

CHAPTER 6 - HASHING

DATE -

LEC 1: INTRODUCTION

* Main usage of Data structure is to store our data.

* Operations performed on data structures?

- Insertion

- Searching (most dominant operation) — Needs to be faster

- Deletion

Time Complexity (to search)

* Array = $O(n)$

* Sorted Array = $O(\log n)$

* Linked List = $O(n)$

* Binary Tree = $O(n)$

* Binary Search Tree = $O(\log n)$

* Balanced Binary Search Tree = $O(\log n)$

* Priority queues (min heap & max heap) = $O(n)$

Note:

Hashing → solution to reduce search time

↳ Time Complexity = $O(1)$ i.e. constant time

Direct Address Table → (II) to arrays, was used before Hashing ↳ 1st data structure with $O(1)$ time complexity for searching

100 records have to be stored, each should have unique tag & in range from 1 to 100³

define array of 100 size, store / insert elements at appropriate index
1, 2, 3, ..., 100



↳ Can be used only if Number of keys used are proportional to size of array.

* If number of keys to be inserted are very less but size of each key is very large, then this method is not useful e.g.: - \$1000, 1000000, 1000000000

even to insert these array index has to be this large

LEC 2: HASHING INTRO

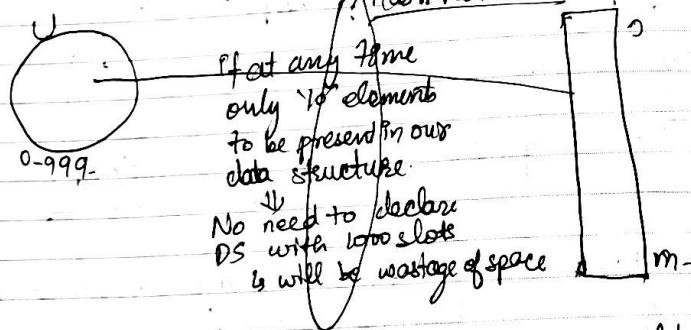
DATE

Problem with DAT:

Problem with DAT:
The size of the array / table should be as big as the
key sizes that are to be stored.
↳ Sometimes we end up declaring size of array as too
large even if number of elements to be stored are
very less.

HASHING means taking a set of very large numbers and mapping it to small numbers which correspond to indices of hash table.

↳ we declare a data structure - large enough to hold elements which might be present at any point in time.



Then we try to map elements into the available slots.

- **Hash Function:** Mapping function that is used to map the 'key' (to be stored from Universal set) to the location at which it will be stored in available Data structures.

→ Also applied at the time of searching & check its availability at the said location in data structure.

eg:- Elements to be inserted are : 121,145, 132,999 (0-999)
we search for 121,145, 132,999 in array N. 121 132 N N 145 N N 999

Flash Function $(0-999) \rightarrow (0-9)$ $\Rightarrow f(v) = v \bmod 10$
 \hookrightarrow which will take values from $(0-999)$, map them to values $(0-9)$

$$h(\text{val}) = \text{mod } 10$$

PROBLEM WITH HASHING : COLLISION

DATE

hash table when more than one elements map to same location
in ? On applying hash function

121	132	145	N	99
0	1	2	3	4

In prey example if we try to store 131 then it will lead to collision

SOLUTION FOR COLLISION RESOLUTION TECHNIQUE

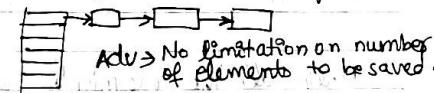
(I) Make Hashing function better.

* We are considering only least significant bit in previous example.

→ Problem:

We can have better hashing function which will decrease the number of collisions but they can't completely vanish.

(II) CHAINING



→ says In Hash-table instead of storing elements directly at its indices, maintain a linked list corresponding to each index in which elements will be stored.

- * More than one elements can be accommodated at same index in linked list.
- Space wasted on pointer. \hookrightarrow Used \Rightarrow when insertion & deletion ^{searching} both involved.

