

COURSENOTES

**CS4104:
Data and Algorithm Analysis**

Clifford A. Shaffer
Department of Computer Science
Virginia Tech
Copyright ©1996, 1998, 2003

CS4014 Prereqs and Major Topics

What you need to already know:

- All that Discrete Math stuff
 - Proof by contradiction and induction
 - Summations
 - Set theory, relations
- The basics of Asymptotic Analysis
 - Big-oh, Big- Ω , Θ
- Most of what was covered in CS2604
 - Basic data structures
 - Algorithms for searching and sorting

What we will do:

- Finally understand upper/lower bounds
- Lower bounds proofs
- Analysis techniques (no hand waving!)
 - Recurrence Relations
- \mathcal{NP} -completeness theory

Process:

- Weekly homework sets (they are hard!)
- Work in pairs

Program Efficiency

[Week 1: Rawlins Sections 1.1–1.4]

Our primary concern is EFFICIENCY.

We want efficient programs. How do we measure the efficiency of a program? (Assume we are concerned primarily with time.)

- On what input?
- How do we speed it up?
- When do we stop speeding it up?
- Should we bother with writing the program in the first place?

Algorithm Efficiency

Since we don't want to write worthless programs, we will focus on algorithm efficiency.

We need a yardstick.

- It should measure something we care about.
- It should be quantitative, allowing comparisons.
- It should be easy to compute (the measure, not the program).
- It should be a good predictor.

We need:

- A measure for problem size.
- A measure for solution effort.
- Use key operations as a measure of solution effort.
- Total cost is a function of problem size and key operations.

Cost Model

To get a measurement, we need a model.

Example:

- Assigning to a variable takes fixed time.
- All other operations take no time.

```
sum = n*n;
```

One assignment was made, so the cost is 1.

```
sum = 0;
for (i=1; i<=n; i++)
    sum = sum + n;
```

Assignments made are $1 + \sum_{i=1}^n 1 = n + 1$.
(Depending on how you want to deal with loop variables, you might want to say it is $2n + 1$.)

```
sum = 0;
for (i=1; i<=n; i++)
    for (j=1; j<=n; j++)
        sum = sum + 1;
```

Assignments made are
 $1 + \sum_{i=1}^n \sum_{j=1}^n 1 = n^2 + 1$.

What makes a model “good”?

- Consider assigning strings by copying. Is this a good model?

Big Issues

How do we create an efficient algorithm?

How do we recognize a “good” algorithm?

How “hard” is a problem?

General Plan:

- Define a PROBLEM.
- Build a MODEL to measure the cost of a solution to the problem.
- Design an ALGORITHM to solve the problem.
- ANALYZE both the problem and the algorithm under the model.
 - Analyze an algorithm to get an UPPER BOUND.
 - Analyze a problem to get a LOWER BOUND.
- COMPARE the bounds to see if our solution is “good enough”.
 - Redesign the algorithm.
 - Tighten the lower bound.
 - Change the model.
 - Change the problem.

Problems

Our problems must be well-defined enough to be solved on computers.

An **problem** is a function (i.e., a mapping of inputs to outputs).

We have different **instances** (inputs) for the problem, where each instance has a **size**.

To **solve** a problem, we must provide an algorithm, a coding of problem instances into inputs for the algorithm, and a coding for outputs into solutions.

An **algorithm** executes the mapping.

- A proposed algorithm must work for ALL instances (give the correct mapping to the output for that input instance).

GOAL: Solve problems with as little computational effort per instance as possible.

Categories of Hard Problems

- A conceptually hard problem.
 - If we understood the problem, the algorithm might be easy. [Natural Language Processing]
 - Artificial Intelligence.
- An analytically hard problem.
 - We have an algorithm, but can't analyze its cost. [Collatz sequence]
 - Complexity Theory.
- A computationally hard problem.
 - The algorithm is expensive.
 - Class 1: No inexpensive algorithm is possible. [TOH]
 - Class 2: We don't know if an inexpensive algorithm is possible. [Traveling Salesman]
 - Complexity Theory
- A computationally unsolvable problem. [Halting problem]
 - Computability Theory.

Towers of Hanoi

Given: 3 pegs and n disks of different sizes placed in order of size on Peg 1.

Problem: Move the disks to Peg 3, given the following constraints:

- A “move” takes the topmost disk from one peg and places it on another peg (the only action allowed).
- A disk may never be on top of a smaller disk.

Model: We will measure the cost of this problem by the number of moves required.

TOH Algorithm

(This is an exercise in the process of problem solving.

Pretend that you have never seen this problem before, and that you are approaching it for the first time.)

Start by trying to solve the problem for small instances.

- 0 disks, 1 disk, 2 disks...
- When we get to 3 disks, it starts to get harder.
- Can we generalize the insight from solving for 3 disks? 4 disks?

Observation: The largest disk has no effect on the movements of the other disks. Why?

Recursive Solutions

When we generalize the TOH problem to more disks, we end up with something like:

- Move all but the bottom disk to Peg 2.
- Move the bottom disk from Peg 1 to Peg 3.
- Move the remaining disks from Peg 2 to Peg 3.

How do we deal with the $n - 1$ disks (twice)?

Forward-backward strategy: Solve simple special cases and generalize their solution, then test the generalization on other special cases.

TOH Solution

```
void TOH(int n, POLE start, POLE goal, POLE temp) {  
    if (n == 0) return;          // Base case  
    TOH(n-1, start, temp, goal); // Recurse: n-1 rings  
    move(start, goal);          // Move one disk  
    TOH(n-1, temp, goal, start); // Recurse: n-1 rings  
}
```

Algorithm Upper Bounds

Worst case cost: Maximum cost for the algorithm over all problem instances of size n .

Best case cost: Minimum cost for the algorithm over all problem instances of size n .

\mathcal{A} : The algorithm.

I_n : The set of all possible inputs to \mathcal{A} of size n .

$f_{\mathcal{A}}$: Function expressing the resource cost of \mathcal{A} .

I is an input in I_n .

$$\text{worst cost}(\mathcal{A}) = \max_{I \in I_n} f_{\mathcal{A}}(I).$$

$$\text{best cost}(\mathcal{A}) = \min_{I \in I_n} f_{\mathcal{A}}(I).$$

Examples:

- Factorial: One input of size n , one cost
- Find: Various models for number of inputs, n different costs
- Findmax: Various models for number of inputs, all cases have same cost

Average Case

We may want the average case cost. For each input of size n , we need:

- Its frequency.
- Its cost.

Given this information, we can calculate the weighted average.

Q: Can the average cost be worse than the worst cost? Or better than the best cost?

Analysis of TOH

There is only one input instance of size n .

How does this affect the decision to measure worst, best, or average case cost?

We want to count the number of moves required as a function of n .

Some facts:

- $f(1) = 1$.
- $f(2) = 3$.
- $f(3) = 7$.
- $f(n) = f(n - 1) + 1 + f(n - 1) = 2f(n - 1) + 1, \forall n \geq 4$.

(Actually, we can simplify our list of facts.)

Recurrence Relation

The following is a recurrence relation:

$$f(n) = \begin{cases} 1 & n = 1 \\ 2f(n-1) + 1 & n > 1 \end{cases}$$

How can we find a closed form solution for the recurrence?

It looks like each time we add a disk, we roughly double the cost – something like 2^n .

If we examine some simple cases, we see that they appear to fit the equation $f(n) = 2^n - 1$.

How do we prove that this ALWAYS works?

Proof for Recurrence

Let's ASSUME that $f(n - 1) = 2^{n-1} - 1$, and see what happens.

From the recurrence,

$$f(n) = 2f(n - 1) + 1 = 2(2^{n-1} - 1) + 1 = 2^n - 1.$$

Implication: if there is EVER an n for which $f(n) = 2^n - 1$, then for all greater values of n , f conforms to this rule.

This is the essence of proof by induction.

Proof by Induction

To prove by induction, we need to show two things:

- We can get started (base case).
- Being true for k implies that it is true also for $k + 1$.

Here again is the proof for TOH:

- For $n = 1$, $f(1) = 1$, so $f(1) = 2^1 - 1$.
- Assume $f(k) = 2^k - 1$, for $k < n$.
 - Then, from the recurrence we have

$$\begin{aligned} f(n) &= 2f(n-1) + 1 \\ &= 2(2^{n-1} - 1) + 1 = 2^n - 1 \end{aligned}$$

- Thus, being true for $k - 1$ implies that it is also true for k .
- Thus, we can conclude that the formula is correct for all $n \geq 1$.

Is this a good algorithm?

Lower Bound of a Problem

[Week 2: Rawlins Sections 1.4–1.8]

To decide if the algorithm is good, we need a lower bound on the cost of the PROBLEM.

We can measure the lower bound (over all possible algorithms) for the worst case, best case, or average case.

Consider a graph of cost for each possible algorithm.

- For a given problem size n , the graph shows the costs for all problem instances of size n .

The worst case lower bound is the LEAST of all the HIGHEST points on all the graphs.

\mathcal{A}_M is the set of algorithms within model M that solve the problem.

Lower Bound on Problem P

$$= \min_{\mathcal{A} \in \mathcal{A}_M} \{ \max_{I \in I_n} f_{\mathcal{A}}(I) \}$$

Lower Bound (cont.)

Lower bounds are harder than upper bounds because we must consider ALL of the possible algorithms – including the ones we don't know!

- Upper bound: How bad is the algorithm?
- Lower bound: How hard is the problem?

Lower bounds don't give you a good algorithm. They only help you know when to stop looking.

If the lower bound for the problem matches the upper bound for the algorithm (within a constant factor), then we know that we can find an algorithm that is better only by a constant factor.

Can a lower bound tell us if an algorithm is NOT optimal?

Lower Bounds for TOH

Try #1: We must move each disk at least twice, except for the largest we move once.

- $f(n) = 2n - 1$.

Is this a good match to the cost of our algorithm?

Where is the problem: the lower bound or the algorithm?

Insight #1: $f(n) > f(n - 1)$.

- Seems obvious, but why?
- Is this true for all problems?

Try #2: To move the bottom disk to Peg 3, we MUST move $n - 1$ disks to Peg 2. Then, we MUST move $n - 1$ disks back to Peg 3.

$$f(n) \geq 2f(n - 1) + 1.$$

Thus, TOH is optimal (for our model).

New Models

New model #1: We can move a stack of disks in one move.

New model #2: Not all disks start on Peg 1.

Problem Solving Algorithm

If the upper and lower bounds match,
then stop,
else if close or problem isn't important,
then stop,
else if model focuses on wrong thing,
then restate it,
else if the algorithm is too fat,
then generate slimmer algorithm,
else if lower bound is too weak,
then generate stronger bound.

Repeat until done.

Factorial Growth

Which function grows faster? $f(n) = 2^n$ or $g(n) = n!$

How about $h(n) = 2^{2n}$?

	n	1	2	3	4	5	6	7	8
$g(n)$	$n!$	1	2	6	24	120	720	5040	40320
$f(n)$	2^n	2	4	8	16	32	64	128	256
$h(n)$	2^{2n}	4	16	64	256	1024	4096	16384	65536

Consider the recurrences.

$$h(n) = \begin{cases} 4 & n = 1 \\ 4h(n-1) & n > 1 \end{cases}$$

$$g(n) = \begin{cases} 1 & n = 1 \\ ng(n-1) & n > 1 \end{cases}$$

I hope your intuition tells you the right thing.

But, how do you PROVE it?

Induction? What is the base case?

Using Logarithms

$n! \geq 2^{2n}$ iff $\lg n! \geq \lg 2^{2n} = 2n$. Why?

$$\begin{aligned} n! &= n \times (n-1) \times \cdots \times \frac{n}{2} \times \left(\frac{n}{2}-1\right) \times \cdots \times 2 \times 1 \\ &\geq \frac{n}{2} \times \frac{n}{2} \times \cdots \times \frac{n}{2} \times 1 \times \cdots \times 1 \times 1 \\ &= \left(\frac{n}{2}\right)^{n/2} \end{aligned}$$

Therefore

$$\lg n! \geq \lg \left(\frac{n}{2}\right)^{n/2} = \left(\frac{n}{2}\right) \lg \left(\frac{n}{2}\right).$$

Need only show that this grows to be bigger than $2n$.

$$\begin{aligned} &\left(\frac{n}{2}\right) \lg \left(\frac{n}{2}\right) \geq 2n \\ \iff &\lg \left(\frac{n}{2}\right) \geq 4 \\ \iff &n \geq 32 \end{aligned}$$

So, $n! \geq 2^{2n}$ once $n \geq 32$.

Now we could prove this with induction, using 32 for the base case.

- What is the tightest base case?
- How did we get such a big over-estimate?

Logs and Factorials

We have proved that $n! \in \Omega(2^{2n})$.

We have also proved that $\lg n! \in \Omega(n \lg n)$.

From here, its easy to prove that $\lg n! \in O(n \lg n)$, so $\lg n! = \Theta(n \lg n)$.

This does **not** mean that $n! = \Theta(n^n)$.

- Note that $\lg n = \Theta(\lg n^2)$ but $n \neq \Theta(n^2)$.
- The log function is a “flattener” when dealing with asymptotics.

A Simple Sum

```
sum = 0; inc = 0;
for (i=1; i<=n; i++)
  for (j=1; j<=i; j++) {
    sum = sum + inc;
    inc++;
  }
```

Use summations to analyze this code fragment.
The number of assignments is:

$$2 + \sum_{i=1}^n \left(\sum_{j=1}^i 2 \right) = 2 + \sum_{i=1}^n 2i = 2 + 2 \sum_{i=1}^n i$$

Give a good estimate.

- Observe that the biggest term is $2 + 2n$ and there are n terms, so its at most:
- Actually, most terms are much less, and its a linear ramp, so a better estimate is:

Give the exact solution.

- Of course, we all know the closed form solution for $\sum_{i=1}^n i$.
- And we should all know how to prove it using induction.
- But where did it come from?

A Problem-Specific Approach

Observe that we can “pair up” the first and last terms, the 2nd and $(n - 1)$ th terms, and so on. Each pair sums to:

The number of pairs is:

Thus, the solution is:

A Little More General

Since the largest term is n and there are n terms, the summation is less than n^2 .

If we are lucky, the solution is a polynomial.

Guess: $f(n) = c_1n^2 + c_2n + c_3$.

$f(0) = 0$ so $c_3 = 0$.

For $f(1)$, we get $c_1 + c_2 = 1$.

For $f(2)$, we get $4c_1 + 2c_2 = 3$.

Setting this up as a system of 2 equations on 2 variables, we can solve to find that $c_1 = 1/2$ and $c_2 = 1/2$.

So, if it truly is a polynomial, it **must** be

$$f(n) = n^2/2 + n/2 + 0 = \frac{n(n+1)}{2}.$$

Use induction to prove. Why is this step necessary?

