

Boosting the probability Cont'd

Analysis

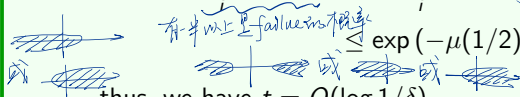
Define

$$Y_i = \begin{cases} 1, & \text{if } |\hat{f}_a - f_a| \geq \epsilon \|f\|_2; \\ 0, & \text{otherwise.} \end{cases}$$

- For $k = O(1/\epsilon^2)$, we have $P(Y_i = 1) < \frac{1}{3}$.
- Note that $\mu = E(\sum_i Y_i) \leq \frac{t}{3}$. Then by the Chernoff bound,

\Rightarrow t 次中估计失败次数

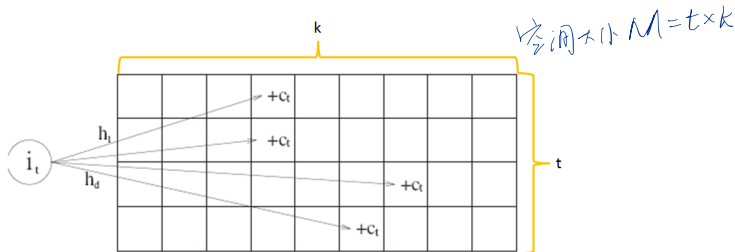
$$P(\text{media far}) \leq P(\sum_i Y_i > \frac{t}{2}) \leq P(\sum_i Y_i > (1 + \frac{1}{2})\mu) \leq \exp(-\mu(1/2)^2/4) < \exp(-t/48) < \delta,$$



thus, we have $t = O(\log 1/\delta)$.

- Finally, we can get an (ϵ, δ) -approximation in space complexity $O(\frac{\log 1/\delta}{\epsilon^2})$ counters.

Count min or Cormode-Muthukrishnan sketch

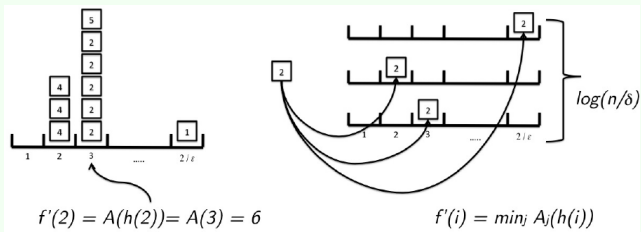


Algorithm

- 1: $C[1 \dots t][1 \dots k] \leftarrow \vec{0}$, where $k = \frac{2}{\epsilon}$ and $t = \lceil \log(1/\delta) \rceil$;
 - 2: Choose t independent hash functions $h_1, h_2, \dots, h_t : [n] \rightarrow [k]$;
- Process** item (j, c) , where $c = 1$:
- 3: for $i = 1$ to t do $C[i][h_i(j)] \leftarrow C[i][h_i(j)] + c$;
- Output:**
- 4: On query a , report $\hat{f}_a = \min_{1 \leq i \leq t} C[i][h_i(a)]$;
- Handwritten blue notes: "没有 $g(x)$ 值只会变大" (no $g(x)$ value will only increase), "发生 collision" (collision occurs), "并非无偏估计" (not an unbiased estimate).

CM sketch analysis

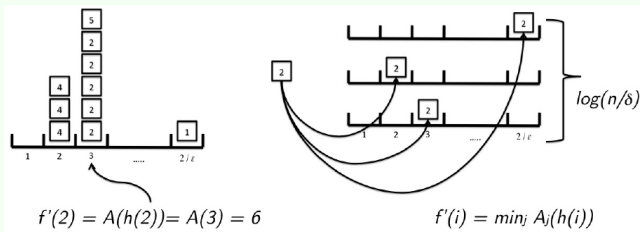
Analysis



- Clearly, for each i , we immediately have $f(a) \leq \text{count}[i, h_i(a)]$. However, the bound may be poor.

CM sketch analysis

Analysis



- Clearly, for each i , we immediately have $f(a) \leq \text{count}[i, h_i(a)]$. However, the bound may be poor.
- To get a better estimator, we will take the minimum over all the rows in count.

CM sketch analysis cont.

Analysis

- For a fixed a , we now analyze the collision in one such counter, say in $\text{count}[i, h_i(a)]$. Let r.v. X_i denote this collision.
- For $j \in [n] \setminus \{a\}$, let

$$Y_{i,j} = \begin{cases} 1, & \text{if } h_i(j) = h_i(a); \\ 0, & \text{otherwise.} \end{cases}$$

be the indicator of the event $h_i(j) = h_i(a)$. Notice that j makes a contribution to the counter iff $Y_{i,j} = 1$ (Note that $E(Y_{i,j}) = \frac{1}{k}$).

- Thus, we have $X_i = \sum_{j \in [n] \setminus \{a\}} f_j Y_{i,j}$. By linearity of expectation,

$$E[X_i] = X_i = \sum_{j \in [n] \setminus \{a\}} \frac{f_j}{k} = \frac{\|f\|_1 - f_a}{k} = \frac{\|f_{-a}\|_1}{k}.$$

高维的 (all elements except a's count) 与 f_a 没有关系 不是无偏估计, 无法用切比雪夫.

- Since each $f_j \geq 0$, we have $X_i \geq 0$, and we can apply Markov's inequality to get (by choosing the value of k)

每一行: $P[X_i \geq \epsilon \|f\|_1] \leq P[X_i \geq \epsilon \|f_{-a}\|_1] \leq \frac{\|f_{-a}\|_1}{k \epsilon \|f_{-a}\|_1} = \frac{1}{2}.$

取适当的 k 使 $\dots = \frac{1}{2}$ $k = \frac{1}{\epsilon}$

CM sketch analysis cont.

Analysis

- The above probability is for one counter. We have t such counters, mutually independent. The excess in the output $\hat{f}_a - f_a$, is the minimum of excesses X_i over all $i \in [t]$. Thus

$$\begin{aligned}
 P[\hat{f}_a - f_a \geq \epsilon \|f\|_1] &\leq P[\hat{f}_a - f_a \geq \epsilon \|f_a\|_1] \\
 &= P[\min\{X_1, \dots, X_t\} \geq \epsilon \|f_a\|_1] = \prod_{i=1}^t P[X_i \geq \epsilon \|f_a\|_1] \leq \frac{1}{2^t} \cdot \delta
 \end{aligned}$$

$\Rightarrow t > \log_2 \frac{1}{\delta}$

- Using our choice of t , this probability is at most δ . Thus, we have shown that, with high probability,

$$f_a \leq \hat{f}_a \leq f_a + \epsilon \|f_a\|_1$$

只增不减, 不用 $f_a - \epsilon \|f_a\|_1$

- Thus the space requirement is therefore $M = O\left(\frac{\log 1/\delta}{\epsilon}\right)$ counters.

Take-home messages

- Data streaming
- Deterministic algorithm
- Randomized algorithm
 - Naive sampling
 - Count sketch
 - Count min sketch