

Divide and Conquer: Counting Inversions¹

1 Counting Inversions

We now look at a closely related problem to merge-sort. Given an array $A[1 : n]$, the pair (i, j) for $1 \leq i < j \leq n$ is called an *inversion* if $A[i] > A[j]$. For example, in the array $[10, 20, 30, 50, 40]$, the pair $(4, 5)$ is an inversion.

COUNTING INVERSION

Input: An array $A[1 : n]$

Output: The number of inversions in A .

Size: n , the size of the array.

There is a naive $O(n^2)$ time algorithm: go over all pairs and check if they form an inversion or not. We now apply the divide-and-conquer paradigm to do better.

If $n = 1$, then the number of inversions is 0. Otherwise, suppose we divide the array into two: $A[1 : n/2]$ and $A[n/2 + 1 : n]$. Recursively, suppose we have computed the number of inversions in $A[1 : n/2]$ and $A[n/2 + 1 : n]$. Let these be I_1 and I_2 , respectively. Note that any inversion (i, j) in $A[1 : n]$ satisfies

- either $i < j \leq n/2$, which implies (i, j) is an inversion in $A[1 : n/2]$, or
- $n/2 + 1 \leq i < j$, which implies (i, j) is an inversion in $A[n/2 + 1 : n]$, or
- $i \leq n/2 < j$, and these are the extra inversions over $I_1 + I_2$ that we need to count.

Let's call any (i, j) of type (c) above a *cross* inversion, and let C denote this number. Then by what we said above, we need to return $I_1 + I_2 + C$. To obtain a “win”, we need to see if we can calculate C “much faster” than $O(n^2)$ time. How do we do that?

After you think about it for a while, there may not seem to be any easy way to calculate C faster than $O(n^2)$. Indeed, there could be $\Theta(n^2)$ inversions in $A[1 : n]$ and so shouldn't it take that much time to count them? How do we get around this? There are two crucial observations that help here.

- The number of cross-inversions between $A[1 : n/2]$ and $A[n/2 + 1 : n]$ is the same as between $\text{sort}(A[1 : n/2])$ and $\text{sort}(A[n/2 + 1 : n])$, where $\text{sort}(P)$ is the sorted order of an array P .
- If $A[1 : n/2]$ and $A[n/2 + 1 : n]$ were sorted, then the cross-inversions can be calculated in $O(n)$ time. This may not be immediate, but if you understand the COMBINE subroutine above, then it should ring a bell. We elaborate it on it below.

Cross-Inversions between Sorted Arrays. Given two sorted arrays $P[1 : p]$ and $Q[1 : q]$, we can count the number of cross-inversion pairs (i, j) such that $P[i] > Q[j]$ in $O(n)$ time as follows. As in COMBINE we start off with two pointers i, j initialized to 1. We also store a counter num initialized to 0 which, at the end, is supposed to contain the answer C . We check if $P[i] > Q[j]$ or not. If it isn't, that is if $P[i] \leq Q[j]$,

¹Lecture notes by Deeparnab Chakrabarty. Last modified : 19th Mar, 2022

These have not gone through scrutiny and may contain errors. If you find any, or have any other comments, please email me at deeparnab@dartmouth.edu. Highly appreciated!

then (i, j) is not a cross-inversion and we simply increment $i = i + 1$. Otherwise, if $P[i] > Q[j]$, then we increment $\text{num} = \text{num} + (p - i + 1)$ and $j = j + 1$. Why do you increment by so much? Didn't you find (i, j) is an inversion and so you should increment by only +1? Well, not only is (i, j) a cross-inversion, so are $(i + 1, j)$, $(i + 2, j)$, and so on till (p, j) . This is crucially using the fact that P is sorted. By doing a single comparison, because of sortedness, we discover a “bunch” of inversions. This gives us the “win” we were looking for. The claim below explains it more formally. We stop when either $i = p + 1$ or $j = q + 1$.

Claim 1. At any stage, suppose the algorithm encounters $P[i] > Q[j]$. Then $\{(i', j) : i' \geq i\}$ are the *only* cross-inversions which involve j .

Proof. Since P is sorted, $P[i'] > Q[j]$ for all $i' \geq i$ and so all such (i', j) are all valid cross-inversions. Now consider any $i'' < i$. Since in the algorithm the pointer is at $i > i''$, at some previous stage the algorithm compared $P[i'']$ and $Q[j'']$ with $j'' \leq j$, and found $P[i''] \leq Q[j'']$. But since Q is sorted, this would imply $P[i''] \leq Q[j]$. This implies for all $i'' < i$, (i'', j) is not an inversion. \square

```

1: procedure COUNTCROSSINV( $P[1 : p], Q[1 : q]$ ):
2:    $\triangleright$   $P$  and  $Q$  are sorted; outputs the number of  $(i, j)$  with  $P[i] > Q[j]$ .
3:    $i \leftarrow 1; j \leftarrow 1; \text{num} \leftarrow 0$ .
4:   while  $i < p + 1$  and  $j < q + 1$  do:
5:     if  $(P[i] > Q[j])$  then:
6:        $\text{num} \leftarrow \text{num} + (p - i + 1)$ 
7:        $j \leftarrow j + 1$ 
8:     else:
9:        $i \leftarrow i + 1$ 

```

Theorem 1. COUNTCROSSINV counts the number of cross inversions between P and Q in time $O(p + q)$.