Why is this not a universal approach to solving summations?

An Even More General Approach

Subtract-and-Guess or Divide-and-Guess strategies.

To solve sum f , pick a known function g and find a pattern in terms of $f(n) - g(n)$ or $f(n)/g(n)$.

Find the closed form solution for $\sum_{i=1}^n i$.
 Examples: $g_1(n) = n$; $g_2(n) = f(n - 1)$.

n	1	2	3	4	5	6	7	8
$f(n)$	1	3	6	10	15	21	28	36
$g_1(n)$	1	2	3	4	5	6	7	8
$f(n)/g_1(n)$	2/2	3/2	4/2	5/2	6/2	7/2	8/2	9/2
$g_2(n)$	0	1	3	6	10	15	21	28
$f(n)/g_2(n)$		3/1	4/2	5/3	6/4	7/5	8/6	9/7

What are the patterns?

$$\frac{f(n)}{g_1(n)} =$$

$$\frac{f(n)}{g_2(n)} =$$

Solving Summations (cont.)

Use algebra to rearrange and solve for $f(n)$

$$\frac{f(n)}{n} = \frac{n+1}{2}$$

$$\frac{f(n)}{f(n-1)} = \frac{n+1}{n-1}$$

Solving Summations (cont.)

$$\begin{aligned}\frac{f(n)}{f(n-1)} &= \frac{n+1}{n-1} \\ f(n)(n-1) &= (n+1)f(n-1) \\ f(n)(n-1) &= (n+1)(f(n)-n) \\ nf(n) - f(n) &= nf(n) + f(n) - n^2 - n \\ 2f(n) &= n^2 + n = n(n+1) \\ f(n) &= \frac{n(n+1)}{2}\end{aligned}$$

Important Note: This is **not a proof** that $f(n) = n(n+1)/2$. Why?

Growth Rates

Two functions of n have different **growth rates** if as n goes to infinity their ratio either goes to infinity or goes to zero.

Estimating Growth Rates

Exact equations relating program operations to running time require machine-dependent constants.

Sometimes, the equation for exact running time is complicated to compute.

Usually, we are satisfied with knowing an approximate growth rate.

Example: Given two algorithms with growth rate c_1n and $c_22^{n!}$, do we need to know the values of c_1 and c_2 ?

Consider n^2 and $3n$. PROVE that n^2 must eventually become (and remain) bigger.

Proof by Contradiction

Assume there are some values for constants r and s such that, for all values of n ,

$$n^2 < rn + s.$$

Then, $n < r + s/n$.

But, as n grows, what happens to s/n ?

Since n grows toward infinity, the assumption must be false.

Some Growth Rates

Since n^2 grows faster than n ,

- 2^{n^2} grows faster than 2^n .
- n^4 grows faster than n^2 .
- n grows faster than \sqrt{n} .
- $2 \lg n$ grows no slower than $\lg n$.

Since $n!$ grows faster than 2^n ,

- $n!!$ grows faster than $2^{n!}$.
- $2^{n!}$ grows faster than 2^{2^n} .
- $n!^2$ grows faster than 2^{2^n} .
- $\sqrt{n!}$ grows faster than $\sqrt{2^n}$.
- $\lg n!$ grows no slower than n .

If f grows faster than g , then

- Must \sqrt{f} grow faster than \sqrt{g} ?
- Must $\lg f$ grow faster than $\lg g$?

Fibonacci Numbers

[Week 3: Rawlins Sections 1.7–2.2]

$$f(n) = f(n - 1) + f(n - 2) \text{ for } n \geq 2;$$
$$f(0) = f(1) = 1.$$

```
int Fibr(int n) {
    if ((n <= 1) return 1;           // Base case
    return Fibr(n-1) + Fibr(n-2);    // Recursive call
}
```

```
long Fibi(int n) {
    long past, prev, curr;
    past = prev = curr = 1;         // curr holds Fib(i)
    for (int i=2; i<=n; i++) {      // Compute next value
        past = prev; prev = curr;   // past holds Fib(i-2)
        curr = past + prev;         // prev holds Fib(i-1)
    }
    return curr;
}
```

The cost of Fibi is easy to compute:

What about Fibr?

Analysis of Fibr

Use divide-and-guess with $f(n - 1)$.

n	1	2	3	4	5	6	7
$f(n)$	1	2	3	5	8	13	21
$f(n)/f(n - 1)$	1	2	1.5	1.666	1.625	1.615	1.619

Following this out, it appears to settle to a ratio of 1.618.

Assuming $f(n)/f(n - 1)$ really tends to a fixed value, let's verify what that value must be.

$$\frac{f(n)}{f(n - 2)} = \frac{f(n - 1)}{f(n - 2)} + \frac{f(n - 2)}{f(n - 2)} \rightarrow x + 1$$

For large n ,

$$\frac{f(n)}{f(n - 2)} = \frac{f(n)}{f(n - 1)} \frac{f(n - 1)}{f(n - 2)} \rightarrow x^2$$

Analysis of Fibr (cont.)

If x exists, then $x^2 - x - 1 \rightarrow 0$.

Using the quadratic equation, the only solution greater than one is

$$x = \frac{1 + \sqrt{5}}{2} \approx 1.618.$$

What does this say about the growth rate of f ?

Order Notation

little oh	$o(f(n))$	$<$	$\lim f(n)/g(n) = 0$
big oh	$O(f(n))$	\leq	
Theta	$\Theta(f(n))$	$=$	$f = O(g)$ and $g = O(f)$
Big Omega	$\Omega(f(n))$	\geq	
Little Omega	$\omega(f(n))$	$>$	$\lim g(n)/f(n) = 0$

I prefer “ f is in $O(n^2)$ ” to “ $f = O(n^2)$ ”

- While n is in $O(n^2)$ and n^2 is in $O(n^2)$,
 $O(n) \neq O(n^2)$.

Note: Big oh does not say how good an algorithm is – only how bad it CAN be.

If \mathcal{A} is in $O(n)$ and \mathcal{B} is in $O(n^2)$, is \mathcal{A} better than \mathcal{B} ?

Perhaps... but perhaps better analysis will show that $\mathcal{A} = \Theta(n)$ while $\mathcal{B} = \Theta(\lg n)$.

Manipulating Order Notation

$$\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} = r \quad \Rightarrow \quad f \text{ is in } O(g)$$

Proof:

$$\begin{aligned} \lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} = r \\ \Leftrightarrow \forall s > 0 \exists c > 0 : \forall n > c, r + s > \frac{f(n)}{g(n)} > r - s \\ \Rightarrow \forall s > 0 \exists c > 0 : \forall n > c, (r + s)g(n) > f(n). \end{aligned}$$

Therefore,

$$\exists c, r_1 > 0 : \forall n > c, f(n) \leq r_1 g(n).$$

Thus, f is in $O(g)$.

n is in $o(n^2)$ since

$$\lim_{n \rightarrow \infty} \frac{n}{n^2} = \lim_{n \rightarrow \infty} \frac{1}{n} = 0.$$

Limitations on Order Notation

Statement: Algorithm \mathcal{A} 's resource requirements grow slower than Algorithm \mathcal{B} 's resource requirements.

Is \mathcal{A} better than \mathcal{B} ?

Potential problems:

- How big must the input be?
- Some growth rate differences are trivial
 - Example: $\Theta(\lg^2 n)$ vs. $\Theta(n^{1/10})$.
- It is not always practical to reduce an algorithm's growth rate
 - Shaving a factor of n reduces cost by a factor of a million for input size of a million.
 - Shaving a factor of $\lg \lg n$ saves only a factor of 4-5.

Practicality Window

In general:

- We have limited time to solve a problem.
- We have a limited input size.

Fortunately, algorithm growth rates are USUALLY well behaved, so that Order Notation gives practical indications.

Searching

Assumptions for search problems:

- Target is well defined.
- Target is fixed.
- Probes are accurate (hit or miss).
- Search domain is finite.
- We (can) remember all information gathered during search.

We search for a record with a key.

A Search Model

Problem:

Given:

- A list L , of n elements
- A search key X

Solve: Identify one element in L which has key value X , if any exist.

Model:

- The key values for elements in L are unique.
- We can compare keys to determine $<$, $=$, $>$.
- Comparison is our only way to find ordering information.

Goal: Solve the problem using the minimum number of comparisons.

Is this a reasonable model and goal?

Linear Search

General algorithm strategy: Reduce the problem.

- Compare X to the first element.
- If not done, then solve the problem for $n - 1$ elements.

```
Position linear_search(L, lower, upper, X) {  
    if L[lower] = X then  
        return lower;  
    else if lower = upper then  
        return -1;  
    else return linear_search(L, lower+1, upper, X);  
}
```

What equation represents the worst case cost?

Worst Cost Upper Bound

$$f(n) = \begin{cases} 1 & n = 1 \\ f(n-1) + 1 & n > 1 \end{cases}$$

Reasonable to guess that $f(n) = n$.

Prove by induction:

Basis step: $f(1) = 1$, so $f(n) = n$ when $n = 1$.

Induction hypothesis: For $k < n$, $f(k) = k$.

Induction step: From recurrence,

$$\begin{aligned} f(n) &= f(n-1) + 1 \\ &= (n-1) + 1 \\ &= n \end{aligned}$$

Thus, the worst case cost for n elements is linear.

Induction is great for verifying a hypothesis.

Approach #2

What if we couldn't guess a solution?

Try: Substitute and Guess.

- Iterate a few steps of the recurrence, and look for a summation.

$$\begin{aligned}f(n) &= f(n - 1) + 1 \\ &= \{f(n - 2) + 1\} + 1 \\ &= \{\{f(n - 3) + 1\} + 1\} + 1\end{aligned}$$

Now what? Guess $f(n) = f(n - i) + i$.

When do we stop? When we reach a value for f that we know.

$$f(n) = f(n - (n - 1)) + n - 1 = f(1) + n - 1 = n$$

Now, go back and test the guess using induction.

Approach #3

Guess and Test: Guess the form of the solution, then solve the resulting equations.

Guess: $f(n)$ is linear.

$$f(n) = rn + s \text{ for some } r, s.$$

What do we know?

- $f(1) = r(1) + s = r + s = 1.$
- $f(n) = r(n) + s = r(n - 1) + s + 1.$

Solving these two simultaneous equations,
 $r = 1, s = 0.$

Final form of guess: $f(n) = n.$

Now, prove using induction.

Lower Bound on Problem

Theorem: Lower bound for the problem is n comparisons.

Proof: By contradiction.

- Assume an algorithm A exists that requires only $n - 1$ (or less) comparisons of X with elements of L .
- Since there are n elements of L , A must have avoided comparing X with $L[i]$ for some value i .
- We can feed the algorithm an input with X in position i .
- Such an input is legal in our model, so the algorithm is incorrect.

Is this proof correct?

Fixing the Proof

Error #1: An algorithm need not consistently skip position i .

Fix:

- On any given run of the algorithm, *some* element i gets skipped.
- It is possible that X is in position i at that time.

Error #2: Must allow comparisons between elements of L .

Fix:

- Include the ability to “preprocess” L .
- View L as initially consisting of n “pieces.”
- A comparison can join two pieces (without involving X).
- The total of these comparisons is k .
- We must have at least $n - k$ pieces.
- A comparison of X against a piece can reject the whole piece.
- This requires $n - k$ comparisons.
- The total is still at least n comparisons.

Average Cost

How many comparisons does linear search do on average?

We must know the probability of occurrence for each possible input.

(Must X be in L ?)

Ignore everything except the position of X in L . Why?

What are the $n + 1$ events?

$$\mathbf{P}(X \notin L) = 1 - \sum_{i=1}^n \mathbf{P}(X = L[i]).$$

Average Cost Equation

Let $k_i = i$ be the number of comparisons when $X = L[i]$.

Let $k_0 = n$ be the number of comparisons when $X \notin L$.

Let p_i be the probability that $X = L[i]$.

Let p_0 be the probability that $X \notin L[i]$ for any i .

$$\begin{aligned} f(n) &= k_0 p_0 + \sum_{i=1}^n k_i p_i \\ &= n p_0 + \sum_{i=1}^n i p_i \end{aligned}$$

What happens to the equation if we assume all p_i 's are equal (except p_0)?

Computation

$$\begin{aligned} f(n) &= p_0 n + \sum_{i=1}^n ip \\ &= p_0 n + p \sum_{i=1}^n i \\ &= p_0 n + p \frac{n(n+1)}{2} \\ &= p_0 n + \frac{1-p_0}{n} \frac{n(n+1)}{2} \\ &= \frac{n+1+p_0(n-1)}{2} \end{aligned}$$

Depending on the value of p_0 , $\frac{n+1}{2} \leq f(n) \leq n$.

Problems with Average Cost

- Average cost is usually harder to determine than worst cost.
- We really need also to know the variance around the average.
- Our computation is only as good as our knowledge (guess) on distribution.

Average Cost Lower Bound

Lower bound on average cost of problem P is:

$$\min_{A \in \mathcal{A}_{\mathcal{M}}} \left\{ \sum_{I \in I_n} P(I) f_A(I) \right\}$$

Average case lower bound for linear search:

- 1) Only compare X against elements in L :
 - The only thing that matters is the position of X in L .
 - We must compare against all elements until we find X , so linear search is optimal.

- 2) Allow comparisons of elements in L :
 - Consider comparison of $L[i]$ and $L[j]$. Do we gain any information?
 - Either we already knew that X is in the range $i..j$ or we didn't.
 - So, comparisons not involving X are redundant.

Practical Considerations

Our analysis only applies to our model.

We may still gain constant factor improvements.

- An iterative solution is usually faster.
- Put X at the end of L as a sentinel. This eliminates a test in the loop.

There may be wider context: special cases.

- Search near the position of last item found.
- Preprocess.

Don't optimize until you have a working algorithm worth keeping!

Sorted List

Change the model: Assume that the elements are in ascending order.

Is linear search still optimal? Why not?

Optimization: Use linear search, but test if the element is greater than X . Why?

Observation: If we look at $L[5]$ and find that X is bigger, then we rule out $L[1]$ to $L[4]$ as well.

More is Better: If we look at $L[n]$ and find that X is bigger, then we know in one test that X is not in L . Great!

- What is wrong here?

Jump Search

What is the right amount to jump?

Algorithm:

- Check every k 'th element ($L[k], L[2k], \dots$).
- If X is greater, then go on.
- If X is less, then use linear search on the k elements.

This is called Jump Search.

Analysis of Jump Search

If $mk \leq n < (m + 1)k$, then the total cost is at most $m + k - 1$ 3-way comparisons.

$$f(n, k) = m + k - 1 = \left\lfloor \frac{n}{k} \right\rfloor + k - 1.$$

What should k be?

$$\min_{1 \leq k \leq n} \left\{ \left\lfloor \frac{n}{k} \right\rfloor + k - 1 \right\}$$

Take the derivative and solve for $f'(x) = 0$ to find the minimum.

