

Contents

- 1 Integer Multiplication 2**
 - 1.1 Recursive Integer Multiplication 2
 - 1.2 Karatsuba Algorithm for Integer Multiplication 3

- 2 $Rank_A(x)$ 4**
 - 2.1 Computing the Rank 4
 - 2.2 Computing Rank in a sorted array 5
 - 2.3 Computing Rank of two elements at once 5
 - 2.4 Computing Rank of n elements at once 5

- 3 Merging two arrays 6**
 - 3.1 Merging by finding ranks 7

- 4 Merge Sort 7**
 - 4.1 Runtime of MERGESORT 8

- 5 Recurrence 9**
 - 5.1 Solving Recurrence Relations 9
 - 5.2 Analysis Of the Karatsuba Algorithm 10
 - 5.3 Analysis Of the Naive Divide & Conquer Multiplication Algorithm 11
 - 5.4 Substitution method for Solving Recurrences 11
 - 5.4.1 Using the substitution method for solving the naive Divide and Conquer multiplication problem 12
 - 5.4.2 Using substitution to analyze Karatsuba Algorithm 13
 - 5.5 Master Theorem 14
 - 5.5.1 Analyzing running times using the Master Theorem 14

- 6 Divide and Conquer Algorithm Design Paradigm 14**

- 7 Counting Inversions 15**
 - 7.1 Collaborative Filtering 18

- 8 Closest Pair 19**
 - 8.1 Naive Approach for Closest Pair 20
 - 8.2 Divide and conquer approach for Closest Pair 20

1 Integer Multiplication

In the first lecture we discussed that we could add or multiply arbitrarily large integers by putting all the digits of the number in an array and then add/multiply digit by digit. So we know how to multiply two integers using the grade school algorithm but we haven't yet discussed if there's a more efficient way to multiply integers. Since we've been using the divide and conquer approach for different problems which had a brute-force running time of n^2 , lets see if we can use a similar idea for multiplication.

1.1 Recursive Integer Multiplication

Let x and y be two $2n$ digit numbers then

$$x = \sum_{i=0}^{2n-1} x_i 10^i, \quad y = \sum_{i=0}^{2n-1} y_i 10^i$$

And we know the following axiom for Real Numbers

$$(a + b)(c + d) = ac + ad + bc + bd$$

Using these two ideas we can come up with an approach for *dividing* the multiplication problem. In particular we can divide x and y into two n digit numbers a, b and c, d respectively. Where

$$a = \sum_n^{2n-1} x_i 10^{i-n} \quad \text{(The left half of } x \text{)}$$

$$b = \sum_0^{n-1} x_i 10^i \quad \text{(The right half of } x \text{)}$$

$$c = \sum_n^{2n-1} y_i 10^{i-n} \quad \text{(The left half of } y \text{)}$$

$$d = \sum_0^{n-1} y_i 10^i \quad \text{(The right half of } y \text{)}$$

Then

$$\begin{aligned}
 xy &= (10^n a + b)(10^n c + d) \\
 &= 10^{2n} \underbrace{(ac)}_{1 \text{ multiplication}} + 10^n \underbrace{(ad + bc)}_{2 \text{ multiplications}} + \underbrace{bd}_{1 \text{ multiplication}}
 \end{aligned}$$

Example : If $x = 2731$ and $y = 1593$, then

$$\begin{array}{ll}
 x = 27 \times 10^2 + 31 & y = 15 \times 10^2 + 93 \\
 a = 27 & b = 31 \\
 c = 15 & d = 93
 \end{array}$$

Giving us:

$$xy = 10^2(27 \times 15) + 10^2(27 \times 93 + 31 \times 15) + 31 \times 93$$

Now each of the numbers a, b, c, d have n digits. The multiplications by powers of 10 is just a shift operation so we won't count those. Leaving that we have four n digit numbers to multiply and then we perform three addition operations. So our recurrence relation looks like

$$T_1(n) = \begin{cases} 1 & n = 1 \\ 4T_1(\frac{n}{2}) + 3n & n > 1 \end{cases}$$

We will later see, when we solve this recurrence that we get a running time that scales quadratically with n .

1.2 Karatsuba Algorithm for Integer Multiplication

But we already have an n^2 algorithm which is much simpler than this. So it would seem that our approach for dividing the problem doesn't really help. However, the four multiplications we had to perform could be reduced to three by using

$$bc + ad = (a + b)(c + d) - ac - bd$$

Taking from our previous example, we can see that this is true:

$$\begin{aligned}
 31 \times 15 + 27 \times 93 &= (27 + 31)(15 + 93) - 27 \times 15 - 31 \times 93 \\
 465 + 2511 &= (58)(108) - 405 - 2883 \\
 2976 &= 2976
 \end{aligned}$$

The Karatsuba algorithm for multiplying two integers uses this improvement. Reducing the number of multiplications by just 1 doesn't seem like much, but as we'll see this gives us a significantly better running time when this is done at every step of the recursion. Our new recurrence relation is

$$T_2(n) = \begin{cases} 1 & n = 1 \\ 3T_2\left(\frac{n}{2}\right) + 3n & n > 1 \end{cases}$$

We will see that this comes out to be about $n^{1.58}$ which is a big improvement over the n^2 algorithm. As an illustration, let's take $n = 1000$: For the brute force n^2 algorithm, we get:

$$1000^2 = 1000000$$

Whereas for the Karatsuba's $n^{1.58}$ algorithm, we get:

$$1000^{1.58} = 54954$$

Comparing these values shows us that Karatsuba reduces the execution time to:

$$\frac{54954}{1000000} \times 100 = 5.94\%$$

of the brute force execution

Note : Since the 3 in $3n$ is just a constant, we will ignore it and just use n from now on. This simplifies our calculations and is asymptotically equal as we will see below (the master theorem). The recurrence for this algorithm that we will refer to from now on would be

$$T_2(n) = \begin{cases} 1 & n = 1 \\ 3T_2\left(\frac{n}{2}\right) + n & n > 1 \end{cases}$$

2 Rank_A(x)

Problem : Given an array A and an element x (which may or may not belong to A), we want to find the number of elements in A that are "less than" x . This is called the Rank of x in A and is written as $Rank_A(x)$

Example : Let $A = \boxed{5 \mid 4 \mid 6 \mid 9 \mid 2 \mid 7 \mid 5 \mid 8}$

Then $Rank_A(5) = 2$, $Rank_A(3) = 1$, $Rank_A(1) = 0$. Notice the rank of the minimum element of A , the maximum element of A and an integer larger than the maximum of A .

