

DATA STREAMING : FLOW OF MASSIVE SEQUENCES OF DATA

- lots of data
 - cannot store it
- n elements
polylog space in n

MOTIVATIONS

- IP traffic in a router
- telephone calls
- query logs

CHALLENGE

- simple problems (statistics) become difficult

example EASY : find MAX : DIFFICULT : find QUANTILES (most frequent top k, etc)

Problems cannot be solved deterministically

FEW EXCEPTIONS: Find the missing element in a permutation of $1, 2, \dots, n$ (just $n-1$ elements arrive)

- SUM $2 \log n$ bit
- XOR $\log n$ bit

COUNT-MIN SKETCH (Cormode - Muthuk)

n items numbered $1 \dots n$

f frequency array $F[i] = \#$ times i appears in the stream

operations: $F[i]++$, $F[i]--$, invariant: $F[i] \geq 0$

▷ cannot store f entirely, have only $O(\log n)$ bits

▷ find approximation \tilde{F} st.

$$\forall i: F[i] \leq \tilde{F}[i] \leq F[i] + \epsilon \|F\|$$

with probability $1 - \delta$

$$\|f\| = \sum_{i=1}^n F[i]$$

Need a couple of notions

- k -wise independence
 - random variables
 - hash functions
- Markov's inequality

K-wise limited independence

X_1, X_2, \dots, X_n random variables with support S_1, S_2, \dots, S_n

K-wise independent if

\forall choice $i_1 < i_2 < \dots < i_k \in [n]$ and $a_{ij} \in S_{ij}$

$$\Pr[X_{i_1} = a_{i_1} \wedge X_{i_2} = a_{i_2} \wedge \dots \wedge X_{i_k} = a_{i_k}] = \Pr[X_{i_1} = a_{i_1}] \times \Pr[X_{i_2} = a_{i_2}] \times \dots \times \Pr[X_{i_k} = a_{i_k}]$$

Implication: $\mathbb{E}[X_{i_1} X_{i_2} \dots X_{i_k}] = \mathbb{E}[X_{i_1}] \times \mathbb{E}[X_{i_2}] \times \dots \times \mathbb{E}[X_{i_k}]$

$\{h\}_{h \in H}$ family of hash functions $[n] \rightarrow [b]$

$\Pr[h \in H] = \frac{1}{|H|}$ uniform distributions

Family H is **K-wise independent** if

$\forall x_1, x_2, \dots, x_k \in [n] \quad b_1, b_2, \dots, b_k \in [b]$

$$\Pr_{h \in H}[h(x_1) = b_1 \wedge h(x_2) = b_2 \wedge \dots \wedge h(x_k) = b_k] = \Pr_{h \in H}[h(x_1) = b_1] \times \Pr_{h \in H}[h(x_2) = b_2] \times \dots \times \Pr_{h \in H}[h(x_k) = b_k]$$

$$\begin{aligned} [x_{i_j} &= h(x_j)] \\ [a_j &= b_j] \end{aligned}$$

Example

$$h(x) = a_0 + a_1 x + a_2 x^2 + \dots + a_{k-1} x^{k-1}$$

$a_i \in \text{finite field } F$

$\triangleright |H| = |F|^k$

\triangleright there $|F|^{k-1}$ solutions to $h(x_i) = b_i$

$$a_0 = b_i - a_1 x_i - a_2 x_i^2 - \dots - a_{k-1} x_i^{k-1}$$

$$\triangleright \Pr [h(x_1)=b_1 \wedge h(x_2)=b_2 \wedge \dots \wedge h(x_k)=b_k]$$

$$\frac{1}{F^k}$$

$$\Pr [h(x_1)=b_1] \wedge \Pr [h(x_2)=b_2] \wedge \dots \wedge \Pr [h(x_k)=b_k]$$

$$\frac{1}{F} \quad \frac{1}{F} \quad \frac{1}{F}$$

$$= \frac{1}{F^k} \quad \text{QED}$$

Space to store $h(x)$ is that of $a_0 a_1 \dots a_{k-1}$
 $O(k \lg |F|)$ bits

Ex. Take prime $p \in [n+1, 2n]$: show $h'(x) = (h(x) \bmod p) \bmod b$ is approximately k -wise indep.

