

CS 368

Canvas

CS368.stanford.edu

Piazza, Gradescope

3 homeworks (60%)

Optional programming assignment, sub. for 1 HW

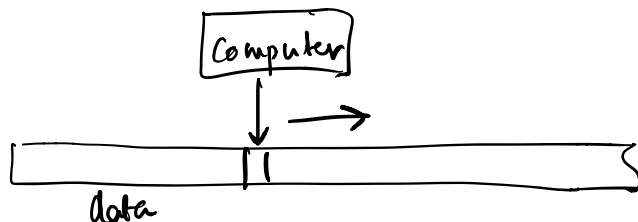
Project (35%)

Piazza participation (5%)

Data Stream Model

- Data does not fit in memory
- polynomial time not good enough!

Google query log.



Warmup: Counter that counts from 1 to n
 $\log n$ bits

exact or deterministic $\rightarrow \log n$ bits needed!

K items: return $[k(1-\epsilon), k(1+\epsilon)]$

Allow failure probability δ (say $\delta < 10^{-6}$)

$O(\log \log n)$ bits suffice!

$\hookrightarrow O_{\epsilon, \delta}$

Instead of k , keep track of $\log k$

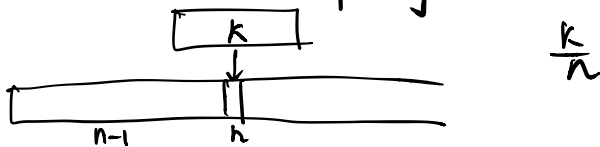


Count Distinct Queries

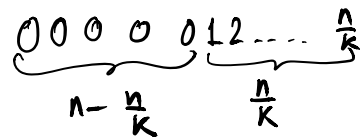
① Hash Table : too large

② Random Sampling

Reservoir sampling



Suppose we pick random sample of k elements out of n



$$\left(1 - \frac{1}{k}\right) \quad \frac{1}{k}$$

$$\left(1 - \frac{1}{k}\right)^k \approx \frac{1}{e}$$

$$\frac{n}{k} \rightarrow \frac{n}{100k}$$

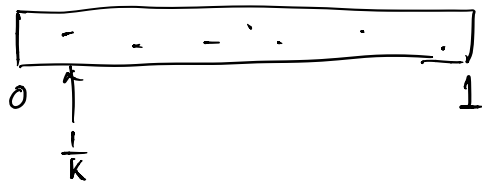
$$\left(1 - \frac{1}{100k}\right)^k \approx \frac{1}{e^{1/100}} \approx 1 - \frac{1}{100}$$

[Flajolet, Martin '85]

$$h: U \rightarrow [0, 1]$$

Assume $h(x) \in_R [0, 1]$

$h(x_1), \dots, h(x_k)$ independent



Claim: If k distinct elev $x_1 \dots x_k$

$$E\left[\min_{i=1..k} h(x_i)\right] = \frac{1}{k+1}$$

$$Y = \min(h(x_1) \dots h(x_k))$$

$$\Pr[h(x_i) \leq t] = t$$

$$\Pr[Y > t] = (1-t)^k$$

$$\Pr[Y \leq t] = 1 - (1-t)^k \quad \text{CDF}$$

$$\Pr[Y \in [t, t+dt]] = k(1-t)^{k-1} dt \quad \text{PDF}$$



$$E[Y] = \int_0^1 t \cdot k(1-t)^{k-1} dt$$

$$E[Y] = \int_0^1 \Pr[Y > t] dt$$

$$= k \int_0^1 (1 - (1-t)) (1-t)^{k-1} dt$$

$$= k \left(\int_0^1 (1-t)^{k-1} dt - \int_0^1 (1-t)^k dt \right)$$

$$= k \left(\frac{1}{k} - \frac{1}{k+1} \right) = \frac{1}{k+1}$$

Is Y close to $E[Y]$?