(III) OPEN ADDRESSING OR Probing
Searching for empty space in which we want to store element

→ If collision occurs, recompute the hash value with some minor modification to hash value - might lead to vacant space (perform recomputation till we get some vacant space)

- a) Linear Probing
 b) Quad " "
 c) Double "

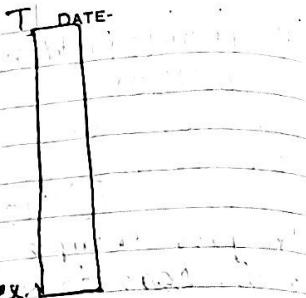
- Used - when no deletion involved
- Disadvantage - No Number of elements to be saved are limited by size of array.

Lec 3: CHAINING

If we want to insert an element or key 'K' then we have to go to location

$T[h(K)]$

Hash function - give index



* Instead of placing element at the index given by hash function.

* Rather element is inserted into the list pointed by the table's location ($T[h(K)]$)

New node is created and inserted at the beginning of the list.

Insertion Time Complexity = $O(1)$ {^{computing Hash func & inserting at beginning of linked list takes constant time}}

Worst case scenarios of Insertion & storing elements? (Search)

When all the elements hash function returns same index
All the element gets stored in linked list present at one index only.

Searching Time Complexity = $O(n)$ {^{of this worst case}}

Worst Case for deletion:

search for the element and then delete it.

Deletion Time Complexity = $O(n)$

Even these worst case complexities of $O(n)$ can be avoided by using Uniform Hash Function.

Uniform Hash Function (Uniform distribution for hashing)
↳ Which uniformly distributes the elements in the array
↳ Does not map all elements to same location

If Uniform Hash Function used: (Avg case for searching)
each location of array will get $\rightarrow \frac{n}{m}$ elements

In every to compute hash function.
 $\Theta\left(1 + \frac{\alpha}{m}\right)$

Average Search Time Complexity $\rightarrow \Theta\left(1 + \frac{\alpha}{m}\right)$

m - size of hash table

Load factor

$$\alpha = \frac{n}{m} \quad \begin{matrix} \text{Total no. of elements} \\ \text{Total size of hash table} \end{matrix}$$

$$\text{If } [n = Km] \Rightarrow \Theta\left(1 + \frac{\alpha}{m}\right) = \Theta(1) \rightarrow \text{Avg case.}$$

↳ Avg Selection $\rightarrow \Theta(1)$

Advantage of chaining : \rightarrow In Avg case \rightarrow most of the operations take constant time.

Deletion is easy \rightarrow In order to delete an element, no need to affect the remaining elements.

Disadvantage :-

↳ Pointers (space is wasted)

LEC 4 : EXAMPLE

DATE -

- Note 96) An advantage of chained hash table over open addressing scheme is worst case complexity of search is less.
- worst case complexity of search is less
 - space used is less
 - deletion is easier
 - None of the above

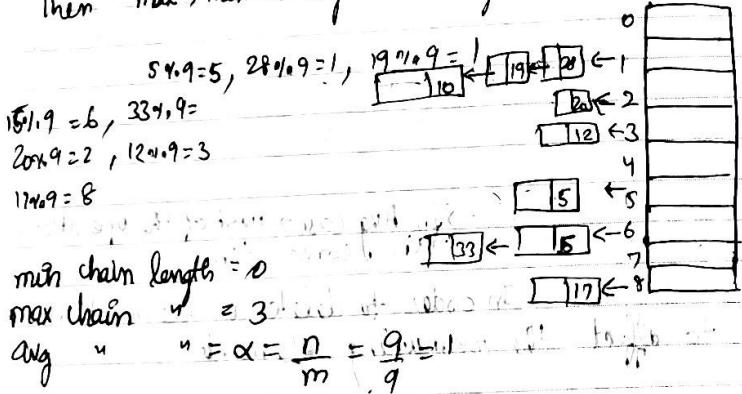
Note 97) Deletion is easier if we use chained hash table.

Note 204

$h(K) = K \bmod 9$, hash table has 9 slots, chaining is used.

Keys: 5, 28, 19, 15, 20, 33, 12, 17 and 10.

