

# COMP 251 Study guide

Francis Piche

January 23, 2018

# Contents

<b>1</b>	<b>Preliminaries</b>	<b>3</b>
<b>I</b>	<b>Recursive Algorithms</b>	<b>3</b>
<b>2</b>	<b>Divide + Conquer Algorithms</b>	<b>3</b>
2.1	MergeSort . . . . .	4
2.2	Binary Search . . . . .	5
2.3	Run Time of Divide + Conquer in General . . . . .	5
2.4	Aside on Recurrences: Domain Transformation . . . . .	6
<b>3</b>	<b>Master Theorem</b>	<b>7</b>
3.1	Tree Method to Prove Master Theorem . . . . .	10
<b>4</b>	<b>Multiplication</b>	<b>11</b>
4.1	Grade School Multiplication . . . . .	11
4.2	Russian Peasant Multiplication . . . . .	11
4.3	Divide + Conquer Multiplication . . . . .	11
4.4	Fast Fourier Transforms . . . . .	12
4.5	Multiplying Matrices . . . . .	13
4.6	Fast Exponentiation . . . . .	14
<b>5</b>	<b>The Median Problem</b>	<b>15</b>
5.1	The Selection Problem . . . . .	15
5.2	Median of Medians . . . . .	16
<b>6</b>	<b>Finding the Closest Pair of Points in the Plane</b>	<b>17</b>
6.1	Exhaustive Search . . . . .	17
6.2	2-D case . . . . .	17
6.3	Widening the Bottleneck . . . . .	18

# 1 Preliminaries

In this course an algorithm is considered **good** if it:

- Works
- Runs in polynomial time. Meaning it runs, in  $O(n^k)$  time. Where  $n$  is (always) the size of the problem. (Number of elements in a list to be sorted etc.)
- Scales multiplicatively with computational power. (If your computer is twice as fast, the problem is solved at least twice as fast)

A bad algorithm is one that:

- Doesn't always work
- Runs in exponential time or greater. Meaning:  $O(k^n)$  time.
- Does not scale well with computational power. (Your computer is twice as fast, but barely any performance boost).

## Part I

# Recursive Algorithms

I won't be going into detail on the specifics of things like how recursion works, MergeSort, BinarySearch, solving recurrences, Big  $O$ , etc. as it's considered prerequisite material. If you need some review, my COMP250 study guide is still publicly available.

## 2 Divide + Conquer Algorithms

Examples:

- MergeSort
- BinarySearch

## 2.1 MergeSort

The MergeSort algorithm involves splitting a list of  $n$  elements in half, sorting each half recursively, and merging the sorted lists back into one. It takes time  $T(\frac{n}{2})$  to sort the list of half size, and time  $O(n)$  to merge the list back together. So the recurrence relation for MergeSort is given by:

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

where  $c$  is some constant.

**Theorem 1.** MergeSort runs in time  $O(n \log(n))$ .

*Proof.* Add **dummy numbers** (extra "padding" to the list), until  $n$  is a power of two.  $n = 2^k$ . We can do this because  $O()$  gives an **upper bound**, and adding numbers will make our solution take longer than the real one. Doing this will make solving the recurrence easier.

Unwinding the formula:

$$\begin{aligned} T(n) &= 2(2(T(\frac{n}{4}) + c\frac{n}{2}) + cn) \\ &= 2^2(T(\frac{n}{4}) + 2cn) \\ &= 2^3(T(\frac{n}{8}) + 3cn) \\ &= 2^4(T(\frac{n}{16}) + 4cn) \end{aligned}$$

Notice we have a pattern emerging.

$$= 2^k(T(1)) + kcn$$

Recall  $2^k = n$ , so  $k = \log_2(n)$  and  $T(1) = 1$  so:

$$= n + n \log_2(n)$$

Which is  $O(n \log n)$ . □

## 2.2 Binary Search

Binary search involves splitting your sorted list into two, and searching that half. So our recurrence is given by:

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

where  $c$  represents the constant work (comparisons, setting new bounds etc.)

