

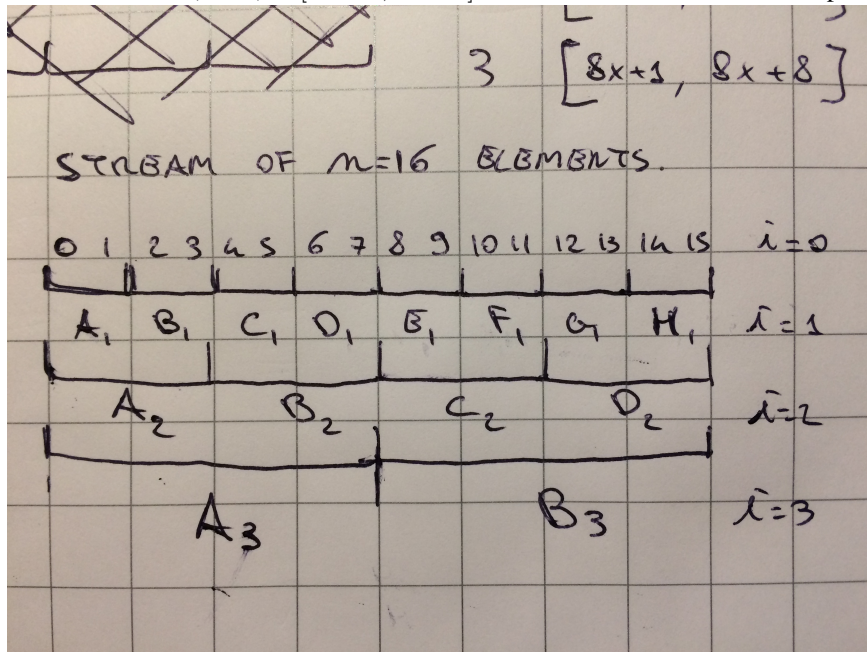
Count-min sketch: range queries

Show and analyse the application of count-min sketch to range queries (i, j) for computing $\sum_{k=i}^j F[k]$. Hint: reduce the latter query to the estimate of just $t \leq 2 \log n$ counters c_1, c_2, \dots, c_t . Note that in order to obtain a probability at most δ of error (i.e. that $\sum_{l=1}^t c_l > \sum_{k=i}^j F[k] + 2\epsilon \log n ||F||$), it does not suffice to say that it is at most δ the probability of error of each counter c_l : while each counter is still the actual wanted value plus the residual as before, it is better to consider the sum V of these t wanted values and the sum X of these residuals, and apply Markov's inequality to V and X rather than on the individual counters.

SOLUTION

We use the Dyadic Interval: a dyadic interval is a range of the form $[x2^i + 1, (x+1)2^i]$ for parameters x and i . Each point in the range $[1 \dots n]$ is a member of $\log_2 n$ dyadic intervals, one for each i in the range $0 \dots \log_2(n) - 1$ (think it as a tree, where we split the interval in 2 at each level).

An example: let's say I have a stream of $N = 16$ different items, this means I need to use $\log N = 4$ CMSs, one for each $i \in \{0, 1, 2, 3\}$. So the first CMS ($i = 0$) will correspond to the interval $[x+1, x+1]$, the second CMS ($i = 1$) to $[2x+1, 2x+2]$ and so on. The elements will be partitioned like in the picture:



The sketch corresponding to $(i = 3)$ will count occurrences of A_3 and B_3 (will count an A_3 every time a number in $\{0, 1, 2, 3, 4, 5, 6, 7\}$ arrives, and a B_3 when a $\{8, 9, 10, 11, 12, 13, 14, 15\}$ arrives.); the $(i = 2)$ sketch will count A_2, B_2, C_2, D_2 and so on. Then for each arbitrary interval (a, b) we just need to sum the corresponding dyadic intervals that makes it. (e.g. if I receive a query for $(7, 13)$, then I can just sum up the occurrences of C_2, G_1 and "7").

A count-min sketch table is kept for each set of dyadic intervals of length 2^i , one for each level in the tree (one for each value of i). Thus we have $\log_2 n$ vectors \tilde{F}_i and, in line of principle, $\log_2 n$ vectors F_i .

