

# Streaming: Frequent Items (etc)

Note Title

2/24/2016

- Count-Min Sketch
- Count Sketch
- Apriori Algorithm (Frequent Itemsets)
- Bloom Filters

Streaming  $A = \langle a_1, a_2, \dots, a_i, \dots, a_m \rangle$

$$a_i \in [n]$$

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

$j \in [n]$

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

$\uparrow$  all ways true  
 $\uparrow$  true w.p.  $\geq 1 - \delta$

Data Structure  
 $\sum_A(j) \rightarrow \hat{f}_j$

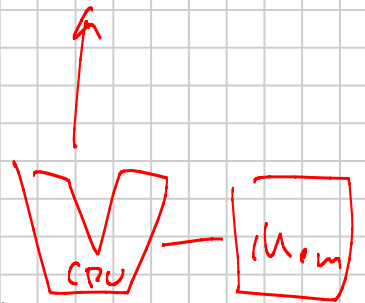
hashing

$$h_r: [n] \rightarrow [k]$$

$$h_r \sim \mathcal{H}$$

$r = \{1, 2, \dots, t\}$

$$h_r(a_i) \rightarrow j \in [k]$$

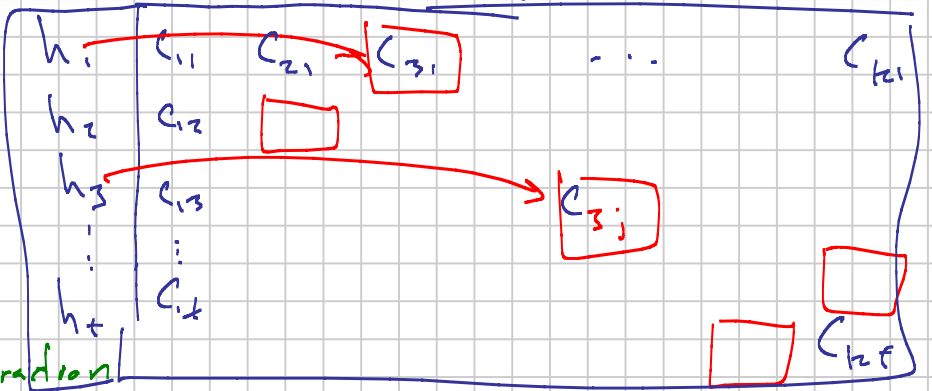


# Count-Min Sketch

Inv:  $(r_j = 0$

## Update

for  $i=1$  to  $m$   
 for  $r=1$  to  $t$   
 $C_{r, h_r(a_i)}++$   
 or -- on subtraction



## Query

$$S_A(g) = \hat{f}_g$$

$$k = \frac{2}{\epsilon}$$

$$t = \log_2 \frac{1}{\delta} \approx 5-10$$

$$\arg \min_{r \in [t]} C_{r, h_r(g)} \rightarrow f_g \leq \hat{f}_g \leq f_g + w$$

↑ How large  $\epsilon$

## Random Variable

$Y_{r,j} \equiv$  Amount of counts from  $g$  in  $C_{r, h_r(g)}$  which hash  $[n]_g$  to  $j$  query

$$Y_{r,j} = \begin{cases} f_j & \text{with prob } 1/k \\ 0 & \end{cases}$$

$$E[Y_{r,j}] = f_j / k \quad E\left[\sum_{j \in [n]} Y_{r,j}\right] = \sum_j E[Y_{r,j}] = \sum_j \frac{f_j}{k} = \frac{m}{k} \frac{\epsilon m}{k}$$

Expected overcount in row  $r$  is  $\frac{m}{k} = \frac{\epsilon m}{2}$

Markov Inequality  $\Pr[X > \alpha] \leq \frac{E[X]}{\alpha} = \frac{1}{2}$

↑  $\alpha = \log_2 \frac{1}{\delta}$  if  $X > 0$

$$\Pr\left[\arg \min_{r \in [t]} \sum_j Y_{r,j} > \epsilon m\right] \leq \left(\frac{1}{2}\right)^t = \left(\frac{1}{2}\right)^{\log_2 \frac{1}{\delta}} = \delta$$

# Space for Count-Min

$$\left( \# \text{ counters} = kt = \frac{2}{\epsilon} \log \frac{1}{\delta} \right) \log m$$

$$+ \left( t = \log \frac{1}{\delta} \right) \frac{\log n}{\text{space of hash}}$$

