

1
21/1

Today lesson is on Hash-tables.

Preface:

Definition : Abstract Data Type - (from geekforgees)

ADT from now on.

ADT is a type (class) for objects whose behaviour is defined by a set of value and a set of operations.

The definition of ADT only mentions what operations are to be performed, but not how these operations will be implemented. It does not specify how data will be organized in memory and what algorithms will be used for implementing the operations. It's called "abstract" because it gives an implementation independent view. The process of providing only the essentials and hiding the details is known as abstraction.

For example : Stack, we can define next operations:

add - element will be added to the set.

remove - in LIFO order.

Etc. There can be some implementation : by array, linked-list etc., but when we look at a stack as ADT we don't care about the implementations, all we know is that there are some operations and that elements will be removed by LIFO order.

Map $\langle K, V \rangle$ myMap = new



ADT

Tree Map

HashMap

LinkedHashMap

Implementations

(7)

Dictionary - ADT

Maintain set of items, each with a key

- insert(item) overwrite any existing key
- delete(item)
- search(key) : return item with given key, or ~~exists~~ report doesn't exist.

Note: It's a little different than BST; in BST if we didn't find a key, we could find the next smaller or larger key: Successor or predecessor, here we are not able to do so.

We are assuming that items have unique keys, no two items have the same key, and one way to enforce that is when you insert an item with an existing key it overwrite whatever key was there (Python, Java..)

One way to solve dictionaries is by balanced binary search trees as AVL - and we can do all of the operations in $\log(n)$ time (we ignore from the fact that doing a search in AVL gives us more information)

This is one solution - it turns out that we can get better.

Like we already saw, in the comparison model the best way to do sort is $n\log(n)$ and search in $\log(n)$, then we saw in the RAM model - if we assume that our items are integers we can sort faster

= sometimes in linear time - here we will see how to search faster than $\log(n)$ - and we are going to get down to constant time. - with no assumptions, except maybe that the keys are integers. $O(1)$ with high probability!

(3)

Python: dict

D[key] ~ search

D[key] = Val ~ insert

del D[key] ~ delete

→ item → (key, value) Java too.

Motivations: (as a result in every modern programming language)

- docdist - comparing between documents. (problem) [C, C++, Ruby, Python, Java]

- databases (hashing, search trees)

- compilers & interpreters

- network router

Server

- Substring search (editors)

- String commonalities

- file\dir synchronization [e.g. dropbox if you send a

- cryptography file, you check through hashing. if the file had changed]

Simple approach:

Direct-access table

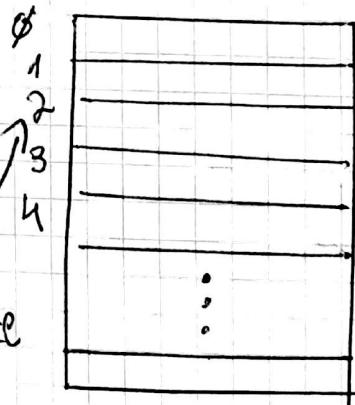
Store items in array indexed by

key [integers keys] - item with

index a - will be in second place
all empty slots will be null

2 problems: Badness

- ① Keys may not be normal integers!
- ② gigantic memory hog



Solution to ①:

