

Hashing

- Direct Address Table =
(DAT)

Searching time of some DS -

- (i) unsorted array — $O(n)$
- (ii) sorted array — $O(\log n)$
- (iii) linked list — $O(n)$.
- (iv) Binary tree(BT) — $O(n)$
- (v) Binary search tree — $O(n)$
- (vi) Balanced BST — $O(\log n)$.
- (vii) priority queue(min & max heap) — $O(n)$.

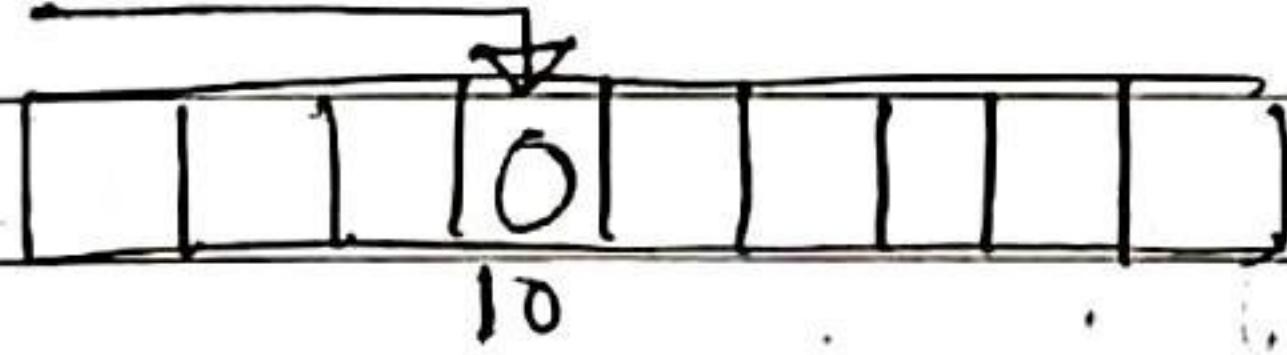
→ If we used hashing technique then on average time taken to search an element is — $O(1)$.

Hashing is used to minimize the search time.

Before using hashing technique people used one more technique called as DAT (direct address table).

DAT: It is similar to Array.

Key — $\{1, 2, 3, \dots, 100\}$



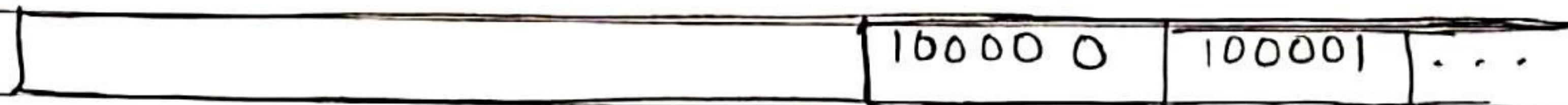
→ The first DS which has supported the searching in — $O(1)$ is called DAT (direct address table). In this table we also have an array and we are going to place an element according to key value.

→ This method is useful only when you know the no. of keys which are going to be in small range and

If the

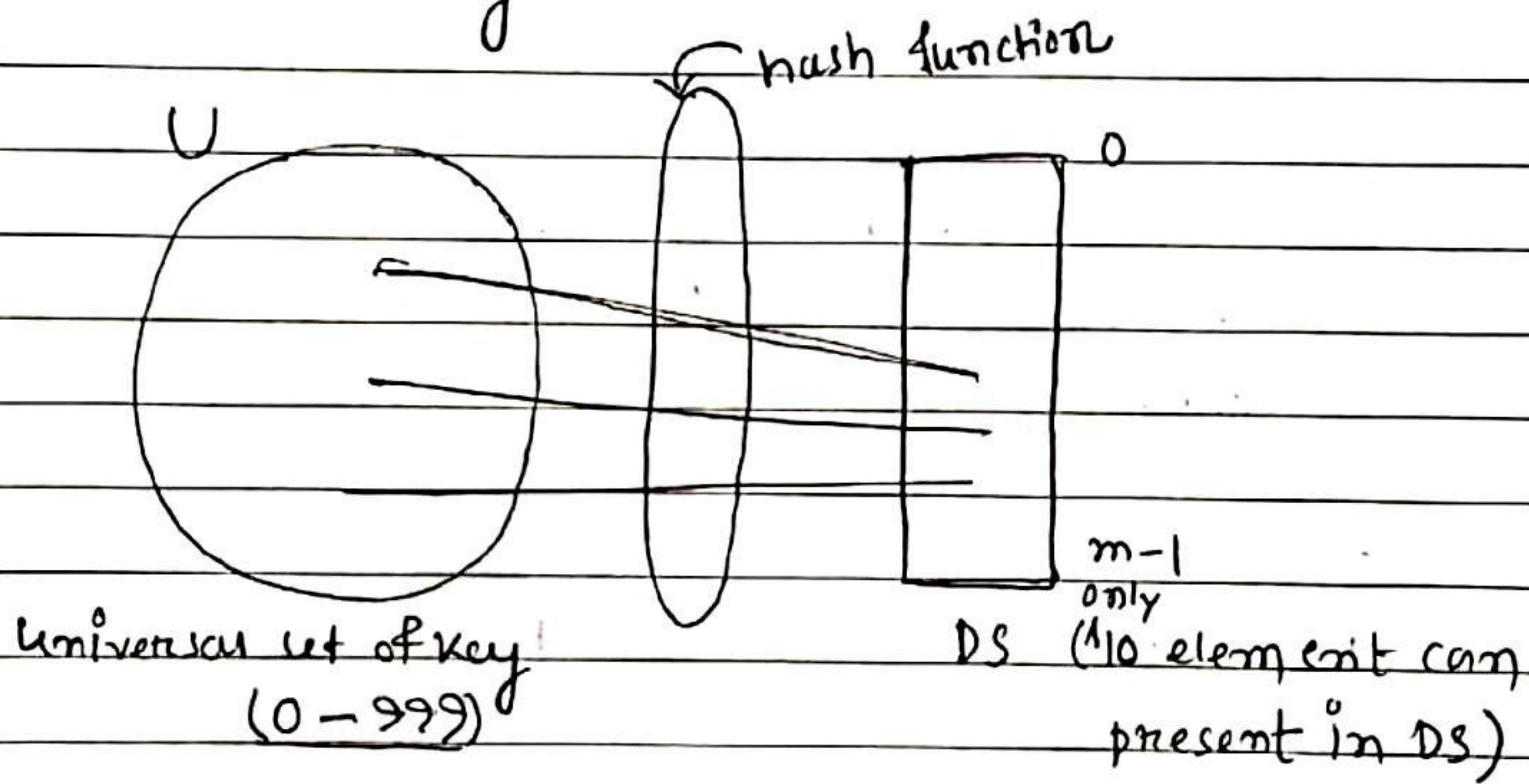
In case keys used are very large numbers, but the total no. of keys are less in number.