Now we are armed to describe the divide-and-conquer algorithm for counting inversions.

```

1: procedure COUNTINV1( $A[1 : n]$ ):
2:    $\triangleright$  Counts the number of inversions in  $A[1 : n]$ 
3:   if  $n = 1$  then:
4:     return 0.  $\triangleright$  Singleton Array
5:    $m \leftarrow \lfloor n/2 \rfloor$ 
6:    $I_1 \leftarrow \text{COUNTINV1}(A[1 : m])$ 
7:    $I_2 \leftarrow \text{COUNTINV1}(A[m + 1 : n])$ 
8:    $B_1 \leftarrow \text{MERGESORT}(A[1 : m])$ 
9:    $B_2 \leftarrow \text{MERGESORT}(A[m + 1 : n])$ 
10:   $C \leftarrow \text{COUNTCROSSINV}(B_1, B_2)$ 
11:  return  $I_1 + I_2 + C$ .

```

Let's analyze the time complexity. As always, let $T(n)$ be the worst case running time of COUNTINV1 on an array of length n . Let $A[1 : n]$ be the array attaining this time, and let's see the run of the algorithm

on this array. The time taken by [Line 6](#) and [Line 7](#) are $T(\lfloor n/2 \rfloor)$ and $T(\lceil n/2 \rceil)$ respectively. The time taken by [Line 10](#) takes $O(n)$ time by what we described above. Furthermore, the [Line 8](#) and [Line 9](#) takes $O(n \log n)$ time. Together, we get the following recurrence

$$T(n) \leq T(\lfloor n/2 \rfloor) + T(\lceil n/2 \rceil) + O(n \log n)$$

The above ‘almost’ looks like the merge-sort recurrence, and indeed the recurrence solves² to $T(n) = O(n \log^2 n)$.

But there is something wasteful about the above algorithm. In particular, if you run it on a small example by hand you will see that you are sorting a lot. And often the same sub-arrays. Whenever you see this, often you can exploit this observation and get a faster algorithm?

Let us use this opportunity to introduce a new idea in the divide-and-conquer paradigm: *get more by asking for more*. This “asking for more” technique is something you may have seen while proving statements by induction where you can prove something you want by actually asking to prove something stronger by induction. In this problem, we ask our algorithm to do more: given an array $A[1 : n]$ it has to count the inversions *and* also has to sort the array too. Now note that in this case [Line 8](#) and [Line 9](#) are not needed any more; this is returned by the new stronger algorithm. We however need to also return the sorted array : but this is what COMBINE precisely does³. So the final algorithm for counting inversions is below.

```

1: procedure SORT-AND-COUNT( $A[1 : n]$ ):
2:   ▷ Returns  $(B, I)$  where  $B = \text{sort}(A)$  and  $I$  is the number of inversions in  $A[1 : n]$ 
3:   if  $n = 1$  then:
4:     return  $(A, 0)$ . ▷ Singleton Array
5:    $m \leftarrow \lfloor n/2 \rfloor$ 
6:    $(B_1, I_1) \leftarrow \text{SORT-AND-COUNT}(A[1 : m])$ 
7:    $(B_2, I_2) \leftarrow \text{SORT-AND-COUNT}(A[m + 1 : n])$ 
8:    $C \leftarrow \text{COUNTCROSSINV}(B_1, B_2)$ 
9:    $B \leftarrow \text{COMBINE}(B_1, B_2)$ 
10:  return  $(B, I_1 + I_2 + C)$ 

```

Now we see that the recurrence for the running time of SORT-AND-COUNT is precisely

$$T(n) \leq T(\lfloor n/2 \rfloor) + T(\lceil n/2 \rceil) + O(n)$$

Theorem 2. SORT-AND-COUNT returns the number of inversions of an array $A[1 : n]$ in $O(n \log n)$ time.

As an application, you are now ready to solve Problem 1 of PSet 0. Go ahead and try it again!

²Warning: Master Theorem doesn’t apply. But if you go back to the kitty method proof, you should be able to recreate the $T(n) = O(n \log^2 n)$. Indeed, every “round” one puts in $\leq Cn \log n$ instead of $\leq Cn$, and possibly even smaller.

³Actually, the COUNTCROSSINV code looks so much like COMBINE, you should not really have two lines ([Line 8](#) and [Line 9](#) in SORT-AND-COUNT), but wrap both of these into a single subroutine. I have separate lines for conceptual clarity at the expense of running time inefficiency (but not in a way that the big-Oh picture is muddled). In your coding assignment, hopefully you will keep this in mind.