prehash

4

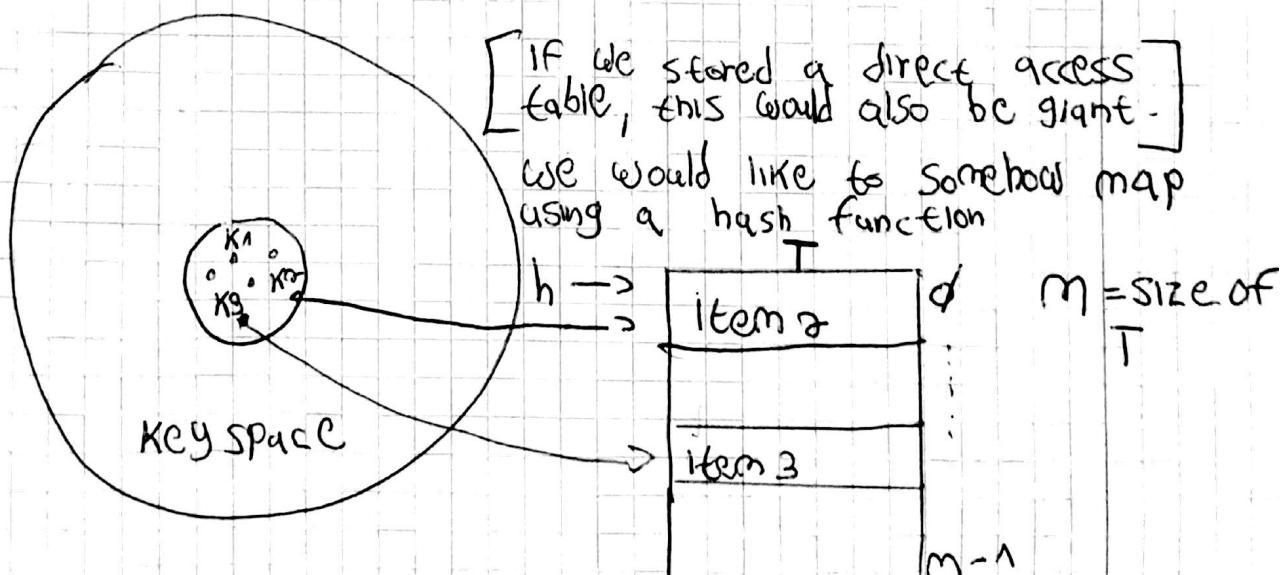
- maps keys to nonneg integers [string, object...]
- in theory keys are finite and discrete.
- In python, `hash(x)` is the prehash of `x` (if you give an integer as a key it returns the integer, and if you give a key it let say `x` so ~~the~~ `hash(x)` is integer)
- $\text{hash('1\phiB')} = \text{hash('1\phiC')} = 64$
- ideally: $\text{hash}(x) = \text{hash}(y) \Leftrightarrow x = y$
- `hash --` can implement `hash` method which tells python what to do, when you call `hash` of your object.
prehash must not change! Don't mess with it.
[for example a list which is a mutable object can't be in a hash table because if you append objects to the list it would get another value]

The more interesting problem is reducing space

Solution to ②: hashing

To hash means cut into pieces and mix around.

- reduce universe of all keys (integers) down to reasonable size m for table



Idea: $m = \Theta(n)$

(5)

n - number of keys in the dictionary (right nom)
That is if we want to store m items we use
 $2m, 3m$ space, but not much more.

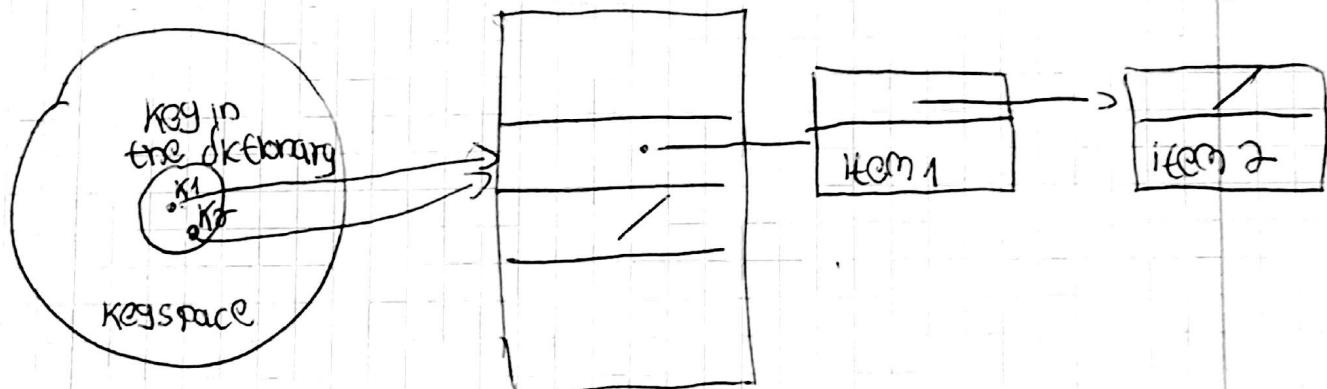
Collision: $h(k_i) = h(k_j)$, but $k_i \neq k_j$

and then there gonna be chaining.

By definition
hashing function
 $h: [0, 1, \dots, m-1]$
hashing function
needs to know what
 m is.

Chaining: linked list. of colliding elements

- In each slot of hash table



- Some slots can be null. Worst case all of the keys are mapped to the same slot by $h(\text{key})$
Worst case $\Theta(n)$ [the reason that it actually work is randomization]

Simple uniform hashing: [false assumption]
each key is equally likely to be hashed to any slot
of the table independent of where other keys
hashing [key one ~~map~~ goes to random place, no matter
what place key two goes to random place.]