Key - {100000, 100001, 100002, ...}



so, that the array size will be 1M, just to store 10 values. In such case DAT fail. (Waste of space)
→ solve this problem we use hashing.

- Introduction to hashing —



hash function :

$$h(U) \rightarrow (0-(m-1))$$

$$h(0-999) \rightarrow (0-9)$$

hash function is any function which will take a key from U and it will map it to one of the values from 0 to m-1.

example:

(0-999)

set of Key - (121, 145, 132, 999)

hash function - (mod 10)

0	1	2	3	4	5	6	7	8	9
NULL	121	132	NULL	N	145	N	N	N	999

(size-10)

When searching for an element - 150, 145

$\begin{array}{|c|} \hline 150 \% 10 \\ \hline = 0 \end{array}$

element is not present.

$\begin{array}{|c|} \hline 145 \% 10 \\ \hline = 5 \end{array}$

element is present.

→ Insertion, ^{deletion and} searching is going to take constant time by using the hashing technique.

→ In hashing technique, there is a problem called Collision.

→ Collision problem =

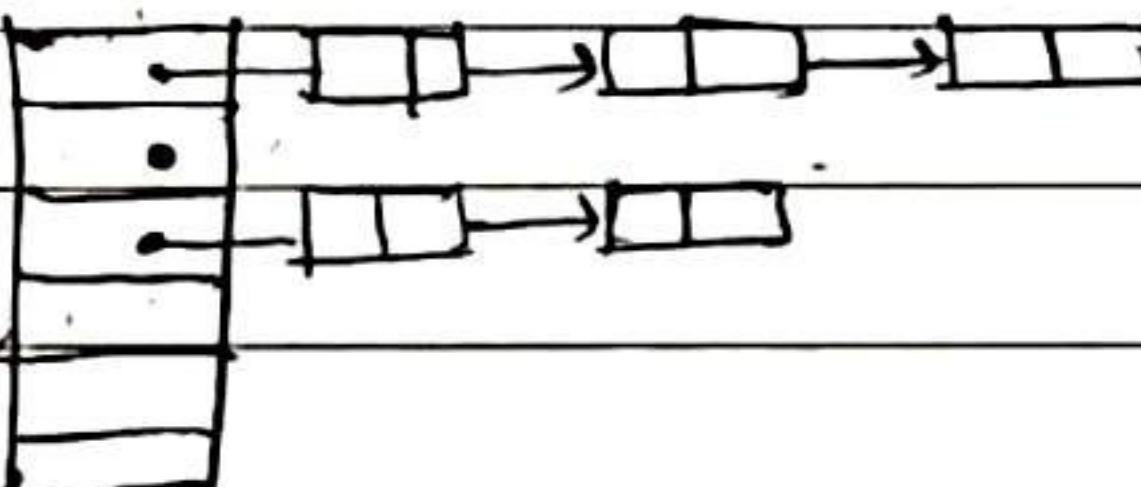
|→ When two key are two element being hash into same cell that condition is called collision.

→ Solving collision problem is -

(i) Better hash function.

(ii) Chaining.

(T, D, S)



(For deletion
chaining is
better)

(iii) Open addressing. (T, S)

(probing - searching for an empty space to insert an element)

