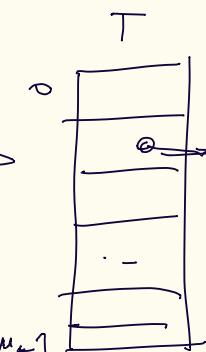


Hash with chaining

$h(x)$



where $m, M \ll U$
 M # of keys

$|T|$

$m-1$

array

Total space: $\frac{m}{\text{pointer}} + \frac{m}{\text{key}} + \text{satellite data}$

pointer

$$m = m \cdot [L + 2 \cdot 64]$$

2 pointers

($m \log_2 m$ bits)

$m \cdot 6$ bits

keys table pointer
Space: $m \lg U + m \lg m + m \lg M$
Bits

• INSERTION = $O(1)$ w.c.

UNREALISTIC simple uniform hash

• SEARCH = AVG length of the list $O(1+\alpha)$ \Rightarrow α

$\forall \text{key } k, P(h(k) = i) = \frac{1}{m}$
 \Rightarrow BUT COMPUTATIONALLY UNFEASIBLE

• DELETION = \propto depends on search $O(1+\alpha)$

hash $h(x)$ that distributes the keys in table entries uniformly distributed at random

Theorem 8.1 Under the hypothesis of simple uniform hashing, there exists a hash table with chaining, of size m , in which the operation Search(k) over a dictionary of n objects takes $\Theta(1 + n/m)$ expected time. The value $\alpha = n/m$ is often called the load factor of the hash table.

Proof:

CONSTANT TIME COMPUTATION OF $h(x)$

$$= 1 + \sum_{i \in D} 1 \cdot P(h(k) = i) = 1 + \frac{m}{m}$$

$$\mathbb{E}[\text{search}] = 1 + \mathbb{E}[\text{length of the list}] = 1 + m \cdot \frac{1}{m} = 1 + \frac{M}{m}$$

$\Rightarrow O(1 + \frac{m}{m})$ time

of keys

$\Rightarrow O(1 + \alpha)$ search

$|T|$

$\frac{m}{m} = \alpha = \text{LOAD FACTOR}$

if size of table $m = n$ cost becomes constant ($\frac{n}{n}$)
but THE PROBLEM is (n) ?

if 1) M is static (fixed # of keys) $\Rightarrow O(k)$

But considering insertion() and deletion()

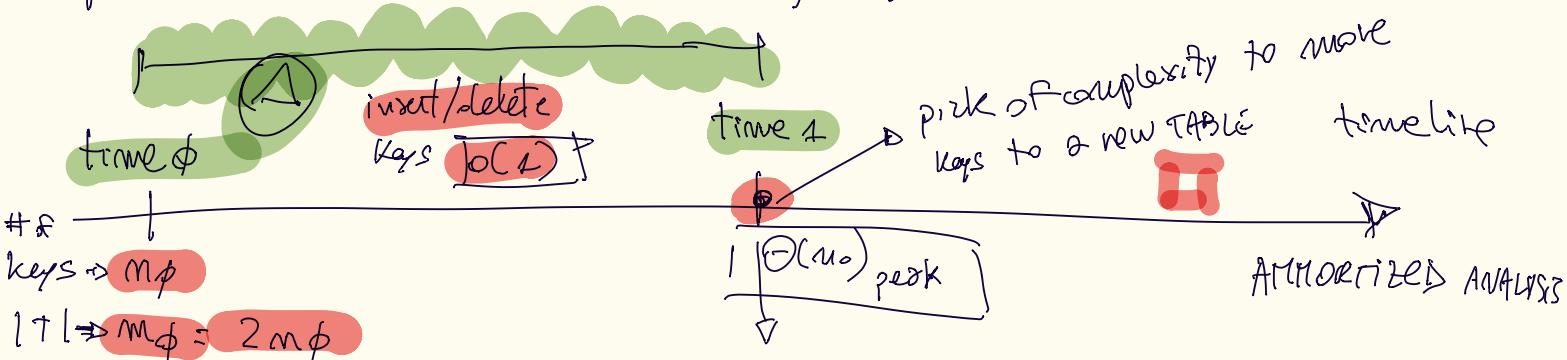
2) M is dynamic (changes) \Rightarrow managed by

GLOBAL REBUILDING

- want: constant time operation when M is dynamic (# of keys)

- at time ϕ , we have some # of keys n_ϕ and size of table $m_\phi = 2m_\phi$

- keys are added to the table, look how complexity changes to the hash approach



\Rightarrow insertion & deletion

$\Rightarrow O(1 + \frac{m}{n}) = \text{CONSTANT TIME}$

$$m_\phi = 6$$

$$m_\phi = 2 \cdot 6 = 12$$

$$\alpha = \frac{m_\phi}{2m_\phi} = \frac{1}{2}$$

$$\frac{n_\phi}{2} = \frac{6}{2} = 3$$

ESCAPE TIME

$\begin{cases} 1) \# \text{ of keys} \geq 2m_\phi \text{ (fill the table)} \\ \text{or} \\ 2) \# \text{ of keys} \leq \frac{m_\phi}{2} \end{cases}$

if ① or ② are true \Rightarrow do something

NOT in ① or ②

$$\frac{m_\phi}{2} < \text{# keys} < \frac{2m_\phi}{m_\phi}$$

$$\frac{m_\phi}{2} < \frac{m_\phi}{4} < \frac{2m_\phi}{m_\phi}$$

\Rightarrow in this range $|T|$ at the beginning is proportional to the # of keys (operation executed between T_1 & T_2 is constant)

In cases ① or ② # keys are too large or too small



in those cases:

table DOWNS

table HALVES

\Rightarrow RESIZE THE TABLE IN ① OR ②

- in case (1) \Rightarrow TABLE size gets $M_1 = 2m\phi = 4m\phi$ (worst)
- in case (2) \Rightarrow TABLE size gets $M_1 = \frac{m\phi}{2} = \frac{2m\phi}{2} = m\phi$ (best)

(4) In that time frame, between T_0 & T_1 is
 IMPORTANT that is guaranteed that # of keys (m) is proportional (for a given constant) to the size of the table
 \Rightarrow operations cost constant time $\Theta(1)$

(5) How much it costs?
 (REBUILDING AT T_1 ?)