2.1 Computing the Rank

One way to compute the rank of an element is just by a simple linear scan of the array A . This is a simple algorithm, taking n steps in the worst-case scenario.

2.2 Computing Rank in a sorted array

Suppose the A is sorted, then can we do any better than n steps? Well we have already studied the binary search algorithm for finding an element in a sorted array. Binary search takes about $\log n$ steps. We can use a slightly modified version of it to find the rank of an element. Recall what binary search returns when the search key x is not in the array.

2.3 Computing Rank of two elements at once

Suppose now that A is sorted and instead of one element x , you're given two elements x and y and you have to compute the rank of both of these elements. Without loss of generality, we can assume $x \leq y$ (if $x > y$ just swap the two). Now we know we can compute the rank of one element using binary search. For two elements we can just repeat the same process for the second element. So the total steps taken would be $\log n + \log n = 2 \log n$. Another thing we could do is that suppose $Rank_A(x)$ turns out to be i . Then instead of using binary search on the whole array, we can use binary search on $B = A[i \dots n]$. Then $Rank_A(y) = Rank_B(y) + Rank_A(x)$. Concretely, we can do the following.

Algorithm 1 : Finding Rank of Two Elements

```
rankXA  $\leftarrow$  getRank( $x, A$ )  
rankYB  $\leftarrow$  getRank( $y, \{A[rankXA] \dots A[n]\}$ )  
return {rankXA, rankYB + rankXA}
```

How much improvement do we get by using this algorithm as compared to just using binary search on A twice? Well let's see, any improvement that we do get is because now we're searching over a smaller array. The size of that array is $n - Rank_A(x)$. So what's the largest possible value for $n - Rank_A(x)$? that would be our worst case. The largest value is when $Rank_A(x) = 0$. So in that case the size of the new array would be n and it would take us $\log n$ steps to compute $Rank_A(y)$ and a total of $\log n + \log n$ steps to compute both the ranks. So in the worst case there is no improvement by this algorithm.

2.4 Computing Rank of n elements at once

Okay so now, instead of two elements x and y we can generalize the problem. Suppose we're given two sorted arrays, A and B . Our task is to find the rank in A of all elements of B . That is we have to compute $\{Rank_A(x) | x \in B\}$. We solve this problem by extending our algorithm for the 2 element case. The algorithm is given below.

Algorithm 2 : Computing $\{Rank_A(x)|x \in B\}$ using Binary Search

```
for  $i = 1$  to  $n$  do
  if  $i == 1$  then
     $rankB[i] \leftarrow getRank(B[i], A)$ 
  else
     $rankB[i] \leftarrow (rankB[i - 1] + getRank(B[i], \{A[rankB[i - 1]] \cdots A[n]\}))$ 
```

Let's analyze the running time for this algorithm. For two elements x and y we calculated that the running time should be $2 \log(n)$. When we extend it to n elements, again at each turn, in the worst case, the size of the remaining array that we have to look at is not reduced at all. So the total running time would be $\log n + \log n + \log n + \dots + \log n = n \log n$.

Taking a step Back : Can we do any better than this ? Let's consider again, the linear scan approach that was discussed in the beginning. As it turns out, in one linear pass through the array A we can determine the rank of all elements of B . The algorithm for that is quite simple, and is described below.

Algorithm 3 : Computing $\{Rank_A(x)|x \in B\}$ using Linear Search

```
 $r \leftarrow 0$ 
for  $i = 1$  to  $n$  do
  if  $A[i] > B[j]$  then
     $rankB[j] \leftarrow r$ 
     $j \leftarrow j + 1$ 
  else
     $r \leftarrow r + 1$ 
```

This concludes our discussion on rank. Now we move on to a seemingly different problem. That of "merging" two sorted arrays. It is quite easy to see that this algorithm takes $2n$ comparisons.

3 Merging two arrays

Problem : The merge problem is as follow. We're given two arrays A and B , each of size n . Both of these arrays are in sorted order. We have to make another sorted array C of size $2n$, that contains all the elements of A and B .

Example : Let $A = \boxed{1 \ 3 \ 4 \ 7}$, $B = \boxed{2 \ 5 \ 6 \ 8}$

Then $C = \boxed{1 \ 2 \ 3 \ 4 \ 5 \ 6 \ 7 \ 8}$

3.1 Merging by finding ranks

One way to do this is to just combine the two arrays into one big array C and then just sort C . Since C is of size $2n$ this would take about $2n \log(2n)$ comparisons, if we use insertion sort with binary search. Using selection sort or bubble sort it would take about $2n(2n-1)/2$ comparisons. But since the two arrays A and B are sorted perhaps we can do better. If you really think about this problem is exactly equivalent to determining the rank of all elements of B in A (or vice versa). If we know the ranks we know in which position the elements of A and B would go in the combined array. Like in the example described above we know

$$\text{Rank}_A(2) = 1, \quad \text{Rank}_A(5) = 3, \quad \text{Rank}_A(6) = 3, \quad \text{Rank}_A(8) = 4$$

