

LEARN CODING

ale66

SORTING

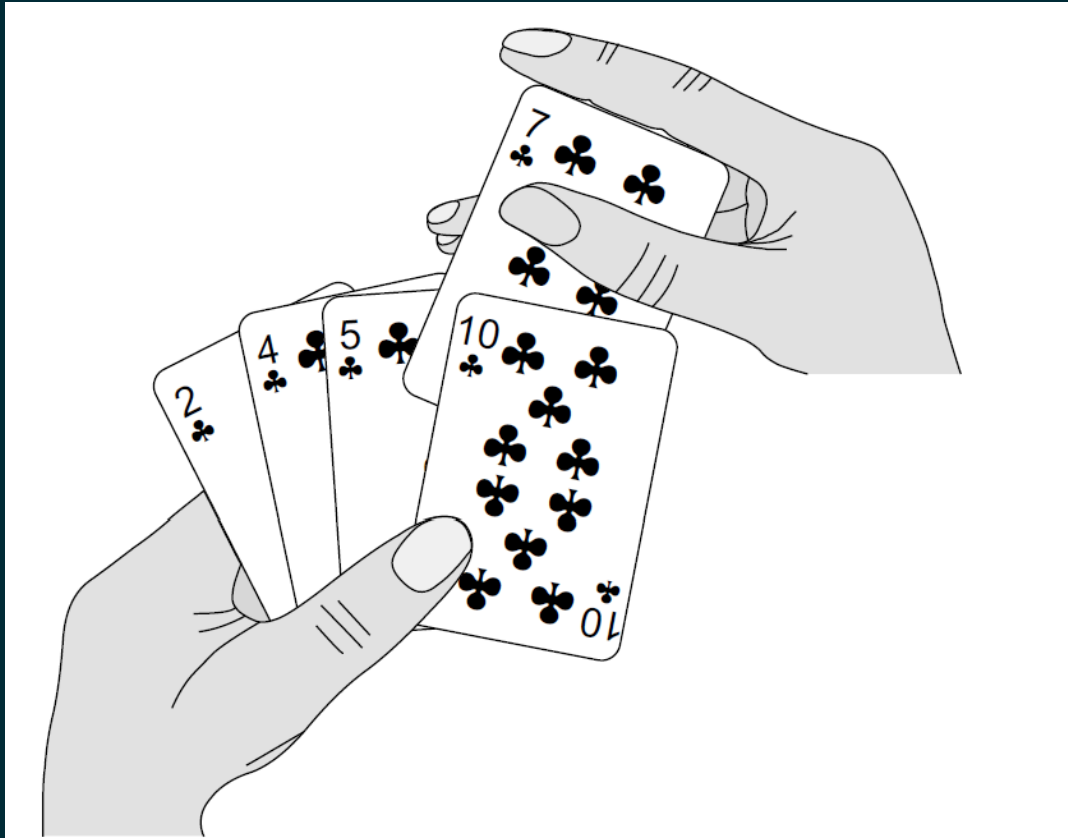


Figure 2.1 Sorting a hand of cards using insertion sort.

STATEMENT

Instance:

a sequence of n integers: $A = a_1, a_2 \dots a_n$

Solution:

A permutation of the values $\pi : [1..n] \rightarrow [1..n]$

Constraint

values never decrease: $a_1 \leq a_2 \dots a_n$

Let's assume that there are no repeated values and Python notation (first elem. is in pos. 0)

```
1 A = [11, 6, 8, 2, 22, 16, 25]  
2  
3 sort(A) = [?, ?, ?, ?, ?, ?]
```

$\pi(0) = ?$

$\pi(1) = ?$

$\pi(2) = ?$

$\pi(3) = ?$

...

Let's assume that there are no repeated values

```
1 A = [11, 6, 8, 2, 22, 16, 25]
2
3 sort(A) = [2, 6, 8, 11, 16, 22, 25]
```

$\pi(0) = 3$: the elem. in position 0 now goes to pos. 3

$\pi(1) = 1$: the elem. in pos. 1 remains there

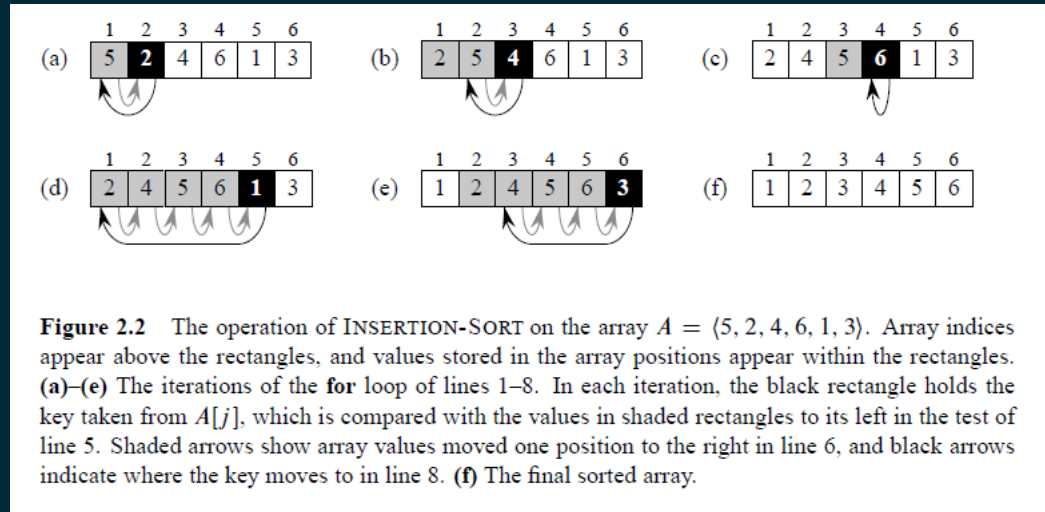
$\pi(2) = 2$: so does the elem. in pos. 2

$\pi(3) = 0$: the elem. in pos. 3 now goes to pos. 0

...

GOOD NEWS ABOUT SORTING

Solvable within $Kn \log_2 n$ steps

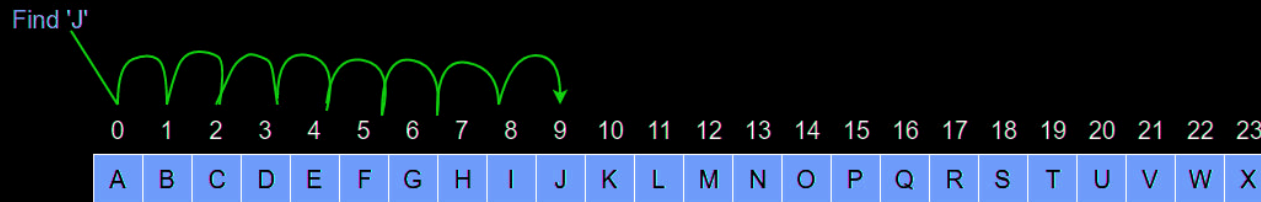


Even quicker when A is already half-sorted

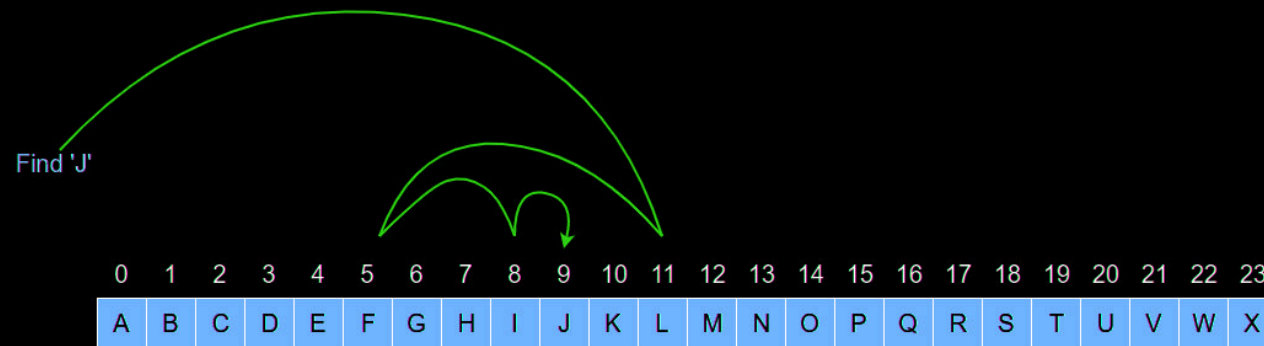
Python runs `powersort()`, an optimised version of TimSort

```
1 a = [11, 6, 8, 2, 22, 16, 25]
2
3 a.sort()
```

SEARCHING



Binary search vs Linear search



Instance:

- a collection of n integers $A = a_1, a_2, \dots, a_n$
- an integer k

Question: does k belong to the collection?

OBSERVATIONS

Strings and texts are searchable: alphanumerical ordering

The textbook algo. takes up to n comparisons

Search is everywhere: implementation must be

- quick, and
- reliable, handles negatives, 0s, NaNs, Nones

CASE STUDY: EGO NETWORKS [FERRARA ET AL, 2011]

- while less than N profiles collected
 - generate random FB ids (a fixed-length, 32-digits integer)
 - test the random id: does it land on an open FB profile?
 - yes: expand the visit to the neighborhood
 - no: go back to generating random ids

Cost: at least N generations+FB accesses; dN comparisons to check for duplicate ids

VANILLA SOLUTION: LINEAR SEARCH

```
1 def linSearch(data_list, search_key):
2     '''Linear search'''
3
4     found = False
5
6     n = len(data_list)
7
8     for elem in data_list:
9
10         if elem == search_key:
11             found = True
12
13     return found
```

EXERCISE

Apply **while** instead of **for** to stop operations as soon as the key value is found

Return the position at which the key was found

Model call:

```
1 my_data = [11, 6, 8, 2, 22, 16, 25]
2
3 my_key = 42
4
5 outcome = yourSolution(my_data, my_key)
```

ORDERED SEQUENCES

Special case: the input sequence is already sorted, either in increasing or decreasing order

a much more efficient algorithm is available: Binary Search

- only works for sorted sequences
- top speed: *at most* $\log_2 n$ comparisons

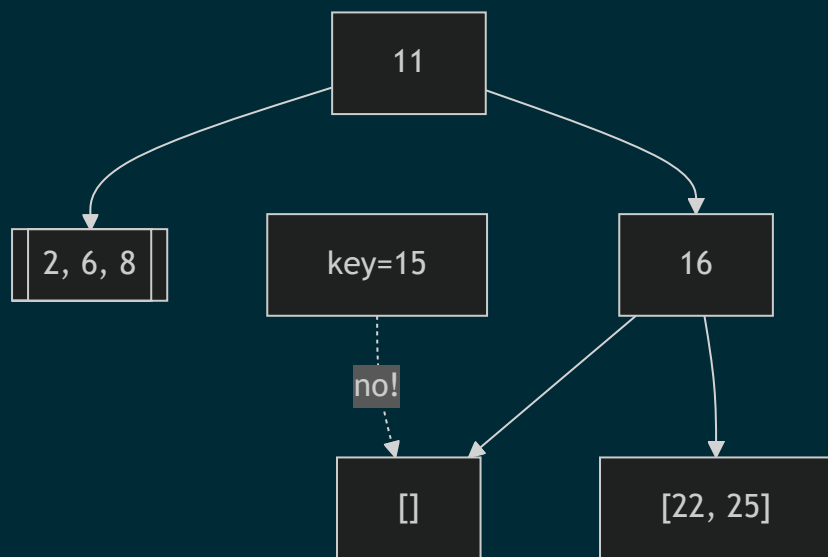
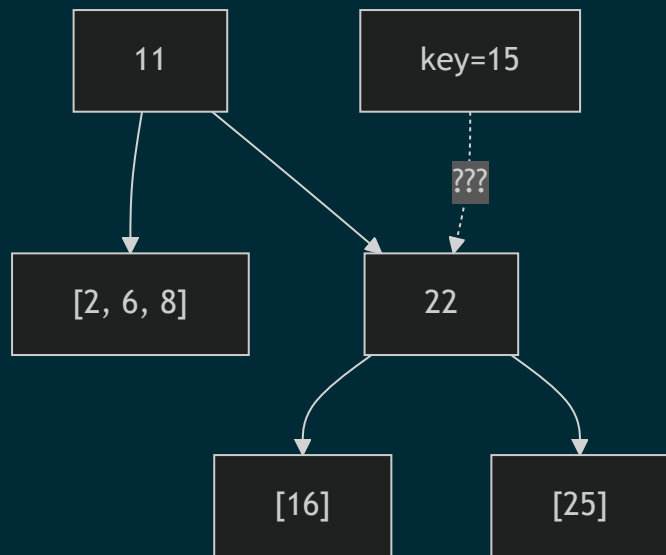
items	Vanilla (linear)	BS
1,000	1,000	10
1,000,000	1,000,000	20
1,000,000,000	1,000,000,000	30