CASE (1): Keys doubles
 \Rightarrow • move keys into the new table \Rightarrow I'm moving $2m\phi$ keys
 • destroy old table

1) Having $2m\phi$ keys \Rightarrow insertion: $\Theta(2m\phi \cdot 1)$

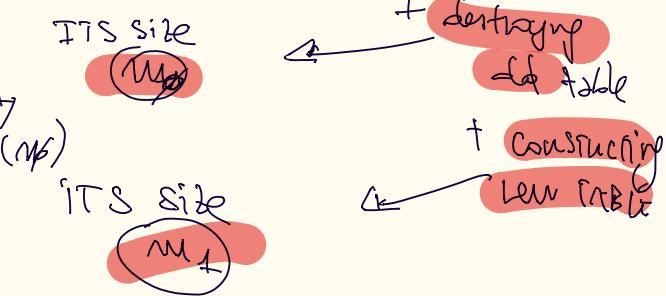
2) Destroy old table $T_0 \Rightarrow \Theta(m\phi) = \Theta(2m\phi)$

3) Create new table $T_1 \Rightarrow \Theta(M_1) = \Theta(4m\phi)$

\Rightarrow AMMORTIZED $\frac{\Theta(m\phi)}{m\phi} = \Theta(1) \Rightarrow$ constant

$$\Rightarrow \Theta(2m\phi \cdot 1 + M_1 + m\phi) = \Theta(2m\phi + 4m\phi + m\phi) = \Theta(7m\phi)$$

$$= \boxed{\Theta(m\phi)}$$
 take a cost proportional to the keys



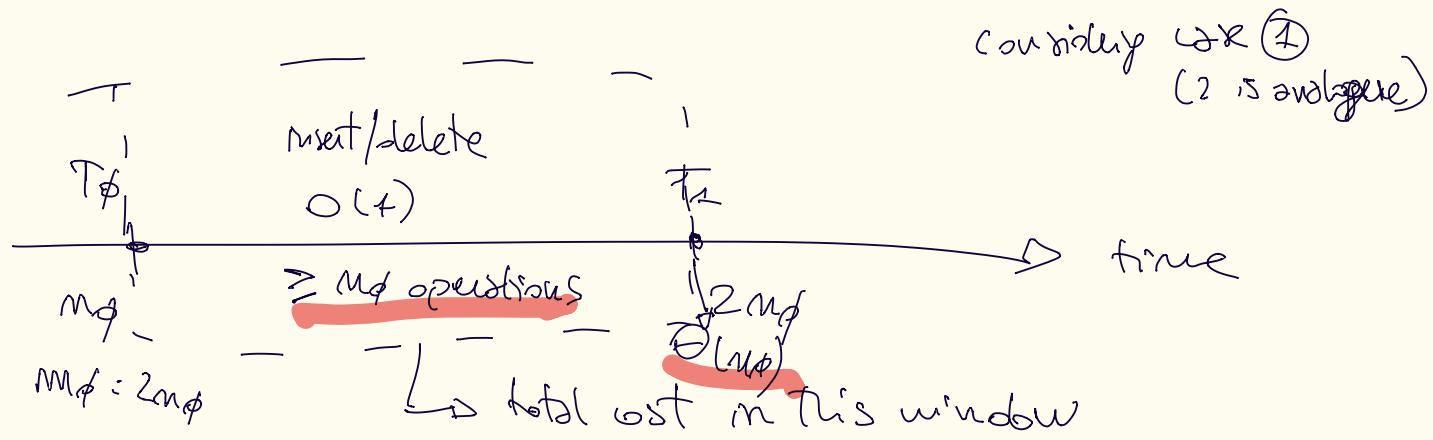
\Rightarrow so at T_1 on the line I'm getting $\Theta(m\phi)$ (peak)
 exactly this cost

point-wise complexity

AMMORTIZED ANALYSIS: not counting (constant, m_0 , constant).

but I'm counting complexity over a sequence of operations

- consider window $T_0 - T_1$



\Rightarrow start with M_0 , when arrive to T_1

keys gets $2M_\phi$

\Rightarrow at least I did M_ϕ operations

$\geq M_\phi$ operations

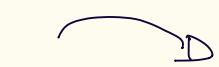
(consider insertion)

(each costs $O(1)$)

\rightarrow so in the window

I pay at least

$$(\Rightarrow) M_\phi + M_\phi = 2M_\phi \text{ so asymptotically } O(M_\phi)$$



\Rightarrow spread total cost to each (M_ϕ operation)

(like a loan)

(involving at least M_ϕ operations)

\Rightarrow worst case cost of a sequence of operations

of operations

like $\frac{\text{WC cost}}{\text{AVG cost}}$

$$\Rightarrow \frac{\Theta(M_\phi)}{M_\phi} \Rightarrow \text{constant } O(1)$$

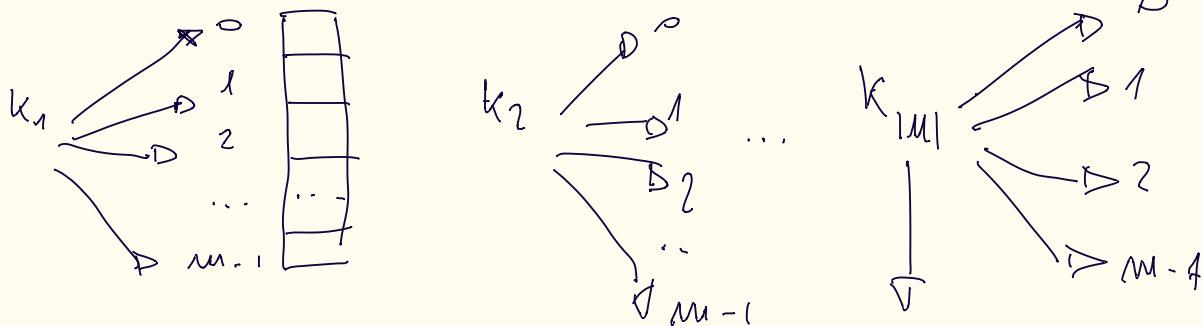
\Rightarrow cost of hashing is amortized constant

\Rightarrow so all SK if I have simple uniform hashing

NOT POSSIBLE

SIMPLE UNIFORM HASHING is not possible!

- ↳ every key can go everywhere in the table (uniformly distributed at random)
- ~ $h: U \rightarrow [m]$
- To have h simple uniform, h must be selected from a class of hash functions ($h \in \mathcal{H}$) such that is able to map every key k in U in every slot of the table T .



