**Last time:**
- matching
- intro hashing
- universal hashing

**Key idea:**
- put the randomness in hash function instead of assuming data is random!
- gives all the nice properties:
  - data distributed randomly throughout table
  - hash fn efficient to store
  
  efficient to compute.

**Today:**

3 applications of hashing & lossy compression
- Bloom filters
- Heavy hitters & count-min sketch
- Distinct elts

with short review of variance & tail bounds

# Heavy hitters

$$3 \quad 5 \quad 7 \quad 3 \quad 4$$

stream of elts $\quad a_1, a_2, a_3, \cdots,$

at any time $t$,

let $\quad f_x^t$ : # times seen element $x$ in

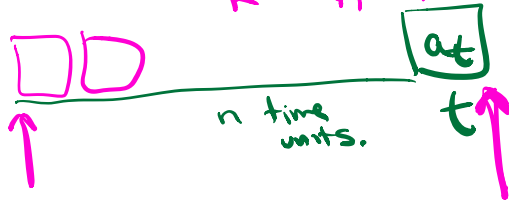$$a_1, a_2, \cdots, a_t$$

Goal: when elt shows up, output that element

if $\quad f_x^t > \dfrac{n}{k}$

$n$ : millions, billions

$k$ : 10's, 100's, 1000's

$a_t$

such an element is
a heavy hitter.

Space used proportional
to # unique elts

$n$ time
units.

$t$

not possible to solve this problem exactly
with sublinear space.

# Modified goal $(\varepsilon, \delta)$

① If $\quad f_x^t > \dfrac{n}{k} \quad$ output $x$.

② If $\quad x$ is output, then with prob at least
$1 - \delta \quad$ it is the case that $\quad f_x^t \geq \dfrac{n}{k} - \varepsilon n$

For example: suppose $\quad k = 25, \quad \varepsilon = 0.01$

$$\delta = \dfrac{1}{2^{10}}$$

① If $\quad f_x^t > \dfrac{n}{25} = 0.04 n \quad$, then output $x$

② If $\quad x$ output then w.prob $\geq \boxed{1 - \dfrac{1}{2^{10}}}$

$$f_x^t \geq \underbrace{0.04 n}_{\frac{n}{k}} - \underbrace{0.01}_{\varepsilon} n = 0.03 n$$

## Count-min sketch:

Designer specifies $n, k, \delta, \varepsilon$
$\implies b, \ell$

keep 2D array called CMS



each row is hash table of size $b$

typical values for $b$ & $\ell$ might be
$b = 1000$
$\ell \approx 5$

when element $x$ shows up

$Inc(x)$: $\forall \; 1 \leq j \leq \ell$ increment $CMS[j][h_j(x)]$

Observe: $\forall$ time $t$, $\forall j$, $\forall x$
$$CMS[j][h_j(x)] \geq f_x^+$$

$Count(x)$: return $\min_{1 \leq j \leq \ell} CMS[j][h_j(x)]$

if this value $\geq \frac{n}{k}$ output $x$ as HH.

by observation: $Count(x) \geq f_x^+$

---

## Construction

① hash fns behave randomly
$\forall \; x, y \atop x \neq y$ $\qquad \forall \; 1 \leq j \leq \ell$ $\quad \Big| \; Pr\Big( h_j(x) = h_j(y) \Big) = \frac{1}{b}$

② hash fns for $1 \leq j \leq \ell$ are indep of each other

$\mathcal{H}$ universal class of hash fns.

