

*Lecture 16*  
*The Count-Min Sketch*

Lecturer: N.R.Aravind

Scribe: N.R.Aravind

In the previous class, we saw the Misra-Gries algorithm for finding heavy hitters (elements with frequency at least  $\epsilon n$  in a stream of length  $n$ ). It used  $O(\frac{1}{\epsilon} \log m)$  space where the domain of the elements is  $\{1, 2, \dots, m\}$ .

We now see the Count-Min algorithm, which is due to Cormode and Muthukrishnan (2003). It gives an estimate of the frequencies of every element, and makes an error of at most  $\epsilon n$  on each estimate.

1. Choose  $t$  hash functions  $h_1, \dots, h_t : \{1, 2, \dots, m\} \rightarrow \{1, 2, \dots, k\}$ , each from a pairwise independent hash family.
2. For  $1 \leq i \leq t$  and  $1 \leq j \leq k$ , initialize counters  $C[i, j]$  to zero. The counters  $C[i, 1], \dots, C[i, k]$  correspond to the  $i$ th hash function.
3. For each element  $x$ , do: set  $i_1 = h_1(x), \dots, i_t = h_t(x)$ . Increment each of  $C[1, i_1], C[2, i_2], \dots, C[t, i_t]$  by one.
4. For each  $y \in \{1, \dots, m\}$ , output  $\min\{C[1, h_1(y)], C[2, h_2(y)], \dots, C[t, h_t(y)]\}$  as the estimated frequency of  $y$ .

**Analysis:** We will choose the values of  $t$  and  $k$ , based on the analysis.

Fix  $y \in \{1, \dots, m\}$  and  $i \in \{1, 2, \dots, t\}$ . We have  $C[i, y] \geq f(y)$  and  $E[C[i, h(y)]] = f(y) + \frac{n - f(y)}{k} \leq f(y) + \frac{n}{k}$ . Thus, by applying Markov's inequality, we get:

$$\Pr(C[i, y] - f(y) \geq \epsilon n) \leq \frac{1}{\epsilon k}.$$

We now choose  $k = \lceil \frac{2}{\epsilon} \rceil$  so that the above probability is at most  $1/2$ . We choose  $t = \lceil (1 + c) \log m \rceil$  so that we get:

$$\Pr(\text{Min}(C[1, h_1(y)], \dots, C[t, h_t(y)]) \geq f(y) + \epsilon n) \leq \frac{1}{2^t} \leq \frac{1}{m^{1+c}}.$$

Finally, by applying the union bound on all  $m$  elements, we can bound the total probability of error by at most  $\frac{1}{m^c}$ . The choice of  $c = \frac{1}{\log m} \log \left( \frac{1}{\delta} \right)$  will bring the error to at most  $\delta$ .

Thus, the total space complexity is  $O \left( \frac{\log n}{\varepsilon} (\log m + \log \left( \frac{1}{\delta} \right)) \right)$ , where  $\log n$  is for the space per counter and the remaining factors are the total number of counters.