In universe U +ive m keys

Mappings: $m^{|U|}$
= $\#$ of simple hash f(x)s $|H| = m^{|U|}$
= $\#$ of mappings: $m^{|U|}$

How big is H ? $\Rightarrow |H|$ must include all $m^{|U|}$ mappings

$$\Rightarrow |H| = \# \text{ mappings} = m^{|U|}$$

To represent each $h \in H$ need at least $\geq \log_2 m^{|U|}$ bits
 \downarrow

bits to represent something: \log_2 something

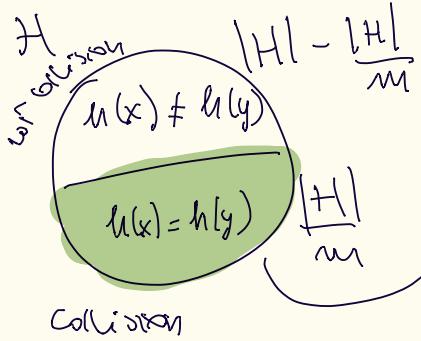
$$\Rightarrow \sum (\log_2 m^{|U|}) = \sum (|U| \log_2 m) \text{ bits}$$

\Rightarrow every function h requires $|U| \log_2 m$ bits

or. Keys of 8 bytes, $m = 2^{8 \text{ bytes}} = 2^{8 \cdot 8 \text{ bits}} = 2^{64}$ Size of the universe
 \Rightarrow NOT AFFORDABLE

UNIVERSAL HASHING

\rightarrow class \mathcal{H} of functions $h: U \rightarrow [m]$ is called universal hashing iff
 & pair of keys $x \neq y \in U$ the $(\# \text{ of } h : h(x) = h(y)) \leq \frac{|\mathcal{H}|}{m}$



such that
you have
a collision

↓
no good
hash fns

↓
represented

so

given $x, y \in U$
 $(x \neq y)$

so $P(\text{given } x, y \in U, x \neq y \text{ s.t. } h(x) = h(y))$

$$= \frac{|\mathcal{H}|}{m} = \frac{|\mathcal{H}|}{m} \cdot \frac{1}{|\mathcal{H}|} = \frac{1}{m}$$

- reasoning with U is too much

\Rightarrow do reasoning with pairs of objects
 (x, y)

- if pair of keys, the
a collision with
probability $\frac{1}{m}$

Universal Hashing with chaining

- pick $h \in \mathcal{H}$ at random
- want to see avg list length

$$\mathbb{E}[\text{list length}] = \sum_{\substack{\text{of cell } i \\ h(k)}} 1$$

$$\leq \sum_{k' \in D} \frac{1}{m} \cdot 1 = \frac{|D|}{m} = \frac{m}{m} = 1$$

$$1 \cdot P(h(k) = h(k')) \leq$$

↑
picked at random

$$\frac{1}{m} = \frac{|\mathcal{H}|}{m} = \frac{|\mathcal{H}|}{m} \cdot \frac{1}{|\mathcal{H}|} = \frac{1}{m}$$

by the definition
of universal hash
function

see as simple uniform hash

BUT EXISTS?

example of universal hash function

- class of universal \mathcal{H}

- $h: U \rightarrow [m] \Rightarrow$ key is $\log_2 |U|$ bits

- every key K taken from U consists $\log_2 U$ bits, $U = |U|$

- take the key, divide it into blocks: $k_0, k_1, k_2, \dots, k_{R-1}$ where $R = \frac{\# \text{ of blocks}}{\log_2 m}$

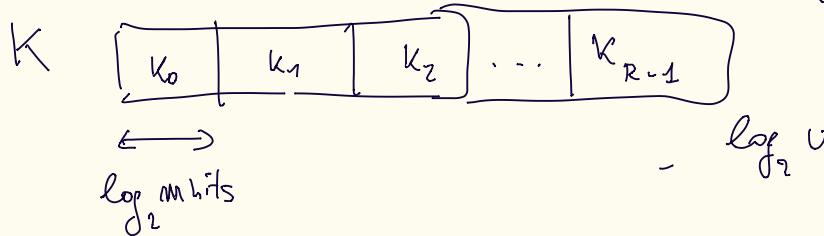
- every block is $\log_2 m$ bits, $\log_2 m$ bits are k_i

$$0 \leq k_i < m$$

opposite

every block is at most the size of domain

key



Picked at random

- construct hash function taken from a class, \Rightarrow I'll have a parameter a that defines the hash functions.

- of size $\log_2 U$

$a \in \{a_0, a_1, a_2, \dots, a_{R-1}\}$

picked at random

$\log_2 m$ bits

$\Omega = \frac{\log_2 U}{\log_2 m}$

$$\Omega = \frac{\log_2 U}{\log_2 m}$$

$$0 \leq a_i < U$$

I multiply
a · key
(but w/out ϕ)

Class \mathcal{H} :

- hash function:

$$h_a(k) = \sum_{i=0}^{R-1} a_i \cdot k_i \bmod m$$

where $a = [a_0, a_1, a_2, \dots, a_{R-1}]$

all blocks sum key prime

- where

$$0 \leq a_i < U$$

$$0 < a < U$$

$$|\mathcal{H}| = \#\alpha = U - 1$$

$\frac{1}{U}$ parameter

at $0 \leq a \leq U$

How many
hash functions
I've?

To be a universal hash function class we have to prove

that: \forall given $x, y; x \neq y$ $\#h_a: [h_a(x) = h_a(y)] \leq \frac{|\mathcal{H}|}{m}$

of hash \rightarrow count this number

\Rightarrow

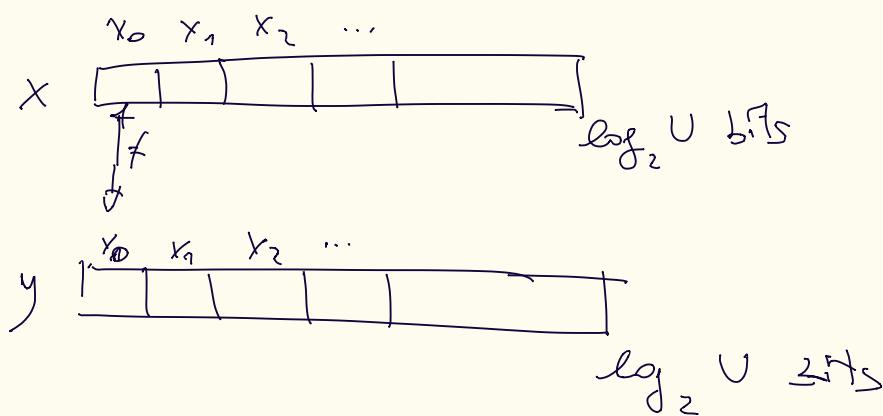
- in order to have a collision

want to count $\#h_a: \sum_{i=0}^{R-1} a_i \cdot x_i \bmod m = \sum_{i=0}^{R-1} a_i \cdot y_i \bmod m$.

so there are ~~number of~~ waste functions such that this equality is true \textcircled{X}

\Rightarrow at least there is 1 pair of blocks that is different

• assume $x \neq y$



\textcircled{X}

$$\Rightarrow \#h_a: \sum_{i=0}^{R-1} a_i \cdot x_i \equiv \sum_{i=0}^{R-1} a_i \cdot y_i \pmod{m}$$

