# Distinct Elements in Vanilla Streaming Model

**Goal:** Let $F_0 = |\{j : F_j > 0\}|$. Given $(\epsilon, \delta)$-approximation algorithm for $F_0$.

**Preliminaries:** Pairwise independent families of hash functions.

A family $H$ of hash functions [mapping $[n] \to [r]$] is said to be <u>pairwise independent</u> if for any <sup>distinct</sup> $X, y \in [n]$ and $v, w \in [r]$, $\Pr_{h \in H}[h(x) = v \text{ and } h(y) = w] = $

Let $X_{(x, v, 1)}, \ldots, X_{(x, v, r)}$ be random variables indicating whether $h(x) = v$, for $h \in H$. $H$ is pairwise-independent hash family iff these indicator variables are pairwise-independent.

$\Pr_{h \in H}[h(x) = v] \cdot \Pr_{h \in H}[h(y) = w] = \frac{1}{r^2}$

**Note:** Implies $h(x)$ uniformly distributed in $[r]$, for any $x$.

**Example:** Assume $r$ is prime. For any $a, b \in [r]$, let $h_{a,b}(x) := ax + b \bmod r$.
Let $H = \{h_{a,b} : a, b \in [r]\}$

**Claim:** $H$ is a pairwise independent hash family.

**Proof:** Fix $x \neq y \in [n]$, $v, w \in [r]$. What is $\Pr_{h \in H}[h(x) = v \text{ and } h(y) = w]$?

It is $\dfrac{|\{h_{a,b} : h_{a,b}(x) = v \text{ and } h_{a,b}(y) = w\}|}{|H| = r^2}$ , **Claim:** numerator $= 1$ ∴ ⌣.

$$ax + b = v \pmod r$$
$$ay + b = w \pmod r$$

Two linearly independent equations in 2 unknowns have a unique solution. ✓

**Note:** $h_{a,b}$ takes only $O(\log r)$ bits to store, and it can be evaluated with $O(1)$ multiplications and additions mod $r$.

**Fun fact:** if $r = 2^{61} - 1$ or $r = 2^{127} - 1$, then reducing mod $r$ is just
[underbrace: mersenne primes]
a bit-wise shift and $O(1)$ machine multiplications & operation, and multiplication mod $r$ can be done with $O(1)$ machine multiplications.

---

**Outline for rest of lecture**
- Recall template for designing $(\epsilon, \delta)$-multiplicative approximation algorithm for $F_0$
  - First, develop "basic estimator" that is unbiased but might have high variance. (say variance $v$)
  - Take the average of "many" copies of the basic estimator to drive down variance $t = \frac{v}{t^2 \epsilon^2}$ [pairwise-independent]
  - Apply Chebyshev to conclude the above is a $(1 \pm \epsilon)$-mult. approx. w.p. $\geq \frac{3}{4}$.
  - Take median of $O(\log(\frac{1}{\delta}))$ copies of the above to get a $(1 \pm \epsilon)$-mult. approx w.p. $\geq 1 - \delta$.

E.g, last time the basic estimator for $F_k$

had variance at most $V \leq k \cdot n^{1-\frac{1}{k}} \cdot F_k^2$.

So took $t = \frac{V}{\epsilon^2 F_k^2} = \frac{k \cdot n^{1-\frac{1}{k}}}{\epsilon^2}$ mutually independent

Copies of the basic estimator and looked at their

average. We then took the median of $O\left(\log\left(\frac{1}{\delta}\right)\right)$

Copies of the above, for total space usage

$$O\left( \frac{k \cdot n^{1-\frac{1}{k}}}{\epsilon^2} \cdot \log\left(\frac{1}{\delta}\right) \cdot \log(n \cdot m) \right)$$

Today: We will give a basic estimator $\hat{F_0}$ that is "close to" unbiased

we will <u>not</u> bound its variance so we won't get an $(\epsilon, \delta)$-mult. approx.

We will use Markov's inequality to show it $\cancel{has\ a}$ satisfies

$$\frac{F_0}{10} \leq \hat{F_0} \leq 10 F_0 \qquad \text{with probability at least } \frac{3}{4}$$

Taking the median of $O(\log(\frac{1}{\delta}))$ copies then gives an estimator $\hat{F_0}'$ s.t,

$$\frac{F_0}{10} \leq \hat{F_0}' \leq 10 F_0 \qquad \text{with probability}$$

at least $1-\delta$.

Tommorow we will see the $\cancel{\text{actual template}}$ $\cancel{\text{for obtaining}}$ an actual $(\epsilon, \delta)$-mult. approx. alg.

same for $F_0$ following the template

↑ (in red) sortof

# A Simple (Suboptimal) Distinct Elements Algorithm.

Given a number $v$, let $zeros(v)$ denote the number of "trailing zeros" in the binary representation of $v$. This is the same as the largest $i$ such that $2^i$ divides $v$.

Assume $r > n$ is a power of two.

