

Part II. Sorting and order statistics

Chapter 6. Heapsort, the use of priority queue

Chapter 7. Quicksort, analysis

Chapter 8. Sorting in linear time, lower bounds

Chapter 9. Medians and order statistics

Chapter 6. Heapsort

heaps

LEFT, RIGHT, PARENT

heapsort

BUILD-MAX-HEAP(A)

MAX-HEAPIFY(A, i)

HEAPSORT(A)

heaps as priority queues

HEAP-MAXIMUM(A)

HEAP-EXTRACT-MAX(A)

HEAP-INCREASE-KEY(a, I, key)

MAX-HEAP-INSERT(A, key)

HEAPSORT(A)

1. BUILD-MAX-HEAP(A)
2. for $i \leftarrow \text{length}[A]$ down to 2
3. do exchange $A[1] \leftrightarrow A[i]$
4. $\text{heap-size}[A] \leftarrow \text{heap-size}[A] - 1$
5. MAX-HEAPIFY(A, 1)

analysis of time

```
BUILD-MAX-HEAP(A)
1. heap-size[A] <-- length[A]
2. for i <-- length[A]/2 downto 1
3.   do MAX-HEAPIFY(A, i)

MAX-HEAPIFY(A, i)
1.1 <-- LEFT[i]
2.r <-- RIGHT[i]
3.if l <= heap-size[A] and A[l] > A[i]
4.  then largest <-- l
5.  else largest <-- i
6.if r <= heap-size[A] and A[r] > A[i]
7.  then largest <-- r
8.if largest != i
9.  then exchange A[i] <-> A[largest]
10.   MAX-HEAPIFY(A, largest)
```

```
HEAP-MAXIMUM(A)
```

```
1. return A[1]
```

```
HEAP-EXTRACT-MAX(A)
```

```
1. if heap-size[A] < 1
```

```
2.   then error "heap underflow"
```

```
3. max <-- A[1]
```

```
4. A[1] <-- A[heap-size[A]]
```

```
5. heap-size[A] <-- heap-size[A] - 1
```

```
6. MAX-HEAPIFY(A, 1)
```

```
7. return max
```

```
HEAP-INCREASE-KEY(A, i, key)
```

1. if key < A[i]
2. then error "new key is smaller than current key"
3. A[i] <-- key
4. while i > 1 and A[PARENT[i]] < A[i]
5. do exchange A[i] <--> A[PARENT[i]]
6. i <-- PARENT[i]

```
MAX-HEAP-INSERT(A, key)
```

1. heap-size[A] <-- heap-size[A] + 1
2. A[heap-size[A]] <-- -oo
3. HEAP-INCREASE-KEY(A, heap-size[A], key)

Chapter 7. Quicksort

the basic idea of "quicksort"

the algorithm details

average case running time (probabilistic analysis)

randomized algorithm

variants

divide-and-conquer for quicksort

divide: partition list $A[p, r]$ into two sublists $A[p, q - 1]$ and $A[q + 1, r]$ such that

- (a) $A[i] \leq A[q]$ for all $i = p, \dots, q - 1$
- (b) $A[i] \geq A[q]$ for all $i = q + 1, \dots, r$

conquer: sort $A[p, q - 1]$ and $A[q + 1, r]$ recursively.

combine: no further work is needed.

QUICKSORT(A, p, r)

1. if $p < r$
2. then $q \leftarrow \text{PARTITION}(A, p, r)$
3. QUICKSORT($A, p, q-1$)
4. QUICKSORT($A, q+1, r$)

```
PARTITION(A, p, r)
```

```
1. x  $\leftarrow$  A[r]
```

```
2. i  $\leftarrow$  p-1
```

```
3. for j  $\leftarrow$  p to r-1
```

```
4.   do if A[j]  $\leq$  x
```

```
5.       then i  $\leftarrow$  i + 1
```

```
6.           exchange A[i]  $\leftrightarrow$  A[j]
```

```
7. exchange A[i+1]  $\leftrightarrow$  A[r]
```

```
8. return i+1
```

example: Figure 7.1 (page 147)

Is there an easier way (not necessarily the most efficient way) to do partition?

analysis of performance

worst case partitioning: $T(n) = T(n - 1) + \Theta(n)$

best case : $T(n) = T(\lfloor n/2 \rfloor) + T(\lceil n/2 \rceil - 1) + \Theta(n)$

balanced partitioning:

$T(n) = T(9n/10) + T(n/10) + \Theta(n)$

using recursive tree method, it can be shown that

$T(n) = O(n \log_{10} n)$ and

$T(n) = \Omega(n \log_{10/9} n)$.