This is a minimum when $k = \sqrt{n}$.

What is the worst case cost?

Roughly $2\sqrt{n}$.

Lessons

We want to balance the work done while selecting a sublist with the work done while searching a sublist.

In general, make subproblems of equal effort.

This is an example of divide and conquer

What if we extend this to three levels?

- We'd jump to get a sublist, then jump to get a sub-sublist, then do sequential search
- While it might make sense to do a two-level algorithm (like jump search), it almost never makes sense to do a three-level algorithm
- Instead, we resort to recursion

Binary Search

[Week 4: Rawlins Sections 2.3–2.5]

```
int binary(int K, int* array, int left, int right) {
    // Return position of element (if any) with value K
    int l = left-1;
    int r = right+1;    // l and r beyond array bounds
    while (l+1 != r) {  // Stop when l and r meet
        int i = (l+r)/2; // Middle of remaining subarray
        if (K < array[i]) r = i;    // In left half
        if (K == array[i]) return i; // Found it
        if (K > array[i]) l = i;    // In right half
    }
    return UNSUCCESSFUL; // Search value not in array
}
```

Worst Case for Binary Search

$$f(n) = \begin{cases} 1 & n = 1 \\ f(\lfloor n/2 \rfloor) + 1 & n > 1 \end{cases}$$

Since $n/2 \geq \lfloor n/2 \rfloor$, and since $f(n)$ is assumed to be non-decreasing (why?), we can use

$$f(n) = f(n/2) + 1.$$

Alternatively, assume n is a power of 2.

Expand the recurrence:

$$\begin{aligned} f(n) &= f(n/2) + 1 \\ &= \{f(n/4) + 1\} + 1 \\ &= \{\{f(n/8) + 1\} + 1\} + 1 \end{aligned}$$

Collapse to

$$f(n) = f(n/2^i) + i = \lg n + 1$$

Now, prove it with induction.

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

Lower Bound

How does n compare to \sqrt{n} compare to $\lg n$?

Can we do better?

Model an algorithm for the problem using a decision tree.

- Consider only comparisons with X .
- Branch depending on the result of comparing X with $L[i]$.
- There must be at least n nodes in the tree. (Why?)
- Some path must be at least $\lg n$ deep. (Why?)

Thus, binary search has optimal worst cost under this model.

Average Cost of Binary Search

An estimate given these assumptions:

- X is in L .
- X is equally likely to be in any position.
- $n = 2^k$ for some non-negative integer k .

Cost?

- One chance to hit in one probe.
- Two chances to hit in two probes.
- 2^{i-1} to hit in i probes.
- $i \leq k$.

What is the equation?

Average Cost (cont.)

$$\frac{1 \times 1 + 2 \times 2 + 3 \times 4 + \dots + \lg n 2^{\lg n - 1}}{n}$$
$$= \frac{1}{n} \sum_{i=1}^{\lg n} i 2^{i-1}$$

$$\begin{aligned} \sum_{i=1}^k i 2^{i-1} &= \sum_{i=0}^{k-1} (i+1) 2^i \\ &= \sum_{i=0}^{k-1} i 2^i + \sum_{i=0}^{k-1} 2^i \\ &= 2 \sum_{i=0}^{k-1} i 2^{i-1} + 2^k - 1 \\ &= 2 \sum_{i=1}^k i 2^{i-1} - k 2^k + 2^k - 1 \end{aligned}$$

Now what? Subtract from the original!

$$\sum_{i=1}^k i 2^{i-1} = k 2^k - 2^k + 1 = (k-1) 2^k + 1.$$

Result

$$\begin{aligned}\frac{1}{n} \sum_{i=1}^{\lg n} i 2^{i-1} &= \frac{(\lg n - 1) 2^{\lg n} + 1}{n} \\ &= \frac{n(\lg n - 1) + 1}{n} \\ &\approx \lg n - 1\end{aligned}$$

So the average cost is only about one or two comparisons less than the worst cost.

If we want to relax the assumption that $n = 2^k$, we get:

$$f(n) = \begin{cases} 0 & n = 0 \\ 1 & n = 1 \\ \frac{\lceil \frac{n}{2} \rceil - 1}{n} f(\lceil \frac{n}{2} \rceil - 1) + \frac{1}{n} 0 + \\ \frac{\lfloor \frac{n}{2} \rfloor}{n} f(\lfloor \frac{n}{2} \rfloor) + 1 & n > 1 \end{cases}$$

Average Cost Lower Bound

Use decision trees again.

Total Path Length: Sum of the level for each node.

The cost of an outcome is the level of the corresponding node plus 1.

The average cost of the algorithm is the average cost of the outcomes (total path length/ n).

What is the tree with the least average depth?

This is equivalent to the tree that corresponds to binary search.

Thus, binary search is optimal.

Changing the Model

What are factors that might make binary search either unusable or not optimal?

- We know something about the distribution.
- Data are not sorted. (Preprocessing?)
- Data sorted, but probes not all the same cost (not an array).
- Data are static, know all search requests in advance.

Interpolation Search

(Also known as Dictionary Search)

Search L at a position that is appropriate to the value of X .

$$p = \frac{X - L[1]}{L[n] - L[1]}$$

Repeat as necessary to recalculate p for future searches.

Quadratic Binary Search

This is easier to analyze:

Compute p and examine $L[\lceil pn \rceil]$.

If $X < L[\lceil pn \rceil]$ then sequentially probe

$$L[\lceil pn - i\sqrt{n} \rceil], i = 1, 2, 3, \dots$$

until we reach a value less than or equal to X .

Similar for $X > L[\lceil pn \rceil]$.

We are now within \sqrt{n} positions of X .

ASSUME (for now) that this takes a constant number of comparisons.

Now we have a sublist of size \sqrt{n} .

Repeat the process recursively.

What is the cost?

QBS Probe Count

Cost is $\Theta(\lg \lg n)$ IF the number of probes on jump search is constant.

Number of comparisons needed is:

$$\sum_{i=1}^{\sqrt{n}} i \mathbf{P}(\text{need exactly } i \text{ probes})$$

This is equal to:

$$\sum_{i=1}^{\sqrt{n}} \mathbf{P}(\text{need at least } i \text{ probes})$$

We require at least two probes to set the bounds, so cost is:

$$2 + \sum_{i=3}^{\sqrt{n}} \mathbf{P}(\text{need at least } i \text{ probes})$$

QBS Probe Count (cont.)

Useful fact (Čebyšev's Inequality):

The probability that we need probe i times (P_i) is:

$$P_i \leq \frac{p(1-p)n}{(i-2)^2n} \leq \frac{1}{4(i-2)^2}$$

since $p(1-p) \leq 1/4$.

This assumes uniformly distributed data.

Final result:

$$2 + \sum_{i=3}^{\sqrt{n}} \frac{1}{4(i-2)^2} \approx 2.4112$$

Is this better than binary search?

What happened to our proof that binary search is optimal?

Comparison

Let's compare $\lg \lg n$ to $\lg n$.

n	$\lg n$	$\lg \lg n$	Diff
16	4	2	2
256	8	3	2.7
64K	16	4	4
2^{32}	32	5	6.4

Now look at the actual comparisons used.

- Binary search $\approx \lg n - 1$
- Interpolation search $\approx 2.4 \lg \lg n$

n	$\lg n - 1$	$2.4 \lg \lg n$	Diff
16	3	4.8	worse
256	7	7.2	\approx same
64K	15	9.6	1.6
2^{32}	31	12	2.6

Not done yet! This is only a count of comparisons!

- Which is more expensive: calculating the midpoint or calculating the interpolation point?

Which algorithm is dependent on good behavior by the input?

Hashing

Assume we can preprocess the data.

- How should we do it to minimize search?

Put record with key value X in $L[X]$.

If the range is too big, then use hashing.

How much can we get from this?

Simplifying assumptions:

- We hash to each slot with equal probability
- We probe to each (new) slot with equal probability
- This is called uniform hashing

Hashing Insertion Analysis

Define $\alpha = N/M$ (Records stored/Table size)

Insertion cost: sum of costs times probabilities for looking at 1, 2, ..., $N + 1$ slots

- Probability of collision on insertion?
 $\alpha = N/M$
- Probability of initial collision and another collision when probing? α^2

$$\sum_{i=0}^{i=N} i \left(\frac{N}{M}\right)^i \frac{M - N}{M}$$

Simpler formulation: Always look at least once, look at least twice with probability α , look at least three times with probability α^2 , etc.

$$\sum_{i=0}^{\infty} \alpha^i = 1 + \alpha + \alpha^2 \dots = \frac{1}{1 - \alpha}$$

How does this grow?

Searching Linked Lists

[Week 5: Rawlins Sections 2.4–3.2]

Assume the list is sorted, but is stored in a linked list.

Can we use binary search?

- Comparisons?
- “Work?”

What if we add additional pointers?

“Perfect” Skip List

Building a Skip List

Pick the node size at random (from a suitable probability distribution).

Skip List Analysis

What distribution do we want for the node depths?

```
int randomLevel(void) { // Exponential distribution
    for (int level=0; Random(2) == 0; level++); // No-op
    return level;
}
```

What is the worst cost to search in the “perfect” Skip List?

What is the average cost to search in the “perfect” Skip List?

What is the cost to insert?

What is the average cost in the “typical” Skip List?

How does this differ from a BST?

- Simpler or more complex?
- More or less efficient?
- Which relies on data distribution, which on basic laws of probability?

Other Types of Search

- Nearest neighbor (if X not in L).
- Exact Match Query.
- Range query.
- Multi-dimensional search.
- Is L static?

Is linear search on a sorted list ever better than binary search?

Selection

How can we find the i th largest value

- in a sorted list?
- in an unsorted list?

Can we do better with an unsorted list than to sort it?

Assumption: Elements can be ranked.

Properties of Relationships

Partial Order: Given a relation on a set S , and a binary operator R ,

- Asymmetric: If aRb , then it is not true that bRa .
- Transitive: If aRb and bRc , then aRc .

Example: S is Integers, R is $<$.

Example: S is the power set of $\{1, 2, 3\}$, R is subset.

A partial order is also called a **poset**.

If every pair of elements in S is relatable by R , then we have a **linear order**.

General Model

For all of our problems on Selection and Sorting:

- The poset has a linear ordering. (Usually natural numbers and a relationship of \leq .)
- Cost measure is the number of 3-way element-element comparisons.

Selection problems:

- Find the max or min.
- Find the second largest.
- Find the median.
- Find the i th largest.
- Find several ranks simultaneously.

Finding the Maximum

```
int Find_max(int *L, int low, int high) {  
    max = low;  
    for(i=low+1; i<= high; i++)  
        if(L[i] > L[max])  
            max = i;  
    return max;  
}
```

What is the cost?

Is this optimal?

Proof of Lower Bound

Try #1:

- The winner must compare against all other elements, so there must be $n - 1$ comparisons.

Try #2:

- Only the winner does not lose.
- There are $n - 1$ losers.
- A single comparison generates (at most) one (new) loser.
- Therefore, there must be $n - 1$ comparisons.

Alternative proof:

- To find the max, we must build a poset having one max and $n - 1$ losers, starting from a poset of n singletons.
- We wish to connect the elements of the poset with the minimum number of links.
- This requires at least $n - 1$ links.
- A comparison provides at most one new link.

Average Cost

What is the average cost for `Find_max`?

- Since it always does the same number of comparisons, clearly $n - 1$ comparisons.

How many assignments to `max` does it do?

Ignoring the actual values in L , there are $n!$ permutations for the input.

`Find_max` does an assignment on the i th iteration iff $L[i]$ is the biggest of the first i elements.

Since this event does happen, or does not happen:

- Given no information about distribution, the probability of an assignment after each comparison is 50%.

Average Number of Assignments

Find_max does an assignment on the i th iteration iff $L[i]$ is the biggest the first i elements.

Assuming all permutations are equally likely, the probability of this being true is $1/i$.

$$1 + \sum_{i=2}^n \frac{1}{i} \times 1 = \sum_{i=1}^n \frac{1}{i}.$$

This sum generates the n th **harmonic number**: \mathcal{H}_n .

Technique

Since $i \leq 2^{\lceil \lg i \rceil}$, $1/i \geq 1/2^{\lceil \lg i \rceil}$.

Thus, if $n = 2^k$

$$\begin{aligned}\mathcal{H}_{2^k} &= 1 + \frac{1}{2} + \frac{1}{3} + \dots + \frac{1}{2^k} \\ &\geq 1 + \frac{1}{2} + \frac{1}{4} + \frac{1}{4} + \frac{1}{8} + \frac{1}{8} + \frac{1}{8} + \frac{1}{8} \\ &\quad + \dots + \frac{1}{2^k} \\ &= 1 + \frac{1}{2} + \frac{2}{4} + \frac{4}{8} + \dots + \frac{2^{k-1}}{2^k} \\ &= 1 + \frac{k}{2}.\end{aligned}$$

Using similar logic, $\mathcal{H}_{2^k} \leq k + \frac{1}{2^k}$.

Thus, $\mathcal{H}_n = \Theta(\lg n)$.

More exactly, \mathcal{H}_n is close to $\ln n$.

Variance

How “reliable” is the average?

- How much will a given run of the program deviate from the average?

Variance: For runs of the program, average square of differences.

Standard deviation: Square root of variance.

From Čebyšev’s Inequality, 75% of the observations fall within 2 standard deviations of the average.

For `Find_max`, the variance is

$$\mathcal{H}_n - \frac{\pi^2}{6} = \ln n - \frac{\pi^2}{6}$$

The standard deviation is thus about $\sqrt{\ln n}$.

- So, 75% of the observations are between $\ln n - 2\sqrt{\ln n}$ and $\ln n + 2\sqrt{\ln n}$.
- Is this a narrow spread or a wide spread?

Finding the Second Best

[Week 6: Rawlins Section 3.3]

In a single-elimination tournament, is the second best the one who loses in the finals?

Simple algorithm:

- Find the best.
- Discard it.
- Now, find the second best of the $n - 1$ remaining elements.

Cost?

Is this optimal?

Lower bound:

- Anyone who lost to anyone who is not the max cannot be second.
- So, the only candidates are those who lost to max.
- Find_max might compare max to $n - 1$ others.
- Thus, we may need $n - 2$ additional comparisons to find second.
- Wrong!

Lower Bound for Second

The previous argument exhibits the necessity fallacy:

- Our algorithm does something, therefore all algorithms solving the problem must do the same.

Alternative: Divide and conquer

- Break the list into two halves.
- Run `Find_max` on each half.
- Compare the winners.
- Run `Find_max` on the winner's half for second.
- Compare that second to second winner.

Cost: $\lceil 3n/2 \rceil - 2$.

Is this optimal?

What if we break the list into four pieces?
Eight?

Binomial Trees

Pushing this idea to its extreme, we want each comparison to be between winners of equal numbers of comparisons.