Then max, min & avg chain lengths in hash table.



Note 204

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 1st 3 insertions?

Uniform distribution of elements using hashing.

Probability that particular key 'K' is mapped to one of the locations in array = $\frac{1}{n}$ {For every slot}

not used real probability = $\left(97 \times \frac{1}{100}\right) \left(97 \times \frac{1}{100}\right) \left(97 \times \frac{1}{100}\right)$
97 slots can be used.

Note 97 Consider a hash table with 'n' 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.

- What is the probability that bucket number '1' is empty after 'K' insertion.
- What is the probability that no collision has occurred in any of the 'K' insertions?
- What is the probability that 1st collision occurs at the K^{th} insertion?

Size of hash table = n

Probability that a key value is hashed to particular bucket = $\frac{1}{n}$

a). Prob that bucket number 1 is empty after 'K' insertion = $\left(1 - \frac{1}{n}\right)^K = \frac{(n-1)}{n}^K$

b). Prob that no collision has occurred in any of the 'K' insertions = $\left(\frac{n}{n}\right) \left(\frac{n-1}{n}\right) \left(\frac{n-2}{n}\right) \dots \left(\frac{n-(k-1)}{n}\right)$

c). Prob that 1st collision occurs at K^{th} insertion

$$= \left(\frac{n}{n}\right) \left(\frac{n-1}{n}\right) \left(\frac{n-2}{n}\right) \dots \left(\frac{n-(K-2)}{n}\right) \left(\frac{K-1}{n}\right)$$

$\underbrace{\quad}_{(K-1)^{\text{th}} \text{ element}}$

K^{th} element has to take a slot taken by $(K-1)$ elements.

LECTURE 7: OPEN ADDRESSING

DATE _____

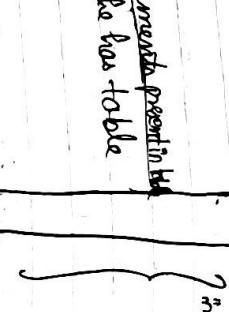
OR CLOSED HASHING

- * All the elements are stored in slots available in the array. Not in linked list outside array.

↳ Consequences - The number of elements stored are restricted by size of array.

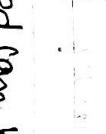
Load factor (α) = $\frac{\text{Total no. of elements present in the hash table}}{\text{Total size of the hash table}}$

$$\text{max load factor} = \frac{m}{m} = 1$$



$$\min \alpha = 0$$

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



- * Number of elements or keys to be stored can be much more than the number of slots available in array.

Collision possible

- ↳ In this case also elements are to be stored in same hash table
- ↳ Collision resolved by again applying a hash function with minor modification

Called 2nd time probing (Finding a location in hash table at which element is to be stored)

$f(K_1) = 1$ If $f(K)$ already filled.
 probing function: $(f(K_1) + 1, 2, \dots, m-1)$ If $f(K)$ already filled.
 (g taking already filled locations into consideration.)

same.

Using key & an second set increase the number by 1 if 1st collision then if 2nd time collision occurs \dots

Hash function is not only a function of key but also a

or probing number

eg of LMS in the function \rightarrow every location of hash table is inspected atleast once.

↳ Probing sequence - The sequence of examination/ inspection of locations of hashtable done in order to insert an element into the hash table & known as probing sequence.

Worst time Complexity for searching :- $O(m)$

↳ In this case, if element not present at some slot then we assume that it can be at some other slot.

We can be lucky if either we find the element on the way or we find a NULL. \rightarrow Searching ends there (and stop on seeing d ('space created' after deletion))

$$\begin{aligned} f(K_1) &= 1 \\ f(K_1) &= 2 \quad (\text{already occupied}) \\ f(K_1, 1) &= 6 \quad ("") \\ f(K_1, 2) &= 3 \quad (\text{null}) \end{aligned}$$

N.	1	X	2	X	3	X	4	X	5	N	6	X
K												

Searching also done in similar sequence.

For searching we will be computing hash func" for K_1 in following sequence.