So we know that 2 would come after one element of A has been inserted (which one?). 5 would come after three elements of A have been inserted into C and so on. And we know that finding the ranks of all elements of B in A takes just n steps. So we should be able to solve this problem in about n steps. Typically though, the way the solution is implemented doesn't require explicitly finding the ranks. We just maintain two running pointers for A and B , compare the elements at those indices, insert the smaller one into C and increment the corresponding pointer. The exact algorithm is given below

Algorithm 4 : Merging A and B

```
 $p \leftarrow 1$ 
 $q \leftarrow 1$ 
for  $i = 1$  to  $2n$  do
  if ( $p > n$  or  $A[p] > B[q]$ ) then
     $C[i] \leftarrow B[q]$ 
     $q \leftarrow q + 1$ 
  else
     $C[i] \leftarrow A[p]$ 
     $p \leftarrow p + 1$ 
```

4 Merge Sort

The above strategy allows us to define another sorting algorithm, called the Merge Sort. Merge Sort works as follows. Given an array, A of size n we divide the array into two halves $A[1 \dots n/2]$ and $A[n/2 + 1 \dots n]$, recursively sort the two halves and then merge the sorted halves using the merge algorithm we just discussed. The pseudocode is given below.

Algorithm 5 : Merge Sort

```
function MERGESORT(A)
  if SIZE(A) = 1 then
    return A
  else
    L ← MergeSort(A[1 ... SIZE(A)/2])
    R ← MergeSort(A[SIZE(A)/2 + 1 ... SIZE(A)])
    return MERGE(L, R)
```

4.1 Runtime of MergeSort

To analyze the running of this algorithm we can set up a recurrence relation. We know that to sort an array of size m we have to sort two arrays of size $m/2$ and then merge the two. We know that the time taken to merge two arrays of size $m/2$ is m . Denote by $T(m)$ the number of comparisons MergeSort takes to sort an array of m elements. It is clear from the algorithm and preceding discussion that

$$T(m) = 2T\left(\frac{m}{2}\right) + m \quad (1)$$

We want to find a closed form solution to this recurrence (rather than this recursive/inductive definition). One way to find such a closed form solution is to just extend the recursive relation. Suppose we want to find what is $T(n)$, then repeatedly applying (1) we get

$$\begin{aligned} T(n) &= 2T\left(\frac{n}{2}\right) + n \\ &= 2\left(2T\left(\frac{n}{4}\right) + \frac{n}{2}\right) + n \\ &= 2\left(2\left(2T\left(\frac{n}{8}\right) + \frac{n}{4}\right)\right) + n + n \\ &= 2\left(2\left(2\left(2T\left(\frac{n}{16}\right) + \frac{n}{8}\right)\right)\right) + n + n + n \end{aligned}$$

In general the k^{th} line of the above sequence of equation is

$$T(n) = 2^k \cdot T\left(\frac{n}{2^k}\right) + \overbrace{n + n + \dots + n}^{k \text{ times}}$$

As we know from the algorithm that $T(1) = 1$ and it is easy to see that the maximum value of k is $\log n$ (assuming n is a power of 2 to rid ourselves from dealing with floors and ceilings), we get

that the last equation will be

$$\begin{aligned}
 T(n) &= 2^{\log n} \cdot T\left(\frac{n}{2^{\log n}}\right) + \overbrace{n + n + \dots + n}^{\log n \text{ times}} \\
 &= n \cdot T\left(\frac{n}{n}\right) + \overbrace{n + n + \dots + n}^{\log n \text{ times}} \\
 &= n \cdot 1 + n \log n \\
 &= n \log n + n
 \end{aligned}$$

We will now discuss recurrence relations, and ways to solve them.

5 Recurrence

The following two paragraphs are mostly copy pasted from <https://courses.engr.illinois.edu/cs573/fa2010/notes/99-recurrences.pdf>

A recurrence is a recursive description of a function, or in other words, a description of a function in terms of itself. Like all recursive structures, a recurrence consists of one or more base cases and one or more recursive cases. Each of these cases is an equation or inequality, with some function value $f(n)$ on the left side. The base cases give explicit values for a (typically finite, typically small) subset of the possible values of n . The recursive cases relate the function value $f(n)$ to function value $f(k)$ for one or more integers $k < n$. For example, the following recurrence describes the identity function $f(n) = n$:

$$f(n) = \begin{cases} 0 & \text{if } n = 0 \\ f(n-1) + 1 & \text{otherwise} \end{cases}$$

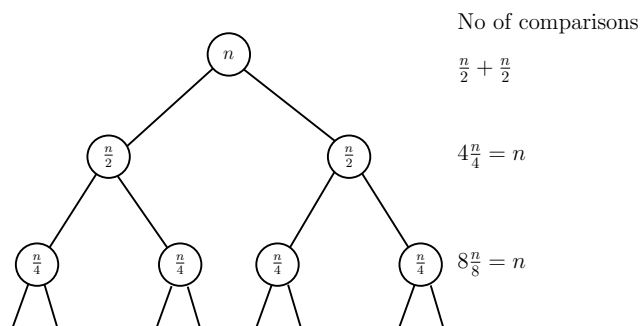
Recurrences arise naturally in the analysis of algorithms, especially in divide and conquer algorithms. In many cases, we can express the running time of an algorithm as a recurrence, where the recursive cases of the recurrence correspond exactly to the recursive cases of the algorithm. Recurrences are also useful tools for solving counting problems — How many objects of a particular kind exist?

By itself, a recurrence is not a satisfying description of the running time of an algorithm. Instead we would like a closed-form solution to the recurrence; this is a non-recursive description of a function that satisfies the recurrence. We will discuss now, how to solve some types of these recurrence relations.