(i) Linear probing

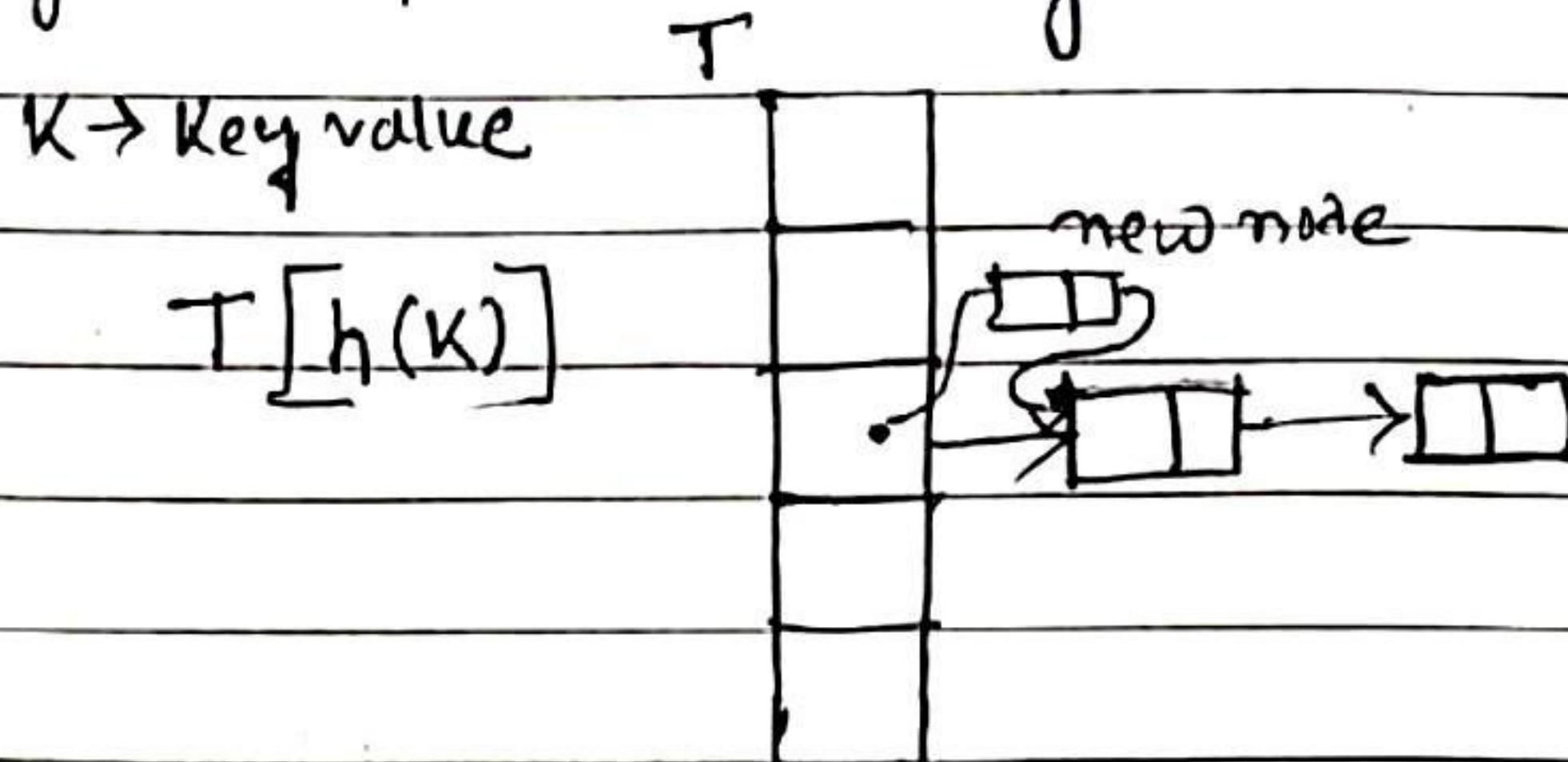


(For deletion open
addressing is not better)

(ii) Quad probing

(iii) Double hashing

- Collision solving technique : [chaining].



~~Time taken~~ hash table
 → Insertion In case of chaining time taken = $O(1)$.

→ In worst case time for searching = $O(n)$.

(n element inserted in the table, all the element might map ∞ (in worst case) to a single entry (cell) and the list of entry might contain as many as n)

→ In worst case Deletion Time = $O(n)$.

$$\alpha = \frac{n}{m}$$

$\alpha \rightarrow \text{load factor.}$

$n \rightarrow \text{total element present in table.}$

$m \rightarrow \text{total size of table}$

→ In this case average search time = $\Theta(1 + \frac{n}{m})$

$$n = K m$$

$$-\Theta(1 + \alpha)$$

$$\frac{n}{m} = K(\alpha)$$

$$-\Theta(1)$$

average deletion time = $\Theta(1)$.

↓ Insertion ↓

~~Adv: Adv:~~

→ deletion is easy in chaining.

disadv:

→ pointers (space wasted).

1996
g - 2014

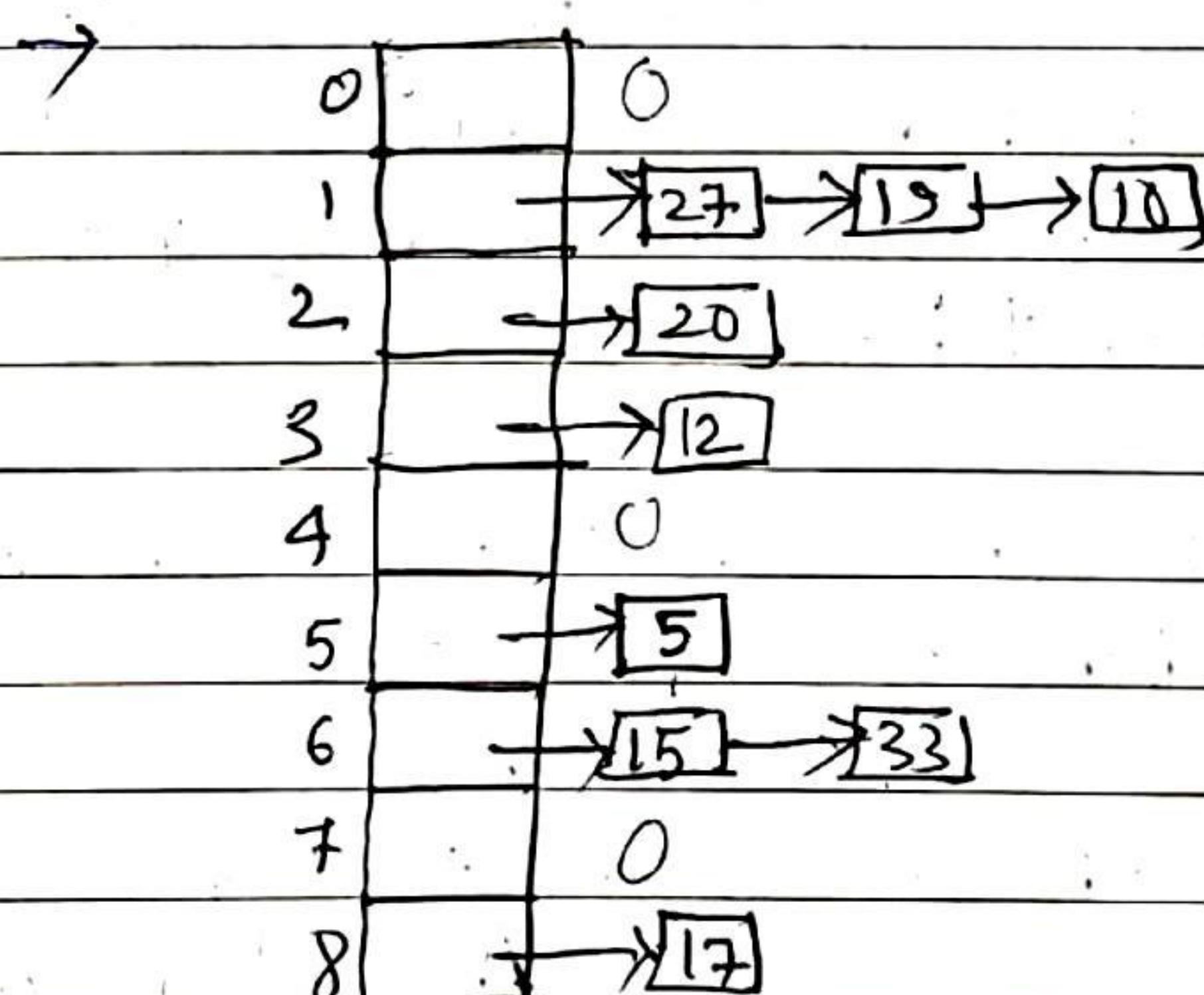
[Question] - ①

An advantage of chain hash table over open addressing scheme is -

- worst case complexity of search is less.
- space used is less.
- deletion is easier.
- none of the above.

g - 2014
[Question] - ②

$h(K) = K \bmod 9$, hash table has 9 slots, chaining is used Keys: 5, 23, 19, 15, 20, 33, 12, 17, and 10. Then, max and min and average chain length in hash table.