(6)

Under this assumption we can analyze hashing with chaining

Analysis:

- expected length of chain for n keys, m slots. = n/m . The probability for each key

is $1/m$, $\underbrace{1/m + 1/m + \dots}_{n \text{ times}} = n/m = \alpha = \text{load}$

factor of
the table

As long as $m = \Theta(n) \rightarrow \alpha$ is constant.

$$= \Theta(1) \text{ running time} = \Theta(1 + |\text{chain}|) = \Theta(1 + \alpha)$$

[We write 1 because α , can be smaller than 1, and we always have the cost of computing the hash function!]

Assume Simple
uniform Hashing \rightarrow

As long as α is Constant
constant $m = \Theta(n) \rightarrow$ running time
all of operations

That is not really why hashing works \rightarrow it's an intuition
why hashing works

How do we construct h ? 3 hash functions
The first and the second are bad and the third
is good, but the first 2 function most of the
time ok!

① DIVISION method $h(k) = k \bmod m$

⊕

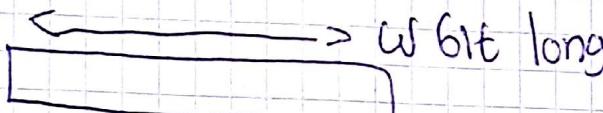
There can be a lot of situations, where it not works good especially when m and k have a common factor e.g. m is even and k always even - we use half of the table. Most of the time works good if m is prime and m not very close to a power of 2 or ten.

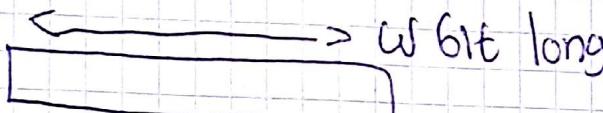
(2) multiplication method

$$h(k) = [(a \cdot k) \bmod 2^w] \Rightarrow (w-r)$$

assuming w bit machine

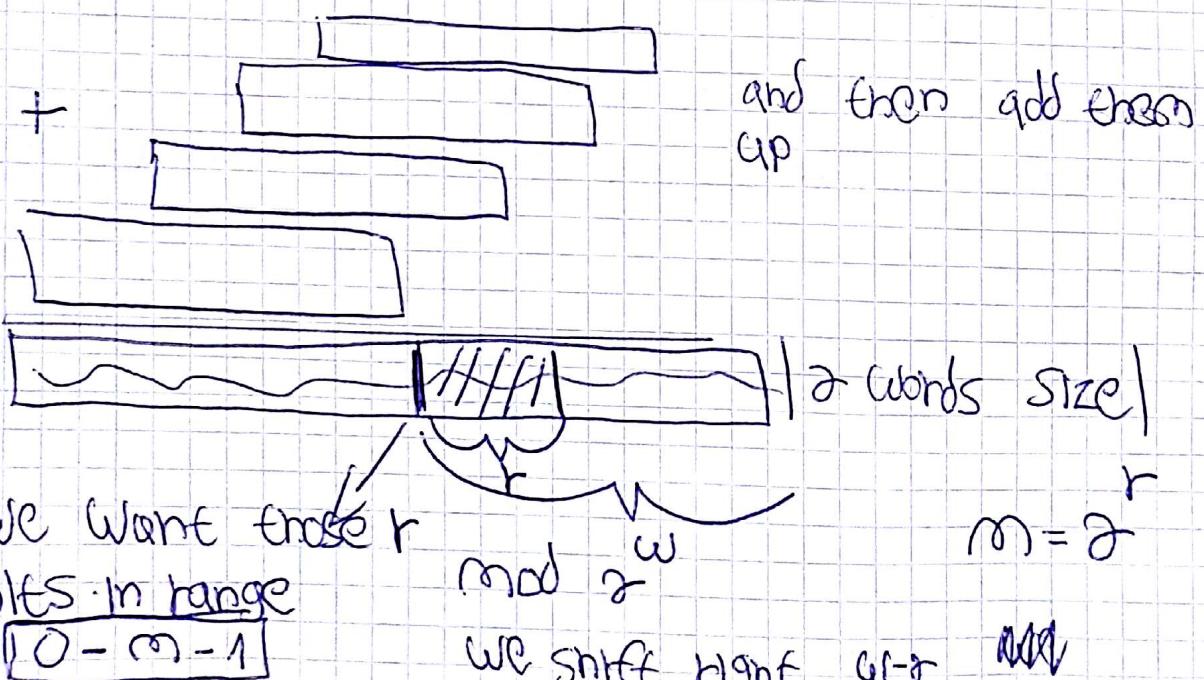
Key K



* X 

$a \cdot 1011010111011010$ we take random integer a among all w bit long integers

we have 4 times 1 in a and thus we have 4 copies of K



In the middle all of the randomization happen ⑧
It works because this is hashing - mixing up numbers
and numbers etc...