5.1 Solving Recurrence Relations

One way to think about what's going on in a recurrence relation is to visualize it as a tree diagram. We will consider the Mergesort algorithm and Karatsuba algorithm and see how we can visualize

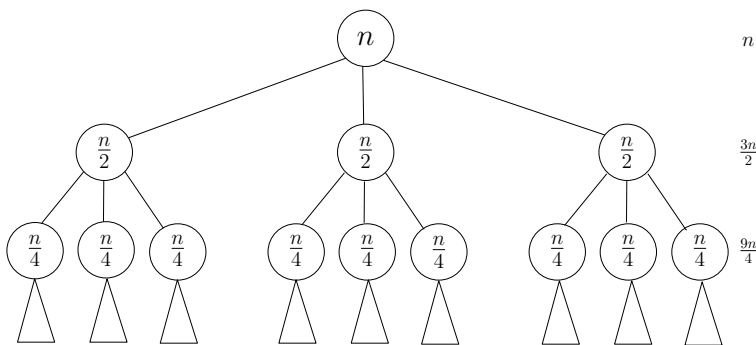
them and see what's happening. This would hopefully provide us with some insight as to how much work our algorithm is doing. We represent each instance of a called function as a node in the tree and calls to a function represent edges. This allows us to figure out how many function calls are made and helps us figure out the amount of work that's being done. Below is a tree diagram for the merge sort algorithm. As a reminder, mergesort splits the array to be sorted into two halves, sorts each half independently and then merges the two halves outside of the recursion. For merging it takes time equal to length of the merged array. For a more detailed description, refer to notes on mergesort.



In the beginning there is a call to mergesort with a problem of size n . That function then calls two other instances of mergesort with problem sizes of $\frac{n}{2}$, each of which make two calls to mergesort with problem sizes of $\frac{n}{4}$ and so on. We know that the calls would stop when the size of the array is 1 (or two.. that's just an implementation issue). So we know there are about $\log n$ levels of the tree. Another thing we can see from this diagram is that at every level n comparisons are made for merging. That means the total time for this algorithm is $n \log n$.

5.2 Analysis Of the Karatsuba Algorithm

Let's now see how the recursion tree for the Karatsuba Algorithm looks like. In the Karatsuba algorithm at each step the problem is divided into 3 subproblems of size $\frac{n}{2}$. So its recursion tree looks as follows



The number of levels of this tree would be the same as in mergesort (i.e $\log n$) since at each level

the problem size is halved, just like in mergesort. However the number of problems grows faster in Karatsuba. And as a result we can see in the diagram that the amount of work being done at each level is not the same. The amount of work at the top (0_{th}) level is n , at the first level $1.5n$ and $2.25n$ at the second level. So it's not immediately clear what the total running time of the algorithm would be. But we can sort of see a pattern. At each level the problem size gets split into 2 and the number of problems grows by a factor of 3. So at level i the number of problems is 3^i and the size of each problem is $\frac{n}{2^i}$. So the total work at each level is $\frac{3^i \times n}{2^i}$.

We can now figure out the total running time that this algorithm would take by summing up the running time over all levels.

$$\begin{aligned} T_2(n) &= \sum_{i=0}^{\log n} \frac{3^i \times n}{2^i} = n \sum_{i=0}^{\log n} \left(\frac{3}{2}\right)^i = n \left(\frac{1 - \frac{3}{2}^{\log n + 1}}{1 - \frac{3}{2}}\right) \\ &\in \mathcal{O}\left(n \times \left(\frac{3}{2}\right)^{\log_2 n}\right) \in \mathcal{O}\left(n \times \left(\frac{3^{\log_2 n}}{n}\right)\right) \in \mathcal{O}\left(3^{\frac{\log_3 n}{\log_3 2}}\right) \\ &\in \mathcal{O}\left(n^{\frac{1}{\log_3 2}}\right) \in \mathcal{O}\left(n^{0.6309}\right) \in \mathcal{O}\left(n^{1.58}\right) \end{aligned}$$

5.3 Analysis Of the Naive Divide & Conquer Multiplication Algorithm

We can similarly imagine a tree for this algorithm. At each step the number of problems would be multiplied by 4 and the problem size would be halved. So at each level the amount of work would be $\frac{4^i \times n}{2^i}$. And the total work would be

$$\begin{aligned} T_1(n) &= \sum_{i=0}^{\log n} \frac{4^i \times n}{2^i} = n \sum_{i=0}^{\log n} 2^i = n \times 2^{\log_2 n + 1} \\ &\in \mathcal{O}\left(n \times 2^{\log_2 n}\right) \in \mathcal{O}(n \times n) \in \mathcal{O}(n^2) \end{aligned}$$

So we see that the running time in this case turns out to be n^2 as was mentioned previously.

5.4 Substitution method for Solving Recurrences

While we have seen that the recursion tree method works pretty well for solving recurrences, we would like to have a method that is simpler and does not require much drawing. Such a method is the substitution method, which simply requires to substitute the value of the function for smaller values from the recurrence and continue until we see a pattern and can't manage it. Let us use the recurrence relation for merge sort, which is:

$$T(n) = 2T\left(\frac{n}{2}\right) + n$$

$$\begin{aligned}
T(n) &= 2T\left(\frac{n}{2}\right) + n = 2\left(2T\left(\frac{n}{4}\right) + \frac{n}{2}\right) + n = 2 * 2 * \left(2T\left(\frac{n}{8}\right) + \frac{n}{4}\right) + n = 2 * 2 * 2\left(T\left(\frac{n}{8}\right) + \frac{n}{8}\right) + n + n + n \\
&\vdots \quad \quad \quad \vdots \\
&= \underbrace{2 * 2 * 2 \dots * 2}_k * T\left(\frac{n}{2^k}\right) + \underbrace{n + n + \dots + n}_k \\
&\vdots \quad \quad \quad \vdots \\
&= \underbrace{2 * 2 * 2 \dots * 2}_{\log n} * T\left(\frac{n}{2^{\log n}}\right) + \underbrace{n + n + \dots + n}_{\log n} \\
&= 2^{\log n} * 1 + n \log n \\
&= n \log n + n
\end{aligned}$$

Hence we get $O(n \log n)$ for the merge sort algorithm.

