

Data Structure

Data $x \in U$

Queries $f_1, f_2, \dots, f_m: U \rightarrow \{0, 1\}$

Goal: Preprocess x and set up memory
so that when given $i \in [m]$,
we can quickly compute $f_i(x)$.

Set Membership.

Set $S \subseteq U$

$$|U| = M$$

$$|S| \leq n$$

Queries: For each $i \in U$, "is $i \in S$ ".

sets of size $\leq n$ is $\binom{M}{\leq n}$.

\Rightarrow Any data structure for this
needs space $\geq \log \binom{M}{\leq n}$ bits.

$$\approx O\left(n \log\left(\frac{M}{n}\right)\right)$$

Word size $\approx O(\log M)$

$$\approx n \log M.$$

(for M big)

$\Rightarrow O(n)$ words.

Going to set:

$O(n)$ ~~size~~ words of memory used

$O(1)$ time to respond to queries.

[FKS '70s]

Pick random map $h: U \rightarrow [n^2]$

Hope: h is one-to-one on S .

Then here is a nice representation:

For each $j \in [n^2]$

store: ① $b_j \in \{0, 1\}$

$b_j = 1$ iff $\exists s \in S$ s.t.

$h(s) = j$

else 0.

② If $b_j = 1$

store the $s \in S$ s.t.

$h(s) = j$.

Pairwise independent hash family

$\mathcal{H} \subseteq \{f: U \rightarrow [K]\}$

s.t. ① For any fixed $u \in U$
for a random $h \in \mathcal{H}$
the distribution of $h(u)$
is uniform over $[K]$.

② For any $u_1, u_2 \in U$, $u_1 \neq u_2$
for random $h \in \mathcal{H}$
the dist. $(h(u_1), h(u_2))$

is uniform over $[K] \times [K]$

\exists such \mathcal{H} with $|\mathcal{H}| = O(|U| \cdot K)$
 $= O(M \cdot K)$

So can specify an elt of \mathcal{H} using $O(\log M + \log K)$ bits.

Claim $\forall S \subseteq [M], |S|=n$

$$P_{h \in \mathcal{H}} \left[\forall x \neq y \in S, h(x) \neq h(y) \right] \geq 1 - \frac{\binom{n}{2}}{K}$$

Proof Fix $x, y \in S$.

$$P_h \left[h(x) = h(y) \right] = \frac{1}{K}$$

Union bound over all $x, y \in S$.

New idea:

Study the number of collisions when hashing a set of size n to a range of size $O(n)$.

For each $x, y \in S, x \neq y$. Let

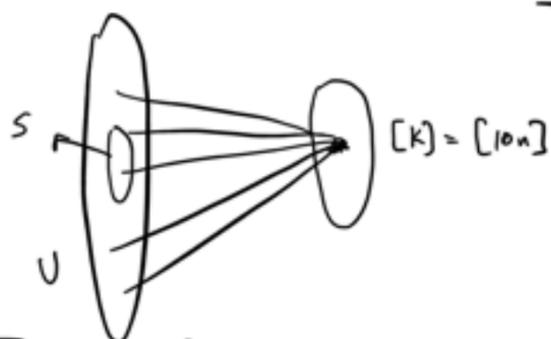
$$Z_{x,y} = 1 \text{ if } h(x) = h(y)$$

$$E[Z_{x,y}] = \frac{1}{K}$$

$$Z = \sum_{\substack{x \neq y \\ x, y \in S}} Z_{x,y} = \# \text{ collisions.}$$

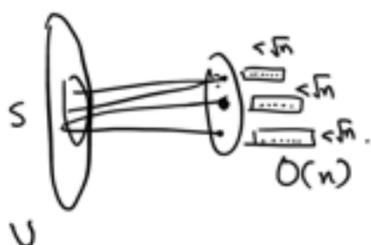
$$E[Z] = \sum_{x \neq y} E[Z_{x,y}] = \frac{\binom{n}{2}}{K}$$

If $K = \frac{10n}{5}$, then $E[Z] \leq \frac{n}{10}$.



If $Z \leq \frac{n}{5}$, then $\forall j \in [K], |h^{-1}(j) \cap S| \leq \sqrt{n}$

$$P_h \left[Z > \frac{n}{5} \right] \leq \frac{E[Z]}{n/5} \leq \frac{5}{10} \leq \frac{1}{2}$$



Final Data Structure

⊗ First hash from U to $[K]$ where $K = O(n)$.

Then for each $i \in [K]$, let

$$S_{i0} = \{s \in S : h(s) = i\} = S \cap h^{-1}(i)$$

Use the $O(n^2)$ space data structure for S_i .

$$\text{Total space used} = \sum_i O(|S_i|^2) + O(n)$$

$$= O(z) + O(n).$$

$$= \underline{O(n)}.$$

For each $i \in [K]$, we write

1. $|S_i|$

2. which $h_i: U \rightarrow [0, |S_i|^2]$

indep. family $\mathcal{H}_i = \{f: U \rightarrow [0, |S_i|^2]\}$
from the canonical p-wise

③ For each $l \in [0, |S_i|^2]$