

Searching Algorithms

* Linear Search

↳ unsorted search list dataset

↳ small dataset

↳ low memory (algorithm has low space complexity [O(1)])

- implemented using a simple for loop

* Binary Search

↳ most popular search

↳ need sorted array as input

↳ not suitable for linked lists (slow middle access)

↳ $T(n) = O(\log_2 n)$

```
def bSearchIterative(arr, target):
```

```
    left = 0
```

```
    right = len(arr) - 1
```

```
    while left <= right:
```

```
        mid = (left + right) // 2      ← floor division
```

```
        if arr[mid] == target:
```

```
            return mid
```

```
        elif arr[mid] > target:
```

```
            right = mid - 1
```

```
        else:
```

```
            left = mid + 1
```

```
    return -1
```

* Jump Search

↳ for sorted arrays

↳ $T(n) = O(\sqrt{n})$

↳ Skips ahead by fixed intervals.

↳ Works better than bSearch for linked lists

↳ bSearch may cause many cache misses, when accessing high latency memory.

↳ jump search is more cache friendly

↳ can be modified for partially sorted arrays

↳ optimum jump size is \sqrt{n} (geeksforgeeks has calculation)

$n = \text{len(arr)}$

step = \sqrt{n}

prev = 0

handles final iteration

while arr[min(step, n) - 1] < target:

 prev = step

 step += \sqrt{n}

 if prev >= n:

 return -1

for i in range(prev, min(step, n)):

 if arr[i] == target:

 return i

return -1

* Interpolation Search

↳ Time = $O(\log(\log(n)))$

↳ effective for uniformly distributed arrays

↳ improves upon binary search by estimating probable position of target

$$\boxed{\text{pos} = \text{low} + \left(\frac{\text{target} - \text{arr}[\text{low}]}{\text{arr}[\text{high}] - \text{arr}[\text{low}]} \times (\text{high} - \text{low}) \right)}$$

`interpolationSearch(arr, low, high, target)`

if ($\text{low} \leq \text{high}$ AND $\text{target} \geq \text{arr}[\text{low}]$ AND $\text{target} \leq \text{arr}[\text{high}]$)

$\text{pos} = \text{low} + ((\text{target} - \text{arr}[\text{low}]) / (\text{arr}[\text{high}] - \text{arr}[\text{low}]))$

* (high - low) / (high - low)

if $\text{arr}[\text{pos}] == \text{target}$

return pos

else if $\text{arr}[\text{pos}] < \text{target}$

return interpolationSearch(arr, pos+1, hi, target)

else

return interpolationSearch(arr, lo, pos-1, target)

Sorting Algorithms

* Stable - maintains relative order of equal elements

* Insertion Sort

↳ split array into 2

↳ assume left half is sorted and add unsorted items to correct position in sorted array

↳ start with 2nd element (take 1st element alone to be sorted part)

↳ compare 1 and 2 and swap if needed

↳ move to 3rd element and compare with 1 and 2 and put in correct position

for $i = 1$ to $n <= \text{start with } 1$ $n = \frac{\text{number of elements}}{\text{size of array}}$

element $\rightarrow \text{key} = \text{arr}[i]$
that we are "inserting" $j = i - 1$

since 0 is already sorted

while ($j \geq 0 \& \text{arr}[j] > \text{key}$)

$\text{arr}[j+1] = \text{arr}[j]$

$j = j - 1$

$\text{arr}[j+1] = \text{key}$

* look for position of key by moving each compared element to the right along the sorted part of the array till we find the correct place

* Shell Sort

- ↳ 1st algo to break quadratic time barrier (i.e. $T(n) \leq O(n^2)$)
- ↳ AKA diminishing increment sort
- ↳ Best $\rightarrow O(n \log n)$ Avg $\rightarrow O(n \log n^2)$ - Worst $O(n^2)$

```

for (int interval = n/2; interval > 0; interval /= 2) {
    for (int i = interval; i < n; i++) {
        int temp = arr[i];
        int j;
        for (j = i; j >= interval && arr[j - interval] > temp; j -= interval)
            arr[j] = arr[j - interval];
        arr[j] = temp;
    }
}

```

- * In worst case, this devolves to insertion sort

* Radix Sort

- ↳ considers the structure of keys
 - ↳ groups 'numbers' by their individual digits (radix)
 - ↳ using radix as key, uses counting sort to sort them
 - ↳ is a stable sort (will not work otherwise)
 - ↳ usually use for a fixed range
 - ↳ has high auxiliary space
 - ↳ $T(n) = O(d \cdot (n + b))$
- n - number of elements
d - max digits in a number
b - base of counting (10 for decimal etc.)

```

int getMax(int arr[], int n) {
    int mx = arr[0];
    for (int i=0; i<n; i++) {
        if (arr[i] > mx) {
            mx = arr[i];
        }
    }
    return mx;
}

```

| | | | | | | | |
|-----|----|----|----|------|----|---|----|
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| 170 | 45 | 75 | 90 | 1800 | 24 | 2 | 66 |

| | | | | | | | | | |
|----|----|---|---|---|---|---|---|---|---|
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
| 12 | 12 | 3 | 4 | 4 | 7 | 2 | 6 | 8 | 1 |