The only candidates for second are losers to the eventual winner.

A binomial tree of height m has 2^m nodes organized as:

- a single node, if $m = 0$, or
- two height $m - 1$ binomial trees with one tree's root becoming a child of the other.

Algorithm:

- Build the tree.
- Compare the $\lceil \lg n \rceil$ children of the root for second.

Cost?

Binomial Tree Representation

We could store the binomial tree as an explicit tree structure.

We can also store the binomial tree implicitly:
In an array.

Assume two trees, each with 2^k nodes, are in the array as:

- First tree in positions 1 to 2^k .
- Second tree in positions $2^k + 1$ to 2^{k+1} .
- The root of a subtree is in the final array position for that subtree.

To join:

- Compare the roots of the subtrees.
- If necessary, swap subtrees so larger root element is second subtree.

Trades space for time.

Adversarial Lower Bounds Proofs

[Week 7: Rawlins Section 3.3-3.5]

Many lower bounds proofs use the concept of an adversary.

The adversary's job is to make an algorithm's cost as high as possible.

The algorithm asks the adversary for information about the input.

The adversary may never lie.

Imagine that the adversary keeps a list of all possible inputs.

- When the algorithm asks a question, the adversary answers, and crosses out all remaining inputs inconsistent with that answer.
- The adversary is permitted to give any answer that is consistent with at least one remaining input.

Examples:

- Hangman.
- Search an unordered list.

Lower Bound for Second Best

At least $n - 1$ values must lose at least once.

- At least $n - 1$ compares.

In addition, at least $k - 1$ values must lose to the second best.

- I.e., k direct losers to the winner must be compared.

There must be at least $n + k - 2$ comparisons.

How low can we make k ?

Adversarial Lower Bound

Call the **strength** of element $L[i]$ the number of elements $L[i]$ is (known to be) bigger than.

If $L[i]$ has strength a , and $L[j]$ has strength b , then the winner has strength $a + b + 1$.

What should the adversary do?

- Minimize the rate at which any element improves.
- Do this by making the stronger element always win.
- Is this legal?

Lower Bound (Cont.)

What should the algorithm do?

If $a \geq b$, then $2a \geq a + b$.

- From the algorithm's point of view, the best outcome is that an element doubles in strength.
- This happens when $a = b$.
- All strengths begin at zero, so the winner must make at least k comparisons for $2^{k-1} < n \leq 2^k$.

Thus, there must be at least $n + \lceil \lg n \rceil - 2$ comparisons.

Find Min and Max

Find them independantly: $2n - 2$.

- Can easily modify to get $2n - 3$.

Should be able to do better(?)

Try divide and conquer.

```
Find_Max_Min(ELEM *L, int lower, int upper) {
    if (upper == lower) return lower, lower;    // n=1
    if (upper == lower+1)                       // n=2
        return max(L[upper], L[lower]),
               min(L[upper], L[lower]);
    mid = (lower + upper)/2;                     // n>2
    max1, min1 = Find_Max_Min(L, lower, mid);
    max2, min2 = Find_Max_Min(L, mid+1, upper);
    return max(L[max1], L[max2]), min(L[min1], L[min2]);
}
```

Cost: $f(n)$ ranges between $3n/2 - 2$ and $5n/3 - 2$.

- For what sort of input does the algorithm work best?

Finding a Better Algorithm

What is the cost with six values?

What if we divide into a group of 4 and a group of 2?

With divide and conquer, we seek to minimize the work, not necessarily balance the input sizes.

When does the algorithm do its best?

What about 12? 24?

Lesson: For divide and conquer, pay attention to what happens for small n .

Algorithms from Recurrences

What does this model?

$$f(n) = \begin{cases} 0 & n = 1 \\ 1 & n = 2 \\ \min_{1 \leq k \leq n-1} \{f(k) + f(n-k)\} + 2 & n > 2 \end{cases}$$

Looking at Table 3.3 in the book, $k = 2$ looks promising.

$$f(n) = \begin{cases} 0 & n = 1 \\ 1 & n = 2 \\ f(2) + f(n-2) + 2 & n > 2 \end{cases}$$

Cost:

What is the corresponding algorithm?

The Lower Bound

Is $\lceil 3n/2 \rceil - 2$ optimal?

Consider all states that a successful algorithm must go through: The state space lower bound.

At any given instant, track the following four categories:

- Novices: not tested.
- Winners: Won at least once, never lost.
- Losers: Lost at least once, never won.
- Moderates: Both won and lost at least once.

Who can get ignored?

What is the initial state?

What is the final state?

How is this relevant?

Lower Bound (cont.)

Every algorithm must go from $(n, 0, 0, 0)$ to $(0, 1, 1, n - 2)$.

There are 10 types of comparison.

Comparing with a moderate cannot be more efficient than other comparisons, so ignore them.

If we are in state (i, j, k, l) and we have a comparison, then:

$$\begin{array}{ll} N : N & (i - 2, j + 1, k + 1, l) \\ W : W & (i, j - 1, k, l + 1) \\ L : L & (i, j, k - 1, l + 1) \\ L : N & (i - 1, j + 1, k, l) \\ & \text{or} \quad (i - 1, j, k, l + 1) \\ W : N & (i - 1, j, k + 1, l) \\ & \text{or} \quad (i - 1, j, k, l + 1) \\ W : L & (i, j, k, l) \\ & \text{or} \quad (i, j - 1, k - 1, l + 2) \end{array}$$

Adversarial Argument

What should an adversary do?

- Comparing a winner to a loser is of no value.

Only the following five transitions are of interest:

$$N : N \quad (i - 2, \quad j + 1, \quad k + 1, \quad l)$$

$$L : N \quad (i - 1, \quad j + 1, \quad k, \quad l)$$

$$W : N \quad (i - 1, \quad j, \quad k + 1, \quad l)$$

$$W : W \quad (i, \quad j - 1, \quad k, \quad l + 1)$$

$$L : L \quad (i, \quad j, \quad k - 1, \quad l + 1)$$

Only the last two types increase the number of moderates, so there must be $n - 2$ of these.

The number of novices must go to 0, and the first is the most efficient way to do this: $\lceil n/2 \rceil$ are required.

Finding the i th Best

We need to find the following poset:

We don't care about the relative order within the upper and lower groups.

Can we do better than sorting? ($\Theta(n \lg n)$)

Can we tighten the lower bound beyond n ?

What if we want to find the median element?

Splitting a List

Given an arbitrary element, split the list into those elements less and those elements greater.

```
int Split(ELEM *L, int lower, int upper, int piv_loc) {
    ELEM pivot = L[piv_loc];
    swap(L[lower], L[piv_loc]);
    piv_loc = lower;
    for (i=lower+1; i<=upper; i++)
        if (pivot > L[i]) {
            piv_loc++;
            swap(L[i], L[piv_loc]);
        }
    swap(L[lower], L[piv_loc]);
    return piv_loc;
}
```

If the pivot is i th best, we are done.

If not, solve the subproblem recursively.

Cost

What is the worst case cost of this algorithm?
Under what circumstances?

What is the average case cost if we pick the pivots at random?

Let $f(n, i)$ be the average time to find the i th best of n elements.

$$f(n, i) = n - 1 + \frac{1}{n} \sum_{k=1}^{n-i} f(n - k, i) + \frac{1}{n} 0 \\ + \frac{1}{n} \sum_{k=n-i+2}^n f(k - 1, i + k - n - 1).$$

Set $j = n - k + 1$.

$$f(n, i) = n - 1 + \frac{1}{n} \sum_{j=i+1}^n f(j - 1, i) \\ + \frac{1}{n} \sum_{j=1}^{i-1} f(n - j, i - j).$$

Let $f(n)$ be the cost averaged over all i .

$$f(n) = \frac{1}{n} \sum_{i=1}^n f(n, i).$$

Technique

$$\begin{aligned}nf(n) &= \sum_{i=1}^n f(n, i) \\ &= n^2 - n + \frac{1}{n} \sum_{i=1}^n \left\{ \sum_{j=i+1}^n f(j-1, i) + \sum_{j=1}^{i-1} f(n-j, i-j) \right\}.\end{aligned}$$

It turns out that the two double sums are the same (just going from different directions).

$$\begin{aligned}nf(n) &= n^2 - n + \frac{2}{n} \sum_{j=1}^{n-1} \sum_{i=1}^j f(j, i) \\ &= n^2 - n + \frac{2}{n} \sum_{j=1}^{n-1} jf(j)\end{aligned}$$

Technique (cont.)

Therefore,

$$n^2 f(n) = n^3 - n^2 + 2 \sum_{j=1}^{n-1} j f(j).$$

This is an example of a full history recurrence.

Solving the Recurrence

If we subtract the appropriate form of $f(n - 1)$, most of the terms will cancel out.

$$\begin{aligned}n^2 f(n) - (n - 1)^2 f(n - 1) &= n^3 - n^2 + 2 \sum_{j=1}^{n-1} j f(j) \\ &\quad - (n - 1)^3 + (n - 1)^2 - 2 \sum_{j=1}^{n-2} j f(j) \\ \Rightarrow n^2 f(n) &= 3n^2 - 5n + 2 + 2(n - 1)f(n - 1) \\ &= (n^2 - 1)f(n - 1) + 3n^2 - 5n + 2.\end{aligned}$$

Estimate:

$$\begin{aligned}n^2 f(n) &= (n^2 - 1)f(n - 1) + 3n^2 - 5n + 2 \\ &< n^2 f(n - 1) + 3n^2 \\ \Rightarrow f(n) &< f(n - 1) + 3 \\ \Rightarrow f(n) &< 3n\end{aligned}$$

Therefore, $f(n)$ is in $O(n)$.

Does this mean that the worst case is linear?

Improving the Worst Case

[Week 8: Rawlins Section 3.5-4.7]

Want worst case linear algorithm.

Goal: Pick a pivot that guarantees discarding a fixed proportion of the elements.

Can't just choose a pivot at random.

Median would be ideal – too expensive.

Choose a constant c , pick the median of a sample of size n/c elements.

Will discard at least $n/2c$ elements.

Selecting an Approximate Median

Algorithm:

- Choose the $n/5$ medians for groups of 5 elements of L .
- Recursively, select the median of the $n/5$ elements.
- Use SPLIT to partition the list into large and small elements around the “median.”

Now, the algorithm for finding the i th element uses the median finding algorithm to recursively reach the goal.

Constructive Induction

Is the following recurrence linear?

$$f(n) \leq f(\lceil n/5 \rceil) + f(\lceil (7n-5)/10 \rceil) + 6\lceil n/5 \rceil + n - 1.$$

To answer this, assume it is true for some constant r such that $f(n) \leq rn$ for all n greater than some bound.

$$\begin{aligned} f(n) &\leq f(\lceil \frac{n}{5} \rceil) + f(\lceil \frac{7n-5}{10} \rceil) + 6\lceil \frac{n}{5} \rceil + n - 1 \\ &\leq r(\frac{n}{5} + 1) + r(\frac{7n-5}{10} + 1) + 6(\frac{n}{5} + 1) + n - 1 \\ &\leq (\frac{r}{5} + \frac{7r}{10} + \frac{11}{5})n + \frac{3r}{2} + 5 \\ &\leq \frac{9r+22}{10}n + \frac{3r+10}{2}. \end{aligned}$$

This is true for $r \geq 23$ and $n \geq 380$.

Thus, we can use induction to prove that,

$$\forall n \geq 380, f(n) \leq 23n.$$

Actually, this algorithm is not practical.

Better to rely on "luck."

Changing the Model

What if we settle for the “approximate best?”

Types of guarantees, given that the algorithm produces X and the best is Y :

1. $X = Y$.

2. X 's rank is “close to” Y 's rank:

$$\text{rank}(X) \leq \text{rank}(Y) + \text{“small”}.$$

3. X is “usually” Y .

$$\mathbf{P}(X = Y) \geq \text{“large”}.$$

4. X 's rank is “usually” “close” to Y 's rank.

We often give such algorithms names:

1. Exact or deterministic algorithm.

2. Approximation algorithm.

3. Probabilistic algorithm.

4. Heuristic.

We can also sacrifice reliability for speed:

1. We find the best, “usually” fast.

2. We find the best fast, or we don't get an answer at all (but fast).

Examples for Findmax

Choose m elements at random, and pick the best.

- For large n , if $m = \lg n$, the answer is pretty good.
- Cost is $m - 1$.
- Rank is $\frac{mn}{m+1}$.

Probabilistic Algorithms

Probabilistic algorithms include steps that are affected by random events.

Problem: Pick one number in the upper half of the values in a set.

1. Pick maximum: $n - 1$ comparisons.
2. Pick maximum from just over $1/2$ of the elements: $n/2$ comparisons.

Can we do better? Not if we want a guarantee.

Probabilistic Algorithm

Pick 2 numbers and choose the greater.

This will be in the upper half with probability $3/4$.

Not good enough? Pick more numbers!

For k numbers, greatest is in upper half with probability $1 - 2^{-k}$.

Monte Carlo Algorithm: Good running time, result not guaranteed.

Las Vegas Algorithm: Result guaranteed, but not the running time.

Sorting

Initial model:

- Sort key has a linear order (comparable).
- We have an array of elements.
- We wish to sort the elements in the array.
- We get information about elements only by comparison of two elements.
- We can preserve order information only by swapping a pair of elements.

To simplify analysis:

- Assume all elements are unique.
- For analysis purposes, consider the input to be a permutation of the values 1 to n .

What if the ALGORITHM could make this assumption?

Swap Sorts

Repeatedly scan input, swapping any out-of-order elements.

Bubble sort: $O(n^2)$ in worst case.

Inversions of an element: the number of smaller elements to the right of the element.

The sum of inversions for all elements is the number of swaps required by bubblesort.

ANY algorithm that removes one inversion per swap requires at least this many swaps.

Worst number of inversions:

Best number of inversions:

Average number of inversions:

- Note that the sum of the total inversions for any permutation and its reverse is $\frac{n(n-1)}{2}$.
- Alternative view: every one of the $\frac{n(n-1)}{2}$ possible inversions occurs in a given permutation or its reverse.

Heap Sort

Heap: complete binary tree with the value of any node at least as large as its two children.

Algorithm:

- Build the heap.
- Repeat n times:
 - Remove the root.
 - Repair the heap.

This gives us list in reverse sorted order.

Since the heap is a complete binary tree, it can be stored in an array.

To delete max element:

- Swap the last element in the heap with the first (root).
- Repeatedly swap the placeholder with larger of its two children until done.

Building the heap

To build a heap, first heapify the two subheaps, then push down the root to its proper position.

