Performance Analysis 2: Asymptotic Analysis

- Choosing abstract performance measures
 - worst case, expected case
- Asymptotic growth rates
 - Why use them? Comparison in the limit. "Big O"
- Analysis of recursive programs

Reading: Weiss, Chapter 5.

Educational Aims

The aims of this topic are:

- 1. to develop a mathematical competency in describing and understanding algorithm performance, and
- 2. to begin to develop an intuitive feel for these mathematical properties.

It is essential for a programmer to be able to understand the capabilities and limitations of different data structures. Asymptotic analysis provides the foundation for this understanding (even though you would not expect to do such analysis on a regular basis).

Worst Case vs Expected Case

Abstract measures of time and space will still depend on actual input data. eg Exhaustive sequential search

Abstract time

- goal is first element in array a units
- goal is last element in array a+bn units

for some constants a and b.

Different growth rates — second measure increases with n.

What measure do we use? A number of alternatives...

Worst Case Analysis

Choose data which have the largest time/space requirements.

In the case of esearch, the worst case complexity is a + bn

Advantages

- relatively simple
- gives an upper bound, or *guarantee*, of behaviour when your client runs it it might perform better, but you can be sure it won't perform any worse

Disadvantages

- worst case could be unrepresentative might be unduly pessimistic
 - knock on effect client processes may perform below their capabilities
 - you might not get anyone to buy it!

Since we want behaviour guarantees, we will usually consider worst case analysis in this unit.

(Note there is also 'best case' analysis, as used by second-hand car sales persons and stock brokers.)

Expected Case Analysis

Ask what happens in the average, or "expected" case.

For eSearch, $a + \frac{b}{2}n$, assuming a uniform distribution over the input.

Advantages

- more 'realistic' indicator of what will happen in any given execution
- reduces effects of spurious/non-typical/outlier examples

For example, Tony Hoare's Quicksort algorithm is generally the fastest sorting algorithm in practice, despite it's worst case complexity being significantly higher than other algorithms.

Disadvantages

- only possible if we know (or can accurately guess) probability distribution over examples (with respect to size)
- more difficult to calculate
- often does not provide significantly more information than worst case when we look at growth rates
- may also be misleading

Asymptotic Growth Rates

We have talked about comparing data structure implementations — using either an empirical or analytical approach.

Focus on analytical:

- independent of run-time environment
- improves understanding of the data structures

We said we would be interested in comparisons in terms of rates of growth.

Theoretical analysis also permits a deeper comparison which the other methods don't — *comparison with the performance barrier inherent in problems*... Wish to be able to make statements like:

Searching for a given element in a block of n distinct elements using only equality testing takes n comparisons in the worst case.

Searching for a given element in an ordered list takes at least $\log n$ comparisons in the worst case.

These are *lower bounds* (in the *worst case*) — they tell us that we are never going to do any better *no matter what algorithm we choose*.

Again they reflect growth rates (linear, logarithmic)

In this section, we formalise the ideas of analytical comparison and growth rates.

Why Asymptopia

We would like to have a *simple description* of behaviour for use in comparison.

- Evaluation may be misleading. Consider the functions $t_1 = 0.002m^2$, $t_2 = 0.2m$, $t_3 = 2\log m$.
- Want a closed form. eg. $\frac{n(n+1)}{2}$ not $n + (n-1) + \dots + 2 + 1$
- Want simplicity.

Difficult to see what $2^{n-\frac{1}{n}}\log n^2 + \frac{3}{2}n^{2-n}$ does. We want to abstract away from the smaller terms...

What simple function does it behave *like*?

Solution

Investigate what simple function the more complex one *tends to* or *asymptotically approaches* as the argument approaches infinity, ie *in the limit*.

Choosing large arguments has the effect of making less important terms fade away compared with important ones.

eg. What if we want to approximate $n^4 + n^2$ by n^4 ? How much error?

n	n^4	n^2	$\frac{n^2}{n^4 + n^2}$
1	1	1	50%
2	16	4	20%
5	625	25	3.8%
10	10000	100	1%
20	160000	400	0.25%
50	6250000	2500	0.04%

Comparison "in the Limit"

How well does one function approximate another? Compare growth rates. Two basic comparisons...

1.

$$\frac{f(n)}{g(n)} \to 0 \quad \text{as} \quad n \to \infty$$

f(n) grows more slowly than g(n).

$$\frac{f(n)}{g(n)} \to 1 \quad \text{as} \quad n \to \infty$$

f(n) is asymptotic to g(n).

In fact we won't even be this picky — we'll just be concerned whether the ratio approaches a constant c > 0.

$$\frac{f(n)}{g(n)} \to c \quad \text{as} \quad n \to \infty$$

This really highlights the distinction between different orders of growth — we don't care if the constant is 0.00000000001 !

'Big O' Notation

In order to talk about comparative growth rates more succinctly we use the 'big O ' notation. . .

Definition

f(n) is O(g(n)) if there is a constant c > 0 and an integer $n_0 \ge 1$ such that, for all $n \ge n_0$,

$$f(n) \le cg(n).$$

— f "grows" no faster than g, for sufficiently large n

— growth rate of f is bounded from above by g

Show (prove) that n^2 is $O(n^3)$.

Proof

We need to show that for some c > 0 and $n_0 \ge 1$,

$$n^2 \le cn^3$$

for all $n \ge n_0$. This is equivalent to

 $1 \leq cn$

for all $n \ge n_0$. Choosing $c = n_0 = 1$ satisfies this inequality.

exercises

Show that 5n is O(3n). Choose eg c=2. Show that 143 is O(1). Choose eg c=143. Show that for any constants a and b, an^3 is $O(bn^3)$. $an^3 \leq cbn^3$ or $1 \leq c\frac{b}{a}$ or $c \geq \frac{a}{b}$.

