has ℓ children, each corresponds to a possible decision of the adversary (delete the first among remaining bins, the second, etc). Each edge from v to a child u is associated with sufficiently many random bits to redistribute the balls in the bin which the adversary deletes when it faces the configuration associated with v . Each leaf corresponds to a configuration with $n/2$ bins containing all the balls. The tree has at most $n!$ distinct leaves. The path from the root to a leaf corresponds to an ordering of the $n/2$ deleted bins (representing a complete deletion order of the bins).

The probability space of this simultaneous game corresponds to sampling all of the random bits on all the edges of the tree T . Once we flip all these bits then the simultaneous game is completely determined and in particular the reward associated with each leaf (deletion ordering) is determined.

The bound we derived in ((*)) bounds the probability that the recourse is high at a specific leaf. By a union bound,

$$\Pr \left[\begin{array}{l} \text{There is an order} \\ \text{of deletion with} \\ \text{reward greater} \\ \text{than } 16n \log n \end{array} \right] \leq n! \cdot \exp(-2 \cdot n \cdot \ln(n)) \leq \exp(n \ln(n) - 2n \ln(n)) = \exp(-n \ln(n))$$

Now, if the reward is small in every leaf then the adversary cannot gain a high reward.