**Theorem 2.** Binary Search is  $O(\log_2(n))$ .

*Proof.* Again we add dummy numbers so that  $n$  is a power of two.  $n = 2^k$

We begin with our recurrence:

$$\begin{aligned} T(n) &= T\left(\frac{n}{2}\right) + c \\ &= T\left(\frac{n}{4}\right) + c + c \\ &= T\left(\frac{n}{8}\right) + c + c + c \\ &= T\left(\frac{n}{2^k}\right) + kc \\ &= T(1) + \log_2(n) \end{aligned}$$

since  $k = \log_2(n)$  which is  $O(\log_2(n))$ . □

## 2.3 Run Time of Divide + Conquer in General

Divide and Conquer is a technique of solving problems that involves taking one large problem of size  $n$ , and breaking it down into  $a$  smaller problems of size  $\frac{n}{b}$ , and solving those problems recursively. They are then combined to produce a solution in time poly-time:  $O(n^d)$ .

So the run-time of a divide and conquer algorithm is:

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

In the case of MergeSort,  $a = 2$ ,  $b = 2$ ,  $d = 1$ .

In the case of BinarySearch,  $a = 1$ ,  $b = 2$ ,  $d = 0$ .

## 2.4 Aside on Recurrences: Domain Transformation

Note that the recurrence for MergeSort is really:

$$T'(n) \leq T'(\lfloor n/2 \rfloor) + T'(\lceil n/2 \rceil) + cn$$

Which we simplified by adding dummy entries. However, we can also say this: Note that this is an informal approximation, since it's really:

$$T'(n) \leq 2T'(\frac{n}{2} + 1) + cn$$

But the  $+1$  doesn't fit with our previous method.

We'll use **domain transformation** to solve this, starting with:

$$\begin{aligned} T(n) &= T'(n + 2) \\ &\leq T'(\frac{n + 2}{2} + 1) + c(n + 2) \end{aligned}$$

plugging in our expression from above

$$\leq T'(\frac{n + 2}{2} + 1) + c'(n)$$

absorbing the  $+2$  into  $c$ .

$$= T'(\frac{n}{2} + 2) + c'(n)$$

simplifying the fraction.

$$= T(\frac{n}{2}) + c'n$$

from our domain transformation at the beginning. Solving this the usual way, we get:

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

But again from our domain transformation:

$$T(n) = T'(n + 2)$$

, so

$$T'(n) = T(n - 2) = O(n \log(n))$$

So we've shown that  $T'(n)$  has the same upper bound as  $T(n)$ .

### 3 Master Theorem

**Theorem 3.** If  $T(n) = aT(n/b) + O(n^d)$  for constants  $a > 0$ ,  $b > 1$ ,  $d \geq 0$ , then:

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

These cases are just a few that occur often in practice when dealing with divide + conquer algorithms.

*Proof.* First we'll need two things. One is the geometric series, and the other is a law of logarithms. Professor Vetta proved them in class, and honestly I doubt you'd be asked to prove them on an exam, but it's good proof practice to go through them so I'll do it here.

$$\sum_{k=0}^l x^k = \frac{1 - x^{l+1}}{1 - x}$$

Proof:

Starting with:

$$(1 - x) \sum_{k=0}^l x^k$$

We can expand it out:

$$= \sum_{k=0}^l x^k - \sum_{k=0}^l x^{k+1}$$

Simplifying the sigma notation:

$$= \sum_{k=0}^l x^k - \sum_{k=1}^{l+1} x^k$$

All terms will cancel except:

$$= x^0 - x^{l+1} = 1 - x^{l+1}$$

Divide through by  $1 - x$

$$= \frac{1 - x^{l+1}}{1 - x}$$

Our second fact to derive is this law of logs:

$$x^{\log_b(y)} = y^{\log_b(x)}$$

Using the power rule of logarithms:

$$\log_b(x)\log_b(y) = \log_b(y^{\log_b(x)})$$

similarly,

$$\log_b(x)\log_b(y) = \log_b(x^{\log_b(y)})$$

so,

$$\log_b(x^{\log_b(y)}) = \log_b(y^{\log_b(x)})$$

Now we're ready for the proof.

