

7. Quicksort

→ Worst-case running time : $\Theta(n^2)$

→ Expected running time : $\Theta(n \lg n)$

→ Constants hidden in $\Theta(n \lg n)$ are small

→ Sorts in place & divide-conquer strategy.

7.1 : Description of quicksort

Quicksort is based on the three-step process of divide-&-conquer.

QUICKSORT(A, p, r)

- ① if $p < r$
 - ② then $q = \text{PARTITION}(A, p, r)$
 - ③ QUICKSORT($A, p, q-1$)
 - ④ QUICKSORT($A, q+1, r$)
- ⑤ Initial call is $\text{QUICKSORT}(A, 0, h)$.

Partitioning

PARTITION(A, p, r)

$x \doteq A[r]$

$i \doteq p-1$

for $j = p$ to $s-1$

do if $A[j] \leq x$

then $i = i + 1$

exchange $A[i] \leftrightarrow A[j]$

return $i + 1$

→ PARTITION always selects the last element $A[r]$ in the subarray

$A[p \dots r]$ as the pivot—the element around which to partition.

→ As the procedure executes, the array is partitioned into four regions, some of which may be empty:

Loop invariant:

- ① All entries in $A[i..j]$ are \leq pivot
- ② All entries in $A[i+1..j-1]$ are $>$ pivot

- ③ $A[n] = \text{pivot}$.

Correctness: Use the loop invariant to prove correctness of PARTITION:

Initialization: Before the loop starts, all the conditions of the loop invariant are satisfied, because n is the pivot & the subarrays $A[i..j-1]$ & $A[i+1..j-1]$ are empty.

Maintenance: While the loop is running, if $A[j] \leq \text{pivot}$, then $A[i..j]$ & $A[i+1..j-1]$ are swapped & then i & j are incremented. If $A[j] > \text{pivot}$, then increment only j .

Termination: When the loop terminates, $j=1$, so all elements in

A are partitioned into one of the three cases: $A[i..j] \leq$ pivot & $A[i+1..n-1] > \text{pivot}$ & $A[n] = \text{pivot}$.

Time for partitioning: $\Theta(n)$ to partition an n -element subarray.

12 Performance of quicksort:

The running time of quicksort depends on the partitioning of the subarrays:

→ If the subarrays are balanced, then quicksort can run as fast as mergesort.

→ If they are unbalanced, then quicksort can run as slowly as insertion sort.

* Worst-case

→ Occurs when the subarrays are completely unbalanced.

→ Have 0 elements in one subarray & $n-1$ elements in the other subarray.

→ Get the recurrence

$$T(n) = T(n-1) + T(0) + \Theta(n)$$

$$= T(n-1) + \Theta(n)$$

$$= \Theta(n^2)$$

→ Same running time as insertion sort.

→ In fact, the worst-case running time occurs when quicksort takes a sorted array as input. but insertion sort runs in $O(n)$ time in this case.

* Best-Case:

→ Occurs when the subarrays are completely balanced every time.

→ Each subarray has $\leq \frac{n}{2}$ elements.

→ Get the recurrence

$$\begin{aligned} T(n) &= 2T(\frac{n}{2}) + \Theta(n) \\ &= \Theta(n \lg n). \end{aligned}$$

* Balanced partitioning

→ Quicksort's average running time is much closer to the best case than to the worst case.

→ Imagine that PARTITION always produces a $1\text{-}1$ split.

→ Get the recurrence

$$T(n) \leq T(3n/10) + T(n/10) + O(n)$$

$$= O(n \lg n)$$

→ Intuition: Look at the recursion tree.

4.3: Randomized Version of quicksort

→ We have assumed that all IP permutations are equally likely.

→ This is not always true.

→ To correct this, we add randomization to quicksort.

→ Instead, we use random sampling, or picking one element at random.

→ Don't always use $A[1]$ as the pivot. Instead, randomly pick an element from the subarray that is being sorted.

RANDOMIZED-PARTITION(A, p, r)
 $i = \text{RANDOM}(p, r)$

exchange $A[i] \leftrightarrow A[1]$
return PARTITION($A[1:r]$).

Randomly selecting the pivot element will, on average, cause the split of the input array to be reasonably well balanced.

RANDOMIZED-QUICKSORT(A, p, r)

if $p > r$

then $q = \text{RANDOMIZED-PARTITION}(A[p:r])$

RANDOMIZED-QUICKSORT($A[p, q-1]$)

RANDOMIZED-QUICKSORT($A[q+1:r]$)

Exercise:

7.1.1: Using figure 7.1 as a model, illustrate the operation of PARTITION on the array $A = \langle 13, 19, 9, 5, 12, 8, 4, 2, 1, 2, 6, 11 \rangle$.