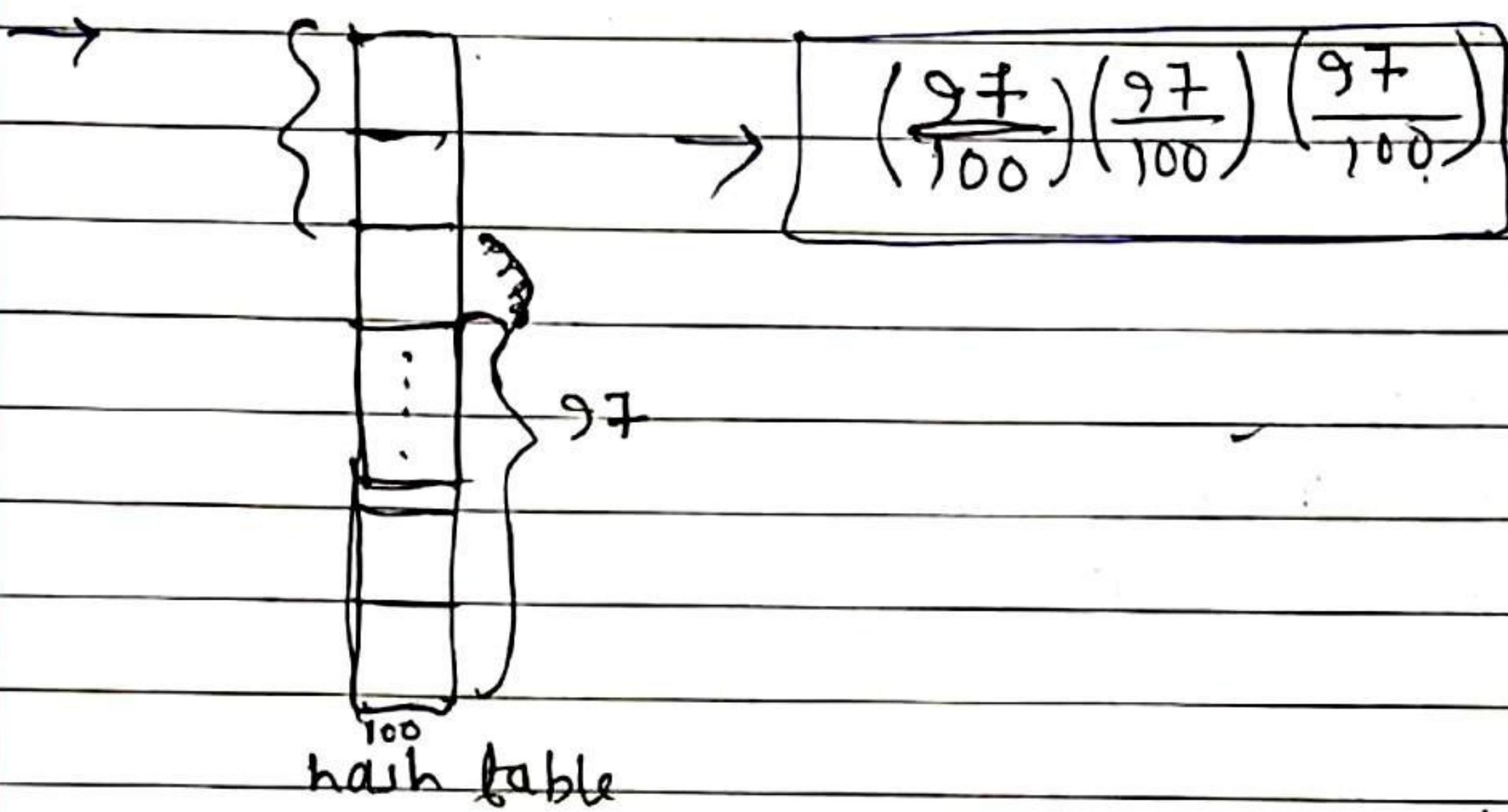
✓ max chain length → 3

✓ min " " → 1

✓ avg " " → $\frac{0+3+1+1+0+1+2+0+1}{9} \rightarrow 1$.

gate-
2019**Question - ③**

Consider a hash table with 100 slots. Collisions are resolved using chaining. Assuming simple uniform hashing, what is the probability that the first 3 slots are unfilled after the first 3 insertions?

gate-
1997**Question - ④**

Consider a hash table with 'n' buckets buckets, where chaining is used to resolve collisions. The hash function is such that the probability that a key value is hashed to a particular bucket is $\frac{1}{n}$. The hash table is initially empty and 'K' distinct values are inserted in the table.

a) What is the probability that bucket no. 1 is empty after K insertions.

b) What is the probability that no collision has occurred in any of the 'K' insertions?

c) What is the prob that first collision occurs at the Kth insertion.

$$\textcircled{a} \rightarrow Y_n$$

$$(1 - Y_n)$$

$$= \left(\frac{n-1}{n}\right)$$

$$\left(\frac{n-1}{n}\right) \left(\frac{n-1}{n}\right) \cdots \left(\frac{n-1}{n}\right) = \left(\frac{n-1}{n}\right)^K.$$

$$\textcircled{b} \quad \left(\frac{n}{n}\right) \left(\frac{n-1}{n}\right) \left(\frac{n-2}{n}\right) \cdots \left(\frac{n-(k-1)}{n}\right)$$