Assume  $n$  is a power of  $b$ , and split up the problem into all its chunks.

$$T(n) = n^d + a\left(\frac{n}{b}\right)^d + a^2\left(\frac{n}{b^2}\right)^d + \dots + a^l\left(\frac{n}{b^l}\right)^d$$

(this is just if you'd "unwound" the whole recursion down to its simplest form like we did in the MergeSort/Binary Search proofs.)

Each term is the amount of work it will take at each level of the recursion.

Notice you can factor out:

$$\begin{aligned} &= n^d \left(1 + a\left(\frac{1}{b}\right)^d + a^2\left(\frac{1}{b^2}\right)^d + \dots + a^l\left(\frac{1}{b^l}\right)^d\right) \\ &= n^d \left(1 + \left(\frac{a}{b}\right)^d + \left(\frac{a}{b^d}\right)^2 + \dots + \left(\frac{a}{b^d}\right)^l\right) \end{aligned}$$

That looks like a geometric series! So let's look at the cases:

Case 1:  $a < b^d$

Applying the geometric series formula:



$$\begin{aligned}
&= n^d \sum_{k=0}^l \left(\frac{a}{b^d}\right)^k \\
&= n^d \frac{1 - \left(\frac{a}{b^d}\right)^{l+1}}{1 - \frac{a}{b^d}}
\end{aligned}$$

we can remove the  $\frac{a}{b^d}^{l+1}$  term with this inequality:

$$\leq n^d \frac{1}{1 - \frac{a}{b^d}}$$

which is  $O(n^d)$ .

Case 2:  $a = b^d$

Since  $\frac{a}{b^d} = 1$ :

$$= n^d(1 + 1 + 1 + \dots + 1)$$

There are  $l + 1$  terms, but we said  $n$  was a power of  $b$ , ( $n = b^l$ ) so,  $l = \log_b(n)$ , thus:

$$= n^d(\log_b(n) + 1)$$

which is  $O(n^d \log_b(n))$

Case 3:  $a > b^d$

Again from geometric series, and multiplying through by -1:

$$n^d \frac{\left(\frac{a}{b^d}\right)^{l+1} - 1}{\frac{a}{b^d} - 1}$$

Again this inequality holds:

$$\leq n^d \frac{\left(\frac{a}{b^d}\right)^{l+1}}{\frac{a}{b^d} - 1}$$

Which is  $O(n^d \left(\frac{a}{b^d}\right)^l)$  which we can simplify:

$$\left(\frac{n}{b^l}\right)^d a^l$$

but  $n = b^l$ , so:

$$\begin{aligned} &= (1)a^l \\ &= a^{\log_b(n)} \end{aligned}$$

now by our second fact:

$$= n^{\log_b(a)}$$

which is  $O(n^{\log_b(a)})$

□

It's **much** more important to understand the proof than it is to memorize the theorem.

### 3.1 Tree Method to Prove Master Theorem

A more intuitive way to think of the proof is with a *Recursion Tree*.

The root node of the tree has label  $n$ , and each node has  $a$  children (except the leaves).  $a$  is called the *branching factor*. Each child is labelled  $\frac{n}{b^d}$  where  $d$  is the depth. The labels represent the size of the sub problems.

The number of nodes at each level is  $a^d$ .

Case 1 is when the root level "dominates" all other levels, so the running time is just  $O(f(n))$  where  $f(n)$  is the amount of work at the root level.

Case 2 is when all levels are roughly the same weight. So the total running time is just  $O(f(n)l)$  where  $l$  is the number of levels.

Case 3 is when the leaves dominate, so the running time is  $O(a^l)$  since the leaves each take time  $O(1)$ , and there are  $a^l$  of them.

## 4 Multiplication

### 4.1 Grade School Multiplication

This takes  $n^2$  multiplications when you multiply two  $n$ -digit numbers. so the runtime is  $\Omega(n^2)$

### 4.2 Russian Peasant Multiplication

Super weird looking algorithm but it works!

