

# SORTING

Algorithm	Worst-case running time	Best-case running time	Average-case running time
Selection Sort	$O(n^2)$	$O(n^2)$	$O(n^2)$
Inception Sort	$O(n^2)$	$O(n)$	$O(n^2)$
Merge Sort	$O(n \log n)$	$O(n \log n)$	$O(n \log n)$
Quick Sort	$O(n^2)$	$O(n \log n)$	$O(n \log n)$

smallest data item bubbles up to top of stack  
top compares first two elements

Date: .....

## BUBBLE SORT

→ simplest sorting algorithm that works by repeatedly swapping the adjacent elements if they are in wrong order.

→ basic for sorting small data

### Advantages.

- ↳ simple to write & easy to understand as it's iterative
- ↳ requires few lines of code
- ↳ little memory ahead

### Disadvantages.

- ↳ time consuming
- ↳ avg. time increases exponentially as number of elements increase
- ↳ not suitable for large sized data.

Best Case:  $O(n)$  (when input data is already sorted)

Worst Case:  $O(n^2)$

Average Case:  $O(n^2)$

## Pseudocode.

**Bubble Sort (array, num\_elems)**

for  $k=1$  until  $num\_elems - 1$  do

loop ( $k = num\_elems - 1$ ) # max passes required =  $n - 1$

$i = 0$

# You don't want to include the last element

because then index  $i + 1$  would go out of bounds

loop ( $i \leq num\_elems - 2$ ):

if ( $array[i] > array[i + 1]$ ):

swap ( $array[i]$ ,  $array[i + 1]$ )

$i = i + 1$

$K = k + 1$

return.

↳ Bubble Sort is slow sorting algorithm

↳ Selection Sort is also quadratic time complexity,  
which is also slow.

## POSSIBLE OPTIMIZATIONS.

- at any point during bubble sort algorithm, you will have a sorted part of your array.
- no need to walk through the sorted part, because it's already sorted.

Optimized Bubble Sort (array, num\_elems):

$k = 1$

loop ( $k \leq \text{num\_elems} - 1$ )

$i = 0$

$\text{is\_swap} = \text{False}$

loop ( $i \leq \text{num\_elems} - k - 1$ ):

if  $\text{array}[i] > \text{array}[i+1]$ :

$\text{swap}(\text{array}[i], \text{array}[i+1])$

$\text{is\_swap} = \text{True}$

$i = i + 1$

if ( $\text{is\_swap} == \text{False}$ ):

$\text{break}$

$k = k + 1$

$\text{return}$

in every iteration  
of selection sort, the min  
element from unsorted array  
is picked & moved to begining  
of unsorted array

algorithm divides the input list  
into two parts.

Date: .....

# SELECTION SORT

→ sorts an array by repeatedly  
finding the minimum element  
from unsorted part. & putting  
at the beginning.

Advantages:

↳ performs well on small sized

↳ in-place algorithm

Disadvantages:

↳ poor efficiency when dealing with huge list of items.

↳ requires  $n^2$  no. of steps.

Best Case:  $O(n^2)$

Worst Case:  $O(n^2)$

Average Case:  $O(n^2)$

- Walk through array

- Find minimum

- Swap this with value located at index = 0

- Repeat

- Walk through array

- Find second minimum

- Swap this with value at index = 1

- ① initialize minimum value to location 0.
- ② traverse the array to find minimum element.
- ③ while traversing if any element smaller than minimum value found then swap

Pseudocode .

Selection Sort ( $A, n$ ):

$i = 0$  index of starting fully inserted in array from

loop  $i \leq n-2$ :

min\_index =  $i$

$j = i+1$

loop  $j \leq n-1$ :

if  $A[j] < A[min\_index]$ :

min\_index =  $j$

$j = j+1$

temp =  $A[i]$

$A[i] = A[min\_index]$

$A[min\_index] = temp$

$i = i+1$

return .

## INSERTION SORT

Advantages .

↳ straightforward implementation

↳ good performance when dealing with smaller lists.

inefficient on large  
data sets

more efficient than others.

Date: \_\_\_\_\_

## Disadvantages:

- not best performing sorting algorithm.
- for every 'n' element that needs to be sorted,  $n^2$  steps required
- time consuming.