1) I know just x_0 and y_0 are different (takes x_0 out and \cancel{y}_0 out)

$$a_0 \cdot x_0 + \sum_{i=1}^{R-1} a_i \cdot x_i \equiv a_0 \cdot y_0 + \sum_{i=1}^{R-1} a_i \cdot y_i \pmod{m}$$

$$\Rightarrow a_0 \cdot x_0 - a_0 \cdot y_0 \equiv \sum_{i=1}^{R-1} a_i \cdot x_i - \sum_{i=1}^{R-1} a_i \cdot y_i \pmod{m}$$

$$a_0 (x_0 - y_0) \equiv \sum_{i=1}^{R-1} a_i (x_i - y_i) \pmod{m}$$

$\neq 0$ since by ①
the assumption
they're different

prime ②

an integer

$$\text{since } ① \& ② \Rightarrow \exists \text{ an inverse } \Rightarrow \exists V (x_0 - y_0)^{-1} \pmod{m}$$

$\Rightarrow \exists$ an integer such that $(x_0 - y_0)^{-1} \cdot (x_0 - y_0) = 1 \pmod{m}$

\Rightarrow multiply all sides by this integer: $(x_0 - y_0)^{-1}$

$$a_0 = (x_0 - y_0)^{-1} \cdot \sum_{i=1}^{R-1} \alpha_i (x_i - y_i) \pmod{m}$$

#a are found by

$$\begin{matrix} \alpha_0 & , & \alpha_1 & , & \alpha_2 & \cdots & \alpha_{R-1} \\ \downarrow m & & \downarrow m & & \downarrow m & & \downarrow m \\ \log_m \text{ bits} & & \text{idem} & & & & \text{(A)} \end{matrix}$$

in this
way
(fixed)

Fixed

\Rightarrow can be chosen ~~anyways~~
in m -ways

To have a collision $\Rightarrow \alpha_0$ must not be chosen as I want it

\Rightarrow How many #a induce the collision? (A)

$$(\text{the: } h_a(x) = h_a(y) \leq \frac{|H|}{m})$$

$$\Rightarrow m^{R-1} - 1$$

but a to

$$2^{\log m^R} = 2^{R \log m}$$

$$m^{R-1} - 1 = m^R : m^1 - 1 = \frac{m^R}{m} - 1 = \frac{2^{R \cdot \log m}}{m} - 1$$

$$\text{NB: } R = \frac{\log v}{\log m}$$

$$= \frac{2^{\frac{\log v}{\log m} \cdot \log m}}{m} - 1 = \frac{2^{\log_2 \frac{v}{m}}}{m} - 1$$

$$= \frac{v}{m} - 1 = \frac{v-m}{m}$$

$$\leq \frac{v-1}{m} = \frac{|H|}{m} \checkmark$$

$\text{avg } H = \# \text{ of hash functions} = d = m$

just another example

$$\Rightarrow \text{so } h_a(k) = \left[\sum_{i=0}^{d-1} a_i \cdot k^i \bmod p \right] \bmod m$$

↑
prime
 $p > m$

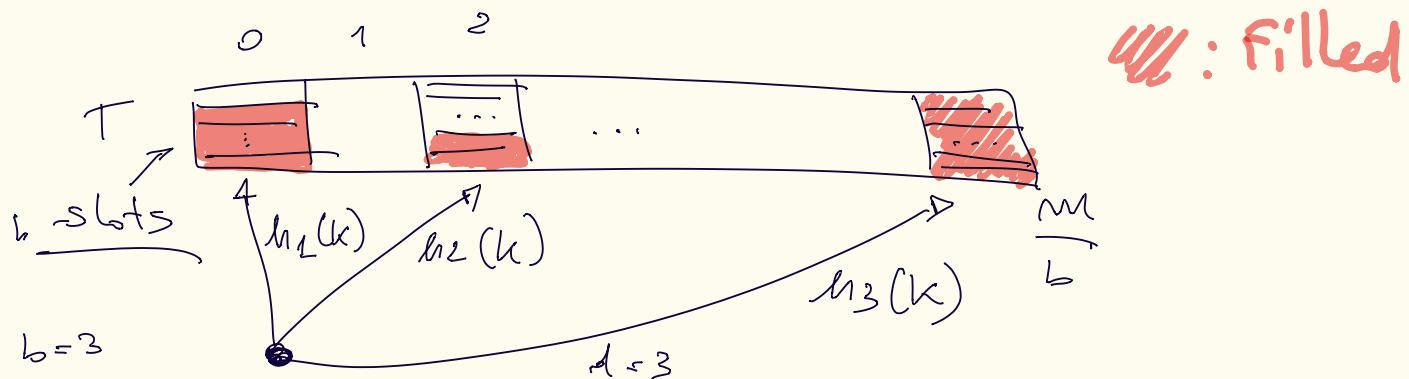
↓
power of 2

D - LEFT HASHING

- Want avoid scanning list } of hash with chaining
- " reduce # of pointers }

\Rightarrow idea

- take the TABLE
- PARTITION IT INTO BEANS
- EVERY BEAN has d -slots \Rightarrow table is $\frac{m}{d}$
- $d \leq \# \text{ of hash functions}$
- ex. $d = 3$



IDEA: Given a key, since I've 3 hash fns, I check 3 slots
 \Rightarrow AMONG SLOTS, TAKE THE MOST EMPTY (2)

- Search is constant $O(3 \cdot d)$

$E[\text{length of LONGEST list}] = \frac{\log \log m}{\log d}$

$d=2$ best case
 $d \geq 2$
 $+ O(1)$

• In $(d=1)$ we don't have choice \Rightarrow we insert directly (sort of chaining)
↳ complexity

$$\Downarrow \quad E[\text{length of longest list}] = \frac{\log M}{\log \log M}$$

Power of 2 choices

\Rightarrow with 4 choice is changing

\Rightarrow switching from 2 choices

\Rightarrow LENGTH OF THE LONGEST LIST DROPS .

CUK20:

Theo: \forall pair of positions i, j , $\forall c > 1$; if $m \geq 2 \cdot c \cdot n$
 $\Rightarrow P(\textcircled{i} \xrightarrow{?} \textcircled{j}) = \frac{1}{m \cdot c^L}$

$$\Rightarrow \alpha = \frac{m}{m} \leq \frac{\frac{m}{m}}{2c} \cdot \frac{1}{\cancel{m}} = \frac{1}{2c} < \frac{1}{2}$$

$\Rightarrow c > 1$

proof by induction

BASE CASE $L=1$

$$P(\textcircled{i} \xrightarrow{?} \textcircled{j}) = P(\exists \text{ a key } k:$$

$$\begin{aligned} h_1(k) &= i \wedge h_2(k) = j \\ h_1(k) &= j \wedge h_2(k) = i \end{aligned} \leftarrow \begin{array}{l} + (v) \\ + (v) \end{array} \right) \leq$$

union bound

$$\leq m \cdot P(k: \text{satisfy })$$

$$= m \cdot \left(P(h_1(k) = i \wedge h_2(k) = j) + P(h_1(k) = j \wedge h_2(k) = i) \right) =$$

$$= m \cdot \frac{2}{m^2} \leq \frac{m}{c} \cdot \frac{1}{m} = \frac{1}{cm} \quad \checkmark$$