---

```
Mult(x,y){  
  if x = 1 then output y  
  if x is odd then output y + Mult(floor(x/2),2y)  
  if x is even then output Mult(x/2, 2y)  
}
```

---

This actually comes from if you take the binary representation of  $x$ : say  $x = 46_{10}$  then  $x = 101110_2$ . The bits that are 1's will have the  $y$  added step, and the zero bits will just have the doubling step. Weird right?

Notice that this means the number of steps is just the number of bits in  $x$ . The number of digits in the result will be at most  $2n$ , so if we need to then add these, we add at most  $n$  numbers of  $2n$  digits so takes time  $O(n^2)$

### 4.3 Divide + Conquer Multiplication

Notice that a number  $x$  can be written as:

$$x = x_n x_{n-1} \dots x_{\frac{n}{2}+1} x_{\frac{n}{2}} \dots x_2 x_1$$

where the  $x_i$  are the digits.

Then we have:

$$x = 10^{\frac{n}{2}} x_L + x_R$$

where  $n$  is the number of digits,  $x_L$  is the first  $\frac{n}{2}$  digits, and  $x_R$  is the last  $\frac{n}{2}$

So by expanding:

$$xy = (10^n x_L y_R + 10^{\frac{n}{2}}(x_L y_R + x_R y_L) + x_R y_R)$$

Notice that this now involves four products of  $\frac{n}{2}$  digit numbers. So the recursion is:

$$T(n) = 4T(\frac{n}{2}) + O(n)$$

We have  $a = 4, b = 2, d = 1$ , which is case 3 of the master theorem.

Which means the running time is:

$$O(n^{\log_2(4)})$$

which simplifies to:

$$O(n^2)$$

Thanks to Gauss, we can actually use this fact:

$$x_L y_R + x_R y_L = x_R y_R + x_L y_L - (x_R - x_L)(y_R - y_L)$$

which is actually only 3 unique products. (adding is cheap)

So our new running time is:

$$T(n) = 3T(\frac{n}{2}) + O(n)$$

which is case 3 of the master theorem, so

$$\begin{aligned} O(n^{\log_2(3)}) \\ = O(n^{1.59}) \end{aligned}$$

## 4.4 Fast Fourier Transforms

These are  $O(n \log(n))$  for multiplying n-bit numbers. They'll be studied more in-depth at the end of the course (time-permitting).

## 4.5 Multiplying Matrices

There are  $n$  multiplications to calculate each entry of the result matrix, and there are  $n^2$  entries, so  $O(n^3)$

Using divide + conquer, divide into 4 sub-matrices:

$$\begin{bmatrix} x_{11} & x_{12} & x_{13} & \dots & x_{1n} \\ x_{21} & x_{22} & x_{23} & \dots & x_{2n} \\ \dots & \dots & \dots & \dots & \dots \\ x_{d1} & x_{d2} & x_{d3} & \dots & x_{dn} \end{bmatrix} = \begin{bmatrix} A & B \\ C & D \end{bmatrix}$$

So if we let:

$$x = \begin{bmatrix} A & B \\ C & D \end{bmatrix} y = \begin{bmatrix} E & F \\ G & H \end{bmatrix}$$

then:

$$XY = \begin{bmatrix} AE + BG & AF + BH \\ CE + DG & CF + DH \end{bmatrix}$$

So multiplying involves eight products with  $\frac{n}{2} \times \frac{n}{2}$  and the recurrence is:

$$T(n) = 8T\left(\frac{n}{2}\right) + O(n^2)$$

which is Case 3 of the master theorem, so runtime is  $O(n^{\log_2 8})$  which is  $O(n^3)$ , no improvement.

There actually is a trick to do better.