$f(K_1, 3) = 2$ Before declaring that element can't be inserted or not present in hash table
 $f(K_1, 1) = 6$ (in case of searching) each index of array is examined except since
 $f(K_1, 2) = 3$ Note: Selection affects searching, but doesn't affect insertion.
 $f(K_1, 3) = 4$
 $f(K_1, 4) = 5$
 $f(K_1, 5) = 1$ whereas in case of searching, 'd' considered as occurred
 $f(K_1, 6) = 7$
 $f(K_1, 7) = 0$

In case of insertion, 'd' considered as NULL like d

like d

like d

like d

Worst case scenario after deletion:

↳ If all the cells are free after deletion or are NULL for only 1 element present in hash table, search will take a lot of time.

If i_k is present at last probe, then irrespective of empty hash table we have to check entire table.

Note:- If Deletion is involved, Chaining is much better than Open Addressing.

Adv over Chaining:-

↳ Space not wasted for pointers.

LEC 8: METHODS OF OPEN ADDRESSING

LINEAR PROBING : probing location in linear manner

* We have to just modify hash function to make it work as linear probing.

* We have following hash function:-

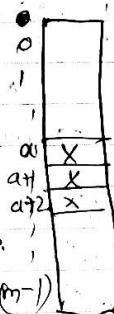
$$h: U \rightarrow \{0, 1, 2, \dots, m-1\} \quad h(K) = a$$

hash function is going to take a key again and apply hash function to get location at which we are in hash table.

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

Why linear?

From 1st location, all the locations - linearly increasing are probed till the end then from 1st location till the 1st probed location.



Hash function for linear probing

$$l_1(K) = a$$

$$\begin{aligned} l_1(K, 1) &= (h(K) + 1) \bmod m \\ &= (a+1) \rightarrow \text{filled} \end{aligned}$$

$$\begin{aligned} l_1(K, 2) &= (h(K) + 2) \bmod m \\ &= a+2 \end{aligned}$$

Probing sequence :-

The sequence of inspecting/examining locations of hash table before finally storing that element.

e.g.: If sequence is:] If 'm' slots in 1 table, every slot $(0, 1, 2, \dots, m-1)$ has to be probed exactly once.

Number of probe required = m

$l_1(K, i) \rightarrow$ ensure probing of all locations

can be :

$$(1, 2, \dots, m-1, 0)$$

$$(2, 3, \dots, m-1, 0, 1)$$

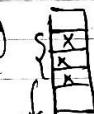
Number of probing sequences = $m^{\frac{m}{2}}$

Note:- If two numbers have initial probe number to be same then they will have same probe sequence. } called Secondary clustering.

All 'm' probes will generate different slots chosen

* Hash function for linear probing should examine each slot of the hash table.

* Problem with Linear Probing: Secondary Clustering Primary Clustering



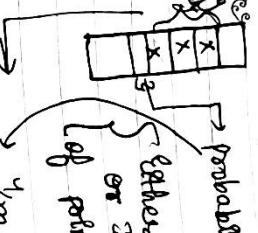
* In case of linear probing, there will be continuous runs of filled slots

* Means many items will fall together in continuous runs; then a gap, again continuous filled slots. \rightarrow increasing search time. } we have to search through entire cluster

$1/m \rightarrow$ probability for each slot in case of uniform distribution

DATE-

primary cluster
secondary cluster
→ probability that this empty slot will get an element
either element directly gets mapped to the location
or if that particular key maps on to the slot
of primary cluster.



$\frac{1}{m}$

Note #
If slots are in primary clustering then
probability that $(q+1)^{\text{th}}$ slot will get filled is
 $\left(\frac{i+1}{m}\right)$

Avg. Search Time $\geq O(1) \rightarrow$ practically as experimentally

* Summary of Linear probing

- Will have 2 hash function → Normal Hashing func & probing
- If slot is full, the sequence which we follows to insert an element is known as probing sequence.

- Probe sequence depends on initial probe number.

- 2 disadvantages → Primary Clustering

Why called linear probing?