**Analysis.** Fixing $t \le n$, elt $x$ that arrives at time $t$.

$$Z_j = CMS[j][h_j(x)] \quad \text{random variable.}$$

$$Z_j = f_x^+ + \sum_{y \ne x} f_y^+ \, W_{xy}$$

$$W_{xy}^j = \begin{cases} 1 & h_j(x) = h_j(y) \\ 0 & o.w. \end{cases}$$

$$E(Z_j) = f_x^+ + \sum_{y \ne x} f_y^+ \underbrace{E[W_{xy}^j]}_{\frac{1}{b}} \quad \text{by linearity of expectation}$$

$$\le f_x + \frac{t}{b} \le f_x^+ + \frac{n}{b}$$

$$\underline{E(Z_j - f_x^+) \le \frac{n}{b}}$$

$$Pr\left(Z_j - f_x^+ > \frac{2n}{b}\right) \le \frac{1}{2}$$

> **Markov's Inequality**
> $X$ is nonnegative r.v.
> $$Pr(X \ge c \cdot E(X)) \le \frac{1}{c}$$

$$Pr\left(\text{Count}(x) - f_x^+ > \frac{2n}{b}\right) \le \frac{1}{2^{\ell}} \quad \Leftarrow \text{bad event}$$

$$Z_1 - f_x^+ > \frac{2n}{b}$$
$$Z_2 - f_x^+ > \frac{2n}{b}$$
$$\vdots$$
$$Z_{\ell} - f_x^+ > \frac{2n}{b}$$

**Conclusion:** $\quad \underline{Pr\left(\text{Count}(x) \ge f_x^+ + \frac{2n}{b}\right)} \le \frac{1}{2^{\ell}} \quad *$

with $\frac{2n}{b} = \varepsilon n$ and $\frac{1}{2^\ell} = \delta$

**Modified goal $(\varepsilon, \delta)$**

$\Rightarrow$ ① If $f_x^+ > \frac{n}{k}$ output $x$.

$\Rightarrow$ ② If $x$ is output, then with prob at least $1-\delta$ it is the case that $f_x^+ > \frac{n}{k} - \varepsilon n$

$\uparrow \frac{n}{2k}$

choose $b$ & $\ell$
so that ⓐ $\frac{2n}{b} = \varepsilon n$ ⓑ $\frac{1}{2^\ell} = \delta$

$$b = \frac{2}{\varepsilon} \qquad \ell = \log_2\left(\frac{1}{\delta}\right)$$

$$\varepsilon = \frac{1}{2k}$$
$$b = 4k$$

$$\delta = \frac{1}{2^{100}}$$
$$\ell = 100$$

Suppose $\boxed{f_x^+ < \frac{n}{k} - \varepsilon n}$

$$Pr\left(\text{Count}(x) \ge \frac{n}{k}\right) \le Pr\left(\text{Count}(x) > f_x^+ + \varepsilon n\right) \le \delta$$
$$\underset{< \frac{n}{k} - \varepsilon n}{\overline{\phantom{xxx}}}$$

# Universal hash family $\mathcal{H}$
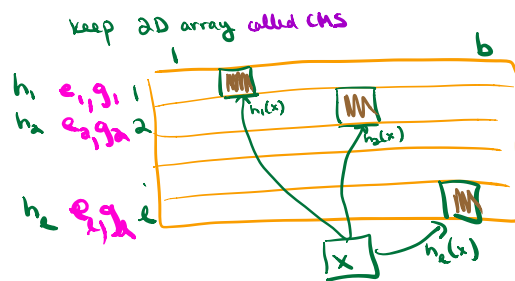
$U = \{0, 1, \ldots, u-1\}$

each $h \in \mathcal{H}$     $h: U \longrightarrow \{0, 1, \ldots, b-1\}$

If $h$ is chosen uniformly at random from $\mathcal{H}$

$$\forall x \neq y \quad \Pr_{h \in \mathcal{H}}\left(h(x) = h(y)\right) \leq \frac{2}{b}$$

Choose any prime # $p > u$

$$\mathcal{H} = \left\{ h(x) = (ex + g) \bmod p \bmod b, \right.$$
$$\text{where} \quad 1 \leq e \leq p-1$$
$$\left. 0 \leq g \leq p-1 \right\}$$

$$|\mathcal{H}| = p(p-1)$$

keep 2D array called CHS



$h_1 \quad e_1, g_1 \quad 1$

$h_2 \quad e_2, g_2 \quad 2$

$h_\ell \quad e_\ell, g_\ell \quad \ell$