

hashing

hashing = doğramak, kıymak

- * inefficient in operations that require any ordering
 - ↳ such as finding min or max, sorting

average time for insert,
delete and find $O(1)$
- worst $O(n)$ -

applications:

- dictionary
- ctrl+f (find command)
- spelling correction
- authentication (doğrulama)

insertion of data with its key

data 1: 22 → $22 \% 10 = 2$

data 2: 35 → $35 \% 10 = 5$

data 3: 27 → $27 \% 10 = 7$

data 4: 32 → $32 \% 10 = 2 \rightarrow$ collision

hash function
 $\text{key of the data} \% \text{table size} = \text{hash value}$

hash table

[0]
[1]
[2] 22
[3]
[4]
[5] 35
[6]
[7] 27
[8]
[9]

collision handling

- 1- separate chaining (linked lists)
 - 2- linear probing (arastırma, search)
 - 3- quadratic probing
 - 4- double hashing
- } open addressing

! important to have
prime table size
(for uniform distribution)

① separate chaining:

- the array elements are pointers to the first node of the lists
- new item is inserted to the front of the list
- load factor should be around 1

[0] 3 → 9 → 6

[1] 7 → 4

[2] 2 → 14

load factor (λ): $\frac{n(\text{elements})}{\text{table size}}$

* cost of searching = constant time + expected chain length = $O(1 + \lambda)$ - average $(1 + \frac{\lambda}{2})$ -

open addressing:

- all the data go inside the table, so bigger table is needed
- the load factor should be below 0.5.
- if a collision occurs, alternative cells are tried until an empty cell is found

② linear probing

hash function

hash table
m

$$\left(\begin{array}{c} \text{key of} \\ \text{the data} \end{array} \% \begin{array}{c} \text{table} \\ \text{size} \end{array} + i \right) \% \begin{array}{c} \text{table} \\ \text{size} \end{array} = \text{hash} \\ \text{value}$$

data 1: 22 \rightarrow 2 ($i=0$)

data 2: 33 \rightarrow 3 ($i=0$)

data 3: 24 \rightarrow 4 ($i=0$)

data 4: 52 \rightarrow 2 + $i \rightarrow i=1$ [3] full
 $i=2$ [4] full
 $i=3$ [5] ✓

[0]
 [1]
 [2] 22
 [3] 33
 [4] 24
 [5] 52
 [6]
 [7]
 [8]
 [9]

! use lazy deletion (if you remove item, algorithm fails. mark it in a special way)

• primary clustering = when the blocks of occupied cells start forming

insertion = $\left(1 + \left(\frac{1}{1-\lambda} \right)^2 \right) \cdot \frac{1}{2}$ (the average cells that are examined)

find = successful is same as insertion, unsuccessful average $\left(1 + \frac{1}{1-\lambda} \right) \cdot \frac{1}{2}$

example of linear probing analysis

	successful	
[0] 9	* 20 20% \rightarrow [5] ✓	1
[1]	* 30 30% \rightarrow [8] ✓	1
[2] 2	* 2 2% \rightarrow [2] ✓	1
[3] 13	* 13 13% \rightarrow [2], [3] ✓	2
[4] 25	* 25 25% \rightarrow [5], [4] ✓	1
[5] 24	* 24 24% \rightarrow [2], [3], [4], [5]	4
[6]	* 10 10% \rightarrow [10] ✓	1
[7]	* 9 9% \rightarrow [5], [10], [0] ✓	3
[8] 30		$\frac{15}{15}$
[9] 20		avg = $\frac{15}{8}$
[10] 10		

	unsuccessful	
* 0	[0], [1] null ✓	2
* 1	[1] null ✓	1
* 2	[3], [5], [4], [5], [6] null ✓	5
* 3	[3], [4], [5], [6] null ✓	4
* 4	[4], [5], [6] null ✓	3
* 5	[5], [6] null ✓	2
* 6	[6] null ✓	1
* 7	[7] null ✓	1
* 8	[8], [3], [10], [5], [1] ✓	5
* 9	[9], [10], [0], [1] ✓	4
* 10	[10], [0], [1] ✓	$\frac{3}{31}$

$$\text{avg} = \frac{31}{8}$$

③ quadratic probing

hash function

$$\left(\begin{array}{c} \text{key of} \\ \text{the data} \end{array} \% \begin{array}{c} \text{table} \\ \text{size} \end{array} + i^2 \right) \% \begin{array}{c} \text{table} \\ \text{size} \end{array} = \text{hash} \\ \text{value}$$

$$H_i = H_{i-1} + 2i - 1 \pmod{M}$$

• eliminates primary clustering problem of linear probing

• there is no guarantee to find an empty cell (problem)

(if size is prime and $\lambda \leq 0.5$, it is okay)

• secondary clustering = when the same alternative cells are occupied.

rehashing: expanding the table, when the load factor reaches 0.5, dynamically.

- double to a prime
- reinsert the new table by using new hash function!

④ double hashing

- eliminates secondary clustering

hash function

$$\text{hash1}(k) + j \cdot \text{hash2}(k) = \text{hash value}$$