Best Case:  $O(n)$

Worst Case:  $O(n^2)$

Average Case:  $O(n^2)$

Pseudocode:

```
insertion_sort(a-list, n):
    loop index : 1 to (n-1):
        current_value = a-list[index]
        free_index = index
        loop free_index > 0 and a-list[free_index-1] > current_value
            a-list[free_index] = a-list[free_index-1]
            free_index = free_index - 1
        a-list[free_index] = current_value
    return
```

55	-2	34	10	0	None	-2	55	34	10	0
↑	↑					↑	↑			

swap(anay, j, j-1)      -2 | 34 | 55 | 10 | 0

BUBBLE SORT

```

def bubble_sort(lst):
    if len(lst) == 1
        print(lst)
        return
    for j in range(len(lst)-1):
        for i in range(len(lst)-j-1):
            if lst[i] > lst[i+1]:
                lst[i], lst[i+1] = lst[i+1], lst[i]
        print(lst)
    pass
bubble_sort(lst)

```

SELECTION SORT

```

def selection_sort(lst):
    for i in range(len(lst)):
        minimum = i
        for j in range(i+1, len(lst)):
            if lst[j] < lst[minimum]:
                minimum = j
        lst[i], lst[minimum] = lst[minimum], lst[i]
    print(lst)

```

Date: .....

## INSERTION SORT

```
def insertion_sort(1st):
```

```
    for i in range(1, len(1st)):
```

```
        element = 1st[i]
```

```
        while i > 0 and 1st[i-1] > element:
```

```
            1st[i] = 1st[i-1]
```

```
            i = i-1
```

```
        1st[i] = element
```

```
    print(1st)
```

\* is selection better than bubble in  
worst case?

# Quick Sort

0	1	2	3	4	5
9	8	10	17	15	14

Note: that the left & right partitions do not need to be sorted for 10 to be sorted.

10 is in its sorted location because all elements smaller than it are on left and all elements bigger than it are on the right.

a	b
85	24
63	45
17	31
31	96
96	50

→ our pivot is 50 for now

→ as soon as l and r switch places, we are done

→ is the value at l and r appropriate as compared to pivot (50 in this case)?

with finding a place for pivot.

↳ since r is bigger than pivot so it's fine, so we update it (go left)

85	24	63	45	17	31	96	50
l					r		

→ "r" sees 31, which is smaller than the pivot (50).  
31 shouldn't be on "r" side of array

31	24	63	45	17	85	96	50
l					r		

→ swap the values of l and r.

\* rcl

→ have rcl so why not swapped-

Advantages

- one of the fastest algorithm on average.
- does not need additional memory.

Disadvantages

- worst-case complexity is  $O(n^2)$ .

Average Case:  $O(n \log n)$

Worst Case:  $O(n^2)$

Best Case:  $O(n \log n)$ .

## QUICK SORT

↳ selection & insertion sort.

→ fast algorithm

→ divide & conquer algorithm.

→ recursive algorithm with base case.

→ works on every object consisting of elements.

→ consider 3 numbers and find minimum out of them.

→ first swap with minimum and then swap with second.

→ divide into two parts.

→ swap with minimum.

→ swap with second.

→ swap with third.

→ swap with fourth.

→ swap with fifth.

→ swap with sixth.

→ swap with seventh.

→ swap with eighth.

→ swap with ninth.

→ swap with tenth.

Date: .....

- Quick sort is an example of a divide and conquer strategy.
- What is "divide & conquer"?
- Strategy in which you break the given problem down into subproblems recursively
  - ↳ these subproblems are similar to original problem
- These subproblems are solved recursively.
- The solutions to these subproblems are combined to form solution for original problem.

- Divide: the problem into number of subproblems.
- Conquer the subproblems by solving them recursively.
- combine the solutions to the subproblems into the solution for the original problem.

\* When is Quick Sort better than Merge Sort?

# MERGE SORT

↪ let's say you have unsorted array.

↪ let's split it up into two equal parts.

↪ it has  $n=8$  elements, so two subarrays of length 4 each.

Advantages:

- \* Merge sort is stable sort
- \* it is easy to understand
- \* it gives better performance

Disadvantages:

- \* it requires extra memory space.
- \* copy of elements to temporary array
- \* requires additional array.
- \* slow process.

Complexity of Merge Sort. The merge sort algorithm passes over the entire list & requires at most  $\log n$  passes and merges  $n$  elements in each pass. The total number of comparisons required by merge sort is given by  $O(n \log n)$ .

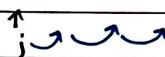
## Merge Sort - Concept.

