

CS 228: Introduction to Data Structures

Lecture 14

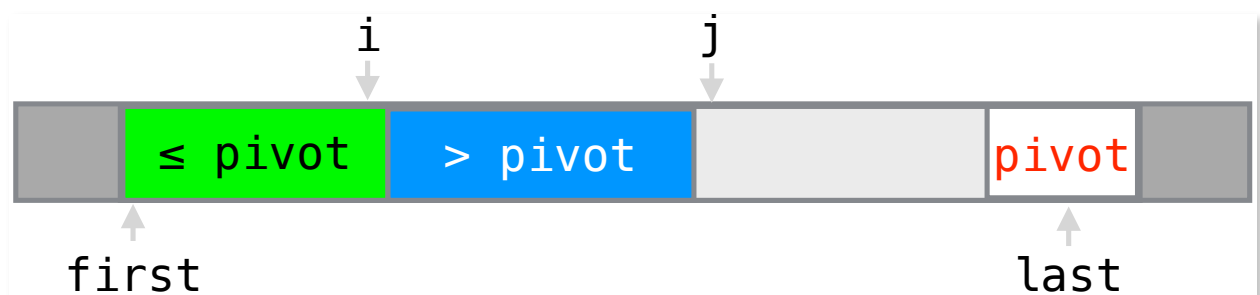
Friday, September 23, 2016

Quicksort (Continued)

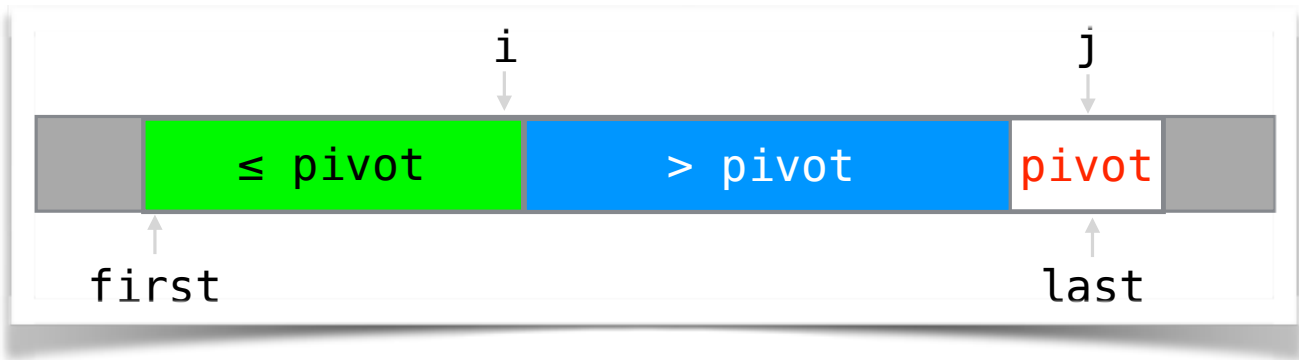
PARTITION maintains the following loop invariant.

After each iteration of the **for** loop,

- if $\text{first} \leq k \leq i$, then $\text{arr}[k] \leq \text{pivot}$,
- if $i+1 \leq k \leq j-1$, then $\text{arr}[k] > \text{pivot}$, and
- if $k = \text{last}$, then $\text{arr}[k] = \text{pivot}$.



At termination, $j == \text{last}$, so the array looks like this:



Therefore, one swap between $\text{arr}[i+1]$ and $\text{arr}[\text{last}]$ gives us the desired postcondition.

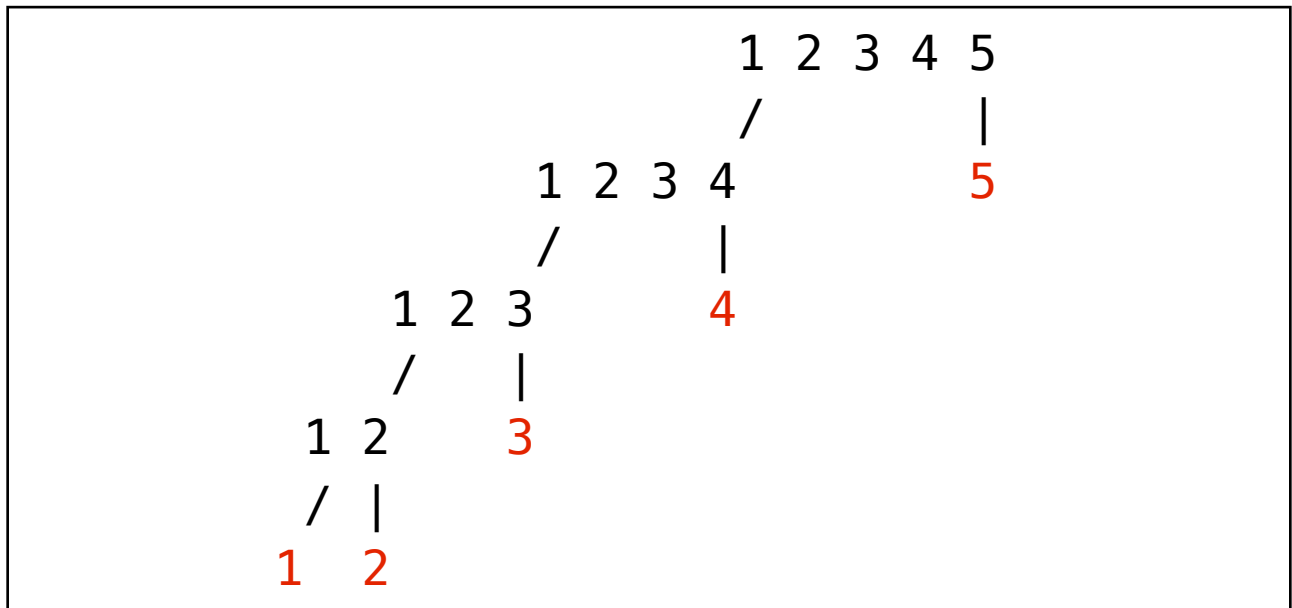
Time complexity of PARTITION. There are $n - 1$ iterations, where $n = \text{last} - \text{first} + 1$. Since each iteration takes $O(1)$ time, PARTITION takes $O(n)$ time.

