

L12: Streaming  
Count-Min Sketch  
(and friends)

Stream  $A = \langle a_1, a_2, \dots, a_n \rangle$   $a_i \in [m]$

- one pass

- small space



$n$  very large  
 $m$  very large

but  $\log n$  or  $\log m$   
count  
label

frequency  $j \in [m]$

$$f_j = |\{a \in A \mid a = j\}|$$

$$F_1 = \sum_j f_j = \text{total } n \text{ count}$$

$$F_2 = \sqrt{\sum_j f_j^2}$$

typically  $F_1 \gg F_2$

$$F_0 = \sum_j f_j^0 = \# \text{ distinct items } \approx F_1 \gg F_2$$

|       |   |   |   |   |   |
|-------|---|---|---|---|---|
| $f_j$ | 3 | 3 | 2 | 1 | 0 |
| $j$   | a | b | c | d | e |

a a b b c a e b c d  
 $n=10$

# Frequency Approximation

$\forall j \in [m] \rightarrow \hat{f}_j$  so

$$|f_j - \hat{f}_j| \leq \epsilon F_1 \\ \leq \epsilon F_2$$

$$F = n \\ \rightarrow \epsilon F_1 = \epsilon n$$

$$\text{size } \frac{1}{\epsilon} \text{ poly}(\log n, \log m)$$

$$\frac{1}{\epsilon^2} \cdot \text{poly}(\log \frac{n}{m})$$

sketch  $S(A)$

data structure

- insert  $(a_i)$

- query  $(g \in [m]) \Rightarrow \hat{f}_g$

tradeoff: space  $(S(A))$  vs accuracy  $\epsilon$

MS:

$$f_j \leq \hat{f}_j \leq f_j + \epsilon n$$

CS:

$$f_j - \epsilon F_2 \leq f_j \leq f_j + \epsilon F_2$$

w.p.

$$> 1 - \delta$$

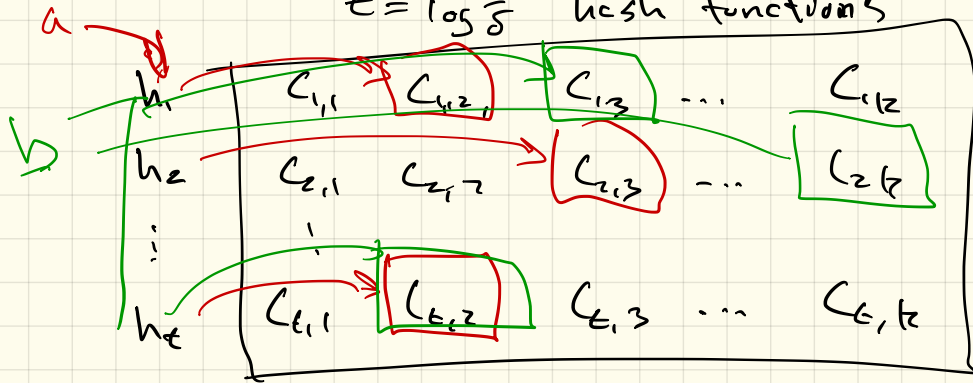
$\delta$   
prob of failure

# Count - Min Sketch

$t = \lceil \log \frac{1}{\epsilon} \rceil$  counters

$t = \log \frac{1}{\delta}$  hash functions

$h_j: [m] \rightarrow [k]$



insert  $a \in A \quad q \in [m]$

for  $j=1$  to  $t$

$C_{j, h_j(a)} ++$

for each row  
hash to counter  
increment

query  $g \in [m]$

$f_g = \min_{j \in [t]} C_{j, h_j(g)}$

•  $f_g \leq \hat{f}_g$  : each  $g$  always hashes to some  $(s, h_g(s))$ , only increments

•  $f_g \leq f_g + W$  : say  $W \leq \epsilon n = \epsilon F_1 = \epsilon \sum f_j$   
 w.p.  $> 1 - \delta$   
 1 row / hash fun

$s \in [m]$   $Y_s = \begin{cases} f_s & \text{if } h(s) = h(g) \\ 0 & \text{otherwise} \end{cases}$  w.p.  $1/k$

$(k = 2/\epsilon)$   $X = \sum_{\substack{s \in [m] \\ s \neq g}} Y_s$

$E[X] = E\left[\sum_s Y_s\right]$

Expected overcount

$\leq F_1/k = \frac{\epsilon n}{2}$

1 hash fun

$\Rightarrow \Pr[\text{one counter } w > \epsilon n] \leq 1/2$

$= \sum_s E[Y_s] = \sum_s \frac{f_s}{k}$

$\leq F_1/k$

Markov Ineq

R.V.  $X \geq 0$   $E[X] = \mu$

$\Pr[X > \alpha] = \frac{\mu}{\alpha} = \frac{\epsilon n / 2}{\epsilon n} = \frac{1}{2}$

## t hash functions

$$\epsilon = \log_2\left(\frac{1}{\delta}\right)$$

1 hash function  $h_i$

$$P_i = \Pr[w_i > \epsilon n] \leq 1/2$$

$h_1, \dots, h_t$  chosen  $h_i \in \mathcal{H}$  independently

t hash functions

$$w = \min_j w_j \quad \text{only } w > \epsilon n \\ \text{if } \underline{\text{all}} \quad w_j > \epsilon n$$

$$\Pr[w > \epsilon n] = P_i^t \leq \frac{1}{2^t} = \frac{1}{2^{\log_2(1/\delta)}} = \delta$$

# Count Sketch

$$k = 4/\epsilon^2$$

$$t = \log_2 \left( \frac{m}{\epsilon} \right)$$

both chopsp  
randomly  
from hash family

$t$ : hash fxn  $h_j: [m] \rightarrow [k]$

$t$ : sign hash fxn  $s_j: [m] \rightarrow \{-1, +1\}$

|          |          |           |           |           |
|----------|----------|-----------|-----------|-----------|
| $s_1$    | $h_1$    | $C_{1,1}$ | $C_{1,2}$ | $C_{1,k}$ |
| $\vdots$ | $\vdots$ |           |           |           |
| $s_t$    | $h_t$    | $C_{t,1}$ |           | $C_{t,k}$ |

because  
of  
median

insert  $a \in A, a \in [m]$

for  $j = 1$  to  $t$

either add  
or subtract

$$C_{j,h_j(a)} = C_{j,h_j(a)} + s_j(a)$$

$$|f_j - \hat{f}_j| \leq \epsilon F_2$$

query  $q \in [m]$

$$\hat{f}_q = \text{median}_{j \in [t]} (C_{j,h_j(q)})(s_j(q))$$

w.p.  $\geq 1 - \delta$

# Bloom Filters

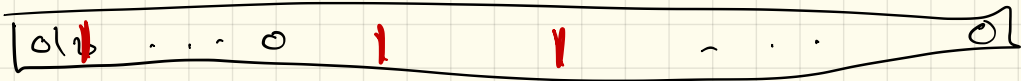
Set data structure  $B$

• if  $g \in [m]$  is in  $B$  then reports  
always

• if  $g \notin B$ , usually says  
not in  $B$

↳ False positives

$k$  bits



$k$  hash functions  $h_j : [m] \rightarrow [k]$

insert

for  $a \in [m]$ : set  $B_{h_j(a)} = 1$

query

if all  $h_j(g) = 1 \rightarrow$  Yes



# Apriori Algo.

Frequent Itemsets

Stream / Set  $A = \{a_1, a_2, \dots, a_n\}$

$a = \text{Set } \{b_1, b_2, \dots, b_t\}$   
 $b_j \in [m]$