a better be odd, and should not be close to power
of two 2^r 2^{r-1} // theoretically this can be bad

③ Universal hashing

$$h(k) = [(ak+b) \text{ mod } p] \text{ mod } m$$

random $\in \{0, \dots, p-1\}$

key

to make it between 0 - $m-1$

prime $> |a|$
universe

for Worst case keys

$$\Pr \{ h(k_1) = h(k_2) \} = 1/m$$

Probability

That is that if we choose to use in Universal hashing
the probability that will be collisions is $1/m$ that is
in the worst case (because a, b are random) \rightarrow that is
the best case for us.

(9)

Crucial thing that is missing is: What should m be?

If it will be too small, eventually we will insert more and more keys and λ will be too big

If it will be big - kind of wasteful the whole idea is to avoid having one slot for every possible key

How to choose m ?

- want $m = \Theta(n)$
 $\lambda(n) \rightarrow \lambda \text{ constant}$ $\Rightarrow \lambda = \Theta(1)$
 $O(n) \rightarrow \text{space linear}$

Idea: Start small: $m = 8$

grow/shrink as necessary

If $n > m$: grow table

Grow table: $m \rightarrow m'$

- make table of size m'
- build new hash h' prime
- rehash for each item in T :
 $T'.insert(item)$ time

$$\Theta(h + m + m') = \Theta(n)$$

\downarrow \downarrow \downarrow
list slots new
table

(10)

The size of the new table $T \rightarrow m'$ have to update in the b' .

$$m' = 2m.$$

$$m' = m+1 \quad \text{Wrong answer}$$



Cost of inserts if we choose to increment m by 1

Cost of n insertions:



$$\Theta(1+2+3+\dots+n) \quad n \text{ insertions}$$

each ins cost of $\Theta(n) = \Theta(n^2)$

$\Rightarrow (1, 2, 4, 8, 16, \dots, n) = \text{geometric series } \Theta(n)$
(can take $O(n)$ time to double m)

This is data structure that suppose to be constant time for operation, here in insertion when we need to add and the table is full $\rightarrow O(n)$

$\log n$ of them are bad, but in the rest you don't do anything \rightarrow this is idea that we call

TABLE DOUBLING

Amortization:

- Operation takes $T(n)$ amortized if K operations take $\leq K \cdot T(n)$ time

think of meaning " $T(n)$ on average", where the average over all operations

In table doubling

So the amortized running is $\Theta(1)$.

k insertions take $\Theta(k)$ time that means

$\Theta(1)$ amortized / insert [you don't care that some insertions are expensive, most of the time is $\Theta(1)$]

Also: k inserts & deletes take $\Theta(k)$

problem: suppose we have n insertions, m

became bigger and then we delete all of the keys \rightarrow So now $n=0$, $m = \max$ value of n how can we fix that?

(1) If $m = n/2$ then shrink $\rightarrow m/2$

SLOW: 2^k insert $\longleftrightarrow 2^k + 1$ = $\Theta(n)$ per operation
Delete

[e.g. let say we have 8 size of m then we insert

1 key, double size of $m=16$, $n=9$ now we delete 1 key $m=16$, $n=8$ and then shrink again]

(2) If $m = \frac{n}{4}$ then shrink $\rightarrow m/2$

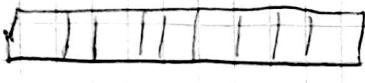
amortized time $\rightarrow \Theta(1)$

maintain the invariant

$n \leq m \leq 4n$

Side effect \rightarrow python lists

python lists also known as resizable arrays

can  random access - constant
append - constant

deleting first item linear time. - moving all other items.

So in fact `list.append` constant amortized!
`list.pop`

There is a way to get rid of the amortized, for people who work in real time systems

String Matching:

Given two strings s & t : does s occur as a substring of t ?

