

Algorithm Design 21/22

Hands On 8 - Count-Min Sketch

Federico Ramacciotti

1 Problem

Consider the counters $F[i]$ for $1 \leq i \leq n$, where n is the number of items in the stream of any length. At any time, we know that $\|F\|$ is the total number of items (with repetitions) seen so far, where each $F[i]$ contains how many times item i has been so far. We saw that CM-sketches provide an FPTAS $F'[i]$ such that $F[i] \leq F'[i] \leq F[i] + \varepsilon \|F\|$, where the latter inequality holds with probability at least $1 - \delta$.

Consider now a range query (a, b) , where we want $F_{ab} = \sum_{a \leq i \leq b} F[i]$. Show how to adapt CM-sketch so that an FPTAS F'_{ab} is provided:

- Baseline is $\sum_{a \leq i \leq b} F'[i]$, but this has drawbacks as both time and error grows with $b - a + 1$.
- Consider how to maintain counters for just the sums when $b - a + 1$ is any power of 2 (less or equal to n):
 - Can we now answer quickly also when $b - a + 1$ is not a power of two?
 - Can we reduce the number of these power-of-2 intervals from $n \log n$ to $2n$?
 - Can we bound the error with a certain probability? *Suggestion*: it does not suffice to say that it is at most δ the probability of error of each individual counter; while each counter is still the actual wanted value plus the residual as before, it is better to consider the sum V of these wanted values and the sum X of these residuals, and apply Markov's inequality to V and X rather than on the individual counters.

2 Solution

2.1 Baseline solution

The baseline solution uses the general Count-Min Sketch of the stream and check the entire interval $[a, b]$ for every query. This solution requires $b - a + 1$ time, so to reduce it to logarithmic time check the solution below.

2.2 Better solution

In a better solution we can store the counters for $\log n$ positions. For example, at position 0, we store the counters for the intervals $[0, 2], [0, 4], [0, 8]$ and so on. As we can see, the space used to store all the counters for all the n positions in the array is $n \log n$. We can reduce the number of counters using a different concept of intervals, as explained below.

2.3 Best solution

The best solution uses ranges with pairs of elements, creating a sort of tree upon the array F . For example, in an array with length 8, we compute the sum of the pairs $[0, 1], [2, 3], [4, 5], [6, 7]$ at the first level, then the sum of the intervals $[0, 3], [4, 7]$ doing the sum of the lower level pairs, and so on, with the last level being the sum of the whole array. The height of the tree is therefore $O(\log n)$.

In this way the error is no more linear in $b - a + 1$, as we walk only the tree in its height in $O(\log n)$ and not the whole interval every time. To store the ranges (namely, the levels of the tree) we use a Count-Min Sketch for each one of them; thus we use $O(\log n)$ sketches for a total of

$O(r * c * \log n)$ space. Updating an element recursively updates the smaller ranges, with the last range being the complete Count-Min Sketch of the actual stream. Querying an interval takes the biggest interval from the first element and reduces the interval iteratively, e.g. with an array of 8 elements and the query $[3, 7]$, we take the first interval $[3, 4]$, then $[5, 6]$ and finally $[7, 7]$ and sum them up.

This solution reduces the error because using only one Count-Min Sketch for the stream and the counters increases the error probability (more elements to map in the same space). We can also use less space by using only one sketch when we get to \sqrt{n} counters, since it is useless for the space to use a new Count-Min Sketch when we have for example 2 or 4 counters.

2.3.1 Error bounding with Markov's inequality

This analysis differs slightly from the solution given above, since we consider using only one Count-Min Sketch for the whole intervals, instead of one per level.

Let's define the set D of dyadic ranges (pair of indices that define an interval), e.g. $[1, n], [1, n/2], [n/2 + 1, n], \dots, [1, 1], [2, 2], \dots, [n, n]$ with $|D| = 2n$.

With D , we have

$$F'_{ab \in D} = F_{ab \in D} + \sum_{gh \in D, gh \neq ab} X_{gh}$$

For notation purposes, we use F_i or $F(i)$ or F_{ab} to refer to F on the intervals i of $[a, b]$.

Given an hash function j in the Count-Min Sketch, define an indicator variable

$$I_{jik} = \begin{cases} 1 & \text{if } h_j(i) = h_j(k) \\ 0 & \text{otherwise} \end{cases}$$

that states if there is a collision on the intervals $i \in D$ and $k \in D$.¹

Now we can see that

$$X_{ji, i \in D} = \sum_{k \in D, k \neq i}^{2n} I_{jik} F_k$$

and its expected value is

$$\begin{aligned} E[X_{ji}] &= E \left[\sum_{k \in D, k \neq i} I_{jik} F_k \right] \\ &= \sum_{k \in D, k \neq i} E[I_{jik} F_k] \\ &= \sum_{k \in D, k \neq i} Pr[I_{jik} = 1] F_k \\ &\leq \sum_{k \in D, k \neq i} \frac{\varepsilon}{e} F_k \\ &\leq \log n \frac{\varepsilon}{e} \|F\| \end{aligned}$$

with $\|F\| = \sum_{ab} F_{ab}$ the sum of all the counters of a level.

In the last pass of the equation above $\log n$ comes from the fact that the sum of all counters is $\|F\|$ and, with the counters "tree" of height $\log n$, it becomes $\log n \|F\|$.

¹Note that here the subscript j refers to the hash function, not the interval: this is the only exception, since every other subscript $x \neq j$ refers to the interval x .

Finally, using Markov's inequality on r hash functions we get

$$\begin{aligned}
Pr[\forall j \in [r] : F'_{ab} \geq F_{ab} + \log n \varepsilon \|F\|] &= \prod_{j=0}^{r-1} Pr\left[F_{ab} + \sum X_{gh} \geq F_{ab} + \log n \varepsilon \|F\|\right] \\
&= \prod_{j=0}^{r-1} Pr\left[\sum X_{gh} \geq \log n \varepsilon \|F\|\right] \\
&\leq \prod_{j=0}^{r-1} \frac{E[\sum X_{gh}]}{\log n \varepsilon \|F\|} \\
&\leq \prod_{j=0}^{r-1} \frac{\log n \frac{\varepsilon}{e} \|F\|}{\log n \varepsilon \|F\|} \\
&= \prod_{j=0}^{r-1} \frac{1}{e} = \left(\frac{1}{e}\right)^r \leq \delta
\end{aligned}$$

In conclusion, we have shown that we provided an FPTAS with a total error probability bounded to δ .