Cost:

$$f(n) \leq 2f(n/2) + 2 \lg n.$$

$$\begin{aligned} f(n) &\leq \sum_{i=0}^{\lg n - 1} 2^{i+1} \lg(n/2^i) \\ &= 2 \sum_{i=0}^{\lg n - 1} 2^i (\lg n - i) \\ &= 2 \lg n \sum_{i=0}^{\lg n - 1} 2^i - 4 \sum_{i=0}^{\lg n - 1} i 2^{i-1} \\ &= 2n \lg n - 2 \lg n - 2n \lg n + 4n - 4 \\ &= 4n - 2 \lg n - 4. \end{aligned}$$

Optimization: When fixing the heap, instead of putting last element in the root, use binary search on path back up to the root.

Quicksort

Algorithm:

- Pick a pivot value.
- Split the array into elements less than the pivot and elements greater than the pivot.
- Recursively sort the sublists.

Worst case:

Pick the pivot at random, so that no particular input has bad performance.

Quicksort Average Cost

$$f(n) = \begin{cases} 0 & n \leq 1 \\ n - 1 + \frac{1}{n} \sum_{i=0}^{n-1} (f(i) + f(n - i - 1)) & n > 1 \end{cases}$$

Since the two halves of the summation are identical,

$$f(n) = \begin{cases} 0 & n \leq 1 \\ n - 1 + \frac{2}{n} \sum_{i=0}^{n-1} f(i) & n > 1 \end{cases}$$

Multiplying both sides by n yields

$$nf(n) = n(n - 1) + 2 \sum_{i=0}^{n-1} f(i).$$

Average Cost (cont.)

Get rid of the full history by subtracting $nf(n)$ from $(n+1)f(n+1)$

$$nf(n) = n(n-1) + 2 \sum_{i=1}^{n-1} f(i)$$

$$(n+1)f(n+1) = (n+1)n + 2 \sum_{i=1}^n f(i)$$

$$(n+1)f(n+1) - nf(n) = 2n + 2f(n)$$

$$(n+1)f(n+1) = 2n + (n+2)f(n)$$

$$f(n+1) = \frac{2n}{n+1} + \frac{n+2}{n+1}f(n).$$

Average Cost (cont.)

Note that $\frac{2n}{n+1} \leq 2$ for $n \geq 1$. Expanding the recurrence, we get

$$\begin{aligned} f(n+1) &\leq 2 + \frac{n+2}{n+1}f(n) \\ &= 2 + \frac{n+2}{n+1} \left(2 + \frac{n+1}{n}f(n-1) \right) \\ &= 2 + \frac{n+2}{n+1} \left(2 + \frac{n+1}{n} \left(2 + \frac{n}{n-1}f(n-2) \right) \right) \\ &= 2 + \frac{n+2}{n+1} \left(2 + \dots + \frac{4}{3} \left(2 + \frac{3}{2}f(1) \right) \right) \\ &= 2 \left(1 + \frac{n+2}{n+1} + \frac{n+2}{n+1} \frac{n+1}{n} + \dots \right. \\ &\quad \left. + \frac{n+2}{n+1} \frac{n+1}{n} \dots \frac{3}{2} \right) \\ &= 2 \left(1 + (n+2) \left(\frac{1}{n+1} + \frac{1}{n} + \dots + \frac{1}{2} \right) \right) \\ &= 2 + 2(n+2) (\mathcal{H}_{n+1} - 1) \\ &= \Theta(n \lg n). \end{aligned}$$

Lower Bound for Sorting

What is the smallest number of comparisons needed to sort n values?

Clearly, sorting is as hard as finding the min and max element: $\lceil 3n/2 \rceil - 2$.

- Why?

Information theory says that, if an algorithm uses only binary decisions to distinguish between n possibilities, then it must use at least $\lg n$ such decisions on average.

How is this relevant?

There are $n!$ permutations to the input array.

So, by information theory, we need at least $\lg n! = \Theta(n \lg n)$ comparisons.

Using the decision tree model, what is the average depth of a node?

This is also $\Theta(\lg n!)$.

Linear Insert Sort

Put the element i into a sorted list of the first $i - 1$ elements.

Worst case cost:

Best case cost:

Average case cost:

What if we use binary search? (This is called binary insert sort.)

Optimal Sorting

[Week 9: Rawlins Section 4.8, 6.1-6.3]

If we count ONLY comparisons, binary insert sort is pretty good.

What is the absolute minimum number of comparisons needed to sort?

For $n = 5$, how many comparisons do we need for binary insert sort?

Binary search is best for what values of n ?

Binary search is worst for what values of n ?

Build the following poset:

Now, put in the fifth element into the chain of 3.

Now, put in the off-element.

Total cost?

Ten Elements

Pair the elements: 5 comparisons.

Sort the winners of the pairings, using the previous algorithm: 7 comparisons.

Now, all we need to do is to deal with the original losers.

General algorithm:

- Pair up all the nodes with $\lfloor \frac{n}{2} \rfloor$ comparisons.
- Recursively sort the winners.
- Fold in the losers.

Finishing the Sort

We will use binary insert to place the losers.

However, we are free to choose the best ordering for inserting.

Recall that binary search is best for $2^k - 1$ items.

Pick the order of inserts to optimize the binary searches.

We can form an algorithm: Binary Merge.

This sort is called merge insert sort

Optimal Sort Algorithm?

Merge insert sort is pretty good, but is it optimal?

It does not match the information theoretic lower bound for $n = 12$.

- Merge insert sort gives 30 instead of 29 comparison.

BUT, exhaustive search shows that the information theoretic bound is an underestimate for $n = 12$. 30 is best.

Call the optimal worst cost for n elements $S(n)$.

- $S(n + 1) \leq S(n) + \lceil \lg(n + 1) \rceil$.
Otherwise, we would sort n elements and binary insert the last.
- For all n and m ,
 $S(n + m) \leq S(n) + S(m) + M(m, n)$ for
 $M(m, n)$ the best time to merge two sorted lists.
- For $n = 47$, we can do better by splitting into pieces of size 5 and 42, then merging.

A Truly Optimal Algorithm

Pick the best set of comparisons for size 2.

Then for size 3, 4, 5, ...

Combine them together into one program with a big case statement.

Is this an algorithm?

Numbers

Types of problems:

- Raise a number to a power.
- Find common factors for two numbers.
- Tell whether a number is prime.
- Generate a random integer.
- Multiply two integers.

These operations use all the digits, and cannot use floating point approximation.

For large numbers, cannot rely on hardware (constant time) operations.

- Measure input size by number of binary digits.
- Multiply, divide become expensive.

Analysis of Number Problems

Analysis problem: Cost may depend on properties of the number other than size.

- It is easy to check an even number for primeness.

If you consider the cost over all k -bit inputs, cost grows with k .

Summary:

- Arithmetical operations are not cheap.
- There is only one instance of value n .
- There are 2^k instances of length k .
- The size of value n is $\log n$.
- The cost may decrease when n increases in value, but generally increases when n increases in length.

Exponentiation

How do we compute m^n ?

We could multiply $n - 1$ times.

Can we do better?

Approaches to divide and conquer:

- Relate m^n to k^n for $k < m$.
- Relate m^n to m^k for $k < n$.

If n is even, then $m^n = m^{n/2}m^{n/2}$.

If n is odd, then $m^n = m^{\lfloor n/2 \rfloor}m^{\lfloor n/2 \rfloor}m$.

```
Power(base, exp) {  
    if exp = 0 return 1;  
    half = Power(base, exp/2);  
    half = half * half;  
    if (odd(exp)) then half = half * base;  
    return half;  
}
```

Analysis of Power

$$f(n) = \begin{cases} 0 & n = 1 \\ f(\lfloor n/2 \rfloor) + 1 + n \bmod 2 & n > 1 \end{cases}$$

Solution:

$$f(n) = \lfloor \lg n \rfloor + \beta(n) - 1$$

where β is the number of 1's in the binary representation of n .

How does this cost compare with the problem size?

Is this the best possible? What if $n = 15$?

What if n stays the same but m changes over many runs?

In general, finding the best set of multiplications is expensive (probably exponential).

Largest Common Factor

The largest common factor of two numbers is the largest integer that divides both evenly.

Observation: If k divides n and m , then k divides $n - m$.

So,

$$f(n, m) = f(n - m, n) = f(m, n - m) = f(m, n).$$

Observation: There exists k and l such that

$$n = km + l \text{ where } m > l \geq 0.$$

$$n = \lfloor n/m \rfloor m + n \bmod m.$$

So, $f(n, m) = f(m, l) = f(m, n \bmod m)$.

$$f(n, m) = \begin{cases} n & m = 0 \\ f(m, n \bmod m) & m > 0 \end{cases}$$

```
int LCF(int n, int m) {
    if (m == 0) return n;
    return LCF(m, n % m);
}
```

Analysis of LCF

How big is $n \bmod m$ relative to n ?

$$\begin{aligned}n \geq m &\Rightarrow n/m \geq 1 \\&\Rightarrow 2\lfloor n/m \rfloor > n/m \\&\Rightarrow m\lfloor n/m \rfloor > n/2 \\&\Rightarrow n - n/2 > n - m\lfloor n/m \rfloor = n \bmod m \\&\Rightarrow n/2 > n \bmod m\end{aligned}$$

The first argument must be halved in no more than 2 iterations.

Total cost:

Matrix Multiplication

Given: $n \times n$ matrices A and B .

Compute: $C = A \times B$.

$$c_{ij} = \sum_{k=1}^n a_{ik}b_{kj}.$$

Straightforward algorithm:

- $\Theta(n^3)$ multiplications and additions.

Lower bound for any matrix multiplication algorithm: $\Omega(n^2)$.

Another Approach

Compute:

$$m_1 = (a_{12} - a_{22})(b_{21} + b_{22})$$

$$m_2 = (a_{11} + a_{22})(b_{11} + b_{22})$$

$$m_3 = (a_{11} - a_{21})(b_{11} + b_{12})$$

$$m_4 = (a_{11} + a_{12})b_{22}$$

$$m_5 = a_{11}(b_{12} - b_{22})$$

$$m_6 = a_{22}(b_{21} - b_{11})$$

$$m_7 = (a_{21} + a_{22})b_{11}$$

Then:

$$c_{11} = m_1 + m_2 - m_4 + m_6$$

$$c_{12} = m_4 + m_5$$

$$c_{21} = m_6 + m_7$$

$$c_{22} = m_2 - m_3 + m_5 - m_7$$

7 multiplications and 18 additions/subtractions.

Strassen's Algorithm

(1) Trade more additions/subtractions for fewer multiplications in 2×2 case.

(2) Divide and conquer.

In the straightforward implementation, 2×2 case is:

$$c_{11} = a_{11}b_{11} + a_{12}b_{21}$$

$$c_{12} = a_{11}b_{12} + a_{12}b_{22}$$

$$c_{21} = a_{21}b_{11} + a_{22}b_{21}$$

$$c_{22} = a_{21}b_{12} + a_{22}b_{22}$$

Requires 8 multiplications and 4 additions.

Strassen's Algorithm (cont)

Divide and conquer step:

Assume n is a power of 2.

Express $C = A \times B$ in terms of $\frac{n}{2} \times \frac{n}{2}$ matrices.

By Strassen's algorithm, this can be computed with 7 multiplications and 18 additions/subtractions of $n/2 \times n/2$ matrices.

Recurrence:

$$\begin{aligned}T(n) &= 7T(n/2) + 18(n/2)^2 \\T(n) &= \Theta(n^{\log_2 7}) = \Theta(n^{2.81}).\end{aligned}$$

Current "fastest" algorithm is $\Theta(n^{2.376})$

Open question: Can matrix multiplication be done in $O(n^2)$ time?

Prime Numbers

How do we tell if a number is prime?

One approach is the prime sieve: Test all prime up to $\lfloor \sqrt{n} \rfloor$.

This requires up to $\lfloor \sqrt{n} \rfloor - 1$ divisions.

- How does this compare to the input size?

Note that it is easy to check the number of times 2 divides n for the binary representation

- What about 3?
- What if n is represented in trinary?

Is there a polynomial time algorithm?

Facts about Primes

Some useful theorems from Number Theory:

Prime Number Theorem: The number of primes less than n is (approximately)

$$\frac{n}{\ln n}$$

- The average distance between primes is $\ln n$.

Prime Factors Distribution Theorem: For large n , on average, n has about $\ln \ln n$ different prime factors with a standard deviation of $\sqrt{\ln \ln n}$.

To prove that a number is composite, need only one factor.

What does it take to prove that a number is prime?

Do we need to check all \sqrt{n} candidates?

Probabilistic Algorithms

Some probabilistic algorithms:

- $\text{Prime}(n) = \text{FALSE}$.
- With probability $1/\ln n$, $\text{Prime}(n) = \text{TRUE}$.
- Pick a number m between 2 and \sqrt{n} . Say n is prime iff m does not divide n .

Using number theory, we can create a cheap test that will determine that a number is composite (if it is) 50% of the time.

Algorithm:

```
Prime(n) {  
    for(i=0; i<COMFORT; i++)  
        if !CHEAPTEST(n)  
            return FALSE;  
    return TRUE;  
}
```

Of course, this does nothing to help you find the factors!

Random Numbers

[Week 10: Rawlins Section 6.6-6.7, Dynamic Programming]

Which sequences are random?

- 1, 1, 1, 1, 1, 1, 1, 1, 1, ...
- 1, 2, 3, 4, 5, 6, 7, 8, 9, ...
- 2, 7, 1, 8, 2, 8, 1, 8, 2, ...

Meanings of “random”:

- Cannot predict the next item:
unpredictable.
- Series cannot be described more briefly than to reproduce it: **equidistribution.**

There is no such thing as a random number sequence, only “random enough” sequences.

A sequence is **pseudorandom** if no future term can be predicted in polynomial time, given all past terms.

A Good Random Number Generator

Most computer systems use a deterministic algorithm to select pseudorandom numbers.

Linear congruential method:

- Pick a seed $r(1)$. Then,

$$r(i) = (r(i - 1) \times b) \text{ mod } t.$$

Resulting numbers must be in range:

What happens if $r(i) = r(j)$?

Must pick good values for b and t .

- t should be prime.

Random Number examples

$$r(i) = 6r(i - 1) \bmod 13 =$$

..., 1, 6, 10, 8, 9, 2, 12, 7, 3, 5, 4,
11, 1, ...

$$r(i) = 7r(i - 1) \bmod 13 =$$

..., 1, 7, 10, 5, 9, 11, 12, 6, 3, 8, 4,
2, 1, ...

$$r(i) = 5r(i - 1) \bmod 13 =$$

..., 1, 5, 12, 8, 1, ...
..., 2, 10, 11, 3, 2, ...
..., 4, 7, 9, 6, 4, ...
..., 0, 0, ...

Suggested generator:

$$r(i) = 16807r(i - 1) \bmod 2^{31} - 1.$$

Introduction to the Slide Rule

Compared to addition, multiplication is hard.

In the physical world, addition is merely concatenating two lengths.