You can see this as associating each new sketch to a reencoding of the original stream, where you have partitioned the alphabet in 2^i subsets of symbols and represented each of them with a single new symbol. Say: ABCD as alphabet, map $AB \rightarrow E, CD \rightarrow F$; each time you witness A or B you count an occurrence of E , each time you witness C or D you count an occurrence of F .

Witnessing a new element in the stream will therefore trigger an update to all the $\log_2 n$ tables.

The idea for a range query is to partition the range into dyadic intervals (of a number of sizes, in general, thus from a number of CMS tables) and return as result the sum of the values stored in the CMS tables for the corresponding intervals [see this [link](#) for further reference].

It can be shown¹ that any range will be split at most into $2\log_2 n$ dyadic intervals.

Therefore for each query we access $t \leq 2 \log n$ counters c_1, c_2, \dots, c_t . We then, want to proof $Pr[\sum_{l=1}^t c_l > \sum_{k=l}^j F[k] + 2\epsilon \log n ||F||] < \delta$.

Firstly, we notice that each counter c_i represents an interval $[a, b]$ and its value is $\sum_{k=a}^b (F[k]) + X_i$, where X_i is the rubbish of that particular interval. Then, give an interval $[l, r]$, and the counters we have $\sum_{i=1}^t c_i = \sum_{k=l}^r F[k] + X$, where X is the total error (notice that this works because the intervals are disjoint). Then we substitute to the previous equation and we have:

$$\begin{aligned} Pr[\sum_{l=1}^t c_l > \sum_{k=l}^r F[k] + 2\epsilon \log n ||F||] &< \delta \\ Pr[\sum_{k=l}^r F[k] + X > \sum_{k=l}^r F[k] + 2\epsilon \log n ||F||] &< \delta \\ Pr[X > 2\epsilon \log n ||F||] &< \delta \end{aligned}$$

Now, we apply Markov inequality and by the linearity of the expectation we have:

$$Pr[X > 2\epsilon \log n ||F||] \leq \frac{E[X]}{2\epsilon \log n ||F||} \leq \frac{\sum_i^{2\log_2(n)} E[X_i]}{2\epsilon \log n ||F||}$$

Now, let $d(i)$ be the depth in the tree to which the dyadic interval i belongs (identifying its CMS table). We see that $\forall i. E[X_i] < \frac{\epsilon}{e} ||F_{d(i)}||_1$, but $\forall x, y. ||F_x||_1 = ||F_y||_1 = ||F||_1 \implies E[X_i] < \frac{\epsilon}{e} ||F||_1$.

(by definition $||A||_1 = \sum_{h=1}^n |A[h]|$, that is in our case we sum the frequencies of symbols, which sum up to the same total amount no matter how we group them with a reencoding)

Thus we have

$$\frac{\sum_i^{2\log_2(n)} E[X_i]}{2\epsilon \log n ||F||} \leq \frac{2\frac{\epsilon}{e} \log n ||F||}{2\epsilon \log n ||F||} = \frac{1}{e}$$

This is the probability of error (that the sum of garbage is more than $2\epsilon \log n ||F||$) in row j (the min, potentially distinct for each table).

Since we choose the row of each table that minimize the sum of the counter (by definition), then there must be an error in all the row r . Thus we have $\frac{1}{e^r} = \delta$ since $r = \ln(\frac{1}{\delta})$.

A little example: suppose we have $n = 16$ (item), then a range $[1 \dots 16]$. Thus, we query $(8, 12)$, this interval is fully included, then we split in $[1 \dots 8]$ and $[9 \dots 16]$ (we go down in the tree). $[1 \dots 8]$ is still too big then we split in $[1 \dots 4]$ and $[5 \dots 8]$, $[1 \dots 4]$ is not included then we take it off. $[5 \dots 8]$ for the same reason we split it in $[5 \dots 6]$ and $[7 \dots 8]$ and so on.

¹First, note that:

- No three intervals of the same length can be contained in the partition of the same query, otherwise you could merge two of them;
- If there are two consecutive intervals of the same length, they must belong to two different larger "parent" intervals (otherwise you could replace them with their parent), i.e. their union cannot belong to the dyadic partition.

Thus, there are at most 2 intervals of each possible size in the partition. As the sizes are $\log_2 n$ in all, the partition is made up of at most $2\log_2 n$ intervals.