# Counting distinct elements in data streams

Elements arrive $(a_i)$ from a domain $[m] = \{1..m\}$   $a_1...a_n$   $n$ elements

Goal: <u>Count # distinct items</u>. $(F_0)$

$$F_k = \sum_{i \in A} f_i^k \qquad \text{where } f_i \text{ is the frequency of an item}$$
$$\text{if } A \text{ is the distinct set of items.}$$

<u>Result</u>: with a small amount of memory we will approximate $F_0$ with high probability.

$$Pr\left[ |F_0 - \tilde{F_0}| \leq \varepsilon.F_0 \right] \geq 1-\delta.$$
     ↗ confidence parameter
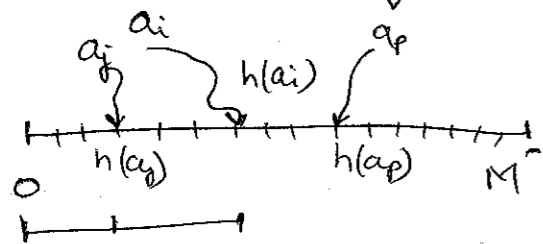     ↳ error parameter

Choose a random hash function $h : [m] \rightarrow [M]$   $M = m^3$.

Note: ensures that the probability of a collision is very small $(\leq \frac{1}{m})$

<u>Basic Idea</u>: Let $t = \frac{c}{\varepsilon^2}$    $c \equiv$ constant. (TBD)

Apply $h(a_i)$, and maintain ~~set~~

$V \equiv$ max over the set of $t$ smallest values in $\{h(a_i)\}$

Output $\tilde{F_0} = \frac{tM}{V}$.



2 smallest values , $V = h(a_i)$

$$\tilde{F_0} = \frac{t.M}{V}$$

Suppose $\tilde{F_0} > (1+\epsilon) F_0$.      Let $b_1, b_2 \ldots, b_{F_0}$ be a listing of

the $F_0$ distinct values

$\Rightarrow$    $h(b_1), h(b_2) \ldots h(b_{F_0})$

contains $\geq t$ elements smaller than $V$.

$$\frac{tM}{V} > (1+\epsilon) F_0 \quad \Rightarrow \quad V < \frac{tM}{F_0(1+\epsilon)}$$

What is $Pr\left[ h(b_i) < \frac{tM}{F_0(1+\epsilon)} \right]$ ?

If $h(b_i)$ is uniformly distributed, then it is      $< \frac{t}{F_0(1+\epsilon)}$

We have $F_0$ (pairwise indep.) events happening. Each event

has prob $p = < \frac{t}{F_0(1+\epsilon)}$.   What is the    chance that $\geq t$ happen ?

Let $X_i = 1$ iff $h(b_i) < \frac{tM}{F_0(1+\epsilon)}$.

$= 0$   otherwise

$$E[X_i] < \frac{t}{F_0(1+\epsilon)} \qquad E\left[ \sum_{i=1}^{F_0} X_i \right] < \frac{t}{1+\epsilon}$$

Let $Y = \sum_{i=1}^{F_0} X_i$ \qquad $E[Y] < \frac{t}{1+\epsilon}$

$$Var[Y] = \sum_{i=1}^{F_0} Var[X_i] \leq \frac{t}{1+\epsilon}.$$

[THIS NEEDS
"PAIRWISE INDEPENDENCE"]
$Pr(X \wedge Y) = Pr(X) \cdot Pr(Y)$
$E[XY] = E[X] \cdot E[Y]$
NOT TRUE FOR
DEPENDENT VARIABLES.

Note that $Var[X_i]$ is actually $p(1-p)$.

Use Chebyshev's Bound

$$\Pr\left[\,|Y - E[Y]| \geq a\right) \leq \frac{Var[Y]}{a^2}$$

$$\Rightarrow \quad \Pr\left[\,|Y - \frac{t}{1+\varepsilon}| \geq a\right] \leq \frac{\frac{t}{1+\varepsilon}}{a^2}$$

$$\Pr[Y \geq t] \leq \Pr\left[\,|Y - \frac{t}{1+\varepsilon}| \geq \frac{t \cdot \varepsilon}{1+\varepsilon}\right] \leq \frac{\frac{t}{1+\varepsilon}}{\frac{t^2 \varepsilon^2}{(1+\varepsilon)^2}} = \frac{(1+\varepsilon)}{t\varepsilon^2}$$

since $t = \frac{c}{\varepsilon^2}$  we get  $\frac{1+\varepsilon}{c}$.

Similar proof when $\tilde{F_0} < (1+\varepsilon) F_0$.

---

**Defn**  $$Var[X] = E\left[(X - \mu_x)^2\right] = E[X^2] - (\mu_x)^2 \qquad \mu_x = E[X]$$

(easy proof)

**FACT**  $$E[X \cdot Y] = E[X] \cdot E[Y] \quad (\text{if } X \text{ \& } Y \text{ are indep})$$

**FACT**  $$Var[X+Y] = Var[X] + Var[Y] \quad \text{if } X \text{ \& } Y \text{ are independent.}$$

$$= E\left[(X+Y - E[X+Y])^2\right]$$

$$= E\left[((X - \mu_x) + (Y - \mu_y))^2\right] = E\left[(X - \mu_x)^2\right] + E\left[(Y - \mu_y)^2\right]$$

$$+ 2E\left[(X - \mu_x)(Y - \mu_y)\right]$$

$$= Var[X] + Var[Y] + 2\left[\underbrace{E[XY - X\mu_y - Y\mu_x + \mu_x\mu_y]}\right]$$

$$= 2\left[\mu_x\mu_y - \mu_x\mu_y - \mu_y\mu_x + \mu_x\mu_y\right]$$

$$= 0$$