Claim:

$$XY = \begin{bmatrix} AE + BG & AF + BH \\ CE + DG & CF + DH \end{bmatrix}$$

is the same as:

$$\begin{bmatrix} S_1 + S_2 - S_4 + S_6 & S_4 + S_5 \\ S_6 + S_7 & S_2 - S_3 + S_5 - S_7 \end{bmatrix}$$

where:

$$S_1 = (B - D)(G + H)$$

$$S_2 = (A + D)(E + H)$$

$$S_3 = (A - C)(E + F)$$

$$S_4 = (A + B)H$$

$$S_5 = A(F - H)$$

$$S_6 = D(G - E)$$

$$S_7 = (C + D)E$$

which is only 7 products! (The additions are negligible)

So we have:

$$T(n) = 7T\left(\frac{n}{2}\right) + O(n^2)$$

Which is Case 3 of the master theorem, so  $O(n^{\log_2(7)})$  which is  $O(n^{2.81})$

## 4.6 Fast Exponentiation

Method of taking exponents in a fast way, since doing:

$$x * x * x * x \dots * x$$

is super slow.

---

```
FastExt(x,n){
  if n=1 output x
  else
    if n is even output FastExp(x, floor(n/2))^2
    if n is odd output FastExp(x, floor(n/2))^2*x
}
```

---

So our recurrence looks like:

$$T(n) = T\left(\text{floor}\left(\frac{n}{2}\right)\right) + O(1)$$

(since we're halving the problem, and doing some constant work at each step)

This is Case 2 of the Master Theorem, so the runtime is  $O(\log_2 n)$

## 5 The Median Problem

### 5.1 The Selection Problem

Want to find the  $k$ th smallest number in a set  $S$ .

Select( $S, k$ )

If  $|S| = 1$  then output  $x_1$ .

Else:

Set  $S_L$  = all numbers less than  $x_1$

Set  $S_R$  = all numbers greater than  $x_1$

If  $|S| = k - 1$  then output  $x_1$  (since if you have  $k-1$  things smaller than  $x_1$ , that can only mean  $x_1$  is the  $k$ th smallest element)

If  $|S| > k - 1$  then output Select( $S_L, k$ ) (since that means the  $k$ th smallest element must be within that set)

If  $|S| < k - 1$  then output Select( $S_R, k-1-|S_L|$ ) (-1 since you know its not  $x_1$ , and -  $|S_L|$  since you know its not in any of those, so you want the  $k-1-|S_L|$ -th element of  $S_R$ .)

The runtime of this algorithm is almost entirely dependent on the choices of pivots, since if you get a "bad" pivot every time, then you would recurse on a set of size  $n-1$ .

$$\begin{aligned} T(n) &= (n-1) + T(n-1) \\ &= . \\ &= . \\ &= . \\ &= 1/2(n(n+1)) \end{aligned}$$

which is  $O(n^2)$ .

We could instead choose our pivot randomly.

The pivot would separate the list into size  $n/4$  to  $3n/4$  with probability  $1/2$ , and so the pivot is good. So the expected running time is:

$$T(n) \leq \frac{1}{2}T\left(\frac{3n}{4}\right) + \frac{1}{2}T(n) + O(n)$$

$$\frac{1}{2}T(n) \leq \frac{1}{2}T\left(\frac{3n}{4}\right) + O(n)$$

$$T(n) \leq T\left(\frac{3n}{4}\right) + O(n)$$

which satisfies Case 1 of the master theorem which is  $O(n)$ .

But what if we want to be **certain** that the worst case will never happen?

## 5.2 Median of Medians

Divide the set  $S$  into groups of size 5. Sort each group and find the median of each group. If you were to find the median of these medians, there would always be less than  $\frac{7}{10}n$  elements in your two groups, which is pretty good. The reason this comes up is:

There's  $\frac{n}{5}$  groups overall. So there's  $\leq \frac{n}{5} \cdot \frac{1}{2} = \frac{n}{10}$  groups to the left. There's 3 elements smaller than the median in its own group, so there's  $\leq \frac{3n}{10}$  elements smaller than the median, which means there's  $\leq \frac{7n}{10}$  elements larger than the median.