Markov's Inequality

$X = \text{random variable} \geq 0, \text{ real } a > 0$

$$\Pr [X \geq a] \leq \frac{\mathbb{E}[X]}{a}$$

Proof -

$$I = \text{indicator variable} = \begin{cases} 0 & \text{if } X < a \\ 1 & \text{if } X \geq a \end{cases}$$

fact $\mathbb{E}[I] = 0 \cdot \Pr[I=0] + 1 \cdot \Pr[I=1]$
 $= \Pr[I=1] = \Pr[X \geq a]$

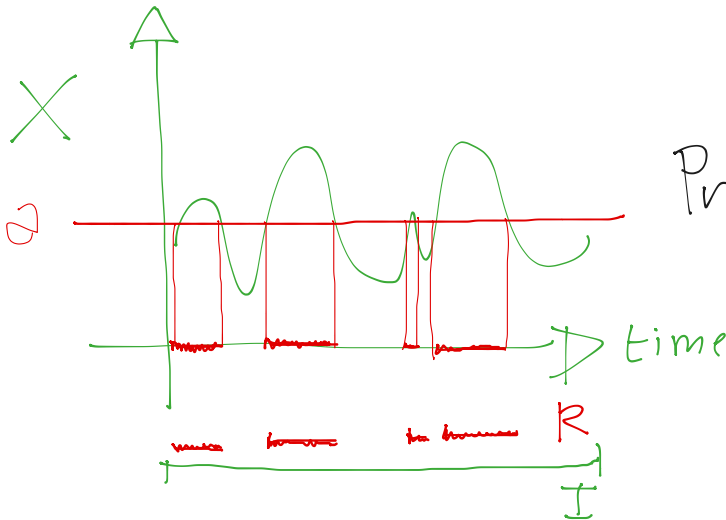
$\triangleright a \leq I \leq X$

$\cdot X < a \Rightarrow I = 0$

$\cdot X \geq a \Rightarrow a \leq I = a \leq X$

$\Rightarrow \mathbb{E}[aI] \leq \mathbb{E}[X]$

$a \cdot \mathbb{E}[I] = \mathbb{E}[X] \quad (\text{see fact})$



$\Pr[X \geq a] = \frac{R}{I}$

COUNT-MIN SKETCH (Cormode - Muthuk)

n items numbered $1 \dots n$

F frequency array $F[i] = \#$ times i appears in the stream

operations: $F[i]++$, $F[i]--$, invariant: $F[i] \geq 0$

▷ cannot store F entirely, have only $O(\log n)$ bits

▷ find approximation \tilde{F} st.

$$\forall i: F[i] \leq \tilde{F}[i] \leq F[i] + \epsilon \|F\|$$

with probability $1 - \delta$

$$\|F\| = \sum_{j=1}^n F[j]$$

Algorithm

1. let $r = \lceil \frac{1}{\delta} \ln \frac{e}{\epsilon} \rceil$ and $c = \frac{e}{\epsilon}$

$e = 2.71828$
Euler's constant

2. let T be a table of $r \times c$ counters initially set to zero

	1	2	...	c
1				
2				
⋮				
r				

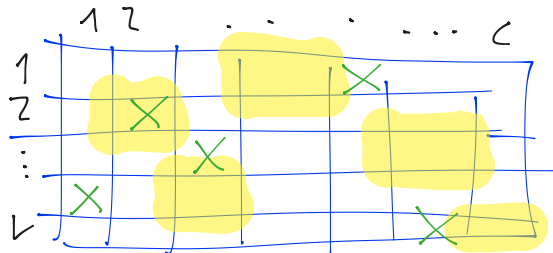
3. Take a family \mathcal{H} of 2-wise independent hash functions
eg $h(x) \iff h(x) = [(ax+b) \bmod p] \bmod c + 1$

4. Choose h_1, h_2, \dots, h_r uniformly and independently from \mathcal{H} eg r pairs $(a_1, b_1) \dots (a_r, b_r)$ $a_i \neq 0$

5. given element i , associate r cells (one per row)

$$T(1, h_1(i)), T(2, h_2(i)), \dots, T(r, h_r(i))$$

$i \neq i'$
might have collisions



6. OPS

$F[i]++ \Rightarrow$ increment by 1
the r cells for i

$F[i]-- \Rightarrow$ decrement, as above

$$7. \tilde{F}[i] = \min_{1 \leq j \leq r} T(j, h_j(i))$$

Fact 1

Storage is $O(v.c) = O(\epsilon^{-1} \log \delta^{-1})$
words of memory, where each
word can store $\|F\|$.

proof

Each entry stores an integer in $[\|F\|]$

Fact 2 $F[i] \leq \tilde{F}[i]$

proof Let $T(j, h_j(i)) = \tilde{F}[i]$

By construction, $\exists i_1, i_2, \dots, i_p$
s.t. one of them is i and

$$T(j, h_j(i)) = \sum_{k=1}^p F[i_k] = F[i] + X_{j,i}$$

where $X_{j,i} \geq 0$ is the excess (s.t. $F[i_k] \geq 0$)

QED

obs To model the excess (due to
hash collisions) use indicator variable

$$I_{jik} = \begin{cases} 1 & \text{if } k \neq i \text{ and } h_j(i) = h_j(k) \\ 0 & \text{otherwise} \end{cases}$$

obs $X_{ij} = \sum_{k=1}^n I_{j,ik} \cdot F[k]$

follows from the fact that $I_{j,ik} = 1$ iff there is a collision in cell $J(j, h_j(i))$

obs $E[I_{j,ik}] = \frac{c}{C}$

proof $E[I_{j,ik}] = 0 \cdot \text{Pr}[I_{j,ik}=0] + 1 \cdot \text{Pr}[I_{j,ik}=1]$
 $= \text{Pr}[I_{j,ik}=1]$

$= \text{Pr}[\exists d \in [C] : k \neq i \wedge h_j(i) = d \wedge h_j(k) = d]$

$= \sum_{d \in [C]} \text{Pr}[k \neq i \wedge h_j(i) = d \wedge h_j(k) = d]$

use the fact that h_j is 2-wise independent

$= \sum_{d \in [C]} \underbrace{\text{Pr}[h_j(i) = d]}_{\frac{1}{C}} \times \underbrace{\text{Pr}[k \neq i \wedge h_j(k) = d]}_{\text{"almost"} \frac{1}{C}}$

since one element in the domain is missing

$= \sum_{d \in [C]} \frac{1}{C^2} = \frac{1}{C} = \frac{c}{C}$

Fact 3 $\tilde{F}[i] \leq F[i] + \varepsilon \|F\|$
 with probability $\geq 1 - \delta$

proof

Let $\tilde{F}[i] = F[i] + X_{ji}$, where $X_{ji} \geq 0$ is the excess
 $\triangleright E[X_{ji}] = E\left[\sum_k I_{jik} F[k]\right] = \sum_k E[I_{jik} F[k]] \leq \sum_k F[k] E[I_{jik}] = \sum_k F[k] \cdot \frac{\varepsilon}{c} = \frac{\varepsilon}{c} \|F\|$

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

$$\begin{aligned} \triangleright \Pr[\tilde{F}[i] > F[i] + \varepsilon \|F\|] &= \\ \Pr[\cancel{F[i]} + X_{ji} > \cancel{F[i]} + \varepsilon \|F\|] &= \\ \Pr[X_{ji} > \underbrace{\varepsilon \|F\|}_{c E[X_{ji}]}] &= \end{aligned}$$

$$\Pr[X_{ji} > c E[X_{ji}]] =$$

i_1, i_2, \dots, i_r chosen uniformly and independently

$$\Pr[X_{i_1} > c E[X_{i_1}]] \times \Pr[X_{i_2} > c E[X_{i_2}]] \times \dots \times \Pr[X_{i_r} > c E[X_{i_r}]]$$

\hookrightarrow Markov's inequality $\leq \frac{E[X_{ji}]}{c E[X_{ji}]} = \frac{1}{c} < \frac{1}{2} \quad 1 \leq i \leq r$

$$< \left(\frac{1}{2}\right)^r = \left(\frac{1}{2}\right)^{\frac{1}{\delta}} = \delta$$

Q.E.D

Roberto Grossi '12