

Algorithms and Data Structures

Advanced sorting algorithms

dr Szymon Murawski

Comarch SA

March 14, 2019

Table of contents I

- 1 Merge sort
- 2 QuickSort
- 3 Counting Sort
- 4 Radix Sort
- 5 Bucket Sort
- 6 Conclusions

Sorting so far

Algorithms

- Selection sort
 - Insertion sort
 - Bubble sort
 - Cocktail sort
-
- All the algorithms we learned so far are quite bad, they have time complexity of $O(n^2)$, which makes them unsuitable for sorting large amounts of data
 - They are however very easy to implement
 - They also require little overhead, which makes them usable for very small arrays
 - Next we will learn about advanced sorting algorithms, which all perform much better than the previous ones!

Course plan

- 1 Merge sort
- 2 QuickSort
- 3 Counting Sort
- 4 Radix Sort
- 5 Bucket Sort
- 6 Conclusions

Divide and conquer sort

Divide and Conquer approach

- **Divide** the problem into smaller subproblems
- **Conquer** the subproblems by solving them recursively
- **Combine** the solutions to the subproblems into the solution to the original problem

Merge sort

- **Divide** the n -element sequence to be sorted into two subsequences of $n/2$ elements each.
- **Conquer** by sorting the subsequences recursively using merge sort.
- **Merge** the two sorted subsequences to produce the sorted answer

Merge sort

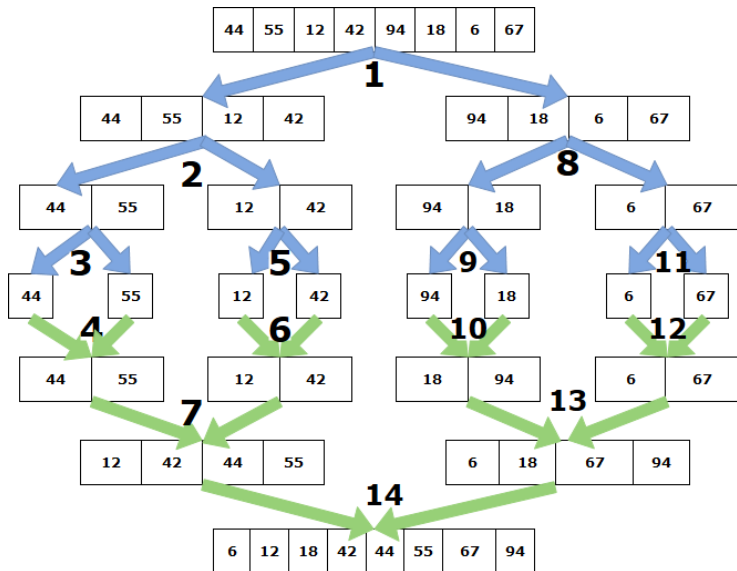
Pseudocode

- Input: Array A , start index p , end index r
- Output: Sorted array A

```
1 MergeSort(A, p, r):  
2   if  $p < r$   
3      $q = \text{floor}((p+r)/2)$   
4     MergeSort(A, p, q)  
5     MergeSort(A, q+1, r)  
6     Merge(A, p, q, r)
```

- Merge sort applies divide and conquer approach to the problem of sorting
- It is recursive algorithm with base case being $p = r$, which means an array of just one number.
- It should be obvious that an array with only one number is already sorted:)
- Pseudocode is quite simple, but the real magic happens inside *Merge* procedure

Merge sort graph



Merge sort call stack

```
1  MergeSort(A, 0, 7)
2  MergeSort(A, 0, 3)
3  MergeSort(A, 0, 1)
4  MergeSort(A, 0, 0)
5  MergeSort(A, 1, 1)
6  Merge(A, 0, 0, 1)
7  MergeSort(A, 2, 3)
8  MergeSort(A, 2, 2)
9  MergeSort(A, 3, 3)
10 Merge(A, 2, 2, 3)
11 Merge(A, 0, 1, 3)
12 MergeSort(A, 4, 7)
13 MergeSort(A, 4, 5)
14 MergeSort(A, 4, 4)
15 MergeSort(A, 5, 5)
16 Merge(A, 4, 4, 5)
17 MergeSort(A, 6, 7)
18 MergeSort(A, 6, 6)
19 MergeSort(A, 7, 7)
20 Merge(A, 6, 6, 7)
21 Merge(A, 4, 5, 7)
22 Merge(A, 0, 4, 7)
```


Merge procedure

- Real magic in merge sort happens in merge procedure
- Merge combines the two sorted sequences into one sorted sequence
- To do that it must create two temporary tables
- The procedure compares the first element in each of those arrays and pulls the smallest in to the output array
- If one array is empty, then procedure just copies all the elements from the second array

Merge procedure pseudocode

```

1  Merge(A, p, q, r):
2      n1 = q - p + 1;
3      n2 = r - q;
4      L[n1] = A[p..q]
5      R[n2] = A[q+1..r] //temp array
6      /* Merge the temp arrays back into A[p..r]*/
7      k = p;
8      while (i < n1 && j < n2)
9          if (L[i] <= R[j])
10             A[k] = L[i];
11             i++;
12         else
13             A[k] = R[j];
14             j++;
15         k++;
16     //Copy the remaining elements of L[], if there are any
17     while (i < n1)
18         A[k] = L[i]; i++; k++;
19     //Copy the remaining elements of R[], if there are any
20     while (j < n2)
21         A[k] = R[j]; j++; k++;

```

Merge sort analysis

- First of the advanced algorithms
- Complexity $O(n \lg n)$
- Complexity is the same in all cases - best, worst or average
- Stable sorting
- Out of place sorting, requires about $O(n)$ additional space
- Perfect, when the data can only be accessed sequentially (linked lists)
- We can use sentinel to skip checking whether a subarray is empty

Course plan

- 1 Merge sort
- 2 QuickSort**
- 3 Counting Sort
- 4 Radix Sort
- 5 Bucket Sort
- 6 Conclusions

Quicksort

- Developed in 1959 by Tony Hoare, while he was a visiting student in Moscow
- Similar to merge sort in that it also divides original array into two subarrays
- Array is divided based on a pivot element - all elements lesser than the pivot form one subsequence, all larger form second subsequence
- Choosing the right pivot has dramatic effect on the whole algorithm!

Pseudocode

- Input: Array A , start index p , end index r
- Output: Sorted array A

```
1  QuickSort( $A$ ,  $p$ ,  $r$ ):  
2    if  $p < r$   
3      pivot = Partition( $A$ ,  $p$ ,  $r$ )  
4      QuickSort( $A$ ,  $p$ , pivot)  
5      QuickSort( $A$ , pivot+1,  $r$ )
```

Partition procedure

Pseudocode

```

1  Partiton(A, p, r)
2      pivot = A[r]
3      i = p - 1
4      for( j = p; j < r - 1; j++ )
5          if A[j] <= x
6              i++
7              swap A[i] with A[j]
8      swap A[i+1] with A[r]
9      return i + 1

```

- i is the last index of array of smaller than pivot numbers
- j is the last index of array of larger than pivot numbers
- $A[j + 1 \dots r - 1]$ is the undivided yet part of array

Partition example

| | | | | | | | | |
|----------|----|----|----|----|----|----|---|----|
| $i = -1$ | 42 | 55 | 12 | 67 | 94 | 18 | 6 | 44 |
|----------|----|----|----|----|----|----|---|----|

|

Partition example

| | | | | | | | | |
|----------|----|----|----|----|----|----|---|----|
| $i = -1$ | 42 | 55 | 12 | 67 | 94 | 18 | 6 | 44 |
| $i = 0$ | 42 | 55 | 12 | 67 | 94 | 18 | 6 | 44 |

|

Partition example

| | | | | | | | | |
|----------|----|----|----|----|----|----|---|----|
| $i = -1$ | 42 | 55 | 12 | 67 | 94 | 18 | 6 | 44 |
| $i = 0$ | 42 | 55 | 12 | 67 | 94 | 18 | 6 | 44 |
| $i = 0$ | 42 | 55 | 12 | 67 | 94 | 18 | 6 | 44 |

|

Partition example

| | | | | | | | | |
|----------|----|----|----|----|----|----|---|----|
| $i = -1$ | 42 | 55 | 12 | 67 | 94 | 18 | 6 | 44 |
| $i = 0$ | 42 | 55 | 12 | 67 | 94 | 18 | 6 | 44 |
| $i = 0$ | 42 | 55 | 12 | 67 | 94 | 18 | 6 | 44 |
| $i = 1$ | 42 | 12 | 55 | 67 | 94 | 18 | 6 | 44 |

|

Partition example

| | | | | | | | | |
|----------|----|----|----|----|----|----|---|----|
| $i = -1$ | 42 | 55 | 12 | 67 | 94 | 18 | 6 | 44 |
| $i = 0$ | 42 | 55 | 12 | 67 | 94 | 18 | 6 | 44 |
| $i = 0$ | 42 | 55 | 12 | 67 | 94 | 18 | 6 | 44 |
| $i = 1$ | 42 | 12 | 55 | 67 | 94 | 18 | 6 | 44 |
| $i = 1$ | 42 | 12 | 55 | 67 | 94 | 18 | 6 | 44 |

|

Partition example

| | | | | | | | | |
|----------|----|----|----|----|----|----|---|----|
| $i = -1$ | 42 | 55 | 12 | 67 | 94 | 18 | 6 | 44 |
| $i = 0$ | 42 | 55 | 12 | 67 | 94 | 18 | 6 | 44 |
| $i = 0$ | 42 | 55 | 12 | 67 | 94 | 18 | 6 | 44 |
| $i = 1$ | 42 | 12 | 55 | 67 | 94 | 18 | 6 | 44 |
| $i = 1$ | 42 | 12 | 55 | 67 | 94 | 18 | 6 | 44 |
| $i = 1$ | 42 | 12 | 55 | 67 | 94 | 18 | 6 | 44 |

|

Partition example

| | | | | | | | | |
|----------|----|----|----|----|----|----|---|----|
| $i = -1$ | 42 | 55 | 12 | 67 | 94 | 18 | 6 | 44 |
| $i = 0$ | 42 | 55 | 12 | 67 | 94 | 18 | 6 | 44 |
| $i = 0$ | 42 | 55 | 12 | 67 | 94 | 18 | 6 | 44 |
| $i = 1$ | 42 | 12 | 55 | 67 | 94 | 18 | 6 | 44 |
| $i = 1$ | 42 | 12 | 55 | 67 | 94 | 18 | 6 | 44 |
| $i = 1$ | 42 | 12 | 55 | 67 | 94 | 18 | 6 | 44 |
| $i = 2$ | 42 | 12 | 18 | 67 | 94 | 55 | 6 | 44 |

|

Partition example

| | | | | | | | | |
|----------|----|----|----|----|----|----|----|----|
| $i = -1$ | 42 | 55 | 12 | 67 | 94 | 18 | 6 | 44 |
| $i = 0$ | 42 | 55 | 12 | 67 | 94 | 18 | 6 | 44 |
| $i = 0$ | 42 | 55 | 12 | 67 | 94 | 18 | 6 | 44 |
| $i = 1$ | 42 | 12 | 55 | 67 | 94 | 18 | 6 | 44 |
| $i = 1$ | 42 | 12 | 55 | 67 | 94 | 18 | 6 | 44 |
| $i = 1$ | 42 | 12 | 55 | 67 | 94 | 18 | 6 | 44 |
| $i = 2$ | 42 | 12 | 18 | 67 | 94 | 55 | 6 | 44 |
| $i = 3$ | 42 | 12 | 18 | 6 | 94 | 55 | 67 | 44 |

|

Partition example

| | | | | | | | | |
|----------|----|----|----|----|----|----|----|----|
| $i = -1$ | 42 | 55 | 12 | 67 | 94 | 18 | 6 | 44 |
| $i = 0$ | 42 | 55 | 12 | 67 | 94 | 18 | 6 | 44 |
| $i = 0$ | 42 | 55 | 12 | 67 | 94 | 18 | 6 | 44 |
| $i = 1$ | 42 | 12 | 55 | 67 | 94 | 18 | 6 | 44 |
| $i = 1$ | 42 | 12 | 55 | 67 | 94 | 18 | 6 | 44 |
| $i = 1$ | 42 | 12 | 55 | 67 | 94 | 18 | 6 | 44 |
| $i = 2$ | 42 | 12 | 18 | 67 | 94 | 55 | 6 | 44 |
| $i = 3$ | 42 | 12 | 18 | 6 | 94 | 55 | 67 | 44 |
| $i = 3$ | 42 | 12 | 18 | 6 | 44 | 55 | 67 | 94 |

Quicksort example

| | | | | | | | |
|----|----|----|----|----|----|----|----|
| 42 | 55 | 12 | 67 | 94 | 18 | 6 | 44 |
| 42 | 12 | 18 | 6 | 44 | 55 | 67 | 94 |
| 6 | 12 | 18 | 42 | 44 | 55 | 67 | 94 |
| 6 | 12 | 18 | 42 | 44 | 55 | 67 | 94 |
| 6 | 12 | 18 | 42 | 44 | 55 | 67 | 94 |
| 6 | 12 | 18 | 42 | 44 | 55 | 67 | 94 |
| 6 | 12 | 18 | 42 | 44 | 55 | 67 | 94 |
| 6 | 12 | 18 | 42 | 44 | 55 | 67 | 94 |

Pivot problem

- Depending on the chosen pivot we can get either balanced or imbalanced split:
 - Balanced split - both subarrays are roughly the same size
 - Imbalanced split - one subarray is much larger than the other
- Worst case scenario is one subarray of size zero
- Choosing the right pivot has dramatic impact on the complexity of the algorithm - in worst case scenario time complexity is $O(n^2)$, compared to $O(n \lg n)$ for average case

Choosing pivot

- First/last element of the array - low cost of computing, but is vulnerable to data distribution
- Random element from the array - quite good, but good random number generator is needed
- Median of three elements - good, also quite fast to compute
- Median of five or more elements - not much better than median-of-3
- BFPRT (Blum-Floyd-Pratt-Rivest-Trajan) - guarantees logarithmic worst case, while also making average case slightly worse. Sometimes called median of medians

Quicksort analysis

- Complexity $O(n \lg n)$ in average case
- In worst case complexity $O(n^2)$
- Worst case is quite rare, also can be mitigated by good pivot select strategy
- Unstable sort
- In-place sort
- Best algorithm for standard sorting, used widely in programming
- About two times faster than merge sort
- Uses tail recursion

Course plan

- 1 Merge sort
- 2 QuickSort
- 3 Counting Sort**
- 4 Radix Sort
- 5 Bucket Sort
- 6 Conclusions

Lower bound for sorting

- All the algorithms described to this point operate on comparing two numbers, to determine the sort order
- It can be shown, that any comparison sort algorithm requires $O(n \lg n)$ comparisons in the worst case.
- Asymptotically, merge and heapsort (will be introduced later) are optimal sorts (quicksort has $O(n^2)$ in the worst case)
- It is however possible to sort in linear time(!), but it requires some knowledge of the data we are about to sort

Counting sort

- Let's assume that we are given an array of numbers to sort. We also know that there are only k unique values among the number we are sorting
- We calculate count c_k how many each of the unique numbers appear in the input array
- Output array is constructed by just placing the k number c_k times, starting from the smallest k with non-zero c_k
- Notice, there are no comparisons here!

Counting sort

Pseudocode

```
1  CountingSort(A, k) :
2  int count[0..k] = 0
3  int output[]
4  for (int i=0; i < length(A); i++)
5      count[A[i]]++;
6  //count[x] now contains the number of elements equal to x
7  for (i=1; i<=k; i++)
8      count[i] += count[i-1]
9  // count[x] now contains the number of elements less than or
   equal to x
10 for (i=length(A) - 1; i >= 0; i--)
11     output[count[A[i]]-1] = A[i]
12     count[A[i]]--
13 // We do it in reverse order for stability
14 return output
```

Counting sort simple example

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| A | 2 | 5 | 3 | 0 | 2 | 3 | 0 | 3 |

| | | | | | | |
|---|---|---|---|---|---|---|
| | 0 | 1 | 2 | 3 | 4 | 5 |
| C | 2 | 0 | 2 | 3 | 0 | 1 |

(a)

| | | | | | | |
|---|---|---|---|---|---|---|
| | 0 | 1 | 2 | 3 | 4 | 5 |
| C | 2 | 2 | 4 | 7 | 7 | 8 |

(b)

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| B | | | | | | | 3 | |

| | | | | | | |
|---|---|---|---|---|---|---|
| | 0 | 1 | 2 | 3 | 4 | 5 |
| C | 2 | 2 | 4 | 6 | 7 | 8 |

(c)

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| B | | 0 | | | | | 3 | |

| | | | | | | |
|---|---|---|---|---|---|---|
| | 0 | 1 | 2 | 3 | 4 | 5 |
| C | 1 | 2 | 4 | 6 | 7 | 8 |

(d)

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| B | | 0 | | | | 3 | 3 | |

| | | | | | | |
|---|---|---|---|---|---|---|
| | 0 | 1 | 2 | 3 | 4 | 5 |
| C | 1 | 2 | 4 | 5 | 7 | 8 |

(e)

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| B | 0 | 0 | 2 | 2 | 3 | 3 | 3 | 5 |

(f)

Counting sort analysis

- Time complexity is $O(k + n)$.
- Usually k is of the order of n , we have linear time complexity!
- Does not use any comparison operators
- Stable sort
- Out of place sort - requires additional $O(n + k)$ space
- If $k \gg n$ then it is unpractical to use it
- Can only work on positive numbers - algorithm described before does not work if the numbers can be negative!

Counting sort with negative numbers

Pseudocode

```
1  CountingSortWithNegativeNumbers(A):
2    max = max_element(A);
3    min = min_element(A);
4    int range = max - min + 1;
5    int count[range];
6    int output[length(A)];
7    for(i = 0; i < length(A); i++)
8        count[A[i]-min]++;
9    for(i = 1; i < count.size(); i++)
10        count[i] += count[i-1];
11    for(i = length(A); i >= 0; i--)
12        output[ count[A[i]-min] - 1 ] = A[i];
13        count[A[i]-min]--;
14    return output
```

Course plan

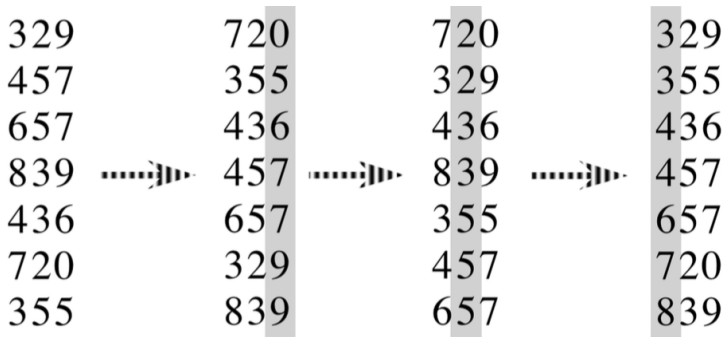
- 1 Merge sort
- 2 QuickSort
- 3 Counting Sort
- 4 Radix Sort**
- 5 Bucket Sort
- 6 Conclusions

Radix sort

- Counting sort is not usable when $k \gg n$.
- Still, we do not need to use comparison sorting algorithms!
- Assume we have to sort n numbers, each with at most d digits
- We can then use counting sort and sort digit by digit, starting with the least significant one!
- The counting sort needs to be stable!

```
1  RadixSort(A,d):  
2    for i=1; i<=d; i++  
3      CountingSort on digit d in A
```

Radix Sort example



Radix sort analysis

- Linear time, when input array is sufficiently larger than number of digits
- Out of place sorting
- Stable sort
- Radix sort can be faster than quicksort if the array is huge, or the keys small
- Still, radix sort is only usable on integer values
- Because of that two points it is often omitted in libraries

Course plan

- 1 Merge sort
- 2 QuickSort
- 3 Counting Sort
- 4 Radix Sort
- 5 Bucket Sort**
- 6 Conclusions