So the max size of the sets is  $\frac{7n}{10}$

Finding the median of the medians is done recursively, by partitioning into 5 groups, and putting a recursive call on finding the pivot.

So the recursive formula is:

$$T(n) \leq T\left(\frac{7n}{10}\right) + T\left(\frac{n}{5}\right) + O(n)$$

Notice the Master Theorem doesn't apply here, instead we need to use the recursion tree method.

First our problem of size  $n$  is broken into two problems, one of size  $\frac{7n}{10}$  and the other size  $\frac{2n}{10}$ . Continuing down recursively, we actually get one side of the tree ending before the other. Namely, the  $7/10$  side will reach the leaves later than the  $2/10$  side.



However, up until the point that this end is reached, we're doing  $(\frac{9}{10})^l n$  work at each level. Beyond this, the work needed at each level only decreases, so it's  $\leq (\frac{9}{10})^l n$ . These terms are geometrically decreasing, so the first term dominates, and we get  $O(n)$ .

## 6 Finding the Closest Pair of Points in the Plane

How fast can we solve this?

### 6.1 Exhaustive Search

Calculate the distance between every pair of points, choose the shortest pairwise distance.  $O(n^2)$ . Is there a faster algorithm?

In one-Dimension, notice that the closest pair of points needs to be next to each other on the line. So we only need to find how far each **pair** is. ( $n - 1$  distances to calculate).

### 6.2 2-D case

Simply taking the closest in their x-coordinate (or y coordinate) doesn't work since they could be close in x but very far in y.

A divide + conquer approach is to separate the points into two groups of size  $\frac{n}{2}$ , so we want our dividing line to pass through the median x-coordinate.

We can now recursively search for the closest pairs in each group.

But what if the closest pair is **between** the two groups?

So we have to check to see if there's a better solution with an endpoint in each group. How can we do this efficiently? (This is the bottleneck step).

### 6.3 Widening the Bottleneck

Notice that by solving the subproblems recursively we can find the smallest distance between two points in both the left and right subproblems call this  $\delta$ . So we know that if a better solution exists, it will be within  $\delta$  from the dividing line.

This seems much better! But what if all the points are within  $\delta$  of the dividing line? Well then this doesn't help much.

There's actually a trick we can do.

We can break up the area into squares of size  $\frac{\delta}{2}$ , and no two points will lie in the same square. This is because if two points are in the same square, then there are on the same side of the dividing line. These points are within  $\delta \frac{\sqrt{2}}{2}$  (by construction of the boxes) from each other, but this is  $< \delta$ , so this contradicts the minimality of  $\delta$ .

We can now use this fact to derive another fact:

Suppose there's a point on either side of the dividing line with distance less than  $\delta$ . We can prove that there will be at most 10 points between them in the y-ordering. (Within the area filled with boxes).

#### Proof

Since the squares are of size  $\frac{\delta}{2}$ , then the two points are either on the same row, or one is within two rows above the other. (or else it would be further than  $\delta$ ) Now, since there can only be one point per-box, there's at most 10 points between them. (count the boxes for yourself!)

Now recall the 1-D case, we can now just look at every pairwise distance on a group where the points are at most 11 apart (rather than the ones that are next to each other as before). So you need to find the distance between a given point, and the next 11 distances.

So at most  $11n$  distances to calculate.

## 6.4 The Finished Algorithm

- Find the point with the median x-coordinate
- Partition using this point
- Recursively find the closest pair of points in each half
- Find the closest pair within the small range given by  $\delta$ , by checking the nearest 11 points (in the y-ordering) for each point.
- Among the three pairs found, (left, right, crossing) output the closest pair.

## 6.5 The Runtime (Enhanced)

Two subproblems of size  $\frac{n}{2}$ , and the work at each level is: finding the median  $O(n)$ , partitioning  $O(n)$ , making the smaller group (within  $\delta$  of dividing line)  $O(n)$ , applying the 1-D algorithm  $O(n)$ . So our recurrence looks like:

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

which is case 2 of the master theorem, so  $n \log(n)$ .