INDUCTIVE STEP (True for $L-1$, proof for L)

$$P(\textcircled{i} \xrightarrow{?} \textcircled{j}) = P(\exists z: (\textcircled{i} \xrightarrow{?} \textcircled{j}) \xrightarrow{L-1} (\textcircled{j} \xrightarrow{?} \textcircled{z})) \leq$$

$$\leq m \cdot P(\textcircled{i} \xrightarrow{?} \textcircled{j} \xrightarrow{?} \textcircled{z}) \stackrel{\text{inductive hypothesis}}{\leq} m \cdot \left(\frac{1}{m \cdot c^{L-1}} \cdot \frac{1}{c \cdot m} \right) =$$

$m = \# \text{ of nodes}$
 $(m = \# \text{ edges})$

$$= \frac{1}{m \cdot c^L} \quad \checkmark$$

REPEAT proof

THEO: $\forall i, j, \forall c > 1$ if $m \geq 2 \cdot c \cdot m$ $P(\textcircled{0} - \textcircled{j}) = \frac{1}{m \cdot c^2}$
inductive step ($L=1$)

$$P(\textcircled{0} - \textcircled{j}) = P(\exists k : h_1(k) = i \wedge h_2(k) = j) \leq \\ \underbrace{\sum_{k=1}^m}_{\text{m choices}} P(\exists k : h_2(k) = j \wedge h_1(k) = i) \stackrel{\text{satisfy that condition}}{\approx} \frac{1}{m}$$
$$\leq m \cdot P(\exists k : \text{satisfy that condition}) = \\ \leq m \cdot \frac{2}{m^2} \leq \frac{m}{c} \cdot \frac{1}{m} = \frac{1}{m \cdot c}$$

inductive step: ok ~~for L=1~~, proof for L

$$P(\textcircled{0} \xrightarrow{L} \textcircled{j}) = P(\exists z : (\textcircled{0} \xrightarrow{L-1} \textcircled{j} \xrightarrow{1} \textcircled{z})) \leq \\ \leq m \cdot P(\textcircled{0} \xrightarrow{L-1} \textcircled{j} \xrightarrow{1} \textcircled{z}) \stackrel{\substack{\text{chosen from} \\ \text{the } m \\ \text{hypothys}}}{\leq} m \cdot \left(\frac{1}{m \cdot c^{L-1}} \cdot \frac{1}{m \cdot c} \right) = \\ = \frac{1}{m \cdot c^L} \quad \checkmark$$

WHERE CLKOP MASTING FAILS \Rightarrow cycle

$$P(\textcircled{i} \textcircled{j}) \leq P(\textcircled{i} \textcircled{j}) =$$

$$= P(\exists i, \exists l : \textcircled{i} \textcircled{j}^l) = \sum_i \sum_{l \geq 1} \frac{1}{m \cdot c^2} = \\ = \frac{1}{m} \sum_{l \geq 1} \frac{1}{m \cdot c^l} = \sum_{l \geq 1} \left(\frac{1}{c} \right)^l = \frac{1}{c-1} \Rightarrow$$

\Rightarrow PROBABILITY OF FAILURE WHILE BUILDING THE GRAPH
 \Rightarrow SURELY KNOWS THAT CONVERGES TO

Conclusion 1: if $c = 3$ R (FTA2S only constraint) $\leq \frac{1}{2}$

Conclusion 2: $P(\textcircled{1} \xrightarrow[\substack{\text{aka} \\ i, j \text{ collide}}]{\substack{\text{is connected}}} \textcircled{2}) = P(\exists \text{ a path from } \textcircled{1} \xrightarrow[\substack{\text{aka} \\ i, j \text{ collide}}]{\substack{\text{is connected}}} \textcircled{2} \text{ to } \textcircled{2}) \leq$

$$\leq \sum_{L \geq 1} \frac{1}{m \cdot c^L} = \frac{1}{m} \sum \frac{1}{c^L} = \frac{1}{m} \cdot \frac{1}{c-1} = \frac{1}{m(c-1)},$$

well known series \Rightarrow

$$= \frac{1}{m} \text{ if } c \geq 2$$

$\Rightarrow P(\text{Collision}) \leq \frac{1}{m}$

Avg # keys colliding with some key k : $\frac{m}{m} < 1$
but $m \leq 6m$

\Rightarrow Avg cost of

insertion $O(1)$
time

MINIMAL ORDERED (perfect) HASHING , Given a Dictionary

if $k_1 < k_2$

then $h(x_1) < h(x_2)$

NO COLLISIONS

- is static
 - cool for strings (for the rock)



→ Containment is $[m]$ ($\#$ of keys)

$$D = \{ AA, BB, BD, CD \} \quad \text{converges}$$

Elements: 0, 1, 2, 3 (rank) \leftarrow want to get

	b_1	b_2	b_1
$A A$	4	2	0
$B B$	1	4	1
$B D$	3	6	2
$C D$	6	0	3

want create an hash function that

- h_1 , h_2 ^{adversaries}
with random

$$l_{n+1}(x_j) \leq 3 \operatorname{rank}(x_j)$$

$$h_2(xg) = \text{rank}(x) + \text{rank}(sg) \pmod{2}$$

$$(m)_7 > 4$$

conform must be \geq than

number
of strings

To GET h;

$$h(t) = g(h_1(t)) + f(h_2(t)) \bmod m$$

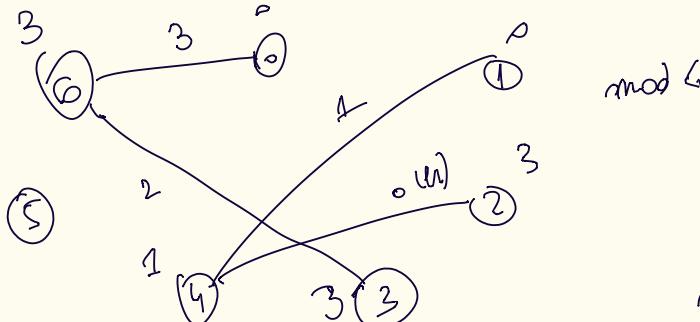
1st example $f(1) + f(2) \bmod 4 = 0$

represent f as an array

A diagram of a linked list with 6 nodes. The nodes are labeled p , 1 , 2 , 3 , 4 , and 5 . The node labeled 5 has a curved arrow pointing to the right, labeled $m-1$.

$$|g| = m = j$$

- constraint graph: nodes = $\lg l = m$



edges: values $\ell_{h_1} \sim \ell_{h_2} = m$

$\approx \#$ of strings

$P(\text{cycle})$ is small

$M = m \log M$, graph is sparse

$$f(6) + f(6) \bmod 4 = 3$$

$$\begin{array}{cc} f(6) + f(3) \bmod 4 = 2 \\ \downarrow & \downarrow \\ 3 & 3 \end{array}$$

$$\begin{array}{cc} f(1) + f(n) \bmod 4 = 1 \\ \downarrow & \downarrow \\ 0 & 1 \end{array}$$