Algorithm: • Choose a random $h$ from a pairwise independent hash family $H$ of hash functions mapping $[n]$ to $[r]$.

• $z \leftarrow 0$

• For each stream update $a_i$:

$$\text{If } zeros(h(a_i)) > z, \text{ then } z \leftarrow zeros(h(a_i))$$

• Output $2^z$.

i.e. the algorithm simply outputs the largest number of trailing zeros it ever saw in a hash value.

Note: Space usage is $O(\log\log m)$ bits for storing $z$, and $O(\log n)$ bits for storing $h$.
Note also that the algorithm implicitly ignores duplicates and its final estimate is independent of the order of the stream.

Intuition: If there are $d$ distinct items, then we expect about $\dfrac{d}{2^i}$ of them to have at least $i$ trailing zeros. When $i = \log d$, this number is $1$. When $i >> \log d$, this number is very close to $0$, so $2^z$ should be fairly close to $2^{\log d + \frac{1}{2}} = \sqrt{2} \cdot d$

**Analysis:** Claim: With probability at least $1 - \text{~~...~~} \geq \frac{5}{8}$, the algorithm will output a number $\hat{d}$ such that

$$\frac{d}{4} \leq \hat{d} \leq 4d.$$

**Proof:** For each $j \in [n]$ and $t \in \mathbb{Z}_+$, let

$$X_{t,j} = \begin{cases} 1 : \text{if } \mathrm{zeros}(h(j)) \geq t \\ 0 : \text{otherwise} \end{cases}$$

and $Y_t = \sum_{j: f_j > 0} X_{t,j}$. Let $z_{max}$ be the final value of $z$ computed by the algorithm. Note:

$$Y_t > 0 \iff z_{max} \geq t \quad \Longrightarrow$$
$$Y_t = 0 \iff z_{max} \leq t - 1$$

> So it's enough to show that
> $$\Pr\left[Y_{\lceil \log(8d) \rceil} \geq \phi\right] \leq \frac{1}{8} \text{ and}$$
> $$\Pr\left[Y_{\lceil \log(\frac{d}{8}) \rceil} = 0\right] \leq \frac{1}{4}$$

**Fact 1:** $\mathbb{E}[Y_t] = \sum_{j: f_j > 0} \mathbb{E}[X_{t,j}] = \sum_{j: f_j > 0} \Pr_{h \in \mathcal{H}}\left[2^t \text{ divides } h(j)\right] = \sum_{j: f_j > 0} \frac{1}{2^t} = \frac{d}{2^t}$

**Fact 2:** ~~...~~

$$\mathrm{Var}[Y_t] \leq \frac{d}{2^t}.$$

**Proof:** $\mathrm{Var}[Y_t] = \mathrm{Var}\left[\sum_{j: f_j > 0} \mathrm{Var}[X_{t,j}]\right] \overset{\text{pairwise independence of } X_{t,j}'s}{=} \sum_{j: f_j > 0} \mathrm{Var}[X_{t,j}]$

$$\leq \sum_{j: f_j > 0} \mathbb{E}[X_{t,j}^2] \overset{X_{t,j} \text{ is } \{0,1\}\text{-valued}}{=} \sum_{j: f_j > 0} \mathbb{E}[X_{t,j}] = \frac{d}{2^t}.$$

- Let $a = \lceil \log_2(8d) \rceil$. ~~...~~ $\Pr[Y_a \geq 1] \overset{\text{Markov}}{\leq} \frac{\mathbb{E}[Y_a]}{1} \leq \frac{d}{2^a} \leq \frac{1}{8}$

- Let $b = \lceil \log_2(\frac{d}{8}) \rceil$. $\Pr[Y_b = 0] \leq \Pr\left[|Y_b - \mathbb{E}[Y_b]| \geq \frac{d}{2^b}\right]$

  $\overset{\text{chebyshev}}{\leq} \underbrace{\frac{\text{Var}(Y_b)}{\left(\frac{d}{2^b}\right)^2}} \overset{\text{Fact 2}}{\leq} \frac{2^b}{d} \overset{\text{definition of } b}{\leq} \frac{2}{8} = \frac{1}{4}$ ∎

To amplify the $\overset{\text{upper bound}}{\text{probability}}$ of a bad estimate from $\frac{3}{8}$ to $\delta$, repeat $O\left(\log\left(\frac{1}{\delta}\right)\right)$ times and take the median.

---

Next time! rather than an estimate $\hat{d}$ s.t. $\frac{d}{8} \leq \hat{d} \leq 8d$ with probability $\geq 1-\delta$, we'll see an actual $(\varepsilon, \delta)$ approximation algorithm using space about $O\left(\log n + \frac{1}{\varepsilon^2}\right)$. ← optimal but we'll miss it by a $\left(\log(\frac{1}{\varepsilon}) + \log\log n \text{ factor}\right)$.

the rough idea will be to run about $\frac{1}{\varepsilon^2}$ copies of the simple algorithm in parallel, and aggregate their estimates into a single, smarter estimate in a clever way.