

Fundamentals of algorithm and asymptotic analysis

① algorithms

1) what is an algorithm?

↳ (computational problem): a (possibly infinite) set of (x, S) pairs, where x is a valid input and S is a set of valid solutions for x

↳ an algorithm provides the "how" — a well defined procedure that outputs a valid solution, given some input x

2) what do we want in an algorithm?

↳ (correctness): an algorithm is said to be correct iff for every valid input, it halts in a valid solution. This is a "worst case" notion of correctness, that it is always correct

↳ (running time): the point of running time analysis is to determine how the number of steps executed by an algorithm scales in input size

3) notions of input size

array → number of array elements
integer → number of bits in binary
graph → $|V| + |E|$

} always important to define

② models of computation (how we define a step)

comparison model

Definition 1.4.1. In the comparison model, the input is an array A of n numbers, and an algorithm can compare two elements ("Is $x > y$, $x < y$, or $x = y$?) in one time unit. No other⁴ operations on the elements are allowed. The array can be manipulated (e.g., permuted or broken into two sub-arrays) without operating on the elements at no cost.

query model

Definition 1.4.9. In the string query model, the input is a string of n bits. In one time unit, an algorithm can query one bit of the string. All other operations on the string can be performed at no cost.

The total number of bits of the input string queried by an algorithm is called its *query complexity*.

Definition 1.4.11. In the graph query model, the input is the symmetric adjacency matrix G of an n -node undirected graph. In one time unit, an algorithm can query one entry in the matrix G . All other operations on G are free.

word RAM model

The framework we mostly use subsequently is a conventional **word RAM** model, where the ‘programming language’ consists of sequential instructions and the running time is the total number of instructions executed. What constitutes an ‘instruction’? We omit a formal specification, but it consists of the typical basic operations found in standard programming languages: arithmetic operations ($+, -, \times, /$, mod, $<$, $>$, $[.]$, ...), memory access operations, and control flow operations (**if**, **return**, ...). Each instruction is allowed to operate on a *word* of data, and the word size is limited. (Clearly, it would be useless to have a model where each instruction could operate on inputs of unrestricted size in one time unit.) We typically assume that numbers (e.g., entries in an array) are each of word-size so that they can be operated on in unit time.

③ analysis of lower bounds / to be vigorous, we should also argue for the existence of the algorithm that runs in N^* steps, vs $< N^* \Rightarrow \min \geq N^*$.

(Adversary arguments) used to prove lower bounds for the running time in concrete computation models. The adversary, by adjusting his responses, makes sure that if an algorithm takes too few steps, then it can come up in at least two different solutions that cannot be differentiated by the algorithm.

input decided on-the-fly by adversary who keeps options open about what actual input is. Adversary ensures that if algorithm takes too few steps, then there are at least two valid inputs which are consistent in the results of the algorithm's steps, yet the solutions for the two inputs are different.

eg **Claim 1.4.3.** Any algorithm solving Max must have running time $\geq n - 1$ in the comparison model.

Proof. Fix any algorithm \mathcal{M} that correctly solves Max on all inputs. The idea of the proof is to show that if \mathcal{M} always makes less than $n - 1$ comparisons, then there are two arrays A and A' so that the two cannot be differentiated based on \mathcal{M} 's comparisons. However, the maximum elements of A and A' occur at different locations, and so, \mathcal{M} must err on either A or A' , a contradiction.

① Take an input A for which \mathcal{M} makes $< n - 1$ comparisons. Construct a graph G on n nodes (indexed by $1, \dots, n$), where nodes i and j are adjacent iff \mathcal{M} compares A_i and A_j . Since G has $< n - 1$ edges, it is disconnected. That is, there exists a partition of the nodes into C_1 and C_2 such that for any $i \in C_1, j \in C_2$, there is no edge between i and j , and so, A_i and A_j are not compared by \mathcal{M} .

② Suppose given A as input, \mathcal{M} outputs A_{i^*} as the maximum. Without loss of generality, let $i^* \in C_1$. Now, consider a new array A' , such that if $i \in C_1$, $A'_i = A_i$, but if $i \in C_2$, $A'_i = A_i + m$ where m is sufficiently large that the maximum element of A' is in C_2 . See an example in Figure 1.1. ③

④ Observe that the comparisons made by \mathcal{M} cannot distinguish between A and A' , as the only comparisons which are different are between $i \in C_1$ and $j \in C_2$, and \mathcal{M} does not make such comparisons. So, \mathcal{M} must give a wrong answer for either A or A' . ⑤ □

- (A)
1. Take an input
 2. Describe some model with which algo M operates on input in $< N^*$ operations, and describe its conclusion
 3. Modify input such that solution is different + proof that solutions different
 4. Show that M gives same solution for both inputs

- (B)
1. Describe algo, and describe how adversary replies
 2. Show that in $< N^*$ steps, there are at least two different inputs that are consistent in algo
 3. Show that inputs (and their solutions) are diff, so one must be wrong

In the context of the graph query model, a problem is evasive if it requires $\binom{n}{2}$ queries.

