

Hashing and Dictionaries

- Ordered dict: under → member, pred, succ, insert, delete } best pretty good
 - ↳ if we reverse pred and succ we can do better
 - keys = $\{0, 1, \dots, n-1\}$ assume usually $n=2^{64}$
- ↳ trivial implementation $\rightarrow |S|=n$
- ↳ Allocate array $T[0, \dots, n-1]$
 - member is array search } unique hash per val
 - insert, del trivial
- but n could bigger than $n!$. Also this is shit

Hash + (main)

Hashing

- $|S| = n$
- Allocate an array $T[0..m-1]$ (m comparable to n)
 - $T[i]$ is a linked list.
- Choose a **hash function** $h : \{0, \dots, u-1\} \rightarrow \{0, \dots, m-1\}$ from some "family" of functions H .
- Insert(x, val) : Append (x, val) to $T[h(x)]$
 - x and y **collide** if $h(x) = h(y)$
- Delete(x) : Remove (x, val) from $T[h(x)]$
- Member(x) : Return true if x is found in $T[h(x)]$

What is a good "family" of functions

- The "gold standard" is h is selected from the set of all functions from $\{0, \dots, u-1\} \rightarrow \{0, \dots, m-1\}$.

– How many functions of this type are there?

$$m^u$$

– How much space does it take to represent such an h ?

$\leftarrow m$ space $\leftarrow u \log m$ bits \leftarrow storing in table bits

- This is called the **IDEAL Hash Model**. Choosing a uniformly random function is completely impractical. We want to find a family that mimics the properties of the Ideal Hash Model.

What is good in a hash func? $h : \{0, \dots, n-1\} \rightarrow \{0, \dots, m-1\}$

→ **Uniformity** $\rightarrow \forall x_i \Pr_{h \in H} (h(x_i) = i) = 1/m$

→ **Collisions** $\rightarrow \forall x \neq y \Pr_{h \in H} (h(x) = h(y)) = 1/m$

→ **Independence** $\rightarrow \forall x_1, \dots, x_k, i_1, \dots, i_k$
 $\Pr_{h \in H} (h(x_1) = i_1, \dots, h(x_k) = i_k) = \prod_{j=1}^k \Pr_{h \in H} (h(x_j) = i_j)$
 $= 1/m^k$

Prob 1

$\Omega \rightarrow$ total set of outcomes, $P : \Omega \rightarrow \mathbb{R}$ so $\sum_{\omega \in \Omega} P(\omega) = 1$

Event $E_i \rightarrow$ subset of $\Omega \rightarrow$ ind of E_i . $E_i = \{\omega \in \Omega \mid P(E_i \cap \{\omega\}) = P(E_i)P(\omega)\}$

Random variable is func $X : \Omega \rightarrow \mathbb{R}$

Exp $\rightarrow E(X) = \sum_{\omega \in \Omega} P(\omega) X(\omega) = \sum_{\omega \in E_i} k \cdot P(X=\omega)$

- random var is iid if $\{x=k_1\}, \{y=k_2\}$ ind & k_1, k_2
- if X, Y ind $E[X \cdot Y] = E[X]E[Y]$
- Lin of Exp.
- if F an event. Indicator rv $I_F = \begin{cases} 0 & \text{if } F \text{ occurs} \\ 1 & \text{if } X \text{ does} \end{cases}$
- $E[I_F] = P(I_F = 1) = \sum_{x \in F} P(x)$

Note: $1+x \leq e^x \Rightarrow \left(1 + \frac{1}{n}\right)^n < e < \left(1 + \frac{1}{n}\right)^{n+1}$

$$\left(1 - \frac{1}{n}\right)^n < e < \left(1 - \frac{1}{n}\right)^{n-1}$$

Hashing with chaining

$S = \{y_1, \dots, y_n\}$ member (x)?