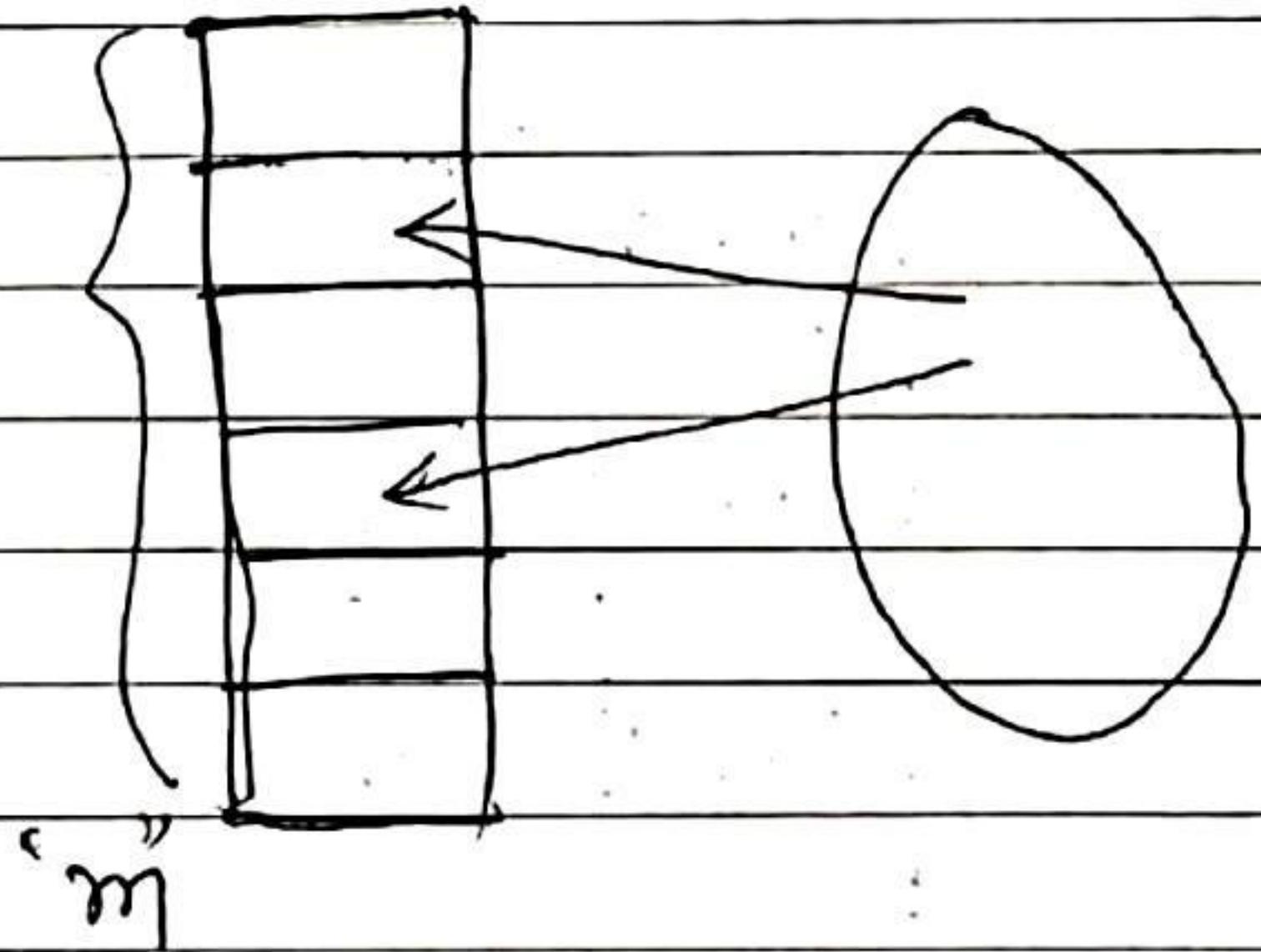
1st ele 2nd ele 3rd ele ... kth element.

$$\textcircled{c} \quad \frac{(k-1)}{n}$$

* Collision Solving Technique: [Open addressing]
or (closed hashing)