Claim 1.4.12. *The problem of determining whether a graph is connected is evasive.*

Proof. We again use an adversary argument. Let \mathcal{M} be an algorithm making $m < \binom{n}{2}$ queries which decides whether its input is connected or not. It issues a sequence of queries $(i_1, j_1), (i_2, j_2), \dots, (i_m, j_m)$, and the adversary replies with b_1, b_2, \dots, b_m .

We now define how the adversary replies to each query. Suppose the adversary receives the ℓ 'th query for $1 \leq \ell \leq m$. The adversary's strategy for choosing b_ℓ is as follows. She considers the graph $G^{(\ell)}$ defined by: $G^{(\ell)}[i_1, j_1] = b_1, \dots, G^{(\ell)}[i_{\ell-1}, j_{\ell-1}] = b_{\ell-1}, G^{(\ell)}[i_\ell, j_\ell] = 0$ and all other entries equal to 1. That is, $G^{(\ell)}$ is the graph which is consistent with the previous query replies, does not have the edge (i_ℓ, j_ℓ) and has all other unqueried edges. If $G^{(\ell)}$ is connected, she replies $b_\ell = 0$; otherwise, she replies $b_\ell = 1$.

② Now, after \mathcal{M} makes all its m queries, we can define two graphs. The first graph, which we call G_0 , is consistent with the adversary's replies but sets all the unqueried edges to 0. The second graph, which we call G_1 , is also consistent with the adversary's replies but sets all the unqueried edges to 1. \mathcal{M} cannot distinguish between G_0 and G_1 because both are consistent with its queries. We argue below that G_0 is disconnected while G_1 is connected, and hence, \mathcal{M} must err on one of them, finishing the proof.

Clearly, G_1 is connected by the definition of the adversary's strategy.

Lemma 1.4.13. G_0 is disconnected.

Proof. Consider a pair of nodes i and j such that (i, j) was not queried by \mathcal{M} . We claim that there is no path between i and j in G_0 .

Let's argue by contradiction. Suppose there is a path between i and j in G_0 . This path only consists of queried edges to which the adversary replied 1. Let (i', j') be the edge on this path that was queried last by \mathcal{M} . Now, observe that according to the adversary's strategy given above, there is a path from i' to j' in the graph G_{con} she considers at this stage. This is because in G_{con} , the path from i' to i , the path from j to j' , and the edge from i to j are all present. Thus, she must reply 0 to the query (i', j') , a contradiction! \square

Decision trees

\hookrightarrow node : comparison
 \hookrightarrow branch : outcome of comparison
 \hookrightarrow leaf : output / solution

we can use a decision tree to model any comparison based algorithm. worst case running time = height

- 1. A height $-h$ binary tree has $\leq 2^h$ leaves.
- 2. A height $-h$ binary tree has $\leq 2^{h+1} - 1$ nodes

$$h = \lfloor \lg \text{leaves} \rfloor + 1$$

$$\text{no. of nodes} = \underbrace{2^0 + 2^1 + 2^2 + \dots}_{h} = 2^h$$

$$\text{GP: } \frac{a(r^h - 1)}{r-1} = \frac{(1)(2^{h+1} - 1)}{(2-1)}$$

Theorem: Any decision tree that can sort n elements must have height $\Omega(n \lg n)$

Claim 1: A height h binary tree has $\leq 2^h$ leaves

Claim 2: The decision tree must contain $n!$ leaves

Recall: Worst-case running time is the **height** of the decision tree.

Ask ourselves: We have $n!$ leaves. What's our *minimum* height? How to relate n and h ? It's $n! \leq 2^h$

$$\begin{aligned}
 h &\geq \lg(n!) \quad (\lg \text{ is monotonically increasing}) \\
 &\geq \lg((n/e)^n) \quad (\text{Stirling's formula}) \\
 &= n \lg n - n \lg e \\
 &= \Omega(n \lg n). \quad \blacksquare
 \end{aligned}$$

④ asymptotic analysis

1) asymptotic bounds

ideq: we can describe how our algorithm's steps scale w input size by using functions to describe their bounds and shape. In particular, we use asymptotic bounds — i.e. as input scales past a certain point, can we bound $f(n)$ by some scaled function $c \cdot g(n)$?