time $\approx |\mathcal{T}[h(x)]|$

↳ linked list in array \mathcal{T} at has value of $h(x)$

What is $E(|\mathcal{T}[h(x)]|) = \sum_{k=0}^m k \Pr(\mathcal{T}[h(x)] = k)$

$$I_i \quad 1 \leq i \leq n = \begin{cases} 1 & \text{if } h(x) = h(y_i) \\ 0 & \text{otherwise} \end{cases}$$

$$\hookrightarrow \Pr(h(x) = h(y_i)) = \Pr(I_i = 1) = \frac{1}{m}$$

$$\begin{aligned} E(|\mathcal{T}[h(x)]|) &= E\left(\sum_{i=1}^n I_i\right) = \sum_{i=1}^n E(I_i) \\ &= \sum_{i=1}^n \Pr\{h(x) = h(y_i)\} \\ &= \begin{cases} \frac{n}{m} & x \notin S \\ 1 + \frac{m-1}{m} & x \in S \end{cases} \end{aligned}$$

Construct a Good Hash Function

- Pick a prime number $p > u$.
- Hash family $H = \{h_{a,b} : a \in \{1, \dots, p-1\}, b \in \{0, \dots, p-1\}\}$.
 - $-g_{a,b}(x) = ax + b \pmod{p}$
 - $-h_{a,b}(x) = g_{a,b}(x) \pmod{m}$
- Claim. For all $x \neq y$, and $r, s \in \{0, \dots, m-1\}$,
 - $\Pr(h(x) = h(y)) \leq 1/m$.
 - $\Pr(h(x) = r \text{ and } h(y) = s) = \Theta(1/m^2)$.

→ let $x \neq y \rightarrow$ let choice of a, b

$$\Pr_{a,b}(g(x) = i \cap g(y) = j)$$

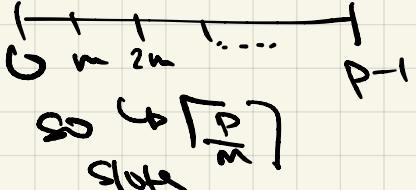
$$\Pr_{a,b}(ax + b = i \pmod{p} \cap ay + b = j \pmod{p})$$

$$\begin{aligned} &\Rightarrow a(x-y) \equiv i-j \pmod{p} \\ &\Rightarrow \begin{cases} a \equiv (i-j)(x-y)^{-1} \pmod{p} \\ b \equiv i - ax \pmod{p} \end{cases} \end{aligned}$$

Mult inverse
as p prime
 $\frac{1}{p-1} \pmod{p}$
 $i=j$

$$\Rightarrow \Pr_{a,b} (g(x)=i \wedge g(y)=j) = \begin{cases} 0 & i=j \\ \frac{1}{P(P-1)} & i \neq j \end{cases} \xrightarrow{\text{diff count both to sum}}$$

$$\rightarrow \Pr(h(x)=h(y)) = \frac{|\{(i,j) \mid i=j \pmod m\}|}{P(P-1)} \xrightarrow{\text{simpl}} \frac{P(\lceil \frac{P}{m} \rceil - 1)}{P(P-1)}$$



 slots $\lceil \frac{P}{m} \rceil$
 $m \quad \dots \quad \lceil \frac{P}{m} \rceil$
 $0 \quad 1 \quad \dots \quad P-1$

$$= \frac{\lceil \frac{P}{m} \rceil - 1}{P-1}$$

$$\leq \frac{\frac{P-1}{m}}{\frac{P-1}{P}} = \frac{1}{m}$$

$$\Pr(h(x)=r \wedge h(y)=s) = \frac{|\{(i,j) \mid i \pmod m=r, j \pmod m=s\}|}{P(P-1)}$$

$P > 3m$

worst case $\lceil \frac{P}{m} \rceil$
 $\leq \frac{\lceil \frac{P}{m} \rceil \lceil \frac{P}{m} \rceil}{P(P-1)} \in O(\frac{1}{m^2})$
 $\Delta \geq \frac{(\lceil \frac{P}{m} \rceil - 1)(\lceil \frac{P}{m} \rceil - 2)}{P(P-1)} \xrightarrow{\text{dif res}} \in O(1)$

choices for $i \pmod m$