$A[1] \leftrightarrow A[5]$

$A = \langle 19, 13, 9, 5, 12, 8, 4, 2, 1, 2, 6, 11 \rangle$

$A[1] \leftrightarrow A[5]$

$\langle 13, 19, 9, 5, 12, 8, 4, 2, 1, 2, 6, 11 \rangle$

$\langle 13, 19, 9, 5, 12, 8, 4, 2, 1, 2, 6, 11 \rangle$

$\langle 13, 19, 9, 5, 12, 8, 4, 2, 1, 2, 6, 11 \rangle$

$\langle 13, 19, 9, 5, 12, 8, 4, 2, 1, 2, 6, 11 \rangle$

$\langle 13, 19, 9, 5, 12, 8, 4, 2, 1, 2, 6, 11 \rangle$

$\langle 13, 19, 9, 5, 12, 8, 4, 2, 1, 2, 6, 11 \rangle$

$\langle 13, 19, 9, 5, 12, 8, 4, 2, 1, 2, 6, 11 \rangle$

$\langle 13, 19, 9, 5, 12, 8, 4, 2, 1, 2, 6, 11 \rangle$

$\langle 13, 19, 9, 5, 12, 8, 4, 2, 1, 2, 6, 11 \rangle$

$\langle 13, 19, 9, 5, 12, 8, 4, 2, 1, 2, 6, 11 \rangle$

$\langle 13, 19, 9, 5, 12, 8, 4, 2, 1, 2, 6, 11 \rangle$

$\langle 13, 19, 9, 5, 12, 8, 4, 2, 1, 2, 6, 11 \rangle$

Randomization of quicksort stops any specific type of array from causing worst-case behavior. For example, an already sorted array cause worst-case behavior in non-randomized Quicksort, but not in RANDOMIZED-Quicksort.

7.1-2:

What value of p does PARTITION return when all elements in the array $A[p..n]$ have the same value? Modify PARTITION so that $q = l \lfloor (p+n)/2 \rfloor$ when all elements in the array $A[p..n]$ have the same value.

7.1-3:

We can modify PARTITION by counting the number of comparisons in which $A[i] = A[j] = A[q]$ & then subtracting that number from the pivot index.

7.1-3': Give a brief argument that the running time of PARTITION

on a subarray of size n' is $\Theta(n')$.

There is a `for` statement whose body executes $n-1-p = \Theta(n)$ times. In the worst-case every time the body of the `if` is executed, but it takes constant time & so does the code outside of the loop. Thus the running time is $\Theta(n)$.

7.1-4: How would you modify QICKSORT to sort into nonincreasing order? We only need to flip the condition on the `if` statement.

7.2: Performance of Quicksort

7.2-1: Use the substitution method to prove that the recurrence $T(n) = T(n-1) + \Theta(n)$ has the solution $T(n) = \Theta(n^2)$, as claimed at the beginning of section 7.2.

We represent $\Theta(n)$ as $C_1 n$ & we guess that $T(n) \leq C_1 n^2$

$$\begin{aligned} T(n) &\stackrel{\text{def}}{=} T(n-1) + C_1 n \\ &\leq C_1(n-1)^2 + C_1 n \\ &= C_1 n^2 - 2C_1 n + C_1 + C_1 n \\ &\leq C_1 n^2, \end{aligned}$$

7.2-2: What is the running time of quicksort when all elements of the array A have the same value?

It is $\Theta(n^2)$, since one of the partitions is always empty (exercise 7.1-2).

7.2-3: Show that the running time of BUCKSORT is $\Theta(n^2)$ when the array contains distinct elements & is sorted in decreasing order.

If the array is already sorted in decreasing order, then, the pivot element is less than all the other other elements. The partition step takes $\Theta(n)$ time, & then leaves you with a subproblem of size $n-1$ & a subproblem of size 0. This gives us the recurrence (considered in 7.2-1 - which we showed has a solution that is $\Theta(n^2)$).

7.2-4: Banks often record transactions on an account in ~~order~~ of the times of transactions, but many people like to receive their bank statements with checks listed in ~~order~~ by check numbers. People usually write checks in ~~order~~ by check number, & merchant usually cash the with reasonable dispatch. The problem of converting time-~~of~~ transactions ordering to check-number ordering is therefore the problem of sorting almost-sorted input. Argue that the procedure

INSERTION-SORT would tend to beat the procedure QUICKSORT on this problem.

The more sorted the array is, the less work insertion sort will do. Namely, INSERTION-SORT is $O(nd)$, where d is the number of inversions in the array. In the above example, the number of inversions tends to be small so insertion sort will be close to linear.