Observation:

$$\log nm = \log n + \log m.$$

Therefore,

$$nm = \text{antilog}(\log n + \log m).$$

What if taking logs and antilogs were easy?

The slide rule does exactly this!

- It is essentially two rulers in log scale.
- Slide the scales to add the lengths of the two numbers (in log form).
- The third scale shows the value for the total length.

Representing Polynomials

A vector \mathbf{a} of n values can uniquely represent a polynomial of degree $n - 1$

$$P_{\mathbf{a}}(x) = \sum_{i=0}^{n-1} \mathbf{a}_i x^i.$$

Alternatively, a polynomial can be uniquely represented by a list of its values at n distinct points.

- Finding the value for a polynomial at a given point is called evaluation.
- Finding the coefficients for the polynomial given the values at n points is called interpolation.

Multiplication of Polynomials

To multiply two $n - 1$ -degree polynomials A and B normally takes $\Theta(n^2)$ coefficient multiplications.

However, if we evaluate both polynomials (at the same points), we can simply multiply the corresponding pairs of values to get the corresponding values for polynomial AB .

Process:

- Evaluate polynomials A and B at enough points.
- Pairwise multiplications of resulting values.
- Interpolation of resulting values.

This can be faster than $\Theta(n^2)$ IF a fast way can be found to do evaluation/interpolation of $2n - 1$ points.

- Normally this takes $\Theta(n^2)$ time. (Why?)

An Example

Polynomial A: $x^2 + 1$.

Polynomial B: $2x^2 - x + 1$.

Polynomial AB: $2x^4 - x^3 + 3x^2 - x + 1$.

Note that evaluating a polynomial at 0 is easy.

If we evaluate at 1 and -1, we can share a lot of the work between the two evaluations.

Can we find enough such points to make the process cheap?

$$AB(-1) = (2)(4) = 8$$

$$AB(0) = (1)(1) = 1$$

$$AB(1) = (2)(2) = 4$$

But: We need 5 points to nail down Polynomial AB. And, we also need to interpolate the 5 values to get the coefficients back.

An Observation

In general, we can write $P_a(x) = E_a(x) + O_a(x)$ where E_a is the even powers and O_a is the odd powers. So,

$$P_a(x) = \sum_{i=0}^{n/2-1} a_{2i}x^{2i} + \sum_{i=0}^{n/2-1} a_{2i+1}x^{2i+1}$$

The significance is that when evaluating the pair of values x and $-x$, we get

$$\begin{aligned} E_a(x) + O_a(x) &= E_a(x) - O_a(-x) \\ O_a(x) &= -O_a(-x) \end{aligned}$$

Thus, we only need to compute the E's and O's once instead of twice to get both evaluations.

Nth Root of Unity

The key to fast polynomial multiplication is finding the right points to use for evaluation/interpolation to make the process efficient.

Complex number z is a **primitive nth root of unity** if

1. $z^n = 1$ and
2. $z^k \neq 1$ for $0 < k < n$.

z^0, z^1, \dots, z^{n-1} are the **nth roots of unity**.

Example: For $n = 4$, $z = i$ or $z = -i$.

Identity: $e^{i\pi} = -1$.

In general, $z^j = e^{2\pi ij/n} = -1^{2j/n}$.

- Significance: We can find as many points on the circle as we need.

Evaluation

Define an $n \times n$ matrix A_z with row i and column j as

$$A_z = (z^{ij}).$$

Example: $n = 4$, $z = i$:

$$A_z = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & i & -1 & -i \\ 1 & -1 & 1 & -1 \\ 1 & -i & -1 & i \end{pmatrix}$$

Let $a = [a_0, a_1, \dots, a_{n-1}]^T$ be a vector.

We can evaluate the polynomial at the n th roots of unity:

$$F_z = A_z a = b.$$

$$b_i = \sum_{k=0}^{n-1} a_k z^{ik}.$$

Another Example

For $n = 8$, $z = \sqrt{i}$. So,

$$A_z = \begin{matrix} & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ & 1 & \sqrt{i} & i & i\sqrt{i} & -1 & -\sqrt{i} & -i & -i\sqrt{i} \\ & 1 & i & -1 & -i & 1 & i & -1 & -i \\ A_z = & 1 & i\sqrt{i} & -i & \sqrt{i} & -1 & -i\sqrt{i} & i & -\sqrt{i} \\ & 1 & -1 & 1 & -1 & 1 & -1 & 1 & -1 \\ & 1 & -\sqrt{i} & i & -i\sqrt{i} & -1 & \sqrt{i} & -i & i\sqrt{i} \\ & 1 & -i & -1 & i & 1 & -i & -1 & i \\ & 1 & -i\sqrt{i} & -i & -\sqrt{i} & -1 & i\sqrt{i} & i & \sqrt{i} \end{matrix}$$

We still have two problems:

1. We need to be able to do this fast. Its still n^2 multiplies to evaluate.
2. If we multiply the two sets of evaluations (cheap), we still need to be able to reverse the process (interpolate).

Interpolation

The interpolation step is nearly identical to the evaluation step.

$$F_z^{-1} = A_z^{-1}b' = a'.$$

What is A_z^{-1} ? This turns out to be simple to compute.

$$A_z^{-1} = \frac{1}{n}A_{1/z}.$$

In other words, do the same computation as before but substitute $1/z$ for z (and multiply by $1/n$ at the end).

So, if we can do one fast, we can do the other fast.

Fast Polynomial Multiplication

An efficient divide and conquer algorithm exists to perform both the evaluation and the interpolation in $\Theta(n \lg n)$ time.

- This is called the Discrete Fourier Transform (DFT).
- It is a recursive function that decomposes the matrix multiplications, taking advantage of the symmetries made available by doing evaluation at the n th roots of unity.

Polynomial multiplication of A and B :

- Represent an $n - 1$ -degree polynomial as $2n - 1$ coefficients:

$$[a_0, a_1, \dots, a_{n-1}, 0, \dots, 0]$$

- Perform DFT on representations for A and B
- Pairwise multiply results to get $2n - 1$ values.
- Perform inverse DFT on result to get $2n - 1$ degree polynomial AB .

Discrete Fourier Transform

```
Fourier_Transform(double *Polynomial, int n) {
    // Compute the Fourier transform of Polynomial
    // with degree n. Polynomial is a list of
    // coefficients indexed from 0 to n-1. n is
    // assumed to be a power of 2.
    double Even[n/2], Odd[n/2], List1[n/2], List2[n/2];

    if (n==1) return Polynomial[0];

    for (j=0; j<=n/2-1; j++) {
        Even[j] = Polynomial[2j];
        Odd[j] = Polynomial[2j+1];
    }
    List1 = Fourier_Transform(Even, n/2);
    List2 = Fourier_Transform(Odd, n/2);
    for (j=0; j<=n-1, J++) {
        Imaginary z = pow(E, 2*i*PI*j/n);
        k = j % (n/2);
        Polynomial[j] = List1[k] + z*List2[k];
    }
    return Polynomial;
}
```

This just does the transform on one of the two polynomials. The full process is:

1. Transform each polynomial.
2. Multiply the resulting values ($O(n)$ multiplies).
3. Do the inverse transformation on the result.

Fibonacci Revisited

Consider again the recursive function for computing the n th Fibonacci number.

```
int Fibr(int n) {
    if (n <= 1) return 1;           // Base case
    return Fibr(n-1) + Fibr(n-2);   // Recursive call
}
```

Cost is Exponential. Why?

If we could eliminate redundancy, cost would be greatly reduced.

- Keep a table

```
int Fibrt(int n, int* Values) {
    // Assume Values has at least n slots, and all
    // slots are initialized to 0
    if (n <= 1) return 1;           // Base case
    if (Values[n] != 0) return Values[n];
    Values[n] = Fibr(n-1) + Fibr(n-2); // Recursive call
    return Values[n];
}
```

Cost?

We don't need table, only last 2 values.

- Key is working bottom up.

Dynamic Programming

The issue of avoiding recomputation of subproblems comes up frequently.

- General solution: Store a table to avoid recomputation.
- Can work bottom up (fill table from smallest to largest)
- Can work top down (recursively), remembering any subproblems that happen to be solved (check table first).

This approach is called

Dynamic Programming

- Name comes from the field of dynamic control systems
- There, the act of storing precomputed values is referred to as “programmng”.

Dynamic Programming is an alternative to Divide and Conquer

- D&C: Split problem into subproblems, solve independently, and recombine.
- DP: Pay bookkeeping costs to remember solutions to shared subproblems.

A Knapsack Problem

Problem: Given an integer capacity K and n items such that item i has integer size k_i , find a subset of the n items whose sizes exactly sum to K , if possible.

Formally: Find $S \subset \{1, 2, \dots, n\}$ such that

$$\sum_{i \in S} k_i = K.$$

Example:

- $K = 163$
- 10 items of sizes 4, 9, 15, 19, 27, 44, 54, 68, 73, 101.

What if K is 164?

Instead of parameterizing problem just by n , parameterize with n and K .

- $P(n, K)$ is the problem with n items and capacity K .

Solving the Knapsack Problem

Think about divide and conquer (alternatively, induction).

What if we know how to solve $P(n - 1, K)$?

- If $P(n - 1, K)$ has a solution, then it is a solution for $P(n, K)$.
- Otherwise, $P(n, K)$ has a solution $\Leftrightarrow P(n - 1, K - k_n)$ has a solution.

What if we know how to solve $P(n - 1, k)$ for $0 \leq k \leq K$?

Cost: $T(n) = 2T(n - 1) + c$.

$T(n) = \Theta(2^n)$.

BUT... there are only $n(K + 1)$ subproblems to solve!

Solution

Clearly, there are many subproblems being solved repeatedly.

Store a $n \times K + 1$ matrix to contain the solutions for all $P(i, k)$.

Fill in the rows from $i = 0$ to n , left to right.

If $P(n - 1, k)$ has a solution, Then

$P(n, k)$ has a solution

Else If $P(n - 1, k - k_n)$ has a solution

Then $P(n, k)$ has a solution

Else $P(n, k)$ has no solution.

Cost: $\Theta(nK)$.

Knapsack Example

$K = 10$.

Five items: 9, 2, 7, 4, 1.

	0	1	2	3	4	5	6	7	8	9	10
$k_1=9$	O	–	–	–	–	–	–	–	–	I	–
$k_2=2$	O	–	I	–	–	–	–	–	–	O	–
$k_3=7$	O	–	O	–	–	–	–	I	–	I/O	–
$k_4=4$	O	–	O	–	I	–	I	O	–	O	–
$k_5=1$	O	I	O	I	O	I	O	I/O	I	O	I

Key:

–: No solution for $P(i, k)$.

O: Solution(s) for $P(i, k)$ with i omitted.

I: Solution(s) for $P(i, k)$ with i included.

I/O: Solutions for $P(i, k)$ with i included AND omitted.

Example: $M(3, 9)$ contains O because $P(2, 9)$ has a solution. It contains I because $P(2, 2) = P(2, 9 - 7)$ has a solution.

How can we find a solution to $P(5, 10)$?

How can we find ALL solutions to $P(5, 10)$?

All Pairs Shortest Paths

For every vertex $u, v \in V$, calculate $d(u, v)$.

Define a **k-path** from u to v to be any path whose intermediate vertices all have indices less than k .

```
void Floyd(Graph& G) {           // All-pairs shortest paths
    int D[G.n()][G.n()];        // Store distances
    for (int i=0; i<G.n(); i++) // Initialize D
        for (int j=0; j<G.n(); j++)
            D[i][j] = G.weight(i, j);
    for (int k=0; k<G.n(); k++) // Compute all k paths
        for (int i=0; i<G.n(); i++)
            for (int j=0; j<G.n(); j++)
                if (D[i][j] > (D[i][k] + D[k][j]))
                    D[i][j] = D[i][k] + D[k][j];
}
```

Reductions

[Week 11: Reductions, \mathcal{NP} -completeness]

A reduction is a transformation of one problem to another.

Purposes: To compare the difficulty of two problems.

- To use one algorithm to solve another problem (upper bound).
- To compare the relative difficulty of two problems (lower bound).

Notation: A problem is a mapping of inputs to outputs.

A definition looks as follows:

SORTING:

- Input: A sequence of integers x_0, x_1, \dots, x_{n-1} .
- Output: A permutation y_0, y_1, \dots, y_{n-1} of the sequence such that $y_i \leq y_j$ whenever $i < j$.

PAIRING

PAIRING:

- Input: Two sequences of integers
 $X = (x_0, x_1, \dots, x_{n-1})$ and
 $Y = (y_0, y_1, \dots, y_{n-1})$.
- Output: A pairing of the elements in the two sequences such that the least value in X is paired with the least value in Y , and so on.

How can we solve this?

One algorithm:

- Sort X .
- Sort Y .
- Now, pair x_i with y_i for $0 \leq i < n$.

Terminology: We say that PAIRING is reduced to SORTING, since SORTING is used to solve PAIRING.

PAIRING Reduction Process

The reduction of PAIRING to SORTING requires 3 steps:

- Convert an instance of PAIRING to two instances of SORTING.
- Run SORTING (twice).
- CONVERT the output for the two instances of SORTING to an output for the original PAIRING instance.

What do we require about the transformations to make them useful?

What is the cost of this algorithm?

PAIRING Lower Bound

We have an upper bound for PAIRING equal to that of SORTING.

What is the lower bound for PAIRING?

Pretend that there is a $O(n)$ time algorithm for PAIRING.

Consider this algorithm for SORTING:

- Transform SORTING to PAIRING with X being the input array for SORTING, and Y an array containing the values 0 through $n - 1$
- Run the $O(n)$ time PAIRING algorithm.
- Take the pairs output by PAIRING and use a simple binsort to order them by the second value of the pair. The first items of the pair will be the sorted list.

What is the cost of this algorithm?

What does this say about the existence of an $O(n)$ time algorithm for PAIRING?

Reduction Process

Consider any two problems for which a suitable reduction from one to the other can be found.

The first problem $P1$ takes input instance \mathbf{I} and transforms that to solution \mathbf{S} .

The second problem $P2$ takes input instance \mathbf{I}' and transforms that to solution \mathbf{S}' .

A reduction is the following three-step process:

- Transform an arbitrary instance \mathbf{I} of problem $P1$ and transform it to a (possibly special) instance \mathbf{I}' of $P2$.
- Apply an algorithm for $P2$ to \mathbf{I}' , yielding \mathbf{S}' .
- Transform \mathbf{S}' to a solution for $P1$ (\mathbf{S}). Note that \mathbf{S} MUST BE THE CORRECT SOLUTION for \mathbf{I} !

Reduction Process (Cont.)

Note that reduction is NOT an algorithm for either problem.

It does mean, given “cheap” transformations, that:

- The upper bound for $P1$ is at most the upper bound for $P2$.
- The lower bound for $P2$ is at least the lower bound for $P1$.