(Big O) $\exists c > 0, n_0 > 0$ s.t. $\forall n > n_0, 0 \leq f(n) \leq c \cdot g(n)$. $f(n) = O(g(n))$

↳ there is at least one (c, n_0) pair where $f(n) \leq c \cdot g(n)$ after n_0 .
i.e. there is at least one way you can tweak $g(n)$ s.t. $g(n)$ overtakes $f(n)$ at some point

(small O) $\forall c > 0, \exists n_0 > 0$ s.t. $\forall n > n_0, 0 \leq f(n) < c \cdot g(n)$. $f(n) = o(g(n))$

↳ for all c , there is a point after which $f(n)$ is strictly bound by $c \cdot g(n)$
i.e. no matter how you tweak the slope, $g(n)$ will overtake $f(n)$ at some point

(Big omega) $\exists c > 0, n_0 > 0$ s.t. $\forall n > n_0, f(n) \geq c \cdot g(n) \geq 0$. $f(n) = \Omega(g(n))$

↳ there is at least one (c, n_0) pair where $f(n) \geq c \cdot g(n)$ after n_0 .
i.e. there is at least one way you can tweak $g(n)$ so that $f(n)$ overtakes $g(n)$ at some point

(small omega) $\forall c > 0, \exists n_0 > 0$ s.t. $\forall n > n_0, f(n) > c \cdot g(n) \geq 0$. $f(n) = \omega(g(n))$

↳ for all c , there is a point after which $f(n)$ is strictly greater than $c \cdot g(n)$
i.e. no matter how you tweak $g(n)$, $f(n)$ will overtake $g(n)$ at some point ↗

(meta) $f(n) = \Theta(g(n))$ iff. $f(n) = O(g(n))$ and $f(n) = \Omega(g(n))$

2) asymptotic bounds as sets

$$O(g(n)) = \{f(n) : \exists c, n_0 > 0 \text{ such that } 0 \leq f(n) \leq c \cdot g(n) \text{ for all } n \geq n_0\}$$

$$\Omega(g(n)) = \{f(n) : \exists c, n_0 > 0 \text{ such that } 0 \leq c \cdot g(n) \leq f(n) \text{ for all } n \geq n_0\}$$

$$\begin{aligned} \Theta(g(n)) &= \{f(n) : \exists c_1, c_2, n_0 > 0 \text{ such that } 0 \leq c_1 \cdot g(n) \leq f(n) \leq c_2 \cdot g(n) \text{ for all } n \geq n_0\} \\ &= O(g(n)) \cap \Omega(g(n)) \end{aligned}$$

$$o(g(n)) = \{f(n) : \forall c > 0, \exists n_0 > 0 \text{ such that } 0 \leq f(n) < c \cdot g(n) \text{ for all } n \geq n_0\}$$

$$\omega(g(n)) = \{f(n) : \forall c > 0, \exists n_0 > 0 \text{ such that } 0 \leq c \cdot g(n) < f(n) \text{ for all } n \geq n_0\}$$

↳ observe that $f(n)$ cannot simultaneously be $O(g(n))$ and $\omega(g(n))$
i.e. small - big cannot occur simultaneously, nor small - small.

3) properties of asymptotic bounds

(L'Hopital's rule)

$$\lim_{x \rightarrow a} \frac{f(x)}{g(x)} = \frac{0}{0} \quad \text{OR} \quad \lim_{x \rightarrow a} \frac{f(x)}{g(x)} = \frac{\pm\infty}{\pm\infty}$$

$$\lim_{x \rightarrow a} \frac{f(x)}{g(x)} = \lim_{x \rightarrow a} \frac{f'(x)}{g'(x)}$$

(reflexivity)

$$\begin{aligned} f(n) &= O(f(n)) \\ f(n) &= \Omega(f(n)) \\ f(n) &= \Theta(f(n)) \end{aligned}$$

(limits) → provide different "views" of f(n)

$$\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} = 0 \implies f(n) = o(g(n))$$

$$\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} < \infty \implies f(n) = O(g(n))$$

$$0 < \lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} < \infty \implies f(n) = \Theta(g(n))$$

$$\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} > 0 \implies f(n) = \Omega(g(n))$$

$$\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} = \infty \implies f(n) = \omega(g(n))$$