Therefore, $T(n) = \theta(n \log n)$.

average case performance

probabilistic analysis: assume that lists to be sorted are random lists.

An intuition that PARTITION should give a balanced partition *with a high probability* on a *random list*:

(1) Consider a list of n elements in which every element has equal chance to be in position i , $i = 1, 2, \dots, n$. That is, for any position j and element $A[i]$,

$$\Pr(\text{rank}(A[i]) = j) = 1/n$$

(2) Then the probability for any chosen pivot to partition the list into two lists of sizes $n/10$ and $9n/10$ is 80%

(3) We can say that most (i.e., 80%) of the time QUICKSORT runs $O(n \log n)$.

Analysis Method 1: taking the average

The pivot could be at any of these positions $1, 2, \dots, n-1, n$. Then

$$T(n) = 1/n * [\sum_{j=2}^{n-1} (T(j-1) + T(n-j)) + 2T(n-1)] + cn$$

Consider two categories of partitions:

“good partition” ($n/4 < j \leq 3n/4$)

“bad partition” ($1 \leq j \leq n/4$ or $3n/4 < j \leq n$)

$T(n)$ is then the sum of two terms (exercise).

One can use the substitution or recursive tree method to prove
 $T(n) = O(n \log n)$.

Analysis Method 2: randomized quicksort

RANDOMIZED-PARTITION(A, p, r)

1. $i \leftarrow \text{random}(p, r)$
2. exchange $A[r] \leftrightarrow A[i]$
3. return PARTITION(A, p, r)

rename elements in A as

z_1, \dots, z_n with z_i being the i th smallest element, and
 $Z_{ij} = \{z_i, \dots, z_j\}$.

Define: $X_{ij} = 1$ iff z_i is compared to z_j

and $X = \sum_{i=1}^{n-1} \sum_{j=i+1}^n X_{ij}$ the number of comparisons.

$$\begin{aligned}
E[X] &= \sum_{i=1}^{n-1} \sum_{j=i+1}^n E[X_{ij}] \\
&= \sum_{i=1}^{n-1} \sum_{j=i+1}^n \Pr(z_i \text{ is compared to } z_j)
\end{aligned}$$

Event z_i is compared to z_j occurs when z_i is chosen as a pivot or z_j is chosen to be a pivot. The chance of each element to be chosen is equally likely.

Note: that z_i or z_j is chosen has nothing do any element outside of Z_{ij} being chosen. Therefore,

$$\begin{aligned}
\Pr(z_i \text{ is compared to } z_j) \\
&= 2/|Z_{ij}| = 2/(j - i + 1)
\end{aligned}$$

Then $E[X] = \sum_{i=1}^{n-1} \sum_{k=1}^{n-1} 2/(k + 1) = O(n \log_2 n)$

Chapter 8 Lower bounds and sorting in linear time

Deriving lower bounds

Comparison-based computational model

- (1) Prove that "finding the max" among n elements needs at least $(n - 1)$ comparisons.
- (2) Prove that "sorting n elements" needs $\Omega(n \log n)$ comparisons.

sorting in linear time

A lower bound for sorting:

decision tree – modeling algorithms/computations

each internal node denotes $(x_i \leq x_j)$, with two outcomes

each leaf denotes a possible output of the algorithm

Claim 1: total number of leaves is $n!$.

Claim 2: the height of the tree is $\Omega(\log n!)$.

Theorem: Sorting needs $\Omega(n \log n)$ comparisons on comparison-based computation models.

Prove: The longest path from the root to a leaf is $\Omega(\log n!)$. I.e., the number of comparisons needed in the worst case is $\Omega(\log(n!))$.

$$\begin{aligned} n! &= n(n-1)(n-n/2)(n-n/2-1)\cdots 2 \times 1 \geq (n/2)^{n/2} 2^{n/2-1} \\ &\geq (n)^{n/2}/2. \end{aligned}$$

or by Stirling's formula:

$$n! = \sqrt{2\pi n} (n/e)^n (1 + O(1/n))$$

$$\Omega(\log(n!)) = \Omega(n \log n)$$

COUNT SORT:

COUNTING-SORT(A, B, k)