Another Reduction Example

How much does it cost to multiply two n -digit numbers?

- Naive algorithm requires $\Theta(n^2)$ single-digit multiplications.
- Faster (but more complicated) algorithms are known, but none so fast as $O(n)$.

Is it faster to square an n -digit number than it is to multiply two n -digit numbers?

- This is a special case, so might go faster.

Answer: No, because

$$X \times Y = \frac{(X + Y)^2 - (X - Y)^2}{4}.$$

If a fast algorithm can be found for squaring, then it could be used to make a fast algorithm for multiplying.

Matrix Multiplication

Standard matrix multiplication for two $n \times n$ matrices requires $\Theta(n^3)$ multiplications.

Faster algorithms are known, but none so fast as $O(n^2)$.

A symmetric matrix is one in which $M_{ij} = M_{ji}$.

Can we multiply symmetric matrices faster than regular matrices?

$$\begin{bmatrix} 0 & A \\ A^T & 0 \end{bmatrix} \begin{bmatrix} 0 & B^T \\ B & 0 \end{bmatrix} = \begin{bmatrix} AB & 0 \\ 0 & A^T B^T \end{bmatrix}.$$

Some Puzzles

1. A hiker leaves at 8:00 AM and hikes over the mountain. The next day, she again leaves at 8:00 AM and returns to her starting point along the same path. Prove that there is a point on the path such that she was at that point at the same time on both days.

2. Take a chessboard and cover it with dominos (a domino covers two adjacent squares of the board). Now, remove the upper left and lower right corners of the board. Now, can it still be covered with dominos?

These puzzles have the quality that, while their answers may be hard to FIND, they are easy to CHECK.

3. Is 667 composite or prime?

Complexity and Computability

Complexity:

- How cheaply can this be computed?
- How hard is this to compute?

Computability:

- When can this be computed?
- Can this be computed at all?

Types of “hard” problems:

- Hard to understand (or specify) the problem
 - Software Engineering
- Hard to design a solution
 - Artificial Intelligence
- Hard to compute in reasonable time
 - Complexity Theory
- Hard (impossible) to do at all
 - Computability Theory

Hard Problems

We say that a problem is computationally “hard” if the running time of the best known algorithm is exponential on the size of its input.

Support:

- Polynomials are closed under composition and addition.
 - Doing polynomial time operations in series is polynomial.
- All computers today are polynomially related.
 - If it takes polynomial time on one computer, it will take polynomial time on any other computer.
- Polynomial time is (generally) feasible, while exponential time is (generally) infeasible.
 - An empirical observation: For most polynomial-time algorithms, the polynomial is of low degree.

Hard Problems (Cont.)

Note that for a faster machine, the size of problem that can be run in a fixed amount of time

- grows by a multiplicative factor for a polynomial-time algorithm.
- grows by an additive factor for an exponential-time algorithm.

Nondeterminism

Imagine a computer that works by guessing the correct solution from among all possible solutions to a problem.

Alternative: Super parallel machine that tests all possible solutions simultaneously.

It might solve some problems more quickly than a regular computer.

Consider a problem which, when given a proposed solution, we can check in polynomial time if the solution is correct.

Even if the number of guesses is exponential, checking (in this case) is polynomial.

Conversely: if you can't guess an answer and check in polynomial time, there can be no polynomial time algorithm!

Nondeterministic Algorithm

An algorithm is **nondeterministic** if it works by guessing the right answer from among a finite number of choices.

Alternatively, imagine a tree of choices, polynomial levels deep.

- A super parallel machine follows all branches of the tree in parallel.
- If any single branch reaches a solution, the problem is solved.

A problem that can be solved in polynomial time by a nondeterministic machine is said to be “in \mathcal{NP} .”

Is Towers of Hanoi in \mathcal{NP} ?

Traveling Salesman Problem

TRAVELING SALESMAN (1):

- Input: A complete, directed graph G with distances assigned to each edge in the graph.
- Output: The shortest simple cycle that includes every vertex.

Problem: How to tell if a proposed solution is *shortest*?

Traveling Salesman (Cont.)

Decision problem: A problem with a YES or NO answer.

TRAVELING SALESMAN (2):

- Input: A complete, directed graph G with distances assigned to each edge in the graph, and an integer K .
- Output: YES if there is a simple cycle with total distance $\leq K$ containing every vertex in G , and NO otherwise.

In \mathcal{NP} : We can guess a cycle, and quickly check if it meets the requirements.

\mathcal{NP} -complete Problems

Many problems are like traveling salesman:

- They are in \mathcal{NP} .
- Nobody knows a polynomial time algorithm.

Is there any relationship between them?

A problem X is said to be \mathcal{NP} -hard if ANY problem in \mathcal{NP} can be reduced to X in polynomial time.

- X is AS HARD AS any problem in \mathcal{NP} .

A problem X is said to be \mathcal{NP} -complete if

1. It is in \mathcal{NP} .
2. It is \mathcal{NP} -hard.

To start the process we need to prove just one problem H is \mathcal{NP} -complete.

- To show that X is \mathcal{NP} -hard, just reduce H to X .
- DON'T GET IT BACKWARDS!

Why Care about \mathcal{NP} -Completeness?

Your boss asks you to write a fast program for TRAVELING SALESMAN.

- Its obviously an easy problem to understand.
- She can easily see some algorithm to solve the problem.
- It must be easy to speed up!

If you can't do the job, what do you tell her?

- I can't do it.
- I can't find evidence that anyone can do it.
- Nobody has been able to do it, despite the fact that many people have tried.

Furthermore, if anyone solved any of this long list of problems, then they would be able to do this problem too.

Satisfiability

Let E be a Boolean expression over variables x_1, x_2, \dots, x_n in Conjunctive Normal form:

$$E = (x_5 + x_7 + \overline{x_8} + x_{10}) \cdot (\overline{x_2} + x_3) \cdot (x_1 + \overline{x_3} + x_6).$$

SATISFIABILITY (SAT):

- INPUT: A Boolean expression E over variables x_1, x_2, \dots in Conjunctive Normal Form.
- OUTPUT: YES if there is an assignment to the variables that makes E true, NO otherwise.

This is the “grand-daddy” \mathcal{NP} -complete problem.

Cook’s Theorem: SAT is \mathcal{NP} -complete.

\mathcal{NP} -completeness Proof Model

[Week 12: \mathcal{NP} -completeness]

Implication: If a polynomial time algorithm can be found for ANY problem that is \mathcal{NP} -complete, then by a chain of polynomial time reductions, ALL \mathcal{NP} -complete problems can be solved in polynomial time.

To show that a decision problem X is \mathcal{NP} -complete:

1. Show that X is in \mathcal{NP} .
 - Give a polynomial-time, nondeterministic algorithm.
2. Show that X is \mathcal{NP} -hard.
 - Choose a known \mathcal{NP} -complete problem, A .
 - Describe a polynomial-time transformation that takes an ARBITRARY instance \mathbf{I} of A to an instance \mathbf{I}' of X .
 - Describe a polynomial-time transformation from \mathbf{S}' to \mathbf{S} such that \mathbf{S} is the solution for \mathbf{I} .

Cook's Proof Outline

1. Any decision problem can be recast as a language acceptance problem:

$$F(I) = \text{YES} \Leftrightarrow L(I') = \text{ACCEPT}.$$

2. Turing machines are a simple model of computation for writing programs that are language acceptors.
3. There is a “universal” Turing machine that can take as input a description for a Turing machine, and an input string, and return the execution of that machine on that string.
4. This in turn can be cast as a boolean expression such that the expression is satisfiable if and only if the Turing machine yields ACCEPT for that string.
5. Thus, *any* decision problem that is performable by the Turing machine is transformable to SAT: This is \mathcal{NP} -hard.

The World of Exponential-time(?) Problems

Question: Does $\mathcal{P} = \mathcal{NP}$?

3-SATISFIABILITY (3 SAT)

Input: Boolean expression E in CNF such that each clause contains exactly 3 literals.

Output: YES if the expression can be satisfied, NO otherwise.

A special case of SAT.

- Is 3 SAT easier than SAT?

Theorem: 3 SAT is \mathcal{NP} -complete.

Proof:

- 3 SAT is in \mathcal{NP} .
 - Guess (nondeterministically) values for the variables.
 - The correctness of the guess can be verified in polynomial time.
- 3 SAT is \mathcal{NP} -hard, by a reduction from SAT to 3 SAT.

3 SAT is \mathcal{NP} -hard

Find a polynomial time reduction from SAT to 3 SAT.

Let $E = C_1 \cdot C_2 \cdot \dots \cdot C_k$ be any instance of SAT.

Strategy: Replace any clause C_i that does not have exactly 3 literals with two or more clauses having exactly 3 literals.

Let $C_i = x_1 + x_2 + \dots + x_j$ where x_1, \dots, x_j are literals.

Replacement

1. $j = 1$, so $C_i = x_1$. Replace C_i with
 $(x_1 + v + w) \cdot (x_1 + \bar{v} + w) \cdot (x_1 + v + \bar{w}) \cdot (x_1 + \bar{v} + \bar{w})$
where v and w are new variables.
2. $j = 2$, so $C_i = (x_1 + x_2)$. Replace C_i with
 $(x_1 + x_2 + z) \cdot (x_1 + x_2 + \bar{z})$
where z is a new variable.
3. $j > 3$. Replace C_i with
 $(x_1 + x_2 + z_1) \cdot (x_3 + \bar{z}_1 + z_2) \cdot (x_4 + \bar{z}_2 + z_3) \cdot \dots$
 $\cdot (x_{j-2} + \bar{z}_{j-4} + z_{j-3}) \cdot (x_{j-1} + x_j + \bar{z}_{j-3})$
where z_1, \dots, z_{j-3} are new variables.

After appropriate replacements have been made for each C_i , a Boolean expression results that is an instance of 3 SAT.

Each replacement is satisfiable if and only if the original clause is satisfiable.

The reduction is clearly polynomial time.

Third Case

If E is satisfiable, then E' is satisfiable:

- Assume x_m is assigned true.
- Then assign $z_t, t \leq m - 2$ as true and $z_k, t \geq m - 1$ as false.
- Then all clauses in Case (3) are satisfied.

If E' is satisfiable, then E is satisfiable:

- Proof by contradiction.
- If x_1, x_2, \dots, x_j are all false, then z_1, z_2, \dots, z_{j-3} are all true.
- But then $(x_{j-1} + x_{j-2} + \overline{z_{j-3}})$ is false, a contradiction.

(Not necessary for proof, but may help insight.)
Conversely, if E is not satisfiable, then E' is not satisfiable.

- E not satisfiable means all x_i are false.
- This leaves E' as

$$(z_1) \cdot (\overline{z_1} + z_2) \cdot \dots \cdot (\overline{z_{j-4}} + z_{j-3}) \cdot (\overline{z_{j-3}})$$

which is NOT satisfiable.

Two Problems

VERTEX COVER:

Input: An undirected graph G and an integer k .

Output: YES if there is a subset of vertices in G of size k or less such that every edge in the graph has at least one of its ends in the subset; NO otherwise.

K-CLIQUE:

Input: An undirected graph G and an integer k .

Output: YES if there is a subset of the vertices of size k or greater that is a complete graph (a clique).

We can reduce either problem to the other by switching G to its inverse G' .

- If edge (i, j) is in G , it is NOT in G' .
- If edge (i, j) is NOT in G , it IS in G' .

K CLIQUE is \mathcal{NP} -Complete

Easy to show that K CLIQUE is in \mathcal{NP} .

Reduce SAT to K CLIQUE.

An instance of SAT is a Boolean expression

$$B = C_1 \cdot C_2 \cdot \dots \cdot C_m$$

where

$$C_i = y[i, 1] + y[i, 2] + \dots + y[i, k_i].$$

Transform this to an instance of K CLIQUE as follows.

$$V = \{v[i, j] \mid 1 \leq i \leq m, 1 \leq j \leq k_i\}.$$

All vertices $v[i_1, j_1]$ and $v[i_2, j_2]$ have an edge between them UNLESS they are two literals within the same clause ($i_1 = i_2$) OR they are opposite values for the same variable.

Set $k = m$.

Example

$$B = (y_1 + y_2) \cdot (\overline{y_1} + y_2 + y_3).$$

B is satisfiable if and only if G has a clique of size $\geq k$.

- B satisfiable implies there is a truth assignment such that $y[i, j_i]$ is true for each i .
- But then, $v[i, j_i]$ must be in a clique of size $k = m$.
- If G has a clique of size $\geq k$, then the clique must have size exactly k and there is one vertex $v[i, j_i]$ in the clique for each i .
- There is a truth assignment making each $y[i, j_i]$ true. That truth assignment satisfies B .

We conclude that K CLIQUE is \mathcal{NP} -hard, therefore \mathcal{NP} -complete.

Coping with \mathcal{NP} -Completeness

1. Organize to reduce costs.
 - Dynamic programming.
 - Backtracking.
 - Branch and Bounds.
2. Find subproblems of the original problem that have polynomial-time solutions.
 - Significant special cases that are useful to answer.
3. Approximation algorithms.
4. Randomized algorithms.
5. Use heuristics.
 - Greedy algorithms.
 - Simulated Annealing.
 - Genetic Algorithms.

Knapsack Analysis Revisited

Fact: Knapsack is \mathcal{NP} -complete.

But we have a $\Theta(nK)$ algorithm!!

Question: How big is K ?

- Input size is typically $O(n \lg K)$ since the item sizes are smaller than K .
- Thus, $\Theta(nK)$ is exponential on input size.

This algorithm is tractable if the numbers are “reasonable.”

- nK can be thousands.
- This is different from TRAVELING SALESMAN which cannot handle $n = 100$.

Such an algorithm is called a pseudo-polynomial time algorithm.

Subproblems and Special Cases

Some restricted cases of \mathcal{NP} -complete problems are useful, and not \mathcal{NP} -complete.

- VERTEX COVER and K CLIQUE have polynomial time algorithms for bipartite graphs.
- 2-SATISFIABILITY has a polynomial time solution.
- Several geometric problems are polynomial-time in two dimensions, but not in three or more.
- KNAPSACK is polynomial if the numbers are “small.”

Approximation Algorithms

Seek algorithms for optimization problems with a guaranteed bound on quality of the solution.

For VERTEX COVER:

- Let M be a maximal (not necessarily maximum) **matching** in G .
 - A matching pairs vertices (with connecting edges) so that no vertex is paired with more than one match.
 - Maximal means pick as many pairs as possible.
- If OPT is the size of a minimum vertex cover, then

$$|M| \leq 2 \cdot OPT$$

because at least one endpoint of every matched edge must be in ANY vertex cover.

BIN PACKING

INPUT: Numbers x_1, x_2, \dots, x_n between 0 and 1, and an unlimited supply of bins of size 1.