On the other hand, if PARTITION does pick pivot that does not participate in an inversion, it will produce an empty partition. Since there is a small number of inversions, QUICKSORT is very likely to produce empty partitions.

7.2-5: Suppose that the splits at every level of quicksort are in proportion $1-\alpha$ to α , where $0 < \alpha < \frac{1}{2}$ is a constant. Show that the minimum depth of a leaf in the recursion tree is approximately $-\lg n / \lg(1-\alpha)$ & the maximum depth is approximately $-\lg n / \lg(1-\alpha)$.

7.3: A randomized version of quicksort

7.3-1: Why do we analyze the expected running time of a randomized algo. & not its worst-case running time?

We analyze the expected run time because it represents the more typical time cost. Also, we are doing the expected run time even the possible randomness underlying computation because it can't be produced adversarially, unlike when doing expected run time over all possible inputs to the algorithms.

7.3-2: When RANDOMIZED-QUICKSORT runs, how many calls are made to the random number generator RANDOM in the worst case? How about in the best case? Give your answer in terms of Θ -notation.

In the worst case, the number of calls to RANDOM is

$$T(n) = T(n-1) + 1 = n = \Theta(n)$$

As for the best case,

$$T(n) = 2T(n/2) + 1 = \Theta(n)$$

This is not too surprising, because each third element (at least) gets picked as pivot.

7.4 Analysis of quicksort

7.4-1: Show that in the recurrence

$$T(n) = \max_{0 \leq q \leq n-1} (T(q) + T(n-q-1)) + \Theta(n),$$

$$T(n) = \Omega(n^2).$$

$$\text{We guess } T(n) \geq cn^2 - 2n,$$

$$(T(n) = \max_{0 \leq q \leq n-1} (T(q) + T(n-q-1)) + \Theta(n))$$

$$\begin{aligned} &\geq \max_{0 \leq q \leq n-1} ((cq^2 - 2q + c(n-q-1)^2) - 2n - 2q - 1) \\ &\geq c \max_{0 \leq q \leq n-1} ((q^2 + (n-q-1)^2 - 2n + 2q - 1) + \Theta(n)) \\ &\geq cn^2 - c(2n-1) + \Theta(n) \\ &\geq cn^2 - 2cn + 2c \\ &\geq cn^2 - 2n. \end{aligned}$$

$$[c \leq 1]$$

7.4-2: Show that quicksort best-case running time is $\Omega(n \lg n)$.

We'll use the substitution method to show that the best-case running time is $\Omega(n \lg n)$. Let $T(n)$ be the best-case time for the procedure QUICKSORT on an input of size n . We have:

f.4-3: Show that the expression $\frac{1}{2}n^2 + \frac{1}{2}n - \frac{1}{2}$ is a maximum over $\varphi = 0, 1, \dots, n-1$ when $\varphi = 0$ or $\varphi = n-1$.

$$T(n) = \min_{1 \leq q \leq n-1} (T(q) + T(n-q-1)) + \Theta(n)$$

Suppose that $T(n) \geq c(n \lg n + 2n)$ for some constant c . Substituting this guess into the recurrence gives

$$f(\varphi) = 2\varphi^2 + (n-\varphi-1)^2$$

$$f''(\varphi) = 4.$$

$$\begin{aligned} T(n) &\geq \min_{1 \leq q \leq n-1} (c\varphi \lg \varphi + 2c\varphi + c(n-q-1) \lg(n-q-1) + \\ &\quad 2c(n-q-1)) + \Theta(n). \\ &= (cn/2) \lg(n/2) + cn + c(n/2-1) \lg(n/2-1) \\ &\quad + cn - 2c + \Theta(n) \\ &\geq (cn/2) \lg n - cn/2 + c(n/2-1)(\lg n - 2) + 2cn \\ &\quad - 2c\Theta(n) \\ &= (cn/2) \lg n - cn/2 + (cn/2) \lg n - cn - c \lg n + 2c \\ &\quad + 2cn - 2c\Theta(n) \\ &= cn \lg n + cn/2 - c \lg n + 2c - 2c\Theta(n). \end{aligned}$$

$f'(q) = 0$ when $q = \frac{1}{2}n - \frac{1}{2}$. $f'(q)$ is also continuous. At q : $f''(q) > 0$, which means that $f'(q)$ is negative left of $f'(q) = 0$ & positive right of it, which means that this is a local minimum. In this case, $f(q)$ is decreasing in the beginning of the interval & increasing in the end, which means that those two points are the only candidates for a maximum in the interval.

$$f(0) = (n-1)^2$$

$$f(n-1) = (n-1)^2 + \Theta(n)$$

7.4-4: Show that RANDOMIZED-QUICKSORT's expected running time is $\Omega(n \lg n)$.