• if I've 2 cycle \Rightarrow review others h_1 & h_2

• SPACE: store h_1, h_2 (universal hashes)

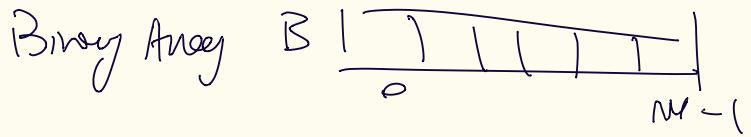
store f \Rightarrow m cells

$$m \approx m \log m$$

Bloom Filter

• r : hash functions

- membership (k) } $O(r)$ time
- insertion (k)
- no deletion



- insertion (k):

$$\forall i = 1, 2, \dots, r \quad B[h_i(k)] = 1$$

- membership (k) = yes iff $\forall i = 1, 2, \dots, r \quad B[h_i(k)] = 1$ (FALSIFYING POSITION)

$$\text{as if } \forall i = 1, 2, \dots, r \quad B[h_i(k)] = 0$$

prob of errors:

→ r -hash functions

$$P(B[j] = 0) = \left(\frac{m-1}{m}\right)^{r \cdot m} = \left(1 - \frac{1}{m}\right)^{r \cdot m} = \left(1 - \frac{1}{m}\right)^{\frac{r \cdot m}{m} \cdot m} = \\ = \left[\left(1 - \frac{1}{m}\right)^m\right]^{\frac{r \cdot m}{m}} \approx e^{-\frac{r \cdot m}{m}}$$

$$\text{so } P(B[j] = 1) = 1 - e^{-\frac{r \cdot m}{m}}$$

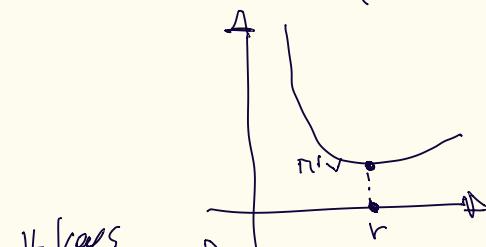
$$\Rightarrow P(\text{errors}) = P(k \notin D, B[h_i(k)] = 1 \forall i = 1, 2, \dots, r) \approx \\ \approx \left(1 - e^{-\frac{r \cdot m}{m}}\right)^r \quad \checkmark$$

• # hash functions that minimize the error: $\frac{m}{m} \ln 2$

• $r = \#$ hash functions

$$r_{\text{optimal}} = \frac{m}{m} \ln 2$$

$$\epsilon_{(\text{error})} = \left(\frac{1}{2}\right)^{r_{\text{optimal}}}$$



$$\# \text{keys} \# h \# \text{sites of BF}$$

$$P(\text{error}) = f(m, m, r) = \left(1 - e^{-\frac{r \cdot m}{m}}\right)^r$$

THEOREM:

SPACE

every DS that makes an error ϵ over m keys (occupying m space) must uses at least m -bits such that:

$$m \geq m \log_2 \frac{1}{\epsilon}$$

$$\epsilon = \left(\frac{1}{2}\right)^{r \text{ optimal}}$$

$$\epsilon = \left(\frac{1}{2}\right)^{\frac{m}{n} \ln 2}$$

$$\frac{1}{\epsilon} = 2^{\frac{m}{n} \ln 2}$$

$$\log \frac{1}{\epsilon} = \log_2 2^{\frac{m}{n} \ln 2}$$

$$\log \frac{1}{\epsilon} = \frac{m}{n} \ln 2 \approx 1.44$$

$$\log \frac{1}{\epsilon} \cdot m \cdot \left(\frac{1}{\ln 2}\right) = m$$

$$\Rightarrow m = \frac{m \cdot \log_2 \frac{1}{\epsilon}}{\ln 2} = 1.44 m \log_2 \frac{1}{\epsilon}$$

SPACES

space
key

pointer

table

HASH WITH CHAINING: $m \log V + m \log m + m \log n$ BITS

•) address :

$\underbrace{m \log V}_{\text{key}} + \underbrace{m \log m}_{\text{pointers}} + \underbrace{m \log n}_{\text{TABLE}}$ BITS

•) Bloom filter :

m bits $m = c \cdot m$

\Rightarrow If key 1 pay c -bits

except of BF

$$\bullet P(\text{error}) = \left(\frac{1}{2}\right)^{\frac{m}{n} \cdot \ln 2}$$

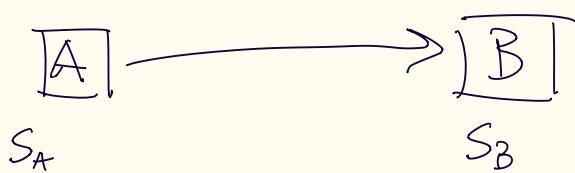
$$\bullet \sigma - \text{optimal} = \frac{m}{n} \cdot \ln 2$$

↓
+ hash fs

• P2P

• A, B want to exchange data (film, movie)

• estimate $A \cap B \Rightarrow S_A \cap S_B$
(with and without errors)



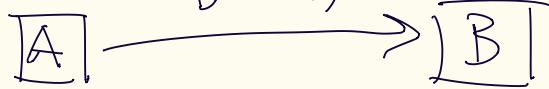
- S_A and S_B are set
of items of the
respectively machines

• A sends to B its set and B computes the intersection
 \Rightarrow too costly!

\Rightarrow send 2 bits really needed

• compute BF for S_A : define size and # functions

$\text{BF}(S_A)$



$\rightarrow S_A$

$\bullet S_B$

\Rightarrow B receives $\text{BF}(S_A)$

$\rightarrow \text{BF}(S_A)$

\bullet

\bullet B computes $S_A \cap S_B$

if size $m = c \cdot n \rightarrow (0, 6195)^c$
error

and $r = \frac{m}{n} \cdot \ln 2$ hash factors

$n = \text{keys}$

\Rightarrow size $\bullet m_A = c \cdot |S_A|$

$\Rightarrow \bullet r_{\text{av}} = \frac{m_A}{M_A} \cdot \ln 2$

\bullet $\text{BF}(S_A)$
if $r > r_{\text{av}}$ \Rightarrow element

is surely
notin the
intersection

\bullet if yes \Rightarrow maybe
we take the element
as belong to intersection
(maybe correct, maybe error)

☒ Yes:

- element such that $s_A \cap s_B$
- extra elements (BF says yes but $\notin s_A \cap s_B$)
- $k \rightarrow \text{BF}(k) = \text{yes}$ but $k \notin s_A \cap s_B$
- CANNOT AVOID \Rightarrow compute the AVG
- \Rightarrow PROBABILITY OF THERE \in