IDEA

Data are sorted: exploit this property to cut down the size of the list *segment* to be checked

```
1 my_sorted_list = [2, 6, 8, 11, 16, 22, 25]
2
3 key = 15
```





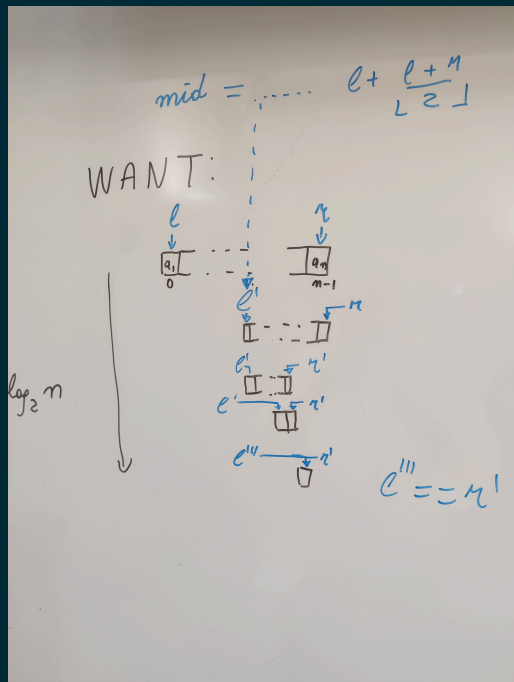
STEPS

- input the ordered list and the key value to be searched
- find the *median* value (here, it's right in the middle!)
- if the median == the key value then stop and say 'yes'
- but if the search key $>$ median then the value, if it exists, can only be in second half of the list
- otherwise, the key value, if it exists, can only be in the first half of the list
- depending on the result of the comparison, continue searching on the 'right' half of the list.

BS halves the searched data at each iteration

Soon, the halving will shrink the list down to just one value, so we check and finish

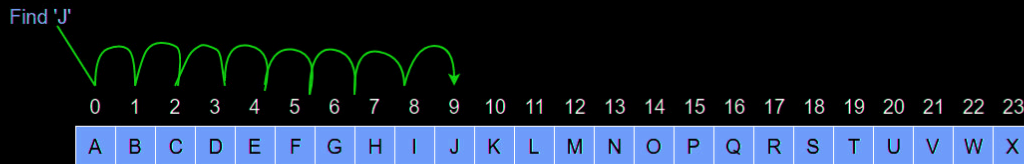
How soon? It will take at most $\log_2 n$ 'cuts' to shrink the list down to 1



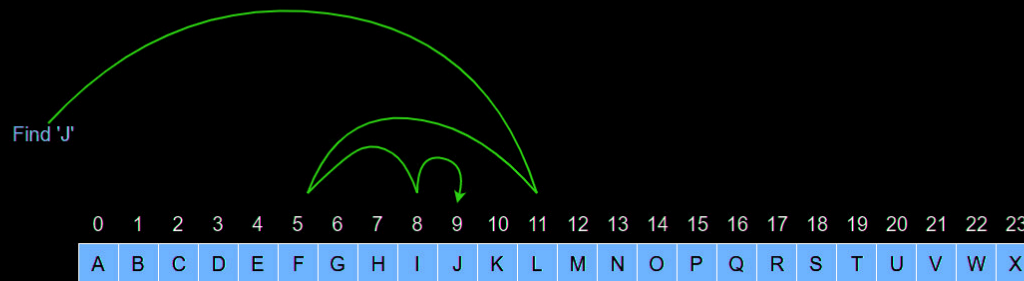
```
1 def bs(a, k):
2     '''Simple Binary search implementation: a is a list of integers, k is a search
3
4     found = False
5     n = len(a)
6
7     # the boundaries of our search
8     l = 0
9     r = n
10
11     while ((found == False) and (l < r)):
12
13         mid = l + int((r-l)/2)
14
15         if k < a[mid]:
16             r = mid
17
18         elif k > a[mid]:
19             l = mid + 1
20
21         else:
22             found = True
23             print('found in position ', mid)
24
25     return found
```

VISUALISATION

Instance: an ordered list of 24 uppercase chars, **key** = 'J'



Binary search vs Linear search



$$\lceil \log_2 24 \rceil = \lceil 4.5849 \rceil = 5 \text{ comparisons}$$