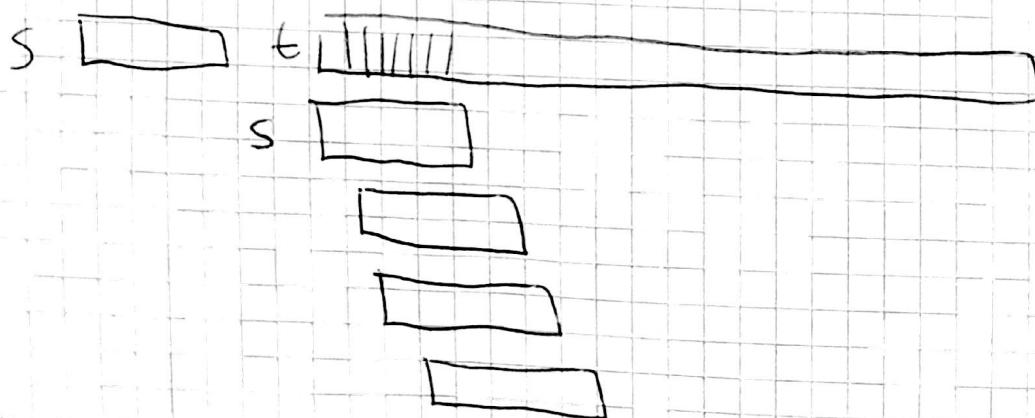
$$s = G_ooG \quad t = \boxed{INBOX}$$

Another application s is what we type in google and t is the entire web.

Simple algorithm

Python code:

```
any (s==t[i:i+len(s)]  
    for i in range(len(t)-len(s)))
```



$$\text{Time: } \Theta(|s| \cdot (|t| - |s|)) = O(|s| \cdot |t|)$$

How can we use Hashing to get linear time
 $\Theta(|s| + |t|)$ guaranteed.

Rolling hash ADT:

- $r.append(c)$: add char. c to the end of x
 - $r.skip(c)$ delete first char of x (assuming it is in c)
- r maintains a string x
- $r()$: hash value of $x = h(x)$

Karp - Rabin algorithm:

for c in s : $rs.append(c)$

rs gives a hash value of s

for c in t [$: \text{len}(s)$]:

$rt.append(c)$

if $rs() == rt()$: ...

for i in $\text{range}(\text{len}(s), \text{len}(t))$:

$rt.skip(t[i - \text{len}(s)])$

$rt.append(t[i])$

if $rs() == rt()$: (Potentially can be equals)

Check whether

$s == t[i - \text{len}(s) + 1 : i + 1]$

If equals found match

$O(s)$

else - happens with probability $\leq \frac{1}{|S|} \Rightarrow$

$O(1)$ expected time

$O(|S| + |E| + \# \text{match} \cdot |S|)$

rs represents hash value of s . $h(s)$

How to build ADT in order to spend constant time for each one of the operations?

division method

$$h(k) = k \bmod m$$



Random
Prime $\geq |S|$



[The bigger m is, the higher probability complexity will be as we analyzed.]

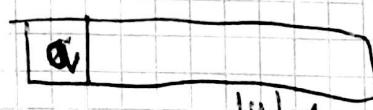
$r.append(c)$:  shift by one = multiply by a

$$a \rightarrow a \cdot a + \text{ord}(c)$$

treat x as a multidigit number (a) in

* base a - alphabet size ($\text{ascII} = 256$)

$r.skip()$



$$a \rightarrow a \overset{-}{\cancel{\cdot}} a \cdot c \cdot a$$

Hashing:

①

Open addressing.

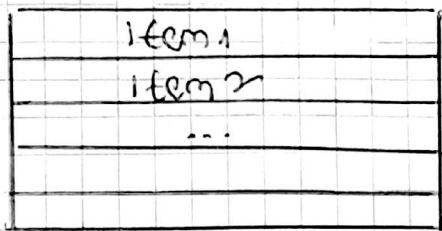
Uniform hashing analysis

Cryptographic hashing

Open addressing:

This is another approach to dealing with collisions.
no chaining

②



At most one per slot.

$m \geq n$



We don't have linked list we can't arbitrarily increase the storage of a slot using a chain

m - number of slots
 n " elements (keys)

③

Probing: (try to insert a key to a computed index of slot and if we fail - try to use a slightly different hash for the key)

Hash function: specified order of slots to probe for a key (for insert/search /delete)

hash function $h : U \times \{0, 1, \dots, m-1\} \rightarrow \{0, 1, \dots, m-1\}$ ↑ arguments
 universe of keys trial count

$h(k_1), h(k_2), \dots, h(k, m-1)$ } Vector
 arbitrary key k ↑ to be a permutation of $\{0, 1, \dots, m-1\}$

(2)

In order to get access for the entirety of our hash table. Even if we have 1 slot left we want to be able to insert to this slot through our hash function.