5.4.1 Using the substitution method for solving the naive Divide and Conquer multiplication problem

As stated previously, the naive divide and conquer algorithm for multiplication has the recurrence relation

$$T_1(n) = 4T_1\left(\frac{n}{2}\right) + n$$

$$\begin{aligned}
T_1(n) &= 4T_1\left(\frac{n}{2}\right) + n = 4\left(4T_1\left(\frac{n}{4}\right) + \frac{n}{2}\right) + n \\
&= 4\left(4\left(4T_1\left(\frac{n}{8}\right) + \frac{n}{4}\right) + \frac{n}{2}\right) + n = 4 * 4 * 4T_1\left(\frac{n}{8}\right) + 4 * 4 * \frac{n}{4} + 4 * \frac{n}{2} + n \\
&\vdots \quad \quad \quad \vdots \\
&= \underbrace{4 * 4 * 4 \dots * 4}_k * T_1\left(\frac{n}{2^k}\right) + \sum_{i=0}^k \left(4^i * \frac{n}{2^i}\right) \\
&\vdots \quad \quad \quad \vdots \\
&= \underbrace{4 * 4 * 4 \dots * 4}_{\log n} * T_1\left(\frac{n}{2^{\log n}}\right) + \sum_{i=0}^{\log n} 2^i n = 4^{\log n} * 1 + n(2^{\log n+1}) \\
&= 2^{2\log n} + n * n = n^2 + n^2 = 2n^2
\end{aligned}$$

Hence we get $O(n^2)$ runtime for the naive divide and conquer multiplication.

5.4.2 Using substitution to analyze Karatsuba Algorithm

Finally, we look into the Karatsuba algorithm and determine its runtime via the substitution method. So we have:

$$\begin{aligned}
T_2(n) &= 3T_2\left(\frac{n}{2}\right) + n = 3\left(3T_2\left(\frac{n}{4}\right) + \frac{n}{2}\right) + n = 3\left(3\left(3T_2\left(\frac{n}{8}\right) + \frac{n}{4}\right) + \frac{n}{2}\right) + n \\
&\vdots \quad \quad \quad \vdots \\
&= \underbrace{3 * 3 * 3 \dots * 3}_k * T_2\left(\frac{n}{2^k}\right) + \sum_{i=0}^k \left(3^i * \frac{n}{2^i}\right) \\
&\vdots \quad \quad \quad \vdots \\
&= \underbrace{3 * 3 * 3 \dots * 3}_{\log n} * T_2\left(\frac{n}{2^{\log n}}\right) + \sum_{i=0}^{\log n} \left(3^i * \frac{n}{2^i}\right) = 3^{\log n} * 1 + n * \sum_{i=0}^{\log n} \left(\frac{3^i}{2^i}\right) \\
&= 3^{\log n} + n * \left(\frac{1 - \left(\frac{3}{2}\right)^{\log n}}{1 - \left(\frac{3}{2}\right)}\right) = 3^{\log n} + n * \left(\frac{1 - \left(\frac{3^{\log n}}{2^{\log n}}\right)}{-\left(\frac{1}{2}\right)}\right) \\
&= 3^{\log n} + n * 2\left(\frac{3^{\log_2 n}}{n}\right) = 3\left(3^{\frac{\log_3 n}{\log_3 2}}\right) = 3\left(3^{\log_3 n}\right)^{\frac{1}{\log_3 2}} = 3n^{\frac{1}{\log_3 2}} = 3n^{1.58}
\end{aligned}$$

Hence, we get $O(n^{1.58})$ as the runtime for the Karasuba algorithm.

5.5 Master Theorem

We've had to analyze the running times of a bunch of divide and conquer algorithms over the past 2,3 classes which had pretty similar recurrence relations. Setting up a recurrence relation, and then figuring out how much work is being done at each level of the recursion tree and then summing up over all levels is a pain and we would like to have some general results that can be used to directly get the running time of a divide and conquer algorithm. As it turns out, most divide and conquer algorithms tend to have the following type of recurrence relation.

$$T(n) = aT\left(\frac{n}{b}\right) + \mathcal{O}(n^d)$$

Where a corresponds to the number of subproblems, n/b corresponds to the size of each subproblem and $\mathcal{O}(n^d)$ is the work performed in the conquer and combine step. There's a theorem called the Master theorem which gives us the following results for these types of recurrence relations

$$T(n) = \begin{cases} \mathcal{O}(n^d) & d > \log_b(a) \\ \mathcal{O}(n^d \log(n)) & d = \log_b(a) \\ \mathcal{O}(n^{\log_b(a)}) & d < \log_b(a) \end{cases}$$

5.5.1 Analyzing running times using the Master Theorem

Let's now apply to master theorem to the various recurrence relations we've discussed up to now. For the merge sort algorithm $a = 2$, $b = 2$ and $d = 1$. So $\log_b(a) = \log_2(2) = 1$ which is the same as d . So case 2 applies here and we get a running time of $n \log n$, the same as when we computed using the recursion tree.

