Streaming Distinct Elements
Sketching Frequencies

Announcements ① Homework 6 to be released Weds,
     due one week later (plus 2-day grace period)

     ② Quiz 4 will be in class, Monday, May 9.

Recap:   $h(x) = ax + b \pmod{M}$    $M \geq m$, prime

     $(a,b) \in [M]^2$ uniformly random.

      Storing a representation of $h$ requires
      only $2 \log_2 (M)$ space to store
      coefficients $a$ and $b$.

   ① $\forall x$   $h(x)$ is unif distrib in $[M]$.

   ② $\forall x \neq y$ then $h(x), h(y)$ are independent
     random variables.

Let   $X_{ik} = \begin{cases} 1 & \text{if} \quad h(a_i) \leq k \\ 0 & \text{if} \quad h(a_i) > k \end{cases}$

   $Y_k = \sum_{i=1}^{d} X_{ik}$ = # of distinct elements in stream
         that hash to $\{1, \ldots, k\}$.

We proved:   ① $\mathbb{E}[X_{ik}] = \dfrac{k}{M}$

     ② $\mathbb{E}[Y_k] = \dfrac{dk}{M}$

     ③ $\text{Var}[Y_k] < \dfrac{dk}{M}$.

# Algorithm for Distinct Elements

Sample random hash function $h$ as above.

Let $t = \left\lceil \dfrac{2(1+\varepsilon)}{\varepsilon^2 \delta} \right\rceil$.

Initialize $(z_1, \ldots, z_t) = 1^t$

    // $z_1, \ldots, z_t$ will store the $t$ smallest distinct
        values in the set $\{ h(a_i) \mid i = 1, \ldots, n \}$
        in increasing order.

for $i = 1, \ldots, n$:
    Observe token $a_i$
    Compute $z = h(a_i)$
    if $z < z_t$:
        update $z_1, \ldots, z_t$ to preserve the
        invariant in the comment above.
    endif
endfor

Output $\dfrac{tM}{z_t}$.

$z_1 \approx \dfrac{M}{d} \quad z_2 \approx \dfrac{2M}{d} \quad \ldots \quad z_t \approx \dfrac{tM}{d}$

# Analysis of the algorithm

$$z_t = k \equiv Y_k = t \text{ and } Y_{k-1} < t.$$

So, we resort to analyzing $Y_k$ for various values of $k$.

Two ways the algo could fail:
- Outputs answer $< (1-\epsilon) d$.

Since answer $= \dfrac{tM}{Z_t}$ this

corresponds to $Z_t > \dfrac{tM}{(1-\epsilon)d}$

$\implies Y_k < t$ for $k = \left\lfloor \dfrac{tM}{(1-\epsilon)d} \right\rfloor$

- outputs answer $> (1+\epsilon)d$

corresponds to $Z_t < \dfrac{tM}{(1+\epsilon)d}$

$\implies Y_\ell \geq t$ for $\ell = \left\lfloor \dfrac{tM}{(1+\epsilon)d} \right\rfloor$

$E[Y_\ell] = \dfrac{d\ell}{m}, \quad \text{Var}[Y_\ell] < \dfrac{d\ell}{m},$

$\Pr\left(Y_\ell \geq t\right) \overset{\color{red}\text{Cheby}}{\leq} \dfrac{d\ell/M}{\left(\frac{d\ell}{M} - t\right)^2}$

$\dfrac{d\ell}{M} = \dfrac{d}{M}\left\lfloor \dfrac{t}{1+\epsilon} \cdot \dfrac{M}{d} \right\rfloor \leq \dfrac{t}{1+\epsilon}$

$t - \dfrac{d\ell}{M} \geq t - \dfrac{t}{1+\epsilon} = \dfrac{\epsilon t}{1+\epsilon}.$

$\left(t - \dfrac{d\ell}{M}\right)^2 \geq \dfrac{\epsilon^2 t^2}{(1+\epsilon)^2}.$ $\qquad \color{teal} t = \left\lceil \dfrac{2(1+\epsilon)}{\epsilon^2 \delta} \right\rceil$

$\leq \dfrac{t/(1+\epsilon)}{\epsilon^2 t^2/(1+\epsilon)^2} = \dfrac{1+\epsilon}{\epsilon^2 t} \leq \dfrac{\delta}{2}.$

Similar calculation using Chebyshev shows
the probability of answer $< (1-\epsilon)d$ is also $\leq \dfrac{\delta}{2}.$

# Sketching Model of Computation

Stream $a_1, \ldots, a_n$

Tokens belong to $[m]$

Also has storage space $s = O(poly(\log n, \log m))$.

After processing stream, the stored representation is used to answer queries from some set $Q$ of potential queries.

$(\varepsilon, \delta)$-PAC property: $\forall g \in Q$ $Pr($answer is $\varepsilon$-accurate$) \geq 1-\delta$.

??

Weds lecture: Estimating frequencies of tokens.

Query $q(x)$ for token $x$ asks,

"How many times did $x$ appear in the stream?"

$\varepsilon$-accurate might mean, for example, answer to $q(x)$ has additive error $\leq \varepsilon n$.