An example

Prove that n^3 is not $O(n^2)$.

Proof (by contradiction) Assume that n^3 is $O(n^2)$. Then there exists *some* c > 0 and $n_0 \ge 1$ such that

$$n^3 < cn^2$$

for all $n \ge n_0$.

That means, $n \leq c$ for all $n \geq n_0$, but that is impossible, as n grows and c is a constant.

From these examples we can start to see that big O analysis focuses on *dominating terms*.

For example a polynomial

$$a_d n^d + a_{d-1} n^{d-1} + \dots + a_2 n^2 + a_1 n + a_0$$

 $-O(n^d)$

 $\begin{array}{l} -\text{ is }O(n^m) \text{ for any } m > d \\ -\text{ is not }O(n^l) \text{ for any } l < d. \\ \text{Here } a_d n^d \text{ is the dominating term, with } degree \ d. \end{array}$

For non-polynomials identifying dominating terms may be more difficult. Most common in CS

- polynomials $1, n, n^2, n^3, \dots$ Usually low-degree
- exponentials $2^n, \ldots$ Corresponds to computational explosion eg Travelling Salesman problem
- logarithmic $\log n, \ldots$

and combinations of these.

Analysis of Recursive Programs

Previously we've talked about:

- The power of recursive programs.
- The unavoidability of recursive programs (they go hand in hand with recursive data structures).
- The potentially high computational costs of recursive programs.

They are also the most difficult programs we will need to analyse.

It may not be too difficult to express the time or space behaviour recursively, in what we call a *recurrence relation* or *recurrence equation*, but general methods for solving these are beyond the scope of this unit.

However some can be solved by common sense!

An example

What is the time complexity of the recursive addition program given below?

```
public static int increment(int i) {return i + 1;}
public static int decrement(int i) {return i - 1;}
public static int add(int x, int y) {
    if (y == 0) return x;
    else return add(increment(x), decrement(y));
}
    if, else, ==, return, etc — constant time
    increment(x), decrement(y) — constant time
```

• add(increment(x), decrement(y))? — depends on size of y

Recursive call is same again, except y is decremented. Therefore, we know the time for add(..., y) in terms of the time for

More generally, we know the time for size n input in terms of the time for size $n - 1 \dots$

$$T(0) = a$$

 $T(n) = b + T(n-1), \quad n > 1$

This is called a *recurrence relation*.

We would like to obtain a closed form — T(n) in terms of n.

If we list the terms, its easy to pick up a pattern. . .

$$T(0) = a T(1) = a + b T(2) = a + 2b T(3) = a + 3b T(4) = a + 4b T(5) = a + 5b \vdots$$

From observing the list we can see that

$$T(n) = bn + a$$

```
public static int multiply(int x, int y) {
   if (y == 0) return 0;
   else return add(x, multiply(x, decrement(y)));
}
```

- if, else, ==, return, etc constant time
- decrement(y) constant time
- add linear in size of 2nd argument
- multiply ?

We use:

const for add terminating case
const for add recursive case
const for multiply terminating case
const for multiply recursive case
for the size of x
for the size of y
time for add with 2nd argument y
time for multiply with arguments \boldsymbol{x} and \boldsymbol{y}

Tabulate times for increasing y...

$$\begin{array}{lll} T(x,0) &=& a'\\ T(x,1) &=& b'+T(x,0)+T_{add}(0)=b'+a'+a\\ T(x,2) &=& b'+T(x,1)+T_{add}(x)=2b'+a'+xb+2a\\ T(x,3) &=& b'+T(x,2)+T_{add}(2x)=3b'+a'+(xb+2xb)+3a\\ T(x,4) &=& b'+T(x,3)+T_{add}(3x)=4b'+a'+(xb+2xb+3xb)+4a\\ &\vdots \end{array}$$

Can see a pattern of the form

$$T(x,y) = yb' + a' + [1 + 2 + 3 + \dots + (y - 1)]xb + ya$$

We would like a closed form for the term $[1 + 2 + 3 + \dots + (y - 1)]xb$. Notice that, for example

$$1 + 2 + 3 + 4 = \frac{4(4+1)}{2}$$
$$1 + 2 + 3 + 4 + 5 = \frac{5(5+1)}{2}$$

In general,

$$1 + 2 + \dots + (y - 1) = \left(\frac{(y - 1)y}{2}\right) = \frac{1}{2}y^2 - \frac{1}{2}y$$

(Prove inductively!)

Overall we get an equation of the form

$$a'' + b''y + c''xy + d''xy^2$$

for some constants a'', b'', c'', d''. Dominant term is xy^2 : linear in x (hold y constant) quadratic in y (hold x constant)

There are a number of well established results for different types of problems. We will draw upon these as necessary.

Summary

Choosing performance measures

- worst case simple, guarantees upper bounds
- expected case averages behaviour, need to know probability distribution more work, may not tell you any more
- amortized case may 'distribute' time for expensive operation over those which must accompany it

Asymptotic growth rates

- compare algorithms
- compare with inherent performance barriers
- provide simple closed form approximations
- $\bullet\,$ big O — upper bounds on growth
- big Ω lower bounds on growth

Analysis of recursive programs

- express as recurrence relation
- look for pattern to find closed form
- can then do asymptotic analysis