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|----|---|
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
| 1 | 1 | 1 | 1 | 2 | 4 | 1 | 2 | 66 | |

```
Void counts(int arr[], int n, int place) {
```

```

    int output[n];
    int count[10] = {0};

```

```

    for (int i=1; i<10; i++) {
        count[i] += count[i-1];
    }

```

```

    for (int i=0; i<n; i++) {
        count[(arr[i]/place)%10]++;
    }

```

```

    for (int i=n-1; i>=0; i--) {
        output[count[(arr[i]/place)%10]-1] = arr[i];
        count[(arr[i]/place)%10]--;
    }

```

```

    for (int i=0; i<n; i++) {
        arr[i] = output[i];
    }

```

```
void radixS(int arr[], int n) {
```

```

    int m = getMax(arr, n);

```

```

    for (int place=1; m/place > 0; place *= 10) {
        counts(arr, n, place);
    }

```

Hashing

- * Have average case time complexity = $O(1)$
Worst case is still $= O(n)$
- * 2 main types
 - ↳ Static hashing
 - ↳ Dynamic hashing
- * Hashing has 3 main parts
 - 1) Key: input given to the hash function
 - 2) Hash function: takes key and returns index of an element in hash table called hash index
 - 3) Hash table: data structures that maps keys to values.
- * Hash table is usually an array of structs that includes a char array to store the key followed by any other required data types
- * Various implementations of hashing
 - ↳ hash list
 - ↳ hash tree
- * Ideal Hashing
 - ↳ Each key is converted into a separate index and stored in home bucket
 - ↳ No collisions
- * But collisions do occur so we consider 2 things.
 - 1) Choose a good hash function to minimize collisions
 - 2) Handle overflow efficiently

- * Clustering - hash function groups multiple key value pairs together
 - not in same bucket but many nearby buckets are filled
- * Collisions - hash function evaluates completely different keys and returns the same hash index thereby storing them in the same bucket.
 - ↳ pigeonhole principle
- * Good hash functions avoids above 2 conditions as much as possible
- * Why can't we use n-sized hash table?
 - ↳ high memory usage
 - ↳ need a perfect hash function that maps each unique key to a unique hash index
 - ↳ very hard to do for dynamic data (data is stored and deleted)

* Characteristics of hash functions

- Efficiency
- Uniform distribution (avoid clustering)
- Determinism (same input key always gives same hash index)
- Minimized collisions

Hash Functions

1) Middle-Square Hash Function

- + Take input key x
- + Square the key (x^2)
- + Choose a specific number of digits from the middle of x^2
- + Use m on extracted digits to fit to table ($m = \text{table size}$)
- Used for small hash tables with integer keys.

2) Multiplicative Method:

$$\begin{aligned} * h(x) &= \text{Floor}(\text{table size} \times (\text{key} \times A \% 1)) \\ &= L(m \cdot (k \cdot A \% 1)) \end{aligned}$$

* A is a constant (usually 0.618 is chosen) [$0 < A < 1$].
Knuth's recommendation

- * Better distribution than division method.
- * Slower than division

4) Folding Hash Function

1) Separate long key into smaller fixed-size parts

2) Combine parts

↳ XOR, AND or Add together

3) Mod by table size

5) Digit Analysis

* Assumes that distribution of keys is known in advance

↳ Analyze digits in key that have skewed distribution

↳ Eliminate skewed digits

↳ Use remaining digits

6) Division Method

* Simplest and most widely used

$$* h(x) = x \% m \quad (m = \text{table size})$$

* need to choose good m

* Generally: a prime number. That is not close to a power of 2.

* For a hash table of size b, & a given bucket will have approximately $2^{32}/b$ integers (when trying to hash all integers)

* If divisor is even, even ints go into even buckets and odd into odd-buckets.

* If divisor is odd, even, odd both will hash into any home.

- + Biased distribution occurs if we use a divisor that is a multiple of prime numbers.
- + Ideally choose a large prime for divisor
- + Or choose divisor such that it has no prime factors smaller than 20.

- + Criterion of a Hash Table
 - ↳ Key density = n/T (n - number of keys in table, T - number of distinct possible keys)
 - ↳ Loading density / factor
 $\alpha = n/(sb)$ (s - number of slots)
 $(b$ - number of buckets)

Overflow Handling

- + When a collision occurs we get an overflow
- + 2 Solutions
 - 1) Search for an empty bucket (open addressing)
 - ↳ linear probing (linear open addressing)
 - ↳ Quadratic Probing
 - ↳ Random Probing
 - 2) Each bucket has a list of all keys indexed to it
 - ↳ Array ~~linear~~ to linear list
 - ↳ Chaining

④ Open addressing - all elements stored directly into hash table

- * ^{Quadratic} Random probing $\Rightarrow h'(k, i) = (h(k) + i^2) \% m \quad \left. \begin{matrix} i = \text{number of collisions} \\ \text{for given } k. \end{matrix} \right\}$
- * Random probing $\Rightarrow h'(k) = (h(k) + r[i]) \% m$