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

Because if collision occurs at location $h(k)$ then the next location to be probed is a linear distance from initially probed location:

Date 2008

keys: 12, 18, 13, 2, 3, 23, 5, 15

DATE

1	6
2	12
3	13
4	2
5	3
6	23
7	5
8	18
9	15

$$h'(1, 0) = (h(1) + 0) \bmod 10$$

$$h(1) = 12 \cdot 1 \cdot 10 = 2$$

$$h(1, 1) = 18 \cdot 1 \cdot 10 = 8$$

$$h(1, 2) = 13 \cdot 1 \cdot 10 = 3$$

$$h(1, 3) = 2 \cdot 1 \cdot 10 = 3 \quad (\text{already occupied})$$

$$h(1, 4) = (2+1) \cdot 1 \cdot 10 = 3$$

$$h(1, 5) = (2+2) \cdot 1 \cdot 10 = 4$$

$$h(1, 6) = 3 \cdot 1 \cdot 10 = 3$$

$$h(1, 7) = 4 \cdot 1 \cdot 10 = 4$$

$$h(1, 8) = 5 \cdot 1 \cdot 10 = 5$$

$$h(2, 0) = 12 \cdot 1 \cdot 10 = 2$$

$$h(2, 1) = 18 \cdot 1 \cdot 10 = 8$$

$$h(2, 2) = 13 \cdot 1 \cdot 10 = 3$$

$$h(2, 3) = 2 \cdot 1 \cdot 10 = 2$$

$$h(2, 4) = (2+1) \cdot 1 \cdot 10 = 3$$

$$h(2, 5) = (2+2) \cdot 1 \cdot 10 = 4$$

$$h(2, 6) = 3 \cdot 1 \cdot 10 = 3$$

$$h(2, 7) = 4 \cdot 1 \cdot 10 = 4$$

$$h(2, 8) = 5 \cdot 1 \cdot 10 = 5$$

$$h(2, 9) = 6 \cdot 1 \cdot 10 = 6$$

$$h(2, 10) = 7 \cdot 1 \cdot 10 = 7$$

$$h(2, 11) = 8 \cdot 1 \cdot 10 = 8$$

$$h(2, 12) = 9 \cdot 1 \cdot 10 = 9$$

$$h(2, 13) = 10 \cdot 1 \cdot 10 = 10$$

$$h(2, 14) = 11 \cdot 1 \cdot 10 = 11$$

$$h(2, 15) = 12 \cdot 1 \cdot 10 = 12$$

$$h(2, 16) = 13 \cdot 1 \cdot 10 = 13$$

$$h(2, 17) = 14 \cdot 1 \cdot 10 = 14$$

$$h(2, 18) = 15 \cdot 1 \cdot 10 = 15$$

$$h(2, 19) = 16 \cdot 1 \cdot 10 = 16$$

$$h(2, 20) = 17 \cdot 1 \cdot 10 = 17$$

$$h(2, 21) = 18 \cdot 1 \cdot 10 = 18$$

$$h(2, 22) = 19 \cdot 1 \cdot 10 = 19$$

BC-9: EXAMPLES ON LINEAR PROBING

DATE

1	6
2	12
3	13
4	2
5	3
6	23
7	5
8	18
9	15

Ques 2010 $h(k) = k \mod 10$
Linear probing

(Date 11-08)
Hash table size = $6 - 10$
Hash probing used.

$$h(k) = k \mod 11, \text{ linear probing}$$

Keys: 43, 36, 92, 87, 11, 47, 13, 14
what is the index into which last record is inserted?

$$h(k, i) = h(k + i) \mod 11$$

$$\begin{cases} h(43, 1) = 43 \\ h(43, 2) = 43 + 11 = 10 \\ h(43, 3) = 43 + 22 = 35 \\ h(43, 4) = 43 + 33 = 76 \\ h(43, 5) = 43 + 44 = 87 \\ h(43, 6) = 43 + 55 = 98 \\ h(43, 7) = 43 + 66 = 109 \\ h(43, 8) = 43 + 77 = 120 \\ h(43, 9) = 43 + 88 = 131 \\ h(43, 10) = 43 + 99 = 142 \end{cases}$$



$$\text{Prob of collision} = \frac{1}{10}$$

Aftee

$$\text{Prob of getting next element} = \frac{1}{11}$$

$$\begin{cases} h(4, 1) = 5 \mod 11 = 5 \\ h(7, 1) = 71 \mod 11 = 5 \\ h(11, 1) = 13 \mod 11 = 2 \\ h(13, 1) = 14 \mod 11 = 3 \\ 15 \mod 11 = 4 \end{cases}$$