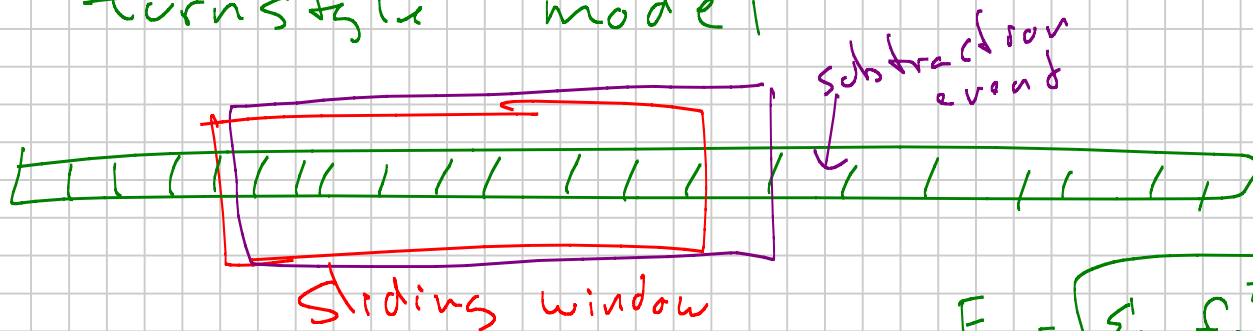
space of hash

$\log \frac{1}{\delta} \approx 5-10$   
x bigger

MG Sketch  $\frac{1}{\epsilon} (\log n + \log m)$

CM-Sketch can do subtraction

turnstile model



$$F_2 = \sqrt{\sum_{j \in [m]} f_j^2}$$

## Count-Sketch

$$|f_g - \hat{f}_g| \leq \epsilon F_2$$

Sign hash function  $S_i$   $h_i$

$$S_i : [n] \rightarrow \{-1, +1\}$$

Update

for  $a=1$  to  $m$

for  $r=1$  to  $t$

$$C_{r,h_r}(a_i) = C_{r,h_r}(a_i) + S_r(a_i)$$

$C_{r1}$	$C_{r2}$	...	$C_{rt}$
$C_{r2}$	$t = 2 \log(1/\delta)$		
$\vdots$	$t = O(1/\epsilon^2)$		
$\vdots$			
$C_{rt}$			$C_{rt}$

Query

$$E[\hat{f}_g] = f_g$$

$$S_A(g) = \hat{f}_g$$

median  $r \in [t]$   $\{ C_{r,h_r}(g) + S_r(g) \}$

# Frequent Items (Apriori Algorithm)

$$A = \langle a_1, a_2, \dots, a_m \rangle$$

$$a_i = \{ b_{i1}, b_{i2}, b_{i3}, \dots \} \quad b_{ij} \in [n]$$

"market basket"

Goal: All subsets of  $[n]$  which co-occur more than  $\phi m$  baskets.

If  $\{milk, eggs\}$  occurs  $\phi m$  times then both  $\{milk\}$  and  $\{eggs\}$  occurs  $\geq \phi m$  times

→ Only Return Maximal Sets

→ Apriori Algorithm.

Round 1 Find all individual items occurring more than  $\phi m$  times

$$M = \{ a_i \mid |a_i| \geq \phi m \} \rightarrow \text{Heavy-Hitters } M, \phi M - \epsilon M$$

$$\hookrightarrow \{ b_1, b_2, b_3, \dots, b_k \} \subset [n]$$

Round 2 Pairs  $\{ \{b_1, b_2\}, \{b_1, b_3\}, \dots, \{b_{k-1}, b_k\} \}$  usually  $\ll n$

Counters for all pairs, (or  $k$  HLL dim- $\epsilon$ )

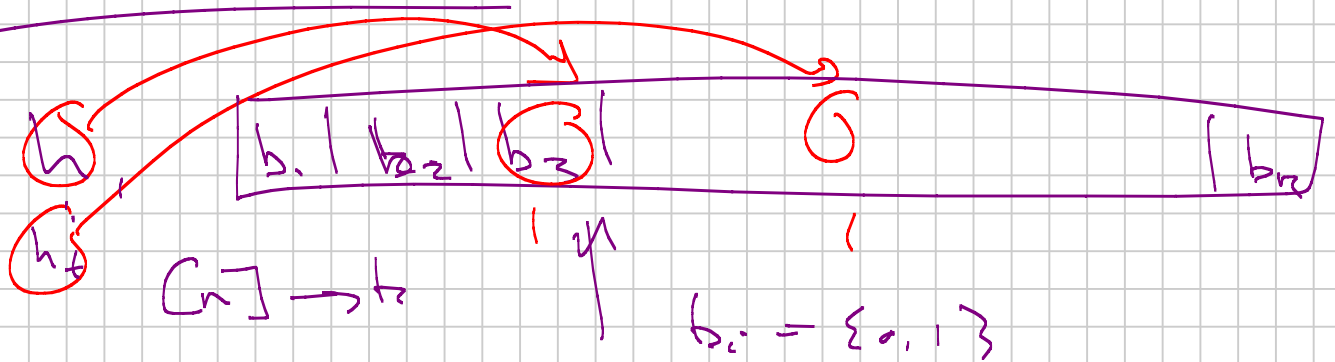
Set of Heavy pairs

$\{\{b_1, b_j\} \dots \{, \}\}$

Round 3

---

Bloom Filter



Update

for  $x \in S$

for  $j = 1$  to  $k$

set  $b_{h_j(x)} = 1$

init:  $b_i = 0$

is  $x$  in  $S$ ?

True iff all  $j \in [k]$

$(b_{h_j(x)} = 1)$