Time Complexity of Quicksort

The running time of QS depends heavily on the sequence of pivots.

Best case: This is achieved when the pivot always splits the array in half. Then, the recursion tree has height $O(\log n)$ with $O(n)$ work per level. The total time is $O(n \log n)$.

Worst case: This is achieved when one of the subarrays to the side of pivot is empty. Then, the recursion tree has height n , with $O(n)$ work per level, for $O(n^2)$ total. With our implementation of partition, this case occurs when the array is sorted:



So, if the array is sorted or nearly sorted, you're better off using insertion sort.

Expected case: A more “typical” case is if array is randomly ordered. Then, odds are good that the pivot will be close to the middle, so each side will contain at most some constant fraction of the elements. Even if the fraction is, say, $0.99n$, the running time is $O(n \log n)$.

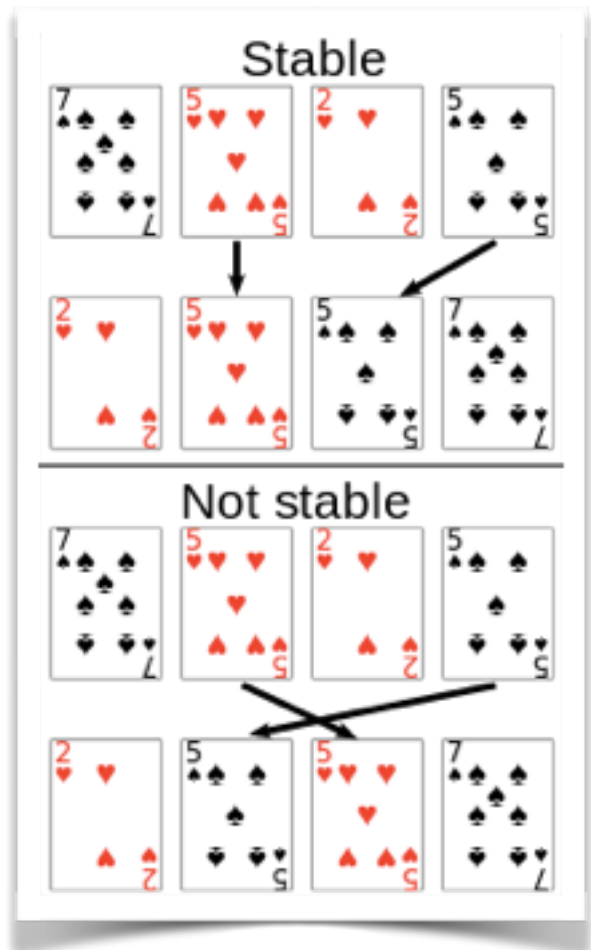
Note. The partition strategy we have studied performs poorly when the data has many duplicates. In particular, if we use it in an array where all elements are *identical*, quicksort will achieve its worst-case running time of $\approx cn^2$, for some constant c . The reason is that the sub-array to the right of the pivot will always be empty (since there are no elements greater than the pivot), while the sub-array to the left will contain $n-1$ elements. There are variants of partition that handle arrays with duplicates better than the version we have studied, but we will not cover those variants here.

Pivot Selection in Practice

Using a randomly selected pivot is equivalent to having randomly ordered data. Thus, we can achieve $O(n \log n)$ expected time via random pivot selection. Random pivot selection is rather impractical, though. A more common approach is to extract a small sample from the array and use its median as the pivot. For instance, in ***median-of-three partitioning*** picks three array elements (e.g., first, middle, and last) and uses the median of the three as the pivot.

Stability

The things we have to sort are usually more complex than plain numbers. This means that there might be multiple different correctly sorted versions of the original input. For instance, consider the card sorting example on the right¹. If we sort the cards by rank, there are two possible orderings, depending on how we arrange the two 5 cards. The first ordering preserves the relative ordering between these two cards (hearts before spades); the second one reverses it.



As another example, suppose we have a list of records, each of which consists of the name of a person and that person's age. Suppose the records are arranged in alphabetical order by name; e.g.,

(Alice, 18), (Chip, 14), (Dan, 14), (Ellie, 18)

¹From http://en.wikipedia.org/wiki/Sorting_algorithm

There are four valid ways to sort by age; two of these are

(Chip, 14), (Dan, 14), (Alice, 18), (Ellie, 18)

and

(Dan, 14), (Chip, 14), (Ellie, 18), (Alice, 18)

The first of these preserves the alphabetical order between Alice and Ellie and between Chip and Dan; the second one does not.

More formally, suppose we are sorting a collection of objects by a certain **key**. For instance, in the card sorting example, the key was the rank; in the person/age example, the key was the age. A sorting algorithm is **stable** if whenever there are two records R and S with the same key, and R appears before S in the original list, then R will always appear before S in the sorted list.

Stability can be an important consideration in practice, since data often has a certain underlying structure that we might want to preserve (e.g., alphabetical order). Not all sorting algorithms are stable, although we can always make them stable with some additional work.

Exercise. Which of the sorting algorithms that we have seen are stable? For any algorithm that is not stable, can you show how you can convert it to a stable one?