85	24	63	45	17	31	96	50
----	----	----	----	----	----	----	----

find the mid value

$$\frac{0+3}{2} \rightarrow \begin{matrix} 24 & 45 & 63 & 85 \\ 0 & 1 & 2 & 3 \end{matrix}$$

$$\begin{matrix} 17 & 31 & 50 & 96 \\ 0 & 1 & 2 & 3 \end{matrix}$$



17	24	31	45	50	63	85	96
----	----	----	----	----	----	----	----

i

k

j

l

\* i and j are compared

24	45	63	88
----	----	----	----

0

1

\* j < i, j+l

\* i < j, i+l

\* j < i, j+l

\* i < j, i+l

\* j < i, j+l

\* i < j, i+l

\* j < i, j+l

\* i < j, i+l

\* j < i, j+l

\* i < j, i+l

\* j < i, j+l

\* i < j, i+l

83	45	63	88
----	----	----	----

24
----

V

63	88
----	----

24	45
----	----

24	45	63	88
----	----	----	----

24	45	163	88
----	----	-----	----

MERGE SORT

```
def merge(L, R, A)
```

lenL = len(L)

lenR = len(R)

i = j = k = 0

while i < lenL and j < lenR:

if L[i] <= R[j]

A[k] = L[i] *(first part of the section = first half)*

i = i + 1 *(first skip) the section = first half*

k = k + 1

else:

A[k] = R[j]

j = j + 1

: (then: k=k+1) until k is less than lenL *(first half) and then*

while i is less than lenL: *(second half) and then*

A[k] = L[i] *(second half) because we have*

i = i + 1

k = k + 1

while j is less than lenR: *(second half)*

A[k] = R[i]

j = j + 1

k = k + 1

def merge\_savt(an):

if len(arr) <= 1:

return an

mid = len(an) // 2

left\_half = an[:mid]

right\_half = arr[mid:]

left\_half = merge\_savt

right\_half = merge\_savt

sorted\_an = []

i = j = 0

→ Why is it important to create two separate copies of subarray?

→ Can we do it without making any copies, ie in-place?

while i < len(left\_half) and j < len(right\_half):

if left\_half[i] <= right\_half[j]:

sorted\_arr.append(left\_half[i])

i = i + 1

else:

sorted\_arr.append(right\_half[j])

j = j + 1

sorted\_arr += left\_half[i:]

sorted\_an += right\_half[j:]

return.

# RECURSION

There are two parts

- ① Base Case → responsible for
- ② Recursive Case → responsible

→ What would have happened if we did not have a base case defined?

↳ if there is no base case use of recursive definition becomes infinitely long and any program will never terminate & produce result.

- When you call a function, that call is saved in a stack, called a stack frame.
- All of variables and other associated information with that particular function call are saved.
- This also means that the stack frame provides a scope for the variables for each function call.
- So for an array of length =  $n$ ,  
• total number of calls =  $n$   
• recursive calls =  $n - 1$

## Recursion in Math

↳ one of the most obvious math definitions that can be stated in a recursive manner is definition of integer factorial.

$$\hookrightarrow 1! = 1 \quad (\text{base case})$$

$$\hookrightarrow N! = N * (N-1)! \quad (\text{recursive call})$$

## Recursion vs Iteration.

↳ we can calculate  $5!$  using a loop.

```
int fiveFactorial = 1
for (int i = 1; i <= 5; i++)
```

    fiveFactorial \*= i

↳ or we can calculate using recursion.

factorial (int n)

    if ( $n == 1$ )

        return 1

    else

        return n \* factorial (n-1)

\* typically, a recursive solution uses both more memory & more processing time than an iterative solution

## Fibonacci Series.

↳ calculation of sequence of Fibonacci numbers  
 $f_n$  can be stated as;

$$\begin{aligned} f_0 &= 0 \\ f_1 &= 1 \end{aligned} \quad \left. \begin{array}{l} \text{base cases.} \\ \text{recursion.} \end{array} \right\}$$

$$f_n = f_{(n-1)} + f_{(n-2)}$$

## Rules of Recursion.

- ↳ Base Case: Always have atleast one case that can be solved without recursion.
- ↳ Make Progress: Any recursive call should make progress toward a base case.
- ↳ Never Duplicate work by solving same instance of a problem in seperate recursive calls.

## Direct vs Indirect Recursion Call.

- ↳ direct recursion is when a method calls itself
- ↳ indirect recursion is when method calls a second method that can call first method.

# HASHING

Unsorted Map - ADT:

- ↳ these are called dictionaries.
- ↳ key-value pairs.
- ↳ keys are unique; values not necessarily unique.
- ↳ order doesn't matter

# query(`key`)

`M[key]` returns the value `v` associated with key `k` in map `M`, if one exists; otherwise raise a `KeyError`.