1. for i \leftarrow 0 to k {k is the largest element}
2. do C[k] \leftarrow 0
3. for j = 1 to length[A]
4. do C[A[j]] \leftarrow C[A[j]] + 1
5. {C[i] contains the number of elements = i}
6. for i \leftarrow 1 to k
7. do C[i]=C[i] + C[i-1]
8. {C[i] contains the number of elements \leq i}
9. for j \leftarrow length[A] downto 1
10. do B[C[A[j]]] \leftarrow A[j]
11. C[A[j]] \leftarrow C[A[j]] - 1

example: A: 2 5 3 0 2 3 0 3 C: 2 0 2 3 0 1

analysis? $T(n) = O(k + n)$

RADIX SORT:

329	720	720	329
457	355	329	355
657	436	436	436
839	457	839	457
436	657	355	657
720	329	457	720
355	839	657	839

RADIX-SORT(A, d)

1. for $i \leftarrow 1$ to d
2. sort A on digit i

Lemma 8.4: Given n b -bit numbers and any positive $r \leq b$.

RADIX-SORT uses $\Theta((b/r)(n + 2^r))$ time.

Proof: each number is of $\lceil b/r \rceil$ digits of r bits each.

BUCKET SORT: assuming uniform distribution of inputs

BUCKET_SORT(A)

1. $n \leftarrow \text{length}[A]$
2. for $i \leftarrow 1$ to n
3. do insert $A[i]$ into list $B[nA[i]]$
4. for $i \leftarrow 0$ to $n-1$
5. do sort list $B[i]$ with insertion sort
6. concatenate the list $B[0], B[1], \dots, B[n-1]$

e.g.: A: .78 .17 .39 .26 .72 .94 .21 .12 .23 .68

B: 0 /
1 -> .12 -> .17
2 -> .21 -> .23 -> .26
3 -> .39
4 /
5 /
6 -> .68
7 -> .72 -> .78
8 /
9 -> .94

analysis? average time

Chapter 9. Medians and order statistics

The selection problem

Input: A set A of n (distinct) numbers and i , $1 \leq i \leq n$;

Output: $x \in A$, the i th smallest element in A .

Selection in expected linear time (but worst case $\Theta(n^2)$)

Selection in worst case linear time

Selection in expected linear time

RANDOMIZED-SELECT(A, p, r, i)

1. if p=r
2. then return A[p]
3. q ← RANDOMIZED-PARTITION(A, p, r)
4. k ← q - p + 1
5. if i = k
6. then return A[q]
7. else if i < k
8. then return RANDOMIZED-SELECT(A, q, q-1, i)
9. else return RANDOMIZED-SELECT(A, q+1, r, i-k)

analysis?

worst case running time $\Theta(n^2)$.

average case (expected time): $E[T(n)]$

Note: $Pr(x_k \text{ is the pivot}) = 1/n$

Define $X_k = 1$ iff $A[p..q]$ has exactly k elements.

Then $E[X_k] = 1/n$ and $X_k = 1$ for exactly one value of k .

$$T(n) \leq \sum_{k=1}^n X_k (T(\max\{k-1, n-k\}) + O(n))$$

$$E[T(n)] = \sum_{k=1}^n E[X_k \cdot T(\max\{k-1, n-k\})] + O(n)$$

$$= \sum_{k=1}^n 1/n \cdot E[T(\max\{k-1, n-k\})] + O(n)$$

$$\max\{k-1, n-k\} = k-1 \text{ if } k > n/2$$

$$\max\{k-1, n-k\} = n-k \text{ if } k \leq n/2$$

$$E[T(n)] \leq 2/n \sum_{k=n/2}^{n-1} E[T(k)] + O(n)$$

Solving the recurrence:

We have $E[T(n)] = O(n)$.

Selection in worst case linear time

deterministically finds a good partition:

- (1) divide n elements into $n/5$ groups of 5 elements
- (2) find the median of each group
- (3) recursively find the median x of medians
- (4) use x as the pivot

why this is a good partition?

Note: the number of elements $\leq x$ is at least:

$$3(1/2\lceil n/5 \rceil - 2) \geq 3n/10 - 6$$

similarly, the number of elements $\geq x$ is at least:

$$3(1/2\lceil n/5 \rceil - 2) \geq 3n/10 - 6$$

$$T(n) \leq T(\lceil n/5 \rceil) + T(\lceil 7n/10 + 6 \rceil) + O(n)$$

when $n \geq 140$