

## SKETCHING ALGORITHM

A stream  $a_1, a_2, \dots, a_m, \dots$  where  $a_i \in U$

$F[a] = \# \text{ times } a \text{ appeared so far}$ ,  $a \in U$  : COUNTER

$\|F\| = \sum_{a \in U} F[a]$  → ISSUE  $U$  is huge, we don't enough memory!

$a$  is the most frequent if  $F[a] \geq F[b] \quad \forall b \in U$

### GOAL

- small space
- two user-defined parameters
  - $\delta$  = error probability
  - $\epsilon$  = we cannot compute  $F[\cdot]$  exactly  
so we give an  $\epsilon$ -approximation

## COUNT-MIN Sketch

- $|U|=n$
- update  $F[i]++$  (in general, you can have  $F[i] += v$ )
- query: return  $F[i]$        $i \in U = [n]$        $v$  can be negative

$$\|F\| = \sum_{i \in U} F[i]$$

Given  $\epsilon, \delta$  we estimate  $\tilde{F}[i]$  such that

$$F[i] \leq \tilde{F}[i] \leq F[i] + \underbrace{\epsilon \|F\|}_{\text{absolute error}}$$

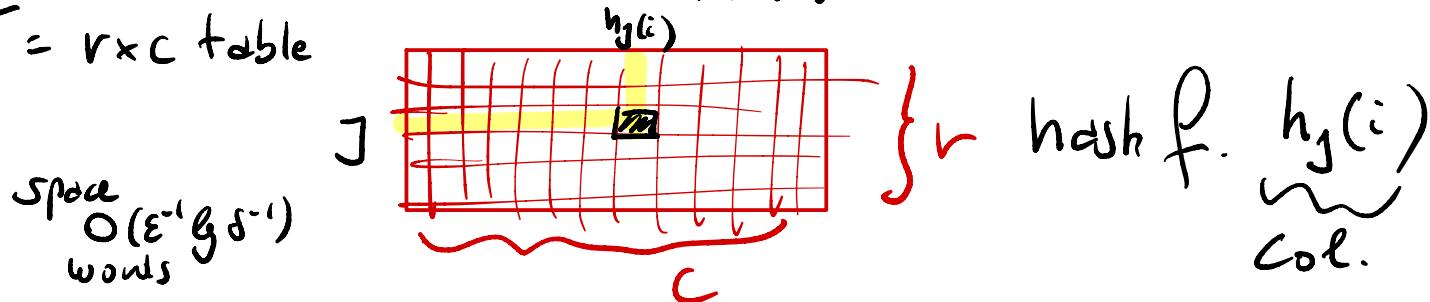
holds with probability  $\geq 1 - \delta$

Limitation: if  $F[i]$  is small compared to  $\epsilon \|F\|$ ,  
the estimation is not useful

CM-stretch is a sort of Bloom filter with counters

Input:  $\epsilon, s$       Output:  $r \times c$  table

- $r = \# \text{rows} = \ln s^{-1}$
- $c = \# \text{columns} = \frac{c}{\epsilon}$  where  $e = \text{the base of the natural log} : \ln = 2.71828\dots$
- $T = r \times c$  table

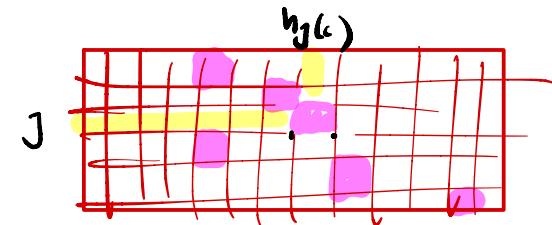


- $h_0, h_1, \dots, h_{r-1} \in \mathcal{F}$  randomly chosen from universal family

- TIME is  $O(r) = O(\ln s^{-1})$

each item  $i \in U \rightarrow$  table cells :

$$h_0(i), h_1(i), \dots, h_{r-1}(i)$$



update :  $F[i]++$  increment by 1 the cells at positions  $\langle j, h_j(i) \rangle$   
for  $0 \leq j \leq r-1$

query  $\tilde{F}[i] \rightarrow \text{return } \min_{0 \leq j \leq r-1} T[j][h_j(i)]$

Issue:  $\tilde{F}[i] \geq F[i]$  by construction,  
but it could be too large!

Let  $j$  be the row that gives the minimum

$$\tilde{F}[i] = T[j][h_j(i)] = F[i] + \underbrace{\dots}_{j_i}$$

"garbage" due to the updates

$F[k]++$  of other items  $k \neq i$

We want to show that  $\Pr[X_{ij} > \varepsilon \|F\|] < S$

Indicator Variable

$$I_{jik} = \begin{cases} 1 & \text{if } h_j(i) = h_j(k) \\ 0 & \text{o.w.} \end{cases} \quad k+i$$

$$\Rightarrow X_{ji} = \sum_{\substack{k \neq i \\ h_j(i) = h_j(k)}} F[k] = \sum_{k=1}^n I_{jik} \cdot F[k] \quad (\text{X})$$

$h: U \rightarrow [c]$   
 $h \in \mathcal{H}$   
 2-way inter

Q.  $E[I_{jik}] = P_r[I_{jik} = 1] = \frac{1}{c} = \frac{\varepsilon}{c}$

$$E[X_{ji}] = \sum_{k=1}^n E[I_{jik}] \cdot F[k] \quad (\text{X})$$

$\underbrace{E[I_{jik}]}_{\text{not a random var.}} \quad = \sum_{k=1}^n \frac{\varepsilon}{c} \cdot F[k] = \frac{\varepsilon}{c} \|F\|$

$\|F\| = \sum_k F[k]$

$c \cdot E[X_{ji}] = \varepsilon \|F\|$

We apply MI ( $\Pr[X > z] \leq \frac{E[X]}{z}$ ):  $z = \epsilon \|F\|$ ,  $X = X_{j,i}$ )

$$\Pr[X_{j,i} > \epsilon \|F\|] \stackrel{\text{MI}}{\leq} \frac{E[X_{j,i}]}{\epsilon \|F\|} = \frac{E[X_{j,i}]}{e E[X_{j,i}]} = \frac{1}{e}$$

We are done :

$$\Pr[\tilde{F}[i] > F[i] + \epsilon \|F\|] =$$

$$\Pr[V_j : \Pr[X_{j,i} > \epsilon \|F\|] =$$

$n_0, n_1, \dots, n_{r-1}$   
are independently  
chosen

$$\prod_{j=1}^r \frac{1}{e} = \frac{1}{e^r} = \delta$$

value  $\frac{1}{\delta}$