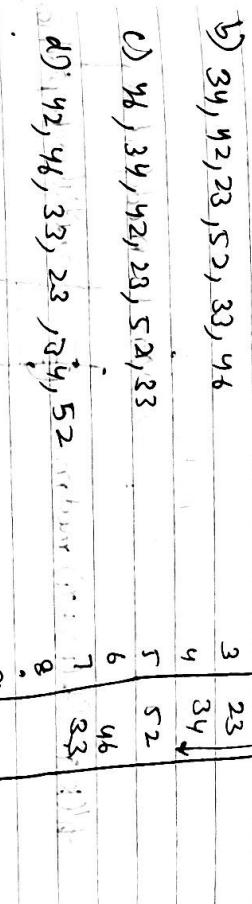
$$h(13, 2) = 18 \mod 11 = 7$$

last record

$$\begin{cases} h(14) = 14 \mod 11 = 3 \\ h(18) = 18 \mod 11 = 7 \end{cases}$$

$$\boxed{\text{last index filled} = 8}$$

- Ques 2010 $h(k) = k \mod 10$
Linear probing
- a) 46, 42, 34, 52, 23, 33
b) 34, 42, 23, 52, 33, 46
c) 46, 34, 42, 23, 52, 33
d) 42, 46, 33, 23, 34, 52



Ques 2010 How many insertion sequences are possible?

$$\begin{array}{l} \text{a) } 46-6, 42-2, 34-4, 52-8 \\ \text{b) } 34-4, 42-2, 23-3, 52-2-3-4-5-6-7 \\ \text{c) } 46-6, 34-4, 42-2, 23-3, 52-2-3-4-5-6-7 \\ \text{d) } 42-2, 46-6, 33-3 \end{array}$$

$$(42, 23, 34) \rightarrow 3! \text{ ways } \Rightarrow 52 \text{ has to come after these}$$

$$\begin{array}{r} 33 \rightarrow \text{after 52} \\ \downarrow \\ \cancel{(42, 23, 34)} \quad \underline{52} - \underline{33} \\ \cancel{33} \end{array}$$

5 places for 46

No. of insertion sequences = $3! \times 5$

= 30 way

(Date 15-07)

The Hash table size is 20. After hashing of how many keys will the probability that any newly inserted key with an existing one exceed at

Prob of collision

= $\frac{1}{20} = \frac{1}{20} = 0.05$

$$\begin{array}{llll} \text{a) } 5 & \text{b) } 6 & \text{c) } 7 & \text{d) } 10 \\ \text{Ans: } 11 & \text{Ans: } 12 & \text{Ans: } 13 & \text{Ans: } 14 \end{array}$$