 $\Leftrightarrow E[\# \text{ errors}] = P(\text{error}) \cdot \# \text{ of checks I do} \sim$
 $= \left(\frac{1}{2}\right)^{\text{round}} \cdot |s_B| =$
 $= \left(\frac{1}{2}\right)^{\frac{m}{m} \ln 2} \cdot |s_B| =$
 $\underbrace{\approx 0.6185}_{\text{.}} \cdot |s_B|$
 \Rightarrow if c is large the prob of error is small

• if I want to be sure, go back the intersection



(SPECTRAL Bloom Filter)

- BF with counters
- supports deletion
- r. hash functions
- insertion(k) \Rightarrow update (increment) the value of the counter
 $\forall i = 1, 2, \dots, r; C[h_i(k)] += 1$
- deletion(k) \Rightarrow update (decrement) the value of the counter
 $\forall i = 1, 2, \dots, r; C[h_i(k)] -= 1$
- freq(k): frequency of a key
 $\min(C[h_i(k)]) \quad \forall i = 1, 2, \dots, r$
- $P(\text{error}) \leq P(\text{overflow}) = \left(\frac{1}{2}\right)^{r-\text{OPT}} = \left(\frac{1}{2}\right)^{\frac{m}{m}(\text{err})} \approx (0, \epsilon)^{\frac{m}{n}}$

$$\begin{aligned} 2^{\frac{m}{n}} &= (2 \cdot 2)^3 \\ 2^6 &= 4^3 \end{aligned}$$

Repeat:

Hash with chaining

- in presence of a simple uniform hash \mathbb{I} the hash with chaining must support search in $O(1 + \frac{m}{n})$, $\frac{m}{n} = 2$ is ad factor

$$\mathbb{E}[\text{list length}] = 1 + \frac{1}{m} \cdot m = 1 + \frac{m}{m}$$

✓

$$\sum_{k=0}^m P(h(k) = i) = m \cdot \frac{1}{m}$$

if $m = n$

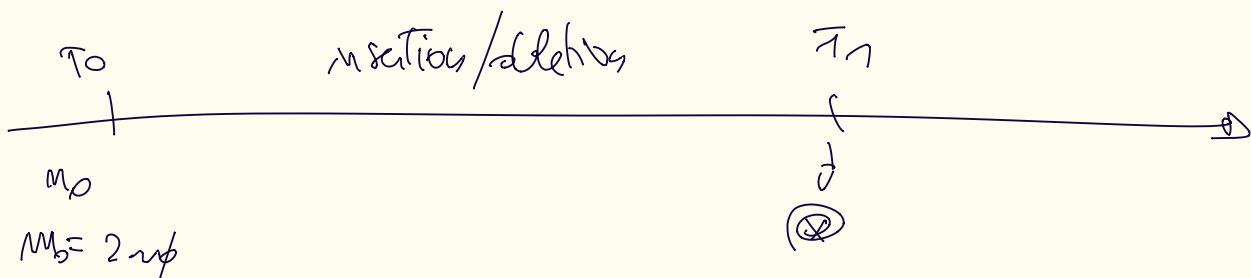
$$\Rightarrow 1 + \frac{1}{1} = \textcircled{1}$$

⇒

constant

ok if m is static

want this when m is dynamic



- could happen 2 cases: ⑦ m doubles $\leq 2m_0$ $m_1 = 2 \cdot 2m_0 = 4m_0$
⑧ m halves $\leq \frac{m_0}{2}$ $m_1 = \frac{2m_0}{2} = m_0$
(or \pm collapse for ⑧)

- ① cost of:
- moving m_0 keys = $O(1 + m_0) = O(m_0) +$
 - build new table = $O(m_1) = O(4m_0) = O(m_0) +$
 - destroy old table = $O(m_0) = O(2m_0) = O(m_0)$
- $\Rightarrow O(m_0)$

The search is constant if we consider the amortized cost:

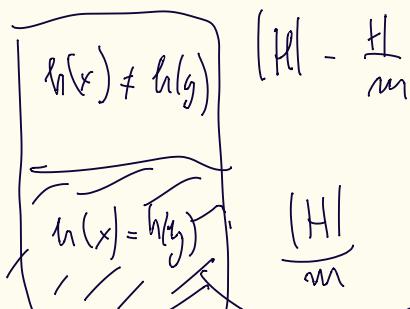
$$\frac{O(m_0)}{m_0} = O(1)$$

Universal Hash

$h: U \rightarrow [m]$ if $h \in H$, H is a class of universal hash functions if $\forall h \in H$ such that $\forall x, y \in U; x \neq y$

$$h(x) = h(y) \leq \frac{|H|}{m}$$

?



$$\Pr(\text{Im } h \text{ is full}) = \frac{|H|}{m} = \frac{|H|}{m} \cdot \frac{1}{|H|} = \frac{1}{m}$$

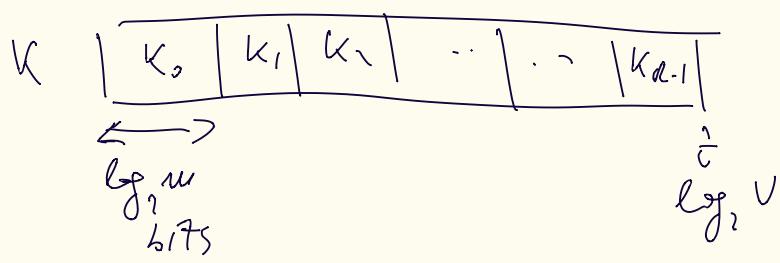
Universal Hash with chaining

pick $h \in H$ at random, AVG list length

$$\begin{aligned} \mathbb{E}[\text{list length}] &= \sum_{\substack{k \in D \\ k \neq k'}} 1 \cdot \Pr(h(k) = h(k')) = \\ &= m \cdot \frac{1}{m} = \frac{m}{m} \end{aligned}$$

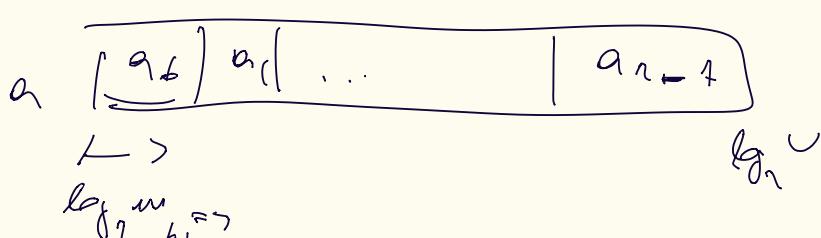
Give an example of Universal Hash function

