

Hashing

Hashing is a technique of mapping a large set of arbitrary data to tabular indexes using a hash function.

It is a method of representing dictionaries for large dataset.

It allows lookups, updating and retrieval operation to occur in a constant time i.e. $O(1)$.

Why it is required?

Linear and Binary search perform lookup/search with complexity of $O(n)$ and $O(\log n)$ resp.

As the size of the dataset increases, these complexities also become significantly high, which is not acceptable.

We need a technique that is independent of the size of data. Hashing allows to do this in constant time $O(1)$.

Hash function -

A hash function is used for mapping each element of a dataset to indexes in the table.

Hash table

Hash table is the data structure that stores elements in key-value pairs -

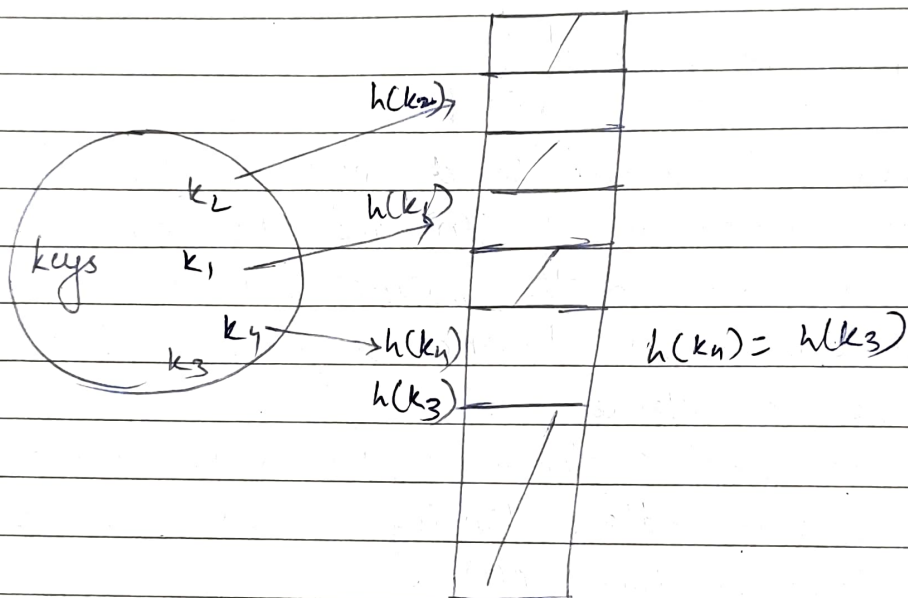
Key - unique integer that is used for indexing the values

Value - data that are associated with keys.

key	data
-----	------

In a hash table, a new index is processed using a key. And, the element corresponding to that key is stored in the index. This is called hashing.

Let 'k' be a key and $h(k)$ be a hash function. Here, $h(k)$ will give us a new index to store the element linked with 'k'.



Search keys - 24, 52, 91, 67, 48, 83

Hash function - $k \bmod 10$

- $k \bmod n$
- Mid square method
- Folding method

$k \bmod 10$

$24 \bmod 10 = 4$

0	
1	91
2	52
3	83
4	24
5	
6	
7	67
8	48

* Division Method

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

$m =$ larger than numbers

to minimize collision, m should be prime or no. with small divisor.

Ex: 10, 15, 12 if $m=5$

$$h(10) = 10 \bmod 5 = 0$$

$$h(15) = 15 \bmod 5 = 0 \Rightarrow \text{collision}$$

$$h(12) = 12 \bmod 5 = 2$$

* Mid-Square Method

$$h(k) = l$$

key is squared (k^2) -

'l' is obtained by deleting the digits from both ends.

$$k - \quad 10 \quad 15 \quad 12$$

$$k^2 - \quad 100 \quad 225 \quad 144$$

$$l - \quad 0 \quad 2 \quad 4$$

* Folding Method

$$h(k) = k_1 + k_2 + k_3 + \dots + k_n$$

- keys are partitioned into parts.

- parts are added together.

$$k - \quad 10 \quad 21 \quad 11$$

$$1, 0 \quad 2, 1 \quad 1, 1$$

$$h(k) = 1 \quad 3 \quad 2$$

* Multiplication Method

- 1) Choose constant 'A' such that $0 < A < 1$
- 2) Multiply 'A' with 'key'.
- 3) Extract the fractional part of $k \cdot A$
- 4) Multiply the result of above by size of hash table i.e. M.
- 5) Resulting hash value is obtained by taking the floor of the result in (4).

$$h(k) = \lfloor M (kA \bmod 1) \rfloor$$

Ex:- $k = 12345$

$$A = 0.35784$$

$$M = 100$$

$$h(12345) = \text{floor} \left[100 (12345 \times 0.35784 \bmod 1) \right]$$

$$= \text{floor} \left[100 (4417.5345 \bmod 1) \right]$$

$$= \text{floor} \left[100 \times 0.5345 \right]$$

$$= \lfloor 53.45 \rfloor$$

$$h(12345) = 53$$

* Collision :-

Since a hash function generates small no. for a given key, there is a possibility that two keys result in the same value. The situation where a newly inserted key maps to an already occupied slot in hash table is called Collision.

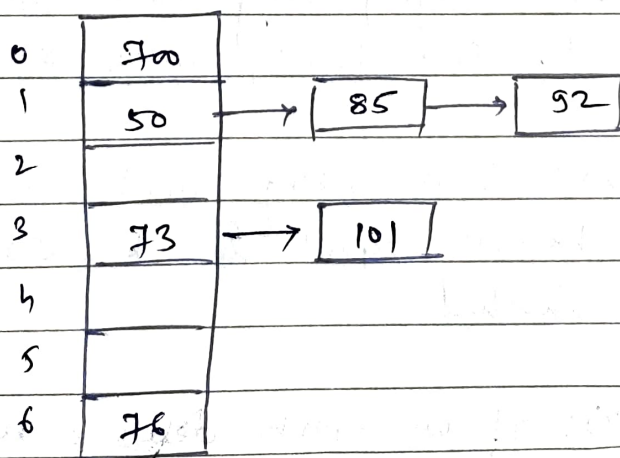
* Collision Resolution Techniques.

- ① Separate Chaining
- ② Open Addressing

① Separate Chaining :-

This is most popular and commonly used technique. Linked list data structure is used to implement this technique. When multiple keys are hashed into same index, then these elements are inserted into a singly-linked list (known as chain).

Ex:- $h(k) = k \bmod 7$ keys - 50, 700, 78, 85, 92, 73, 101



Advantages :-

- * Easy to implement
- * Hash table never fills up. We can always add more elements to chain.
- * Less sensitive to hash function
- * Mostly used when it is unknown how frequently keys are inserted or deleted.

Disadvantages-

- * Wastage of space
- * If chain becomes long, then searching can become $O(n)$ in worst case
- * Use extra spaces for links.

Q

Open Addressing-

Here, all elements are stored in hash table itself. So, at any time, size of hash table must be greater than or equal to the no. of keys to be stored. This is based on probing.

- Insert (k) :- keep probing until an empty slot is found. Once found, insert k.
- Search (k) :- keep probing until slot's key doesn't become equal to 'k' or an empty slot is reached.
- Delete (k) :- If we simply delete a key, then search may fail. So, deleted key's slots are marked as "deleted". Insert(k) can insert at deleted slot but search(k) doesn't stop at deleted slot.



A) Linear probing

Here, hash table is searched sequentially, that starts from the original location of hash.

If location that we get is occupied, then we check for the next location.

Ex: $h(k) = k \bmod 7$ keys are - 50, 700, 76, 85, 92, 73, 101.

0		700	700	700	700	700
1	50	50	50	(50)	(50)	50
2				85	(85)	85
3					92	(92)
4						73
5						101
6			76	76	76	76
	50	700	76	(85)	(92)	73 & 101

Collision

Collision

B) Quadratic Probing

Here, we look for the $(i^{th})^2$ slot in i^{th} iteration.

We always start from original hash location.

If location is occupied, then we check for other slots.

Ex: $h(k) = k \bmod 7$ Insert - 22, 30, 50.

0	
1	22
2	30
3	
4	
5	
6	

22	
30	
50	

$50 \% 7 = 1$, occupied

so search - $1 + 1^2 = 2$

Again occupied, so

search $1 + 2^2 = 5$, unoccupied

Hence, Insert at 5th index

$\rightarrow 1 + 2^2$



c) Double Hashing

In this technique, the increment for probing sequence are computed by using another hash function.

If $h_1(x)$ is full, we check $[h_1(x) + i * h_2(x)] \% m$.
where 'm' is table size. $i = 1, 2, \dots$

Ex: - Hash table size = 7

$$h_1(x) = x \bmod 7 \quad h_2(x) = 1 + x \bmod 5$$

Insert - 27, 43, 92, 74

① $27 \% 7 = \underline{6}$ - Empty so insert 27 in 6th slot

② $43 \% 7 = \underline{1}$ - Empty so insert 43 in 1st slot

③ $92 \% 7 = 1$ - occupied so check -

$$\begin{aligned} & [h_1(x) + i * h_2(x)] \% 7 \\ &= [92 \% 7 + 1 * (1 + 92 \% 5)] \% 7 \\ &= [1 + 1 * (1 + 2)] \% 7 = 4 \% 7 = 4. \\ & \text{Insert } 92 \text{ at } \underline{4^{\text{th}}} \text{ slot} \end{aligned}$$

④ $74 \% 7 = 4$ - occupied -

$$\begin{aligned} & [74 \% 7 + 1 * (1 + 74 \% 5)] \% 7 \\ &= [4 + 1 * (1 + 4)] \% 7 = 9 \% 7 = 2 \\ & \text{Insert } 74 \text{ at } \underline{2^{\text{nd}}} \text{ slot} \end{aligned}$$

0	
1	43
2	74
3	
4	92
5	
6	27