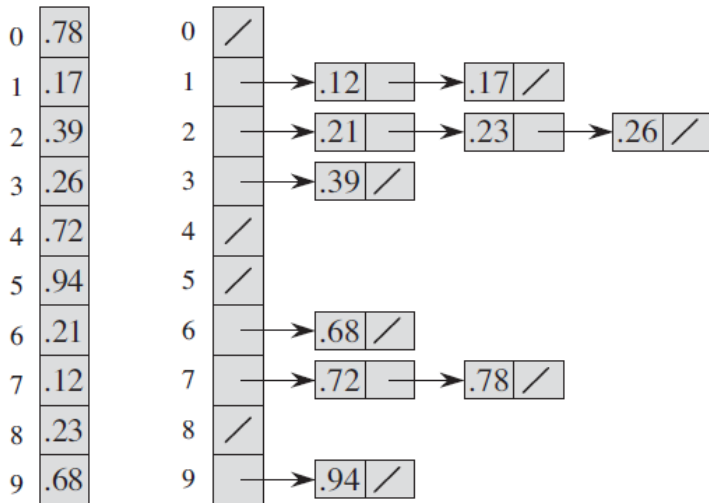
Bucket Sort

- Counting sort assumed that input is an array of small integers
- Bucket sort assumes that input data is generated by a random process, that distributes elements uniformly and independently over the interval $[0, 1)$
- Bucket sort divides the interval $[0, 1)$ into n equal-sized subintervals, called buckets, and then distributes the n input numbers into the buckets
- Since the inputs are uniformly and independently distributed over $[0, 1)$ we expect each bucket roughly the same size
- We sort the numbers in buckets using any other simple algorithm (like insertion sort), as each bucket should be fairly small

Bucket sort pseudocode

```
1  BucketSort(A):
2      n=length(A)
3      int [][] B[n-1] //This is array of arrays – our buckets
4      for (i=0; i < n; i++)
5          B[i] = empty array //initially every bucket is empty
6      for (i=1; i <= n; i++)
7          index = floor(nA[i])
8          B[index].push(A[i])
9      for (i=0; i<n; i++)
10         sort(B[i])
11     //array concatenation
12     output = B[0] + B[1] + ... B[n-1]
13
```

Bucket sort example



Input array

Buckets created for input array

Bucket sort analysis

- Complexity $O(n + \frac{n^2}{k} + k)$, where k is number of buckets
- Out of place sorting
- Stability of sorting depends on algorithm used to sort individual buckets
- Bucket sort is heavily dependent on the uniform distribution of data. In best case scenario it runs in linear time, but worst case is $O(n^2)$

Course plan

- 1 Merge sort
- 2 QuickSort
- 3 Counting Sort
- 4 Radix Sort
- 5 Bucket Sort
- 6 Conclusions**

Summary of sorting algorithms

| Name | Best | Average | Worst | Memory | Stable | In-place |
|----------------|--------------|--------------|--------------|------------|--------|----------|
| Insertion sort | $O(n)$ | $O(n^2)$ | $O(n^2)$ | $O(1)$ | Yes | Yes |
| Selection sort | $O(n^2)$ | $O(n^2)$ | $O(n^2)$ | $O(1)$ | No | Yes |
| Bubble sort | $O(n)$ | $O(n^2)$ | $O(n^2)$ | $O(1)$ | Yes | Yes |
| Cocktail sort | $O(n)$ | $O(n^2)$ | $O(n^2)$ | $O(1)$ | Yes | Yes |
| Merge sort | $O(n \lg n)$ | $O(n \lg n)$ | $O(n \lg n)$ | $O(n)$ | Yes | No |
| Quicksort | $O(n \lg n)$ | $O(n \lg n)$ | $O(n^2)$ | $O(\lg n)$ | No | No |
| Heap sort | $O(n \lg n)$ | $O(n \lg n)$ | $O(n \lg n)$ | $O(1)$ | No | Yes |
| Counting sort | $O(n + k)$ | $O(n + k)$ | $O(n + k)$ | $O(k)$ | Yes | No |
| Radix sort | $O(nk)$ | $O(nk)$ | $O(nk)$ | $O(n + k)$ | Yes | No |
| Bucket sort | $O(n + k)$ | $O(n + k)$ | $O(n^2)$ | $O(n)$ | Yes | No |

<https://www.youtube.com/watch?v=QOYcpGnHH0g>

Choosing the right algorithm

- Choose insertion sort if the input is small, or elements are roughly sorted, or you need very simple code
- Choose heapsort if worst-case scenario is your main concern
- Choose quicksort in general case

General remarks

- In typical applications quicksort behaves better than heapsort
- Simple methods are simple, but not very effective, they should be used only in very specific cases
- For practical use it is best to combine two methods - quicksort for large input (with monitoring and switching to heapsort in worst-case scenario) and insertion sort for smaller subproblems (threshold chosen empirically for a given platform).