\therefore already inserted

9

already inserted

10

already inserted

11

already inserted

12

already inserted

13

already inserted

14

already inserted

15

already inserted

16

already inserted

17

already inserted

18

already inserted

19

already inserted

20

already inserted

21

already inserted

22

already inserted

23

already inserted

24

already inserted

25

already inserted

26

already inserted

27

already inserted

28

already inserted

29

already inserted

30

already inserted

31

already inserted

32

already inserted

33

already inserted

34

already inserted

35

already inserted

36

already inserted

37

already inserted

38

already inserted

39

already inserted

40

already inserted

41

already inserted

42

already inserted

43

already inserted

44

already inserted

45

already inserted

46

already inserted

47

already inserted

48

already inserted

49

already inserted

50

already inserted

51

already inserted

52

already inserted

53

already inserted

54

already inserted

55

already inserted

56

already inserted

57

already inserted

58

already inserted

59

already inserted

60

already inserted

61

already inserted

62

already inserted

63

already inserted

64

already inserted

65

already inserted

66

already inserted

67

already inserted

68

already inserted

69

already inserted

70

already inserted

71

already inserted

72

already inserted

73

already inserted

74

already inserted

75

already inserted

76

already inserted

77

already inserted

78

already inserted

79

already inserted

80

already inserted

81

already inserted

82

already inserted

83

already inserted

84

already inserted

85

already inserted

86

already inserted

87

already inserted

88

already inserted

89

already inserted

90

already inserted

91

already inserted

92

already inserted

93

already inserted

94

already inserted

95

already inserted

96

already inserted

97

already inserted

98

already inserted

99

already inserted

100

already inserted

LECTURE 13: QUADRATIC PROBING

DATE -

$$\text{eg: } m = 10, a_1 = 1, a_2 = 1, h(k) = k^2 \mod 10$$

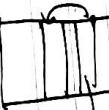
DATE -

- * Similar to linear probing which means in case of collision we will again probe.
- * For second probe different hash function will be used.

Hash function for linear probing: \rightarrow can be modified as

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

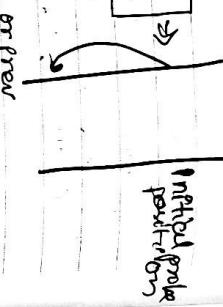
$$h'(k, 1) = (h(k) + 1) \mod m \rightarrow$$



Initial probe position

Hash Function for Quadratic Probing:

$$h'(k) = (h(k) + c_1 k + c_2 k^2) \mod m$$



Probe

- * Distance of probe position from initial probe position increases quadratically.

- * Advantage of quadratic probing over linear probing:

- ↳ Doesn't suffer with primary clustering. Probability that elements will be stored contiguously is reduced in this case.

Secondary clustering possible

↳ If initial probe position is same for any 2 keys then the entire probing sequence will also be same for both.

- * By doing only 'm' probe sequences we should be able to examine the entire table.
- * If we choose a_1, a_2, m randomly then problem may arise.

↳ Problem in quadratic probing
we might not be able to examine 'm' nodes/slots of hash table using 'm' probing -
we might end up assuming that the table is full even without examining some of the locations.

Note: c_1, c_2, m should be chosen in such a way that whatever the initial probe position is, if 'm' probes are done - one should be able to probe/examine the entire table.

↳ Coming up with such numbers is challenging

↳ Note - Due to above mentioned constraints choosing c_1, c_2, m wisely - Quadratic probing is not preferred.
More commonly used method is linear probing.

Note: - Secondary clustering is possible if only initial probe is dependent on key, later probes are independent of key
↳ problem with both linear and quadratic probing

5	X
4	
3	X
2	
1	

Q How many different probe sequences are possible using quadratic probing? $\Rightarrow m^m$ m - initial probe sequences are possible
 Note Number of probe sequences depend on number of initial probe sequence.

Note Numbers of collisions can't be reduced unless we have lot of probe sequences.

LECTURE 14: DOUBLE HASHING

* Instead of having one hash function & then constant displacement, here we will have 2 hash functions.

Hash Function

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

Advantage of new hash func'

→ No primary clustering ∵ probe displacement is no way linear ∵ next probe position decided by $h_2(k)$

Note:- In Double hashing, successive probes do not depend on initial probe as probe sequence can vary even in case of same initial probe.

→ No secondary clustering.

Number of probe sequences possible $\approx m^2$.

Reason: every probe can start at some position from $(0 - (m-1))$ can have displacement varying from $(0 - (m-1))$

Number of hash functions to be used to decide displacement can be decreased, decreasing the chances of collision

Not widely used.

Note:- Double hashing is best known technique. $n = h_1(k)$, m should be decided carefully so that all the hash table slots are examined

→ try to make them relatively prime $h_2(k), m$ or odd $\frac{m}{2}$ prime

Note:- Symbol table in computers is a hash table

Note:- Hashing preferred if major operation is searching, data is static