$$\text{load factor, } \alpha = \frac{n}{m}$$

$$0 \leq \alpha \leq 1$$



$$h: U \rightarrow \{0, 1, \dots, m-1\}$$

$$h: (U \times \{0, 1, 2, \dots, m-1\}) \rightarrow \{0, 1, \dots, m-1\}$$

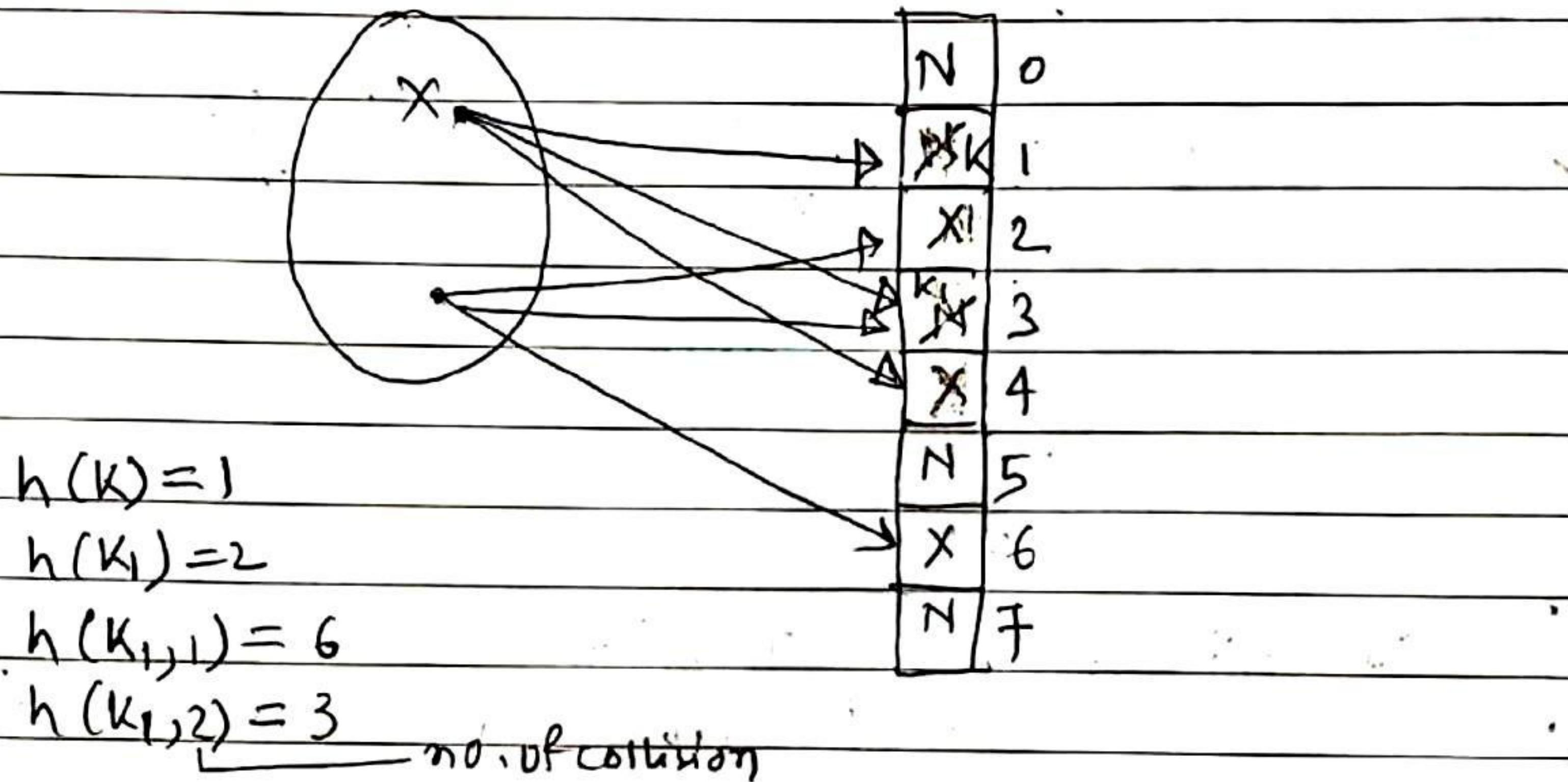
probe Sequence:

try to

Whenever I insert a key, then follow a sequence of examinations, In case if you find any empty cell in this sequence you are going to insert the element there, if you don't find any sequence empty cell in that sequence, then declared

that entire ^{table} is full , you are not able able to insert inserted.

worst time take to searching - $O(m)$



Searching for K_1

$$h(K_1, 0) = 2$$

$$h(K_1, 1) = 6$$

$$h(K_1, 2) = 3$$

$$h(K_1, 3) = 4$$

$$h(K_1, 4) = 5$$

$$h(K_1, 5) = 1$$

$$h(K_1, 6) = 7$$

$$h(K_1, 7) = 0$$

→ In case of full hash table and you are going to unsuccessful search the time complexity for a search is - $O(m)$.

* In Average ~~case~~ time complexity - $O(1)$.

Insert an new element K_1 ,

$$h(K_1, 0) = 2 \checkmark$$

$$h(K_1, 1) = 6 \checkmark$$

$$\cancel{h(K_1, 2) = 3} \checkmark$$

$$h(K_1, 3) = 4 \checkmark$$

$$h(K_1, 4) = 5 \checkmark$$

$$h(K_1, 5) = 1$$

$$h(K_1, 6) = 7$$

$$h(K_1, 7) = 0$$

0	N
1	N
2	a
3	c
4	d
5	K ₁ a N
6	b
7	N

→ We are going to stop & either when we find the searching element or empty cell.

Then Delete an element 'c' from table.

after deleting c search for

K_1 .

0	N
1	N
2	a
3	N
4	d
5	K ₁
6	b
7	N

→ Stop searching here.

and declare
that K_1 is not
present even
when K_1 is
present.

→ In case of deletion open addressing
Create such a problem.

To solve this problem =

We are going to mark deleted place in such a way that it will indicate that some element are earlier there but now it is deleted.

0	N
1	N
2	N
3	D
4	d
5	K ₁
6	b
7	N

D → indicate that, earlier there may be an element but now it deleted.

- Various technique which used in Open addressing =

(i) Linear probing :