For the naive divide and conquer multiplication algorithm we have $a = 4$, $b = 2$ and $d = 1$. Which means $\log_b(a) = \log_2(4) = 2$ which is greater than d . So case 3 applies and we get a running time of

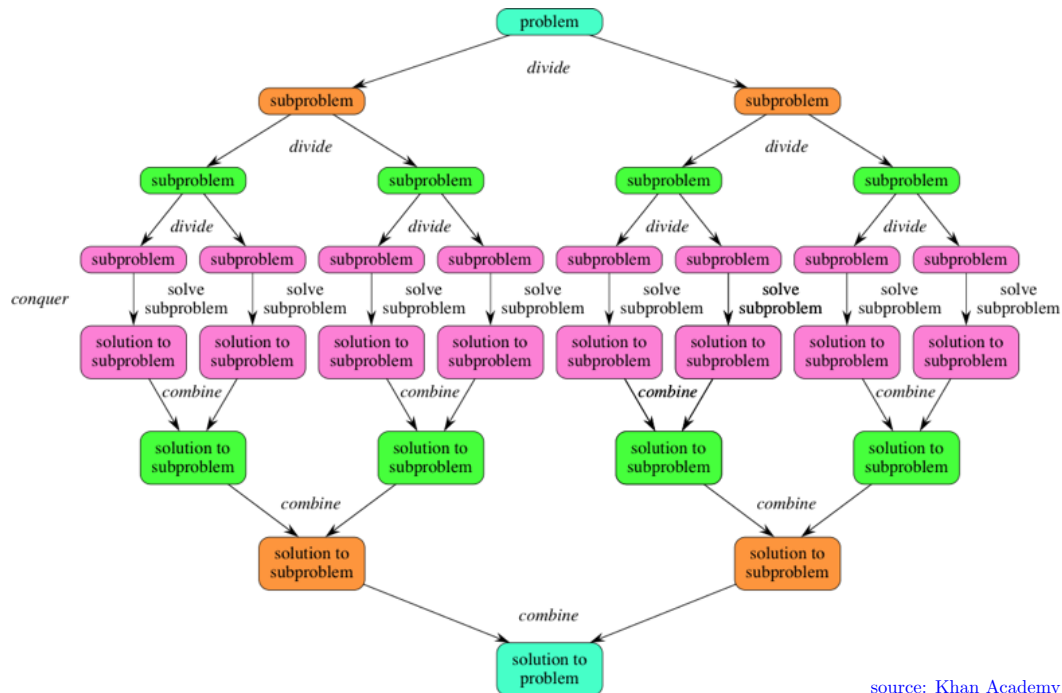
$$n^{\log_b(a)} = n^{\log_2(4)} = n^2$$

For Karasuba too, case 3 applies. You can put in the relevant values and see that the running time comes out to be about $n^{1.58}$ same as we calculated earlier.

6 Divide and Conquer Algorithm Design Paradigm

The merge sort algorithm is the canonical example given for the Divide and Conquer design paradigm. The divide and conquer approach requires recursively breaking a problem into smaller subproblems (the divide part) until they become easy enough to be solved directly. In merge sort,

for example we keep dividing the array into two parts until the size of the array is 1 and we know how we can sort an array with just one element (It's already sorted). The next step is to somehow combine the solutions of the smaller subproblems to obtain the solution of the original problem (the conquer part). This step is usually tricky. In merge sort we used the “merge” algorithm to sort the larger array from two sorted halves.

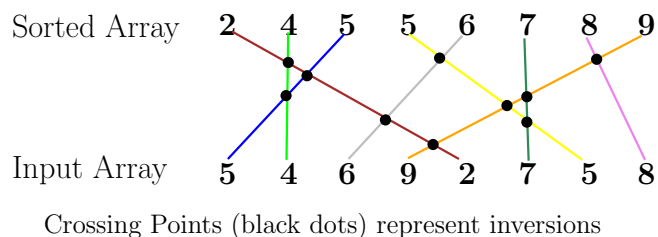


source: Khan Academy

Next, we will discuss two more Divide and Conquer algorithms. The first is the problem of counting inversions, as follows. The second is the problem of finding closest pair of points among a set of points in 2-dimensional Euclidean space. We present the brute-force approach to solve this problem and improve it using the divide and conquer approach.

7 Counting Inversions

Problem : Given an array A an inversion is defined as a pair i, j such that $i < j \wedge A[i] > A[j]$. We want to find how many such pairs are in A .



Example : Let $A = \boxed{5 \mid 4 \mid 6 \mid 9 \mid 2 \mid 7 \mid 5 \mid 8}$. Then the inversions for this array are

$\{(1, 2), (1, 5), (2, 5), (3, 5), (3, 7), (4, 5), (4, 6), (4, 7), (4, 8), (6, 7)\}$

Solution : One way to solve is problem is to simply check all pairs i, j and see how many of them are inversions. There are $\binom{n}{2}$ pairs so that would take about $n^2/2$ steps. That algorithm would work as follows

Algorithm 6 : Counting Inversions Using Brute Force

```

count ← 0
for i = 1 to n do
  for j = i + 1 to n do
    if A[i] > A[j] then
      count ← count + 1

```

As computer scientists, however, we always want to come with better, more efficient algorithms. But why should we think that a better solution should exist for this problem? Well one reason is that this problem is somewhat related to sorting an array (Infact the number of inversions is used as a measure for the sorted-ness of an array). When we're sorting an array we're basically removing all the inversions. And we know we can sort an array in $n \log(n)$ steps. The algorithm for that was discussed in the last class. So it makes sense that there should be a similar algorithm that can count inversions in about the same time too. We used the divide and conquer approach to obtain a better running time for sorting. So let's try to do something similar for counting inversions too. How should be break this problem?

If we divide an array into two halves, then we can split the inversions into 3 types: Left-Left, Right-Right inversions and Left-Right inversions. For example, consider the following array A :

$\boxed{2 \mid 3 \mid 8 \mid 5 \mid 4 \mid 10 \mid 9 \mid 12 \mid 7 \mid 18 \mid 15 \mid 25}$

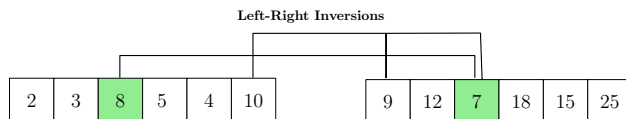
Divide the array into two.

$\boxed{2 \mid 3 \mid 8 \mid 5 \mid 4 \mid 10}$ $\boxed{9 \mid 12 \mid 7 \mid 18 \mid 15 \mid 25}$

Recursively count inversions in each half.

