

27 Apr 2022

# Sketching: Count-Min & Count Sketches

Stream of tokens  $a_1, \dots, a_n$ .

- Each  $a_i \in [m]$

Algorithm maintains data structure of size  $S$  bits.

Goal:  $S = O(\text{poly}(\log n, \log m))$ .

Frequency vector of the stream:  $\vec{f} \in \mathbb{R}^m$

$f_j = \#$  of times token  $j$  appears in stream

$\in \{0, \dots, n\}$ .

$$\|\vec{f}\|_1 = n.$$

Two algorithms today.

Name	Space Required	Approximation Error	Probability
CountMin	$O\left(\frac{\log(m) \log(1/\delta)}{\epsilon}\right)$	$\leq \epsilon n = \epsilon \ \vec{f}\ _1$ <i>one-sided</i>	$1 - \delta$
CountSketch	$O\left(\frac{\log(m) \log(1/\delta)}{\epsilon^2}\right)$	$\leq \epsilon \ \vec{f}\ _2$ <i>two-sided</i>	$1 - \delta$

## Count Min

Given positive integers  $B, t$ .

Sample  $t$  independent hash functions  $h_1, \dots, h_t: [m] \rightarrow [B]$

from a 2-universal hash distribution.

Initialize a 2-D array  $C$ , size  $B \times t$ ,

$$C[k, l] = 0 \quad \forall k, l$$

*# of hash buckets*

*# of hash functions*

$h$  is 2-universal if

$$\forall x \neq y \in [m]$$

$$\forall \bar{x}, \bar{y} \in [B]$$

$$P_r(h(x) = \bar{x} \wedge h(y) = \bar{y}) = \frac{1}{B^2}.$$

//  $C[k, l]$  counts # stream elements  $a_i$   
s.t.  $h_l(a_i) = k$ .

for each  $a_i$  in the stream:

for  $l \in [t]$ :

let  $k = h_l(a_i)$

increment  $C[k, l]$ .

end for

end for

// Done processing stream.

To answer Freq Query( $x$ ):

return  $\min_{l \in [t]} \{ C[h_l(x), l] \}$ .

Analysis. Space Requirement

Store  $t$  hash function descriptions each  
requiring  $O(\log m)$  bits.

$$O(t \log m)$$

Store  $B \times t$  array of counters in  $\{0, \dots, n\}$ .

$$O(Bt \log n)$$

Approximate Correctness

$$\text{For any given } l, \quad E[C[h_l(x), l]] = f_x + \sum_{y \neq x} f_y \cdot P_l(h_l(y) = h_l(x))$$

$$= f_x + \frac{1}{B} \sum_{y \neq x} f_y$$

$$= f_x + \frac{1}{B} (n - f_x).$$

The random variable  $C[h_e(x), l] - f_x$  is  $\geq 0$   
and has exp val.  $\leq \frac{n}{B}$ .

Markov:  $\Pr\left(C[h_e(x), l] - f_x \geq \frac{2n}{B}\right) \leq \frac{1}{2}$ .

Independence of  $h_1, \dots, h_t$ :  $\Pr\left(\min_l \left\{C[h_e(x), l]\right\} - f_x \geq \frac{2n}{B}\right) \leq \frac{1}{2^t}$ .

To get error  $\leq \epsilon n$  w/ prob  $\geq 1 - \delta$

Set

$$\frac{2n}{B} \leq \epsilon n$$

↓

$$B = \left\lceil \frac{2}{\epsilon} \right\rceil$$

$$\frac{1}{2^t} \leq \delta$$

↓

$$t \geq \log\left(\frac{1}{\delta}\right)$$

space required is  $O(t \log m + Bt \log n)$   
 $= O(\log(m) \log\left(\frac{1}{\delta}\right) + \frac{1}{\epsilon} \log\left(\frac{1}{\delta}\right) \log(n))$ .

Count sketch: similar to Count Min but each token is added to the counter with a random  $+/-$  sign.

Consistently use same sign for  $x$  whenever it appears  
 $\text{sign}(x), \text{sign}(y)$  indep. when  $x \neq y$ .

This could help because if many other tokens collide in some hash bucket as  $x$ , their counters may cancel each other out.

## Count Sketch

$$(B = \lceil \frac{3}{\epsilon^2} \rceil, t = (8 \ln(\frac{1}{\delta})))$$

Given  $b, t$  pos. integers.

Sample indep  $Z$ -univ hash functions

$$h_1, \dots, h_t: [m] \rightarrow [B]$$

$$g_1, \dots, g_t: [m] \rightarrow \{-1, 1\}$$

Initialize  $C[k, l] = 0 \quad \forall k \in [B], l \in [t]$

for each element  $a_i$  in stream:

    for each  $l \in [t]$ :

$$C[h_l(a_i), l] \leftarrow C[h_l(a_i), l] + g_l(a_i)$$

    end for

end for

To answer  $\text{FreqQuery}(x)$ :

return median element of  $\{g_l(x) \cdot C[h_l(x), l] \mid l \in [t]\}$

Analysis. Space is  $O(t \log(m) + Bt \log(n))$

$$= O(\log(\frac{1}{\delta}) \log(m) + \frac{1}{\epsilon^2} \log(\frac{1}{\delta}) \log(n))$$

Correctness:

$$E[g_l(x) \cdot C[h_l(x), l]]$$

$$= f_x + \sum_{y \neq x} f_y \cdot E[g_l(x) \cdot g_l(y) \cdot \mathbb{1}_{\{h_l(y) = h_l(x)\}}]$$

$$\begin{cases} 1 & \text{if } h_l(y) = h_l(x) \\ 0 & \text{if not} \end{cases}$$

0"

$$= f_x.$$

$$\text{Var}\left(g_x(\alpha) \cdot C[h_x(\alpha), l]\right) \leq \frac{1}{B} \|f\|_2^2$$

To finish up, Chebyshev  $\Rightarrow$

$$\begin{aligned} \Pr\left(\left|g_x(\alpha) \cdot C[h_x(\alpha), l] - f_x\right| > \varepsilon \|f\|_2\right) \\ &\leq \frac{1}{\varepsilon^2 B} \quad \left(B \geq \left\lceil \frac{3}{\varepsilon^2} \right\rceil\right) \\ &\leq \frac{1}{3}. \end{aligned}$$

Finish up using Hoeffding.