(symmetry)

$$f(n) = \Theta(g(n)) \iff g(n) = \Theta(f(n))$$

(transitivity)

$$f(n) = \mathcal{O}(g(n)) \& g(n) = \mathcal{O}(h(n))$$



$$f(n) = \mathcal{O}(h(n))$$

(complementarity)

$$f(n) = O(g(n)) \iff g(n) = \Omega(f(n))$$

$$f(n) = o(g(n)) \iff g(n) = \omega(f(n))$$

5) Lemmas, formulae and techniques

(upper bound thin) commonly done also to relate terms

↪ substitute something that is bigger to get a more uniform expression

$$\begin{aligned} \lg \lg n! &= (\lg(n)(n-1)\dots(1)) \text{ eq. } n^2 + 3n \lg n \\ &= \lg(n) + \lg(n-1) \dots \leq n^2 + 3n^2 \\ &\leq \lg n + \lg n \dots \leq 4n^2 \\ &= n \lg n \end{aligned}$$

(arithmetic series)

$$1+2+3\dots+n = \frac{n(n+1)}{2}$$

$$a_1 + a_2 + \dots + a_n = \frac{n(a_1 + a_n)}{2}$$

(geometric series)

$$\sum_{i=0}^n a \cdot r^i = a \left(\frac{1 - r^{n+1}}{1 - r} \right)$$

(harmonic series)

$$\sum_{k=1}^n \frac{1}{k} = \frac{1}{1} + \frac{1}{2} + \dots + \frac{1}{n} = 2n \ln n$$

(sum of squares)

$$\sum_{k=1}^n k^2 = 1 + 4 + 9 + \dots + n^2 = \frac{n(n+1)(2n+1)}{6}$$

sum of numbers

$$\sum_{k=1}^n k^3 = 1^3 + 2^3 + \dots + n^3 = \frac{n^2(n+1)^2}{4}$$

integral summation trick

↳ when sum of functions, integrate bound by 1, use definite integrals.

$$\text{eg. } \sum_{i=a}^n f(i) = f(a) + f(a+1) + \dots + f(n) \\ \leq \sum_{i=a}^n 1 = \int_a^{n+1} f(u) du$$

(Factorial)

$$\begin{aligned} \lg(n!) &= \lg(1) + \lg(2) + \lg(3) + \dots + \lg(n) \\ &= \lg(1) + \lg(\lg(2)) + \lg(\lg(3)) + \dots + \lg(\lg(n)) \\ &= \lg(\lg(1) \times \lg(2) \times \lg(3) \times \dots \times \lg(n)) \Rightarrow \lg(\lg(n)!) \end{aligned}$$

⑥ correctness of iterative algorithms

- ↳ in simple and deterministic algorithms, proving correctness is often trivial.
- ↳ loops are a common construct in most algorithms — how can we prove their correctness?

(loop invariant)

for a loop running for N iterations, a loop invariant is a sequence of assertions S_1, S_2, \dots, S_N such that each S_i holds at the start of the i th iteration and such that S_N implies the desired correctness property.

1. initialisation : it is true prior to the first iteration of the loop
2. maintenance : if it is true before the start of the i th iteration, it remains true before the start of the $i+1$ th iteration
3. termination : when the loop terminates, the invariant gives us a useful property that helps show the algorithm is correct

0. define some property

1. base case : loop initializes n _____. show that invariant is true.
2. induction step : assuming invariant is true for i th loop, show that with values carried over, invariance holds for $i+1$ th loop
3. termination : algorithm terminates when _____ (cases?) , show that invariant is still correct after termination. summarize result.

```

INSERTION-SORT( $A$ )
1  for  $j = 2$  to  $A.length$ 
2     $key = A[j]$ 
3    // Insert  $A[j]$  into the sorted sequence  $A[1..j - 1]$ .
4     $i = j - 1$ 
5    while  $i > 0$  and  $A[i] > key$ 
6       $A[i + 1] = A[i]$ 
7       $i = i - 1$ 
8     $A[i + 1] = key$ 

```

Claim 3.1.1 (Inner loop invariant). Let A' be the array A before line 5 in an outer loop iteration with $j = J$. At the start of an inner loop iteration with $i = I$, the following hold:

1. $A[1..I] = A'[1..I]$.
2. $A[I + 2..J] = A'[I + 1..J - 1]$.
3. All elements in $A[I + 2..J]$ are greater than $A'[J] = key$.

Proof. The proof is by induction. The claim is trivial for the base case $I = J - 1$. Suppose the claim holds for I , and the inner loop continues to $I - 1$, we want to argue for $I - 1$. By the induction hypothesis, at the start of the run with $i = I$,

- (1) $A[1..I] = A'[1..I]$
- (2) $A[I + 2..J] = A'[I + 1..J - 1]$
- (3) $A[I + 2..J] > