

sample  $\leftarrow$  fixed portion fixed size. Fitting (Blum). Counting  $\leftarrow$  distinct element distinct moments. Sliding window

- sampling

1. fixed portion sampling.

(1) naive solution.  $S$  queries once,  $d$  queries twice, truly fraction of query with duplicates  $= \frac{d}{d+S}$

however, we have

$$\frac{\frac{d}{2S}}{\frac{S}{2S} + \frac{d}{2S} + \frac{d}{2S}} = \frac{d}{10+9d} \neq \frac{d}{d+S}$$

2. solution: hash key and select 1/k

2. fixed size. sampling

sample size:  $S$

$i \leq S$  时  $i$  必在 sample

$i > S$  时, 以  $\frac{S}{n}$  的概率  $i$  在 sample  $\rightarrow$  if keep, kick off one with prob  $\frac{1}{S}$

with arrive. ele  $x$  already is  $S$ .  $\text{prob}(\text{keep}) = (1 - \frac{S}{n}) + \frac{S}{n} \cdot \frac{S-1}{S} = \frac{n-1}{n}$

2. bloom Filtering.

1. check whether have seen new arrival before  $\rightarrow$  no false negative  
is BF says no, new item not in Set.  
but many false positive, if say yes, not in Set.

2. component.  $\left\{ \begin{array}{l} A: \text{array of } n \text{ bits } [a_1, \dots, a_n] \\ \text{set of hash functions} \rightarrow \text{input element, output } 0, 1, \dots, n-1. \\ \text{set of objects} \end{array} \right.$

3. construction.

用 hash function  $h_j$  将  $S$  中每个 object 映射到  $A$  的  $n$  个 buckets 中. 则  $\forall x \in A, [x] \neq \emptyset$ .

对新 object, 用  $h_1, \dots, h_j$  映射 hash. 若有一个  $h_i(o) = i$ , 且  $A[i] = 0$ , 则  $o \notin S$ .

4. false positive.  $x$  not seen before, but identified as in  $S$ .

fpr. upper bound.  $n$  bits array.  $k$  hash functions.  $m$  elements inserted.  
 $f$ : fraction of 1s in array

$$\text{FPR} = f^k \text{ 其中 } f \leq \frac{mk}{n}$$

5. estimation of  $f$ .

$$\lim_{n \rightarrow \infty} (1 + \frac{1}{n})^n = e \Rightarrow (1 - \frac{1}{n})^n = e^{-1} \Rightarrow (1 - \frac{1}{n})^{n \cdot \frac{d}{k}} = e^{-\frac{d}{k}} \Rightarrow f = 1 - e^{-\frac{km}{n}}$$

$$\text{fpr} = f^k = (1 - e^{-\frac{km}{n}})^k$$

$$\text{optimal } k = \frac{n}{m} \ln 2$$

3. counting: distinct element + distinct moments

1. distinct element.

For Bf, hash element to a long str.  $\Rightarrow$  get  $R$  = longest length of trailing zero  
 $\Rightarrow$  estimation  $= 2^R$

$$1 - e^{-\frac{m}{2^R}} \sim \frac{m}{2^R} \text{ if } 2^R \gg m$$

has for element  $a$  has at least  $r$  trailing 0's :  $p = 1 - e^{-\frac{m}{2r}}$   
 no element  $a$  has  $\sim$  :  $p' = e^{-\frac{m}{2r}}$

if  $\Sigma \gg m$ ,  $p \rightarrow 0, p' \rightarrow 1 \Rightarrow r$  不会太大  
 if  $\Sigma \ll m$ ,  $p \rightarrow 1, p' \rightarrow 0 \Rightarrow r$  不会太小  $> 2^k$  is around  $m$ .

2. distinct moments.

(1)  $k$ th-moments of stream  $S = \sum_{i=1}^n (m_i)^k \rightarrow n = \#$  of distinct value in  $S$ .  
 $m_i = \#$  of occurrences of  $v_i$  in  $S$ .

(2) 0-th-moments :  $\#$  of distinct element

1st -moments : length of  $S$ .

2nd  $\sim$  : Skewness number, 越大越 imbalance.

(3) AMS: estimate 2nd-moments of  $S$ .

随机选  $k$  个数.

$X_k \text{ value} = X$  在被选  $k$  个位置上的 occurrences.

$$\text{estimation} = \frac{n}{k} \sum_{i=1}^k (2X_k \cdot \text{value} - 1).$$

(4) estimation of 3rd-moments  $= \frac{n}{k} \sum (3X_k^3 - 3X_k \cdot v + 1).$

IV. sliding window: counting  $\#$  of 1's in window.