hash function =

$$h: U \rightarrow \{0, \dots, m-1\}$$

$$h(K) = a$$

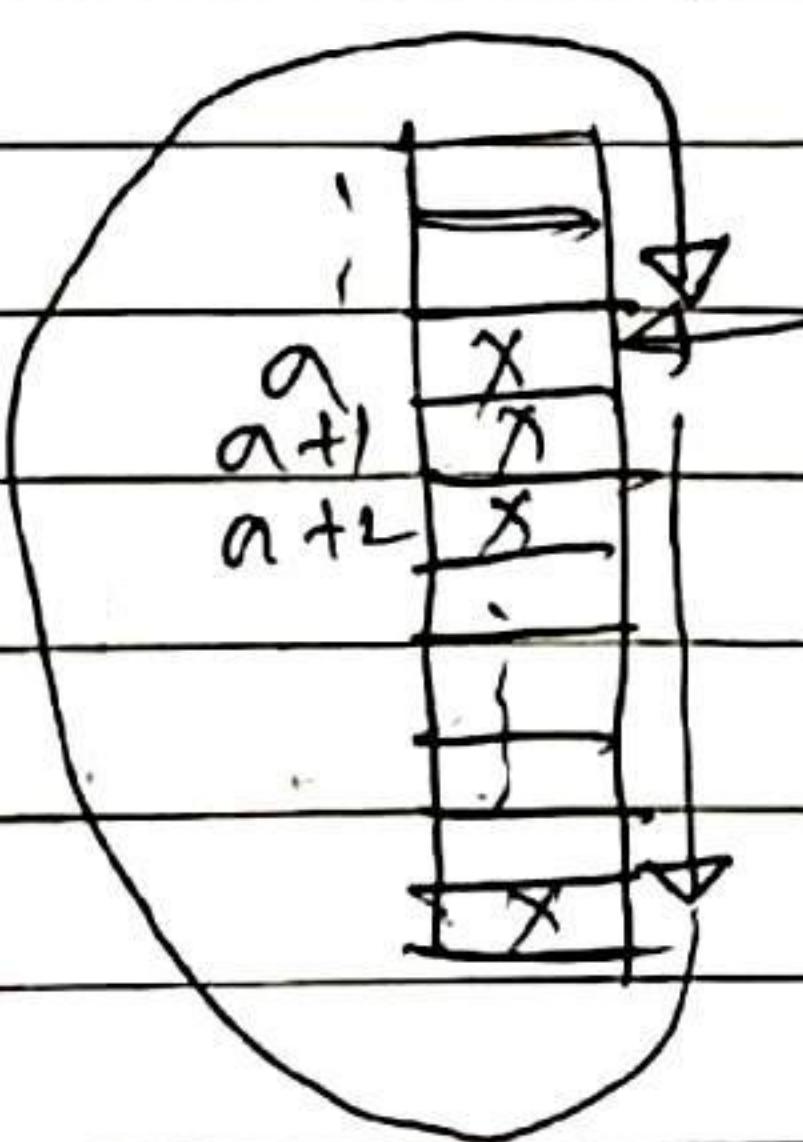
$$h'(K, i) = (h(K) + i) \bmod m$$

$$\begin{aligned} h'(K, 1) &= (h(K) + 1) \bmod m \\ &= (a+1) \bmod m. \end{aligned}$$

0	
1	X
2	X
3	X
a+1	X
a+2	X
a+3	K

$$\begin{aligned} h'(K, 2) &= (h(K) + 2) \bmod m \\ &= (a+2) \bmod m. \end{aligned}$$

$$h'(K, 3) = (h(K) + 3) \bmod m = (a+3) \bmod m$$



Probe sequence →

(0, 1, 2, 3, ..., m-1) always.

(1, 2, 3, ..., m-1, 0) search for

(2, 3, 4, ..., m-1, 0, 1) "m" element

(3, 4, 5, ..., m-1, 0, 1, 2) to insert

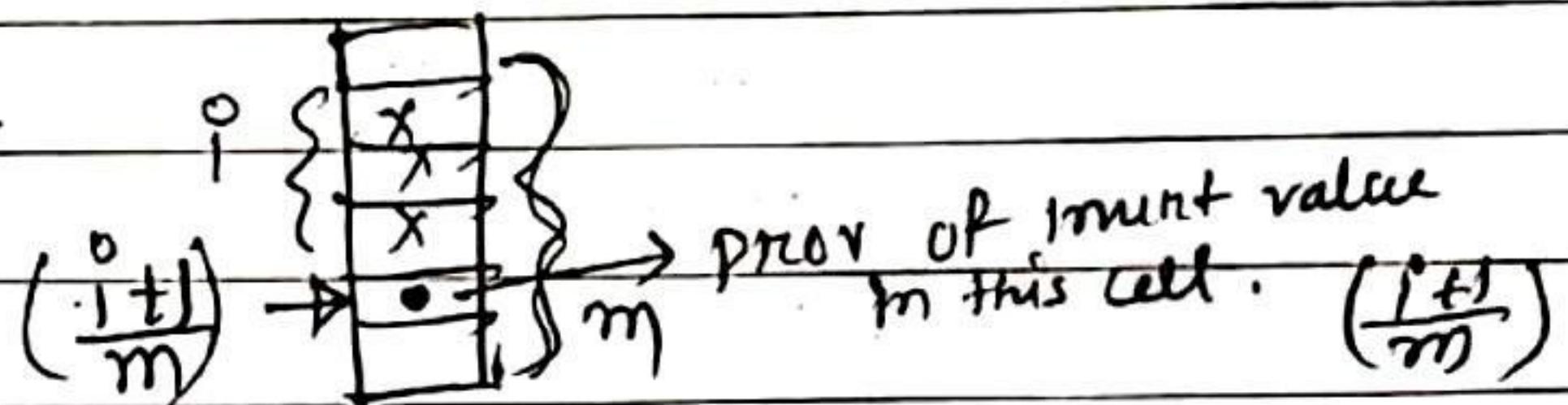
element.

so, probe sequence is possible is "m".

Secondary Clustering =

In case or two elements first probe is found in same value then both of them follow going to follow same probe sequence, this call secondary clustering.

Primary clustering =



~~disadv~~ ^{avg} Linear probing suffer two things (secondary and primary clustering)

Average search time taken to search an element = $O(2.5)$.
Worst case = $O(n)$.

gate
2008

[Question] - ①

Key: 12, 18, 13, 2, 3, 23, 5 and 15.

HST = 10 (hash table size)

$h(K) = K \bmod 10$ and linear probing. What is resultant hash table?

→

0		\leftarrow	get into this $\left(\frac{9}{10} \text{ (prev)}\right)$
1			
2	12	↑	$h(k,i) h(12,0) = (h(k) + i) \bmod 10$
3	13	↑	$h(13,1) = (h(12) + 1) \bmod 10$
4	2		$= (2 + 1) \bmod 10$
5	3		$= 3$
6	23		$h(23,2) = (2 + 2) \bmod 10 = 4$
7	5		
8	8		
9	15		

Primary cluster.

gate
2008Question - ② $HTS = 11 \ (0-10)$ using $h(k) = k \bmod 11$, linear probing method.

Keys : 43, 36, 92, 87, 11, 4, 71, 13, 14.

What is the index into the last record is inserted?

