# DISTINCT ELEMENTS

INPUT:  A stream $s_1 \ldots s_n \in \{1 \ldots N\}$

GOAL:  Estimate # of distinct elements in the stream.

$\{0, \frac{1}{2^k}, \frac{2}{2^k}, \frac{3}{2^k} \ldots\}$

1) Pick a random hash function (pairwise independent) $h: \{1 \ldots N\} \to [0, 1]$

2) Compute $\alpha = $ minimum of $\{h(s_1) \ldots h(s_n)\}$

implement in small space

3) Output $\frac{1}{\alpha} - 1$

$k+1$

$\{ s_1 \cdots s_n \} \rightarrow$ K distinct elements in stream.

$\{ h(s_1) \cdots h(s_n) \} \rightarrow$ K distinct real #s in $[0,1]$



$\approx \frac{1}{k}$

$\min \{ h(s_1), h(s_2) \cdots h(s_n) \}$

$\approx \frac{1}{k}$

**Thm:**

$$\mathbb{E}_h \left[ \min \{ h(s_1) \ldots h(s_n) \} \right] = \frac{1}{K+1}$$

where $\quad K = \#$ of distinct elements in $\{s_1, \ldots s_n\}$

$$h : \{1 \ldots N\} \to [0,1]$$

$\uparrow$
random !!

**Proof:**

Lemma: If $\quad r_1 \ldots r_K \in [0,1]$ uniformly random

the $\quad \mathbb{E}[\min(r_1, \ldots r_K)] = \frac{1}{K+1}$

Truly random hash function $h: \{1..N\} \rightarrow [0,1]$

Let $s_1, \ldots s_k$ be distinct elements in the stream.

$s_{i_1}, s_{i_2}, \ldots s_{i_k}$

1) $\Pr_h \left[ \min(h(s_1), \ldots h(s_k)) \leq \frac{1}{4k} \right] \leq \frac{1}{4}$.

$\Updownarrow$

estimate $\geq 4k$

Proof: $= \Pr_h \left[ \left( h(s_1) < \frac{1}{4k} \right) \vee \left( h(s_2) < \frac{1}{4k} \right) \ldots \vee \left( h(s_k) < \frac{1}{4k} \right) \right]$

$\leq \sum_{i=1}^{k} \Pr \left[ h(s_i) < \frac{1}{4k} \right]$ (union bound)

$= \sum_{i=1}^{k} \left( \frac{1}{4k} \right) = \frac{1}{4}$

2) $\Pr\left[\min\left(h(n_1)\ldots h(s_n)\right) > 4/k\right] \leq e^{-4}$

$\Downarrow$

estimate $< k/4$

truly random function.

$= \Pr\left[(h(s_1) > 4/k) \wedge (h(s_2) > 4/k)\ldots \wedge (h(s_k) > 4/k)\right]$

(independence)

$= \overset{k}{\underset{i=1}{\prod}} \Pr\left[h(s_i) > 4/k\right]$
$\qquad\qquad$ (truly random hash function)

$= \left(1 \overset{||}{-} 4/k\right)^k$



$\approx e^{-4}$

hash family

$$\mathcal{H} = \{ h_1 \ldots h_M : \{1 .. N\} \to [R] \}$$

$$\mathcal{H} = \text{set of all possible functions}$$

To remember a $h \in \mathcal{H}$
need $\log |\mathcal{H}|$ bits.

$\mathcal{H}$ must be small yet some how random

# Pairwise Independent Hash family

A hash family $\mathcal{H}$ is pairwise independent if

$\forall x, y \quad x \neq y \in$ Domain $\{1 \ldots N\}$

$\forall \alpha, \beta \in$ Range $\{1 \ldots R\}$

$$\Pr_{h \sim \mathcal{H}} \left[ (h(x) = \alpha) \wedge (h(y) = \beta) \right] = \Pr \text{ under a completely random function}$$

$$= \frac{1}{R} \cdot \frac{1}{R} = \frac{1}{R^2}$$

Example:

$p$ - prime

$$\mathbb{Z}_p = \{0, 1 \dots p-1\}$$

integers modulo $p$

$$\mathcal{H} = \{ h_{a,b}(x) = ax + b \ (\text{mod } p) \}$$

$a, b \in \mathbb{Z}_p$

$$|\mathcal{H}| = p^2$$

$h_{a,b}(1)$

$h_{a,b}(2)$

$\Downarrow$

$a, b$