# add\_modify(`M`, `key`, `val`)

`M[key] = val`: Associate value `v` with key `k` in Map `M` replacing the existing value if map `M` already contains an item with key equal to `k`.

# delete(`M`, `key`)

`del M[k]`: Remove from map `M` the item with key equal to `k`; if `M` has no such item, then Raise a `KeyError`.

- \* Time complexity (average & worst) is  $O(n)$ .
- \* Binary search not possible  
map isn't sorted.

len(M): return the number of items in map M

iter(M): the default iteration for a map generates a sequence of keys in both ends of the map.

## HASH TABLES

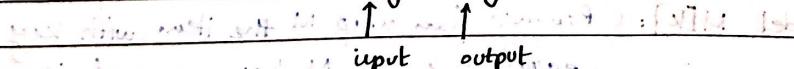
- ↳ faster search approach
- ↳ for unsorted data, we relied on linear search  $\rightarrow O(n)$
- ↳ for data arranged in sorted way, Binary Search would be better  $\rightarrow O(\log n)$

0	1	2	3	4	5	6	7	8	9	10
(1:B)	(2:2)	(3:B)		(5:H)				(9:G)		

searching in a hash table is faster than in an array

- is  $O(1)$
- ↳ use the keys as indices in the array.

$$h(\text{key}) = \text{key} \% (\text{size of array})$$



- Hashing is the process of generating an address/index from a key
- This is done through hash function.
- The generated address/index is referred to as the home address.
- The area in the memory where all the home addresses are located is called the prime area.
- A hash table is a collection of items which are stored in such a way as to make it easy to find them later.
- Each position of hash table, is often called slot.
- Consider,  $h(\text{key}) = \text{key \% table\_size}$
- Assume table\_size to be 7 (slightly greater than total entries)
- Conflict or collision may occur
- To prevent these collisions, we can use good hash function.

## Hashing Methods

Direct

Modulo Division

MidSquare

Rotation

Subtraction

Digit Extraction

Folding

Pseudorandom  
generation

### Direct Method

A direct hash function uses the item's key as bucket index. E.g. if the key is 937, the index is 937.

$$h(\text{key}) = \text{key}$$

#### Limitations;

For example, we cannot have Social Security number as the key using this method because Social Security numbers are nine digits, and we cannot have an array as large as nine digits.

## Subtraction Method

### Modulo-Division Method

→ if key < size then return  
key as it is.

- divides the key by hash table size & uses the remainder as index.
- Use Prime numbers as length of hash tables would work best.

$$h(\text{key}) = \text{key \% table\_size.}$$

\* 1a

- Significance of having a hash table with length that is a Prime Number

↳ When keys are randomly distributed, it probably won't matter what the size of your table is.

↳ However, in the case of non-random keys, a hash table of prime number length will produce the most wide spread distribution of integers to indices.

## Digit Extraction Method

↳ selected digits are extracted from the key and used as the address.

↳ what our hash function/mapping does is that it,

- picks the first, third and fourth digits (starting from the left)

- the resultant three-digit number is our index for that key-value pair.

\* all of these weird operations are there to

somehow reduce the number of collisions and get unique home addresses.

```
def digit_extraction(key):
```

$$d_3 = (\text{key} // 100) \% 10$$

$$d_2 = (\text{key} // 1000) \% 10$$

$$d_1 = (\text{key} // 10000) \% 10$$

$$\text{home\_address} = (d_1 * 100) + (d_2 * 10) + (d_3 * 1)$$

```
return home_address
```

so this home address is the home address of the key.

## Midsquare Method

to begin from our address length you can either take first half or last half off

def midsquare(key):

- calculate the square of your key

→ select your address from the middle of squared number.

$$9452^2 = 89340304 \text{; address is } 3403$$

Folding Methods: same with another methods

→ two methods: ① Fold Shift ② Fold Boundary

① Fold Shift

- the key is divided into parts whose size matches the size of required address
- then either left & right parts are shifted and added with middle part
- any extra digit is dropped off or take mod here with table size.

## (2) Fold Boundary

- you can think of the key is being folded along certain boundaries.
- this results in the key being divided into parts.
- but alternate sets of numbers are reversed.

## Rotation Method

- ↳ works well with keys that are in serial/sequence e.g.
  - student/employee ID's etc: 101, 102, 103, 104, ...
  - 1001, 1002, 1003, 1004, ...
- ↳ a simple hashing algorithm tends to create synonyms when hashing keys like these.
  - rotating the last character to front of the key minimizes this effect.

## Pseudorandom Hashing

- ↳ the key is used as seed in a pseudorandom-number generator, and the resulting random number is then scaled into possible address range using modulo-division.