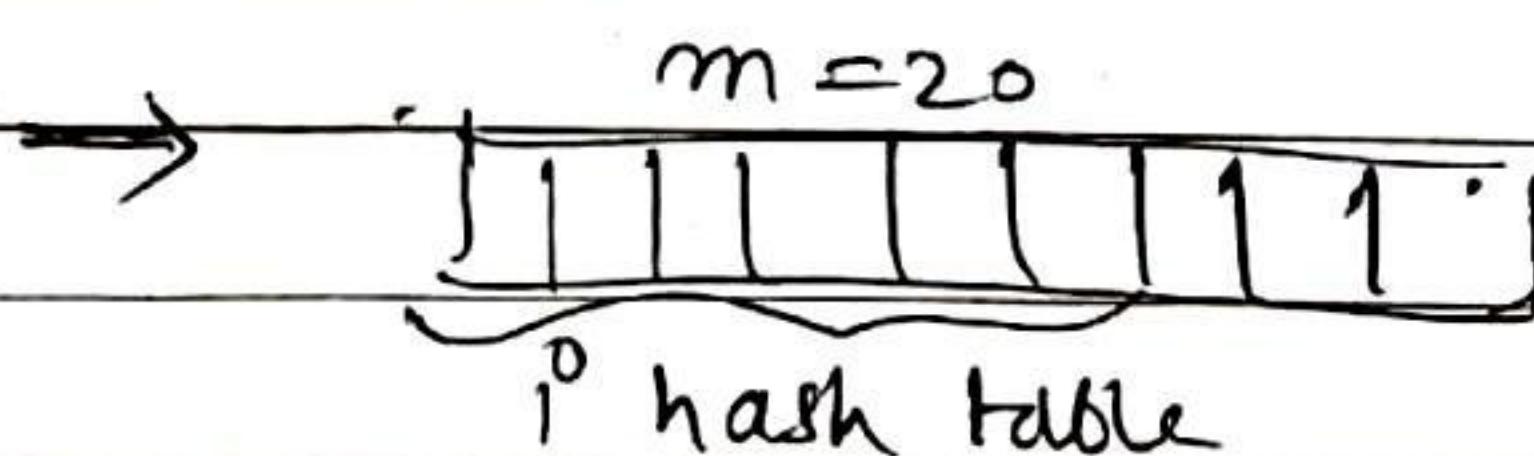
→ in i^{th} index last record is inserted.

0	87	← 11
1	11	
2	13	
3	36	← 14
4	92	← 9
5	4	← 71
6	71	
7	14	
8		← probability that this cell get next insertion
9		$\frac{8}{10} \left(\frac{9}{11} \right)$
10	43	← 87

gate
2007Question - ③

consider a hash function that distributes keys uniformly. The hash table size is 20. After hashing of how many key, will the probability that any new key hashed collides with an existing one exceed 0.5.

- a) 5 b) 6 c) 7 d) 10.



$$\frac{i}{m} > 0.5$$

$$\frac{i}{20} > 0.5 \Rightarrow i > 10$$

gaur
2010Question - ④ $h(k) = k \bmod 10$, linear probing.

0	
1	
2	42
3	23
4	34
5	52
6	46
7	33
8	
9	

a)

- A) After insertion, which will be look like give sequences -
- (x) a) 46, 42, 34, 52, 23, 33
 - (x) b) 34, 42, 23, 52, 33, 46
 - (✓) c) 46, 34, 42, 23, 52, 33
 - (x) d) 42, 46, 33, 23, 31, 52.

- B) How many insertion sequences are possible, to get the give table sequence -

$$\rightarrow \underbrace{(42-23-34)}_{3!} - 52 - 33$$

46 can place take place any of the blank space.

$$3! \times 5 = 30 \text{ (insertion) sequence possible.}$$

(ii) Quadratic probing =

$$h'(k, i) = (h(k) + c_1 i + c_2 i^2) \bmod m.$$

Ex:

$m = 10$ (size of hash Table).

(0 --- 9)

$$c_1 = 1, c_2 = 1$$

$$h(k) = k \bmod 10.$$

$$h(k_1) = 1$$

$$\begin{aligned} h'(k_1, 1) &= 1 + 1 \cdot 1 + 1 \cdot 1^2 \\ &= 3 \end{aligned}$$

$$\begin{aligned} h(k_1, 2) &= 1 + 2 + 4 \\ &= 7 \end{aligned}$$

$$\begin{aligned} h(k_1, 3) &= 1 + 3 + 9 \\ &= (13) \bmod 10 = 3 \end{aligned}$$

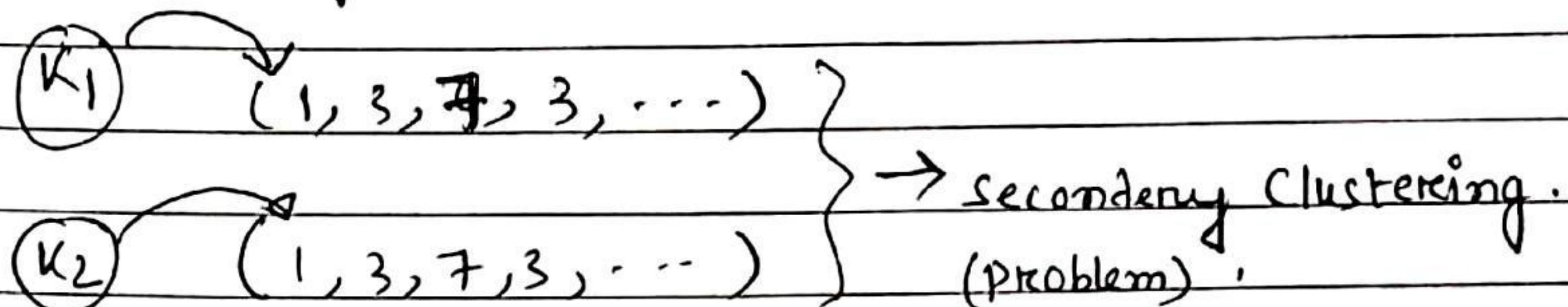
$$h(k_1, 4) = 1 + 4 + 16 = 21 \bmod 10 \quad "m"$$

after 10 probes = 1

In this case, you are not able to examine all the location by using 10 probes of hash table. This is the problem with quadratic probing.

You should ~~carefully~~ choose c_1, c_2 and m value ~~very well~~ carefully. Otherwise it ~~is~~ going to lead problems.

probing sequence -



→ If two keys happen to have the same initial probe location then they are going to have similar probe sequence.

two key's

→ A probe sequence will same when Initial probe is same.

different

→ The no. of probe sequence possible is - "m" -

| every probe sequence depends on initial probe sequence
↳ "m" probe sequences possible.

Best probing technique. (iii) Double hashing

In double hashing no. of probe sequence is possible is m^2
(where "m" → size of hash table) $= O(m^2)$

$$h'(k) = (h_1(k) + i h_2(k)) \bmod m.$$

→ there are no secondary and primary clustering problem.

$\leq m$ prime.

$$(i) \left(h_2(k), m \right)$$

\downarrow \downarrow

Odd 2^k

→ hashing is more better compare to other DS. When the most popular operation is ^{only} searching and data is static).