$\boxed{2 \mid 3 \mid 8 \mid 5 \mid 4 \mid 10}$ $\boxed{9 \mid 12 \mid 7 \mid 18 \mid 15 \mid 25}$
Left-Left = 8-5, 8-4, 5-4 *Right-Right* : 9-7, 12-7, 18-15

Count inversions where a_i and a_j are in different halves and return total inversions count.



We'll refer to the left and right halves of the array as L and R . Then Left-Left inversions are those which only involve pairs from L . Similarly the Right-Right inversions are the inversions that involve pairs from only R . Left-Right inversions are pairs that involve one element from L and one element from R . The Left-Left and Right-Right inversions will be found by recursively dividing the array into even smaller parts. The problem is finding the Left-Right inversions. If we can find the Left-Right inversions in n steps then we can solve the problem in $n \log(n)$ steps. The algorithm would have the same recurrence relation as the merge sort algorithm and hence the same running time.

The problem of finding the left-right inversions is equivalent to finding the ranks in R of all the elements of L . We know that if the two arrays R and L are sorted then we can find the ranks in n steps. So while counting the inversions in the left and right halves of the array we'll do some extra work and also sort them. The idea is that even though we do the some extra work in sorting, the overall time will be reduced because now we can find the inversions much faster. Let's see how that algorithm would work.

Algorithm 7 : Counting Inversions using Divide and Conquer

```

function COUNTINVERSIONS( $A$ )
  if SIZE( $A$ ) = 1 then return ( $A, 0$ )
   $L \leftarrow A[1, 2, \dots, n/2]$ 
   $R \leftarrow A[n + 1, n + 2, \dots, n]$ 
  ( $sortedL, inv_{l,l}$ )  $\leftarrow$  COUNTINVERSIONS( $L$ )
  ( $sortedR, inv_{r,r}$ )  $\leftarrow$  COUNTINVERSIONS( $R$ )
   $inv_{l,r} \leftarrow sum(FINDRANKS(L, R))$  ▷ takes  $n$  steps
  return (MERGE( $L, R$ ),  $inv_{l,l} + inv_{r,r} + inv_{l,r}$ ) ▷ merge takes  $n$  steps

```

So our recurrence relation turns out to be

$$T(n) = 2T\left(\frac{n}{2}\right) + 2n$$

And solving this we get

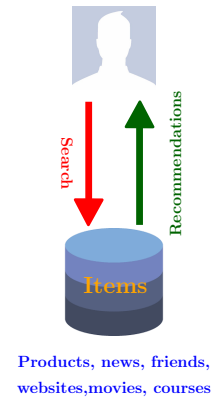
$$T(n) = 2n \log(n)$$

Which is a huge improvement over the n^2 brute force algorithm.

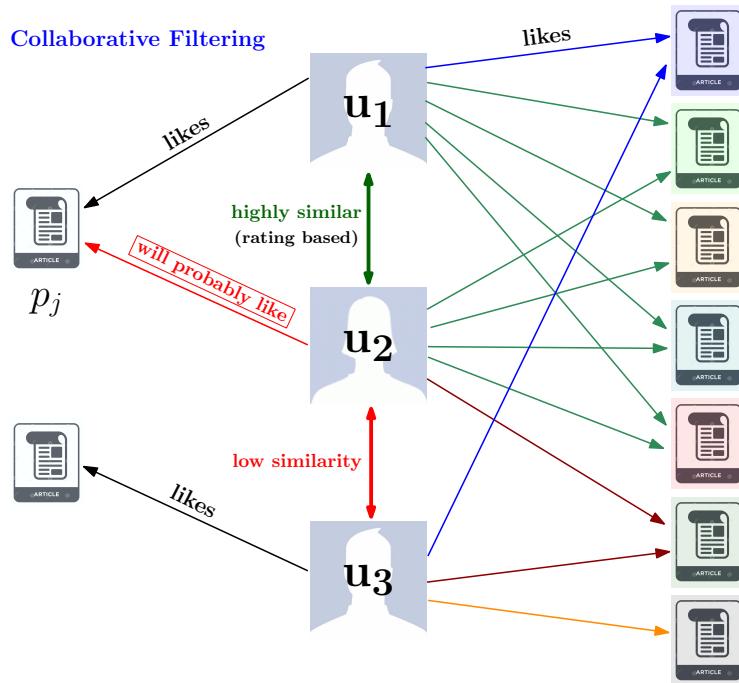
7.1 Collaborative Filtering

One of the areas in which the problem of counting inversions comes up is in recommender systems. A recommender system tries to predict the 'rating' a user gives to a particular item. Like how YouTube tries to recommend you videos based on your viewing history, or how IMDb recommends movies.

	p_1	p_2	p_3	p_j					p_m			
u_1	1	2	1	4		2	3	2	5			2
u_2		1		2	1		2		1			3
u_3	1	1	2		1					1	2	
			3	2		5	2			3	4	
	1		2						5			
u_i		3	2	1	4	5	?	1	3	1	2	1
		4									4	
		5		1							5	
	1	4						1	3	5	1	2
u_n		3	1	1	2	1			4			5



One technique used in recommender systems is that of “*collaborative filtering*”. Collaborative filtering is a method of making automatic predictions (filtering) about the interests of a user by collecting preferences or taste information from many users (collaborating). The underlying assumption of the collaborative filtering approach is that if a person A has the same opinion as a person B on an issue, A is more likely to have B 's opinion on a different issue x than to have the opinion on x of a person chosen randomly” (Wikipedia).



The way this is done is as follows. Based on your activity on the website (How many times you view a particular item, the review you've given some item, etc.) the website builds a list containing your relative preference to various items.

Table 1: User 1

$item_A$	$item_B$	$item_C$	$item_D$	$item_E$
1	2	3	4	5

It then builds a similar list for other users. It then finds the users most similar to you by counting inversions between your list and theirs. Once the website has a list of similar users it can then suggest items to you based on what these similar people like.