We use the same reasoning for the expected number of comparisons, we just take in a different direction.

$$\begin{aligned} E[X] &= \sum_{i=1}^{n-1} \frac{n}{\sum_{j=i+1}^n j - i + 1} \cdot 2 \\ &= \sum_{k=1}^{n-1} \frac{n-i}{k+1} \cdot \frac{2}{n-i} \quad (k \geq 1) \end{aligned}$$

$$E[X] \geq \sum_{i=1}^{n-1} \sum_{k=1}^{n-i} \frac{2^k}{2^k} \cdot \Omega(n \lg n)$$

$$\geq \sum_{i=1}^{n-1} \Omega(n \lg n)$$

$$E[X] = \Omega(n \lg n)$$

Using the master method, we get the solution $\Theta(n \lg n)$.

4.4-5: We can improve the running time of quicksort in practice by taking advantage of the fast running time of insertion sort when its $\Theta(n)$ "nearly" sorted. Upon calling quicksort on a subarray with fewer than k elements, let it simply return without sorting the subarray. After the top-level call to quicksort returns, sum insertion sort on the entire array to finish the sorting process. Argue that this sorting algorithm runs in $\mathcal{O}(nk + n \lg(n/k))$ expected time. How

should we pick k , both in theory & practice?

In the quicksort part of the proposed algorithm, the recursion stops at level $\lg(n/k)$, which makes the expected running time $\mathcal{O}(n \lg(n/k))$. However, this leaves n/k non-sorted, non-interacting subarrays of maximum length k .

Because the nature of the insertion sort algorithm, it will first sort fully one such subarray before considering the next one. Thus, it has the same complexity as sorting each of these arrays, that is $\frac{n}{k} \mathcal{O}(k^2) = \mathcal{O}(nk)$.

In theory, if we ignore the constant factors, we need to solve

$$n \lg n \geq nk + n \lg(n/k)$$

$$\Rightarrow \lg n \geq k + \lg n - \lg k$$

$$\Rightarrow \lg k \leq k.$$

Which is not possible.
If we add the constant factors, we get,

$$c_1 n \lg n \geq c_2 nk + c_3 n \lg(n/k)$$

$$\Rightarrow c_1 \lg n \geq c_2 k + c_3 \lg n - c_3 \lg k$$

$$\Rightarrow \lg k \geq \frac{c_1}{c_3} k.$$

which indicates that there might be a good candidate. Furthermore, the lower-order terms should be taken into consideration too.

In practice, k should be chosen by experiments.

* 4.4-6: Consider modifying the PARTITION procedure by

randomly picking three elements from array A & partitioning about their median (the middle value of the three elements). Approximate the probability of getting at least an α -to $(1-\alpha)$ split, as a function of α in the range $0 < \alpha < 1$.

First, for simplicity's sake, let's assume that we can pick the same element twice. Let's also assume that $0 < \alpha \leq \frac{1}{2}$. In order to get such a split, two out of three elements need to be in the smallest αn elements. The probability of having one is $\alpha^n = \alpha$. The probability of having exactly two is $\alpha^2(1-\alpha) = \alpha$. There are three ways in which all three can be in the smallest αn so the probability of getting such a median is $3\alpha^2 - 2\alpha^3$. We will get the same split if the median is in the largest αn . Since the two events are mutually exclusive, the probability is

$$P_A(\text{OK split}) = (\alpha^2 - 4\alpha^3) = 2\alpha^2(1-2\alpha)$$

* Proof: $\rightarrow l \geq n!$

→ By lemma: $n! \leq l \leq 2^n$ and $2^n \geq n!$

Take log both sides $h \geq \lg(n!)$

Use Stirling approximation: $n! > (n/e)^n$

$$h \geq \lg(n/e)^n$$

$$\begin{aligned} h &\geq n \lg(n/e) \\ &= n \lg n - n \lg e \\ &= \Omega(n \lg n). \end{aligned}$$

* Proof by induction on h .

Basis: $h=0$. Tree is just one node, which is a leaf. $2^h=1$

Inductive step: Assume true for height = $h-1$.

Extend tree of height $h-1$ by making as many new leaves as possible. Each leaf becomes parent to two new leaves.

of leaves for height $h=2$. (# of leaves for height $h-1$)

8: Sorting in Linear Time

Theorem: Any decision tree that sorts n elements for height $\Omega(n \lg n)$.

$$\# \text{leaves} = 2^h$$

$$\rightarrow 2^h \geq n! \quad \rightarrow \text{By lemma: } n! \leq 2^h \text{ and } 2^h \geq n!$$