## Perfect Hash Function

↳ a hash function that maps each item into a

unique slot is referred to as a perfect hash function.

↳ you can increase the size of table to reduce collisions

→ more slots

• more vacant slots means that probability of collision is lower.

• not always feasible though.

↳ so we also need a good hash function to ensure less collisions.

→ so the function is balanced.

## Characteristics of "Good" Hash Functions

- Fast to compute so that querying, addition and deletion are done typically in  $O(1)$ .
- Minimizes collisions.
- Outputs are uniformly distributed along hash table.
- Function is deterministic for a given input, the same output index is returned.

## Collision Resolution.

↳ two or more unique keys hash to same slot in hash table is called collision.

↳  $h(key_1) = h(key_2) = j$  for some static entries.

↳ Load factor of a hash table may indicate collisions.

•  $\lambda = \frac{\text{no of items}}{\text{table size}}$  . measure of how much of the hash table is occupied- bigger load factor, the more hash

• i.e. when the hash table is filled with some values

↳ part 2 on Collision Resolution

Open Addressing      linked List | PITA      Buckets

↳ linear probe

↳ quadratic probe

↳ pseudorandom

↳ key offset

## Open Addressing

- ↳ collisions are resolved in prime area.
- ↳ in case of collisions, you look for **vacant space** within prime area where you can store incoming key.

- ↳ unlike chaining, you do not introduce an additional data structure.

### ① Linear Probing

- ↳ you perform hashing on a key  $k$ , and get an address  $j$ .
- ↳ That slot,  $A[j]$  is already occupied - collision.

- ↳ start at an original hash value position & then move in a sequential manner through slots until we encounter first empty slot.

↳  $h(k) = j$ ,  $A[j]$  is occupied

↳ try looking in,  $A[(j+1) \bmod N]$

↳ if collision, then,  $A[(j+2) \bmod N]$

↳  $A[(j+3) \bmod N] \dots$

↳ keep looking until you found empty slot

↳ so you search until you either,

- find the key you are looking for
- reach an empty slot
- return back to index from where you started.

↳ if you want to remove a key from hash table.

- you can pass your key in hash function
- get index
- and set the value at that index to None.

None broke the search path to get

↳ and instead of placing None at the deleted key's index, place a special marker called "available" or "tombstone".

↳ this indicates that a certain slot wasn't always empty (or None) & there was something here before.

↳ this indicates that a certain slot wasn't always empty (or None) & there was something here before.

Disadvantage: → clustering

- if many collisions occur at same hash value, a number of surrounding slots will be filled by linear probing resolution.

→ impact on other items that are being inserted

and lead to conflicts

42	13	38	49	16	51	27
----	----	----	----	----	----	----

↳ how clustering can happen

variable length

Primary clustering

around the common base.

address.

- can we make collision path independent of have address?

how to know jump size?

Date: .....

is bigger than size me... may fit  
N], where skip is an integer.  
or clustering. point for  
the following note has been.

#  $A[(\text{arg\_ind} + \text{probe\_step}) \% N]$

↳ linear probing,  $f(i) = i$ , for  $i = 0, 1, 2, \dots, m-1$

↳ in quadratic probing,  $f(i) = i^2$ , for  $i = 0, 1, 2, \dots$

•  $A[h(k)+0], A[h(k)+1], A[h(k)+4], A[h(k)+9], \dots$

↳ this avoids primary clustering and instead creates secondary clustering.

↳ linear probing leads to primary clustering.

↳ quadratic solves this, but creates secondary clustering.

↳ secondary clustering is problematic because collision path is dependent on have address.

↳ but on open addressing

• on have address. collision will

be handled by some other stuff. \* this better than primary

(1) we want to come up with clustering because there are a collision path that is unique gaps in the middle but for each key hopefully you understand why

(2) Each key is unique, so if our it's a problem.

path is dependent on key, their

path is also unique.

## 8) Pseudorandom Numbers.

import random

key = int(input("Enter key = ?"))

random.seed(key)

for x in range(5):

print(random.randint(1,19), end=" ")

random.randint(a,b) returns a random integer  
N such that  $a \leq N \leq b$ .

→  $A[h(k)+f(i) \% N]$ , where  $f(i)$  is the pseudorandom number sequence produced using  $k$  as seed.

4)

## Separate Chaining

↪ simple and efficient way of resolving collisions.

↪ each slot in hash table, instead of storing keys/items directly would possess a secondary container, which instead of would contain the keys/items.

↪ more memory required (hash table + secondary containers)