

COMP 251 Study guide

Francis Piche

January 16, 2018

Contents

1	Preliminaries	3
I	Recursive Algorithms	3
2	Divide + Conquer Algorithms	3
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
3	Master Theorem	7
3.1	Tree Method to Proof Master Theorem	10

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 "unwind" 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.