- suppose K divided into blocks, ideally a power of 2 from 0 to $n-1$



$$R = \frac{\log_2 U}{\log_2 m}$$

$$\Rightarrow h_a = \sum_{i=0}^{R-1} k_i \cdot a_i \bmod m$$



$$\log_2 U$$

$$\begin{aligned} 0 \leq a_i &< U \\ 0 < a &< U \end{aligned}$$

m prime

To be universal hash function, given $x, y \in V$, $x \neq y$

$$\#_{\text{coll}} : (h(x) = h(y)) \leq \frac{|H|}{m}$$

that is

- to have no collisions

$$\#_{\text{coll}} : \sum_{i=0}^{n-1} x_i \cdot a_i \bmod m = \sum_{i=0}^{n-1} y_i \cdot a_i \bmod m$$

$$\#_{\text{coll}} : \sum_{i=0}^{n-1} x_i \cdot a_i \equiv \sum_{i=0}^{n-1} y_i \cdot e_i \bmod m$$

$$a_0 \cdot x_0 + \sum_{i=1}^{n-1} (x_i \cdot a_i) \equiv a_0 \cdot y_0 + \sum_{i=1}^{n-1} y_i \cdot e_i \bmod m$$

$$a_0 (x_0 - y_0) \equiv \sum_{i=1}^{n-1} a_i (x_i - y_i) \bmod m \quad \text{multiply for the inverse } (x_0 - y_0)^{-1}$$

$$a_0 \equiv \sum_{i=1}^{n-1} a_i (x_i - y_i) \cdot (x_0 - y_0)^{-1} \bmod m$$

To have no collisions

$$a_0 \not\equiv a_1 + a_2 + \dots + a_{n-1}$$

m -ways

$$m^{R-1} - 1 \leq \frac{m^R}{m} - 1 = \frac{2^{\log_2 m^R}}{m} - 1 =$$

$$= \frac{2^{\log_2 m^R / \log_2 m}}{m} - 1 = \frac{2^{\log_2 m^R / \log_2 m}}{m} - 1 =$$

$$= \frac{2^{\log_2 v}}{m} - 1 = \frac{v}{m} - 1 = \frac{v-m}{m} \leq \frac{v-1}{m} = \frac{|H|}{m} \quad \checkmark$$

cycles within

theo \forall position i, j , $\forall c > 1$ if $m \geq 2 \cdot c \cdot m$

$$P(\overset{c}{\overbrace{i \rightarrow j}}) = \frac{1}{m \cdot c^L}$$

case base $L=1$

$$P(i \rightarrow j) = P(\exists k : h_1(h) = i \wedge h_2(k) = j) \\ h_1(h) = i \wedge h_2(k) = j \quad \text{as.} \\ P(h_1(h) = i) = \frac{1}{m}$$

\Rightarrow

$$\leq m \cdot P(\text{selects } k) = m \cdot \frac{2}{m^2}$$

$$\frac{2m}{m^2} \leq \frac{m}{\frac{c}{m^2}} \quad \frac{2m}{m^2} \leq \frac{m}{c} \cdot \frac{1}{m^2} = \frac{1}{mc}$$

\Rightarrow Inductive step, true for $L=1$

$$P(i \overset{c}{\overbrace{\rightarrow}} j) = P(\exists z : i \overset{L-1}{\overbrace{\rightarrow}} j \overset{L-1}{\overbrace{\rightarrow}} z) \leq \text{upper bound if we use } m$$
$$\leq m \cdot \left(\frac{1}{mc^{L-1}} \right) \cdot \frac{1}{mc} = \frac{1}{mc^L} \quad (\checkmark)$$

Evaluate the probability that cycles having k buildings

$$P(i \overset{c}{\overbrace{\rightarrow}} j) \leq P(i \overset{c}{\overbrace{\rightarrow}}) = \text{re graph}$$

$$= P(\exists \text{ node } i : P(i \overset{c}{\overbrace{\rightarrow}}) = \frac{1}{mc^L}) =$$

$$\stackrel{\text{upper bound}}{\leq} m \cdot \sum_{L \geq 1} \frac{1}{mc^L} = \sum_{L \geq 1} \frac{1}{c^L} = \frac{1}{c-1}$$

Bloom Filter:

Insert(k)
Membership(k)

$R = \# \text{ of hash functions}$

\Rightarrow
Membership(k) $\Rightarrow \forall i=1, \dots, R \quad h_i(k) = 1$

\Rightarrow if yes $\xrightarrow{\text{has } b} \text{ false positive}$
no $\Rightarrow \text{ (prob)}$

Probability

$P[\text{error}]$

$$\cdot P(B[j] = 0) = \left(\frac{m-1}{m}\right)^{R \cdot m} = \left(1 - \frac{1}{m}\right)^{R \cdot m} \cdot \frac{m}{m} =$$

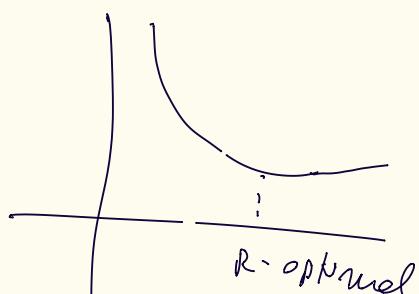
$$= \left[\left(1 - \frac{1}{m}\right)^m \right]^{\frac{R \cdot m}{m}} = e^{-\frac{rm}{m}}$$

$$\cdot P[B[j] = 1] = 1 - e^{-\frac{rm}{m}}$$

$$P[\text{error}] = P(k \notin D \text{ but } B[h_i(k)] = 1) =$$

$$= \left(1 - e^{-\frac{rm}{m}}\right)^R$$

$$R_{\text{optimal}} = \frac{m}{n} \ln 2$$



$$E_{\text{opt}} = \left(\frac{1}{2}\right)^{R_{\text{optimal}}} = \left(\frac{1}{2}\right)^{\frac{m}{n} \ln 2} = \ln 2 \left(\frac{1}{2}\right)^{\frac{m}{n}} =$$
$$= (0, 6)^{\frac{m}{n}}$$

Trees every DS that makes an error ϵ over m keys (of size n) must uses at least

Proof

$$m \geq m \log \frac{1}{\epsilon} \quad (6.7)$$

$$\epsilon = \left(\frac{1}{2}\right)^n \text{opt}$$

$$\epsilon = \left(\frac{1}{2}\right)^{\frac{m}{n} \ln 2}$$

$$\frac{1}{\epsilon} = 2^{\frac{m}{n} \ln 2}$$

$$\log_2 \frac{1}{\epsilon} = \log_2 2^{\frac{m}{n} \ln 2}$$

$$\log_2 \frac{1}{\epsilon} = \frac{m}{n} \ln 2$$

$$\left(\log_2 \frac{1}{\epsilon}\right) n = m \ln 2$$

$$m = \log_2 \frac{1}{\epsilon} \cdot n \cdot \left(\frac{1}{\ln 2}\right) \rightarrow 1, \text{ up}$$