Table 2: User 2

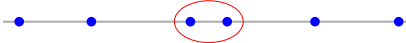
$item_A$	$item_B$	$item_C$	$item_D$	$item_E$
3	1	2	4	5

8 Closest Pair

The second problem is as follows: Given an array of points in the plane, find the pair of points which is closest with respect to Euclidean distance.

8.1 Naive Approach for Closest Pair

Notice that if points are in $1-D$, (an array of real numbers), with sorting it is easy to find the closest pair. In $2-D$, a simple way to solve the problem would be to compare every point with every other point and find minimum distant pair. This takes $O(n^2)$ time.



8.2 Divide and conquer approach for Closest Pair

Input: $2-D$ array P of points.

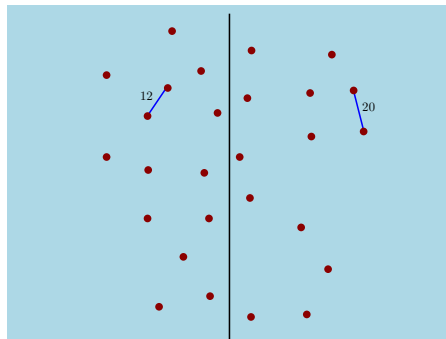
Step 1. Sort the array P with respect to x -coordinates to get $2-D$ array P_x and sort with respect to y -coordinates to get $2-D$ array P_y .

Divide Part:

Step 2. Consider the mid point (median) m of array P_x . Let the set of points on the left side of m be S_1 and on right of m be S_2 . We make array of points in S_1 and call it P_{1x} which is in fact left sub-array of P_x and P_{2x} with points in S_2 . Note that P_{1x} and P_{2x} are sorted in x -direction. We also compute arrays P_{1y} and P_{2y} , which are essentially points in S_1 and S_2 respectively but sorted in y -direction. We compute the arrays P_{1y} and P_{2y} in the following way.

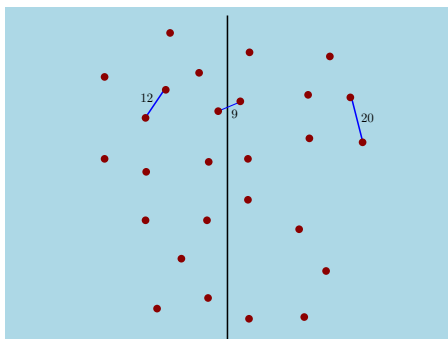
1. Pick first point in P_y and compare its x -coordinate with m .
2. If $x \leq m$, insert the point in P_{1y} , otherwise insert in P_{2y}
3. continue inserting elements in both arrays.

Note that points are inserted in sorted order in both arrays since P_y is sorted. **Analysis** It takes $O(n)$ steps to make the arrays P_{1y}, P_{2y}, P_{1x} and P_{2x} .



Conquer Part:

Step 3. Recursively find closest pairs in S_1 with input data P_{1x}, P_{1y} , and S_2 with input data P_{2x}, P_{2y} , with shortest distance δ_1 and δ_2 respectively.



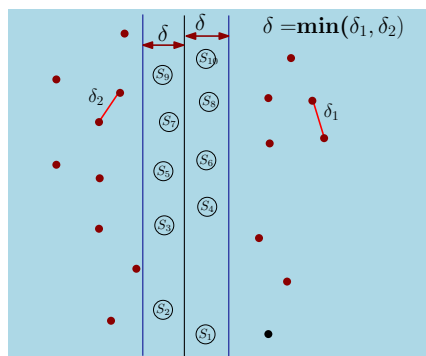
Consider smaller distance of both δ_1 and δ_2 . Let $\delta = \min\{\delta_1, \delta_2\}$. Only thing left to consider in the problem is, what if closest pair has one point in S_1 and other in S_2 ?

Combine Part:

Step 4. Consider $N_1 \subseteq S_1$ and $N_2 \subseteq S_2$ such that points in N_1 and N_2 are in δ distance strip of m . Sort $N_1 \cup N_2$ with respect to Y -coordinates of points. We compute and sort points in $2 - \delta$ strip of m in the following way.

1. Consider first point in P_y array and check if its x -coordinate is less than δ distance from m .
2. if x -distance from m is less than δ then put this point in array $N = N_1 \cup N_2$.
3. check this for every point in P_y .

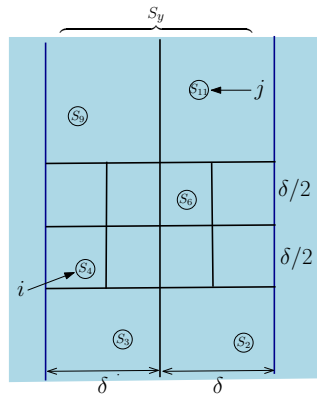
Analysis This step takes $O(n)$ time.



Note that if the closest pair has one point in S_1 and one in S_2 , then that pair must be in N .

For each point x of N , there can be atmost 7 points in N such that distance of x with those 7 points is less than δ . Since if x has distance less than δ with some point, it will be in one of its neighboring squares, shown in diagram below. And there can not be more than one point in any of the small squares otherwise it will contradict the minimality of δ . We are not considering

squares below the point x because these points are sorted and points below have been processed (assuming that traversing is in ascending order).



Step 5. Compute the distance of each point in N to next 7 points in the array and compare it with δ . If distance is smaller than δ then update the closest pair and continue.

Analysis This step takes $O(n)$ time to execute since each point is being compared to constant number of other points.

8.3 Runtime Analysis

In each divide part, we get two half size problems and combining takes $O(n)$ time. So run time for the whole algorithm is

$$T(n) = \begin{cases} 1 & \text{if } n < 3 \\ 2T(\frac{n}{2}) + O(n) & \text{if } n \geq 3 \end{cases}$$

Solving this recursively, $T(n) = O(n \log n)$