OUTPUT: An assignment of numbers to bins that requires the fewest possible number of bins (no bin can hold numbers whose sum exceeds 1).

This problem is \mathcal{NP} -complete.

Heuristic: First fit

- Place a number in the first bin that fits.
- The number of bins used is no more than twice the sum of the numbers.
- First fit can be much worse than optimal.
- Consider 6 of $1/7 + \epsilon$, 6 of $1/3 + \epsilon$, 6 of $1/2 + \epsilon$.

Better Heuristic: Decreasing first fit

- Sort the items, then use first fit.
- This can be proven to yield no more than $11/9$ the optimal number of bins.

Summary

The theory of \mathcal{NP} -completeness gives a technique for separating tractable from (probably) untractable problems.

When faced with a new problem, we might alternate between:

- Check if it is tractable (find a fast solution).
- Check if it is intractable (prove the problem is \mathcal{NP} -complete).

If the problem is in \mathcal{NP} -complete, then use one of the “coping” strategies.

A General Model

[Week 13: Turing Machines]

Want a general model of computation that is as simple as possible.

- Wish to be able to reason about the model.
- “State machines” are simple.

Necessary features:

- Read
- Write
- Compute

Turing Machines

A tape, divided into squares.

A single I/O head:

- Read current symbol
- Change current symbol

Control Unit Actions:

- Put the control unit into a new state.
- Either:
 1. Write a symbol in current tape square.
 2. Move I/O head one square left or right.

Tape has a fixed left end, infinite right end.

- Machine ceases to operate if head moves off left end.
- By convention, input is placed on left end of tape.

A halt state (h) signals end of computation.

“#” indicates a blank tape square.

Formal definition of Turing Machine

A Turing Machine is a quadruple (K, Σ, δ, s) where

- K is a finite set of states (not including h).
- Σ is an alphabet (containing $\#$, not L or R).
- $s \in K$ is the initial state.
- δ is a function from $K \times \Sigma$ to $(K \cup \{h\}) \times (\Sigma \cup \{L, R\})$.

If $q \in K$, $a \in \Sigma$ and $\delta(q, a) = (p, b)$, then when in state q and scanning a , enter state p and

1. If $b \in \Sigma$ then replace a with b .
2. Else (b is L or R) Move head.

Turing Machine Example 1

$M = (K, \Sigma, \delta, s)$ where

- $K = \{q_0, q_1\}$,
- $\Sigma = \{a, \#\}$,
- $s = q_0$,

	q	σ	$\delta(q, \sigma)$
• $\delta =$	q_0	a	$(q_1, \#)$
	q_0	$\#$	$(h, \#)$
	q_1	a	(q_0, a)
	q_1	$\#$	(q_0, R)

Turing Machine Example 2

$M = (K, \Sigma, \delta, s)$ where

- $K = \{q_0\}$,
- $\Sigma = \{a, \#\}$,
- $s = q_0$,

- $\delta = \begin{array}{ccc|c} & q & \sigma & \delta(q, \sigma) \\ \hline & q_0 & a & (q_0, L) \\ & q_0 & \# & (h, \#) \end{array}$

Notation

Configuration: $(q, aaba\underline{\#}\underline{\#}a)$

Halted configuration: q is h .

Hanging configuration: Move left from leftmost square.

A computation is a sequence of configurations for some $n \geq 0$. Such a computation is of length n .

$$\begin{array}{l} (q_0, \underline{a}aaa) \vdash_M (q_1, \underline{\#}aaa) \\ \vdash_M (q_0, \underline{\#}aaa) \\ \vdash_M (q_1, \underline{\#}\underline{\#}aa) \\ \vdash_M (q_0, \underline{\#}\underline{\#}aa) \\ \vdash_M (q_1, \underline{\#}\underline{\#}\underline{\#}a) \\ \vdash_M (q_0, \underline{\#}\underline{\#}\underline{\#}a) \\ \vdash_M (q_1, \underline{\#}\underline{\#}\underline{\#}\underline{\#}) \\ \vdash_M (q_0, \underline{\#}\underline{\#}\underline{\#}\underline{\#}\underline{\#}) \\ \vdash_M (h, \underline{\#}\underline{\#}\underline{\#}\underline{\#}\underline{\#}) \end{array}$$

Computations

Convention:

M is said to halt on input w iff $(s, \#w\#)$ yields some halted configuration.

M is said to hang on input w if $(s, \#w\#)$ yields some hanging configuration.

Turing machines compute functions from strings to strings.

Formally: Let f be a function from Σ_0^* to Σ_1^* . Turing machine M is said to compute f if for any $w \in \Sigma_0^*$, if $f(w) = u$ then

$$(s, \#w\#) \vdash_M^* (h, \#u\#).$$

f is said to be a Turing-computable function.

Multiple parameters: $f(w_1, \dots, w_k) = u$,
 $(s, \#w_1\#w_2\#\dots\#w_k\#) \vdash_M^* (h, \#u\#)$.

Functions on Natural Numbers

Represent numbers in unary notation on symbol I (zero is represented by the empty string).

$f : \mathbb{N} \rightarrow \mathbb{N}$ is computed by M if M computes $f' : \{I\}^* \rightarrow \{I\}^*$ where $f'(I^n) = I^{f(n)}$ for each $n \in \mathbb{N}$.

Example: $f(n) = n + 1$ for each $n \in \mathbb{N}$.

q	σ	$\delta(q, \sigma)$
q_0	I	(h, R)
q_0	$\#$	(q_0, I)

$$(q_0, \#II\#) \vdash_M (q_0, \#III\#) \vdash_M (h, \#III\#).$$

In general, $(q_0, \#I^n\#) \vdash_M^* (h, \#I^{n+1}\#)$.

What about $n = 0$?

Turing-decidable Languages

A language $L \subset \Sigma_0^*$ is Turing-decidable iff function $\chi_L : \Sigma_0^* \rightarrow \{\boxed{Y}, \boxed{N}\}$ is Turing-computable, where for each $w \in \Sigma_0^*$,

$$\chi_L(w) = \begin{cases} \boxed{Y} & \text{if } w \in L \\ \boxed{N} & \text{otherwise} \end{cases}$$

Ex: Let $\Sigma_0 = \{a\}$, and let $L = \{w \in \Sigma_0^* : |w| \text{ is even}\}$.

M erases the marks from right to left, with current parity encode by state. Once blank at left is reached, mark \boxed{Y} or \boxed{N} as appropriate.

Turing-acceptable Languages

M **accepts** a string w if M halts on input w .

- M accepts a language iff M halts on w iff $w \in L$.
- A language is **Turing-acceptable** if there is some Turing machine that accepts it.

Ex: $\Sigma_0 = \{a, b\}$,

$L = \{w \in \Sigma_0^* : w \text{ contains at least one } a\}$.

q	σ	$\delta(q, \sigma)$
q_0	a	(h, a)
q_0	b	(q_0, L)
q_0	$\#$	(q_0, L)

Every Turing-decidable language is Turing-acceptable.

Combining Turing Machines

Lemma: If

$$(q_1, w_1 \underline{a_1} u_1) \vdash_M^* (q_2, ww_2 \underline{a_2} u_2)$$

for string w and

$$(q_2, w_2 \underline{a_2} u_2) \vdash_M^* (q_3, w_3 \underline{a_3} u_3),$$

then

$$(q_1, w_1 \underline{a_1} u_1) \vdash_M^* (q_3, ww_3 \underline{a_3} u_3).$$

Insight: Since $(q_2, w_2 \underline{a_2} u_2) \vdash_M^* (q_3, w_3 \underline{a_3} u_3)$, this computation must take place without moving the head left of w_2

- The machine cannot “sense” the left end of the tape

Combining Turing Machines (Cont)

Thus, the head won't move left of w_2 even if it is not at the left end of the tape.

This means that Turing machine computations can be combined into larger machines:

- M_2 prepares string as input to M_1 .
- M_2 passes control to M_1 with I/O head at end of input.
- M_2 retrieves control when M_1 has completed.

Some Simple Machines

Basic machines:

- $|\Sigma|$ symbol-writing machines (one for each symbol).
- Head-moving machines R and L move the head appropriately.

More machines:

- First do M_1 , then do M_2 or M_3 depending on current symbol.
- (For $\Sigma = \{a, b, c\}$) Move head to the right until a blank is found.
- Find first blank square to left: $L_{\#}$
- Copy Machine: Transform $\#w\#$ into $\#w\#w\#$.
- Shift a string left or right.

Extensions

The following extensions do not increase the power of Turing Machines.

- 2-way infinite tape
- Multiple tapes
- Multiple heads on one tape
- Two-dimensional “tape”
- Non-determinism

Partial Recursive Functions

[Week 14: Partial recursive functions, uncomputability]

A function is Partial Recursive if it can be generated from the following operations.

- $zero() = 0$.
- $succ(x) = x + 1$ for all $x \in \mathbb{N}$.
- $p_i(x_1, \dots, x_n) = x_i$ for $1 \leq i \leq n$.
- Composition: for partial recursive functions h, g_1, \dots, g_m , we can form a new partial recursive function f as follows:

$$f(x) = h(g_1(x), \dots, g_m(x)).$$

We can use any number of arguments in place of x .

- Primitive Recursion: We can construct new function f from partial recursive functions h and g as follows:

$$\begin{aligned} f(x, 0) &= h(x), \\ f(x, succ(y)) &= g(x, y, f(x, y)). \end{aligned}$$

Parital Recursive Functions (Cont.)

unbounded minimalization: for g a partial recursive function,

$$f(x) = \min_y (g(x, y) = 0)$$

that is, the least value y that satisfies $g(x, y) = 0$.

Examples of Partial Recursive Functions

All constants: Apply the successor function to the zero function a constant number of times.

Addition:

$$n + 0 = n.$$

$$n + (m + 1) = (n + 1) + m.$$

Multiplication:

$$n \cdot 0 = 0.$$

$$n \cdot (m + 1) = n \cdot m + n.$$

Non-total function:

$$h(x, y) = \min_z(x = y + z).$$

Gödelization

We can convert numbers to strings by using Unary notation.

We can convert strings to unique numbers as well.

Assign each character in Σ a unique number.

- The i th character has value i .

Denote the i th prime number as p_i .

Represent string $S = s_{i_1}s_{i_2}\dots s_{i_k}$ of length k as:

$$g(s_{i_1}s_{i_2}\dots s_{i_k}) = p_1^{i_1} p_2^{i_2} \dots p_k^{i_k}.$$

Example: CAT is $2^3 3^1 5^{20}$.

Now we can use partial recursive functions to compute functions from strings to strings.

Church's Thesis

We now have two “general purpose” forms of computation: Turing machines and partial recursive functions.

Either can be transformed into the other. Thus, they are “equally powerful.”

We can create a Universal Turing machine:

- input is a specification for a Turing machine M and a string w .
- output is the output from M on w .

Church's Thesis (also **Church-Turing Thesis**): Turing machines are formal versions of algorithms, and no computational procedure will be considered an algorithm unless it can be presented as a Turing machine.

A “thesis” because it is unprovable: asserts that a certain informal concept (algorithm) corresponds to a mathematical concept (Turing machine).

Church's Thesis (cont)

How to disprove Church's Thesis: Devise another model of computation that does finite labor at each step, and which can compute functions not computable by any Turing machine.

Uncomputable Functions

The number of programs is countably infinite

- Can associate a natural number with each possible string, thus each possible program.

Consider the set of functions $f(x) = y$ for x, y natural numbers.

The set of such functions is uncountably infinite.

Diagonalization argument:

Thus, not every function on natural numbers has a program that computes it.

Turing Decidable vs. Turing Acceptable Languages

Any Turing Decidable language is Turing Acceptable.

Problem: Given a Turing Machine M and a string w , will M halt when run on w ?

If there is a Turing machine M_0 that solves this problem, then an algorithm to convert a Turing-accepting machine M_1 to a Turing-deciding machine is:

- Feed M_1 and its input w to M_0 .
- If the answer is “halts” then output \boxed{Y} .
- Otherwise, output \boxed{N} .

This crucial question is called the Halting Problem.

Halting Problem for Programs

Can a **C** program be written to solve the following problem?

Halting Problem:

- Input: A program P and input X .
- Output: “Halts” if P halts when run with X as input. “Does not Halt” otherwise.

Halting Problem Proof

Theorem: There is no program to solve the Halting Problem.

Proof: (by contradiction).

Assumption: There is a **C** program that solves the Halting Problem.

```
bool halt(char* prog, char* input)
{
    Code to solve halting problem
    if (prog does halt on input) then
        return(TRUE);
    else
        return(FALSE);
}
```

Two More Procedures

```
bool selfhalt(char *prog) {
    // Return TRUE if program halts
    // when given itself as input.
    if (halt(prog, prog))
        return(TRUE);
    else
        return(FALSE);
}

void contrary(char *prog) {
    if (selfhalt(prog))
        while(TRUE); // Go into an infinite loop
}
```

The Punchline

What happens when function `contrary` is run on itself?

Case 1: `selfhalt` returns `TRUE`.

- `contrary` will go into an infinite loop.
- This contradicts the result from `selfhalt`.

`selfhalt` returns `FALSE`.

- `contrary` will halt.
- This contradicts the result from `selfhalt`.

Either result is impossible.

The only flaw in this argument is the assumption that `halt` exists.

Therefore, `halt` cannot exist.

Computability Reduction Proof

Turing machine form of Halting Problem: Given Machine M and string w , does M halt on w ?

- Accept proof for programs that this is noncomputable.

Given arbitrary machine M , does it halt on the EMPTY input? (Λ is the EMPTY string.)

This is uncomputable. Proof:

- Suppose that machine M_0 decided the language

$$\{L(M) : M \text{ accepts } \Lambda\}.$$

- Given arbitrary machine M and string w , we can create a new machine M_w that operates as follows on empty input:
 - Write w on the tape.
 - Simulate the execution of M .
- So, we can take arbitrary machine M and string w , create M_w , and invoke M_0 on M_w (with empty input) to solve the original halting problem.
- Thus, M_0 must not exist.

Another Reduction Proof

Does there exist ANY input for which an arbitrary program (machine) halts?

Proof that this is uncomputable:

- Suppose that machine M_0 decided the language

$$\{L(M) : M \text{ accepts some string } w\}.$$

- We can take an arbitrary machine M and string w , and modify it so that it erases its input before proceeding.
- Thus, arbitrary machine M is modified to be M' that effectively is M operating on the empty input.
- Thus, we can take arbitrary machine M and string w , modify it to become M' and feed that to M_0 to solve the problem of deciding if M halts on the empty input.
- We already know that is undecidable.
- Thus, M_0 cannot exist.

Other Noncomputable Functions

Does a program (machine) halt on EVERY input?

Do two programs (machines) compute the SAME function?

Does a particular line (transition) in a program (machine) get executed?

Does a program (machine) compute a particular function?

Does a program contain a “computer virus”?