

## Heavy Hitters (HH)

There is no alg to solve HH in one pass with sublinear space.

We relax our requirements  $\Rightarrow$   $\epsilon$ -HH

Suppose at time  $t \quad \forall i \in U$

$x_i^t = \#$  times seen  $i$  up to time  $t$

Will require estimates

$$\tilde{x}_i^t \in x_i^t \pm \epsilon t \quad (*)$$

that uses space  $\ll \min(t, n)$   
 $\uparrow$   
(# distinct elts)

Idea: counting Bloom filter

$h: U \rightarrow [m]$  random universal       $A$  array with  $m$  elts

when elt  $i$  arrives

$A[h(i)]++;$

$$\tilde{x}_i^+ = \underbrace{A[h(i)]}_{\text{r.v.}} = x_i^+ + \sum_{j \neq i} x_j^+ \mathbb{1}_{h(j)=h(i)}$$

$$E[A[h(i)]] = x_i^+ + \sum_{j \neq i} \frac{x_j^+}{m} \leq x_i^+ + \frac{n}{m}$$

with  $m = \frac{1}{\epsilon}$   
satisfy (\*) in  
expectation

Want low error w/ high probability

$\Rightarrow$  amplify success probability w/ indep repetitions

Do same thing independently  $l$  times w/  $l$  independently

$h_1, h_2, \dots, h_l$   
 $\downarrow \quad \downarrow \quad \dots \quad \downarrow$   
 $A_1, A_2, \dots, A_l$

selected hash fns

Called CountMin sketch

How should we combine entries?

We have overestimates

$$\Rightarrow \tilde{z}_i^+ = \min_{1 \leq j \leq 2} A_j [h_j(i)]$$

$$E(x_i^+ - \tilde{x}_i^+) \leq \frac{n}{m}$$

$$\Rightarrow \text{Markov Inequality} \quad \Pr(x_i^+ - \tilde{x}_i^+ > \frac{2n}{m}) \leq \frac{1}{2}$$

$$\Pr(\tilde{z}_i^+ > \frac{2n}{m}) = \Pr(\text{all estimates too high}) \leq \frac{1}{2^2}$$

How to choose parameters?

$$\tilde{z}_i^+ \in [x_i^+, x_i^+ + \frac{2h}{m}] \quad \text{with prob} \geq 1 - \frac{1}{2^e}$$

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

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

$$\tilde{z}_i^+ \in [x_i^+, x_i^+ + \epsilon n]$$

$$\text{with prob} \geq 1 - \delta$$

useful when  $x_i^+ > 2\epsilon n$

Space usage:  $lm$  counters of  $\log_2 T$  bits each

$$\Rightarrow O\left(\frac{\log \frac{1}{\delta} \log T}{\epsilon}\right)$$

also space for hash fns:  $O(l \log U)$

Example: want to output all items that constitute  $\geq 4\%$  of

flow and output no items that constitute  $< 2\%$