T-Table

0
1
2
3
4
5
6
7



(2)

Insert 586 $h(586, 1) = 1$

As long as the table is empty the insertion is successful.

$h(481, 1) = 6 \dots$ etc

Insert (496)

fail $\swarrow h(496, 1) = 4$ check of T index 4 != null
and first probe is a fail

$h(496, 2) =$ check = 586 = fail

$h(496, 3) = 3$ success

None: empty slot

(flag)

Insert (K, v)

Keep probing until an empty slot is found.

Insert item when found.

(3)

Search (k): As long as the slot encounters are occupied by keys $\neq k$.
 Keep probing until you either encounter k or find an empty slot (you use in the same probe sequence that used for insertion.)

Delete 586 :

0	
1	586
2	133
3	496
4	204
5	481
6	
7	
8	
9	

Problem:

None → According to the search ~~algo~~ algorithm

Search 496 is not exist!

How can we fix it?

Replace deleted item with with DeleteMe flag (diff from none)

0	
1	586
2	
3	
4	
5	
6	

And the psudo code still correct.

Insert treats deleteMe the same as None but
 Search keep going

We can code this up with an array structure, is fairly straightforward

Complexity analyze

Probing strategies:

Ordinary hash function

Linear probing

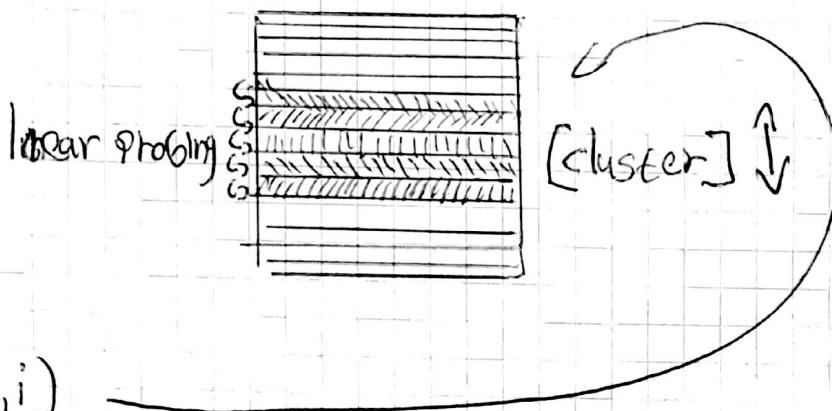
$$h(k, i) = (h'(k) + i) \bmod m$$

That is satisfy the permutation argument?

Yes Permutation ✓

Cluster - notion of clustering is a possibility of start getting consecutive groups of occupied slots, with keep growing

Assumption using linear probing can lead to losing average constant time



In terms of $0.01 < d = \frac{n}{m} < 0.99 \quad \Theta(\log n)$ size probability using linear probing can lead to $\Theta \log n$ size of insert, search which means losing the constant time advantage of hash tables [the problem is not the hash function the problem is the aspect of linearity.

(5)

Double Hashing:

$$h(k, i) = (h_1(k) + i \cdot h_2(k)) \bmod m$$

First we have to figure out how we guarantee a permutation?

If $h_2(k), m$ relatively prime \Rightarrow Permutation

$m = 2^r$, $h_2(k)$ for all k is odd

Uniform Hashing Assumption X Simple Uniform Hashing

Each key is equally likely to have any one of the $m!$ permutations as its probe sequence.

[~~No one~~ found an exact hash function that guarantees exact probability, but very close to it]

$d = \frac{n}{m}$ cost of operations such as search, etc

$\leq \frac{1}{1-d}$ We have to ensure that d is no more than 0.5

number of empty slots

$\frac{m-n}{m}$

$1-d \rightarrow$
in term of prob

$\frac{1}{1-d}$

It means that in order to preserve the size of d (no more than 0.5) we have to implement a resize function

In another way, than in the chaining approach

Open addressing takes less memory \rightarrow because we don't use pointers!

Cryptographic hashes

[Bonus]

(6)

Password Storage

One-way

Open

$h(x)$

It is very hard
to find x

$$\text{s.t } h(x) = a$$

/etc/passwd : devadas : x12uG7

↓
Login name

↓ :

↓ :

If we want to make sure that even the sys-admin
doesn't know our password,

So in the file that have the password there is
the login name and besides to it the hashed
value of the password

So when I log in and type the password the
system calculate through the hash function
that the password is right.

$$h(x') \approx h(x)$$

?

typed in password

problems collisions...