$$\text{Var}[Y] = E[Y^2] - E[Y]^2$$

$$E[Y^2] = \int_0^1 t^2 k(1-t)^{k-1} dt$$

$$= k \int_0^1 (1 - (1-t))^2 (1-t)^{k-1} dt$$

$$= k \left[\int_0^1 (1-t)^{k-1} dt - 2 \int_0^1 (1-t)^k dt + \int_0^1 (1-t)^{k+1} dt \right]$$

$$\begin{aligned}
&= k \left[\frac{1}{k} - \frac{2}{k+1} + \frac{1}{k+2} \right] \\
&= k \left[\left(\frac{1}{k} - \frac{1}{k+1} \right) - \left(\frac{1}{k+1} - \frac{1}{k+2} \right) \right] \\
&= k \left[\frac{1}{k(k+1)} - \frac{1}{(k+1)(k+2)} \right] = \frac{2}{(k+1)(k+2)} \leq 2 E[Y]^2
\end{aligned}$$

$$\text{Var}[Y] = E[Y^2] - E[Y]^2 \leq E[Y]^2$$

Chebyshev's inequality:

$$\Pr[|Y - E[Y]| > \varepsilon E[Y]] \leq \frac{\text{Var}[Y]}{(\varepsilon \cdot E[Y])^2} \leq \frac{1}{\varepsilon^2}$$

Take mean of multiple iid copies.

Y_1, \dots, Y_t be independent copies of Y

$$Z = \frac{Y_1 + \dots + Y_t}{t}$$

$$E[Z] = E[Y]$$

$$\text{Var}[Z] = \frac{1}{t^2} \sum_1^t \text{Var}[Y_i] = \frac{\text{Var}[Y]}{t}$$

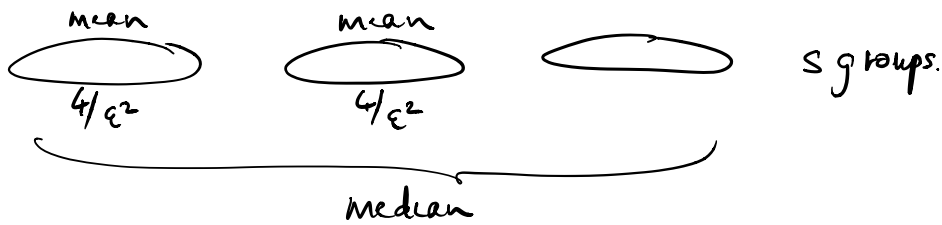
$$\text{Var}[Z] \leq \frac{(E[Z])^2}{t}$$

$$\Pr[|Z - E[Z]| > \varepsilon \cdot E[Z]] \leq \frac{\text{Var}(Z)}{(\varepsilon \cdot E[Z])^2} \leq \frac{1}{\varepsilon^2 t}$$

$$t = \frac{4}{\varepsilon^2}, \text{ failure probability } \leq \frac{1}{4}$$

Median of Means

- ① Z_1, \dots, Z_s independent copies of Z
- ② return median (Z_1, \dots, Z_s)



Median bad \Rightarrow at least half of $z_1 \dots z_s$ bad.

A_i indicator random variable = $\begin{cases} 1 & \text{if } z_i \text{ bad} \\ 0 & \text{otherwise} \end{cases}$

$$E[A_i] \leq 1/4 \quad A = \sum_{i=1}^s A_i$$

$$E[A] \leq s/4$$

$$\begin{aligned} \Pr[\text{median}(z_1 \dots z_s) \text{ bad}] &\leq \Pr[\text{at least } \frac{s}{2} \text{ of } z_1 \dots z_s \text{ bad}] \\ &= \Pr[A \geq s/2] \\ &\leq e^{-(2 \ln 2 - 1) s/4} \\ &\leq e^{-s/11} \end{aligned}$$

Set $s = 11 \ln(1/\delta)$ to get failure prob. $< \delta$

Chernoff Bound: X sum of independent 0-1 random var.
 $\mu = E[X]$

$$\Pr[X \geq (1+\delta)\mu] \leq \left(\frac{e^{-\delta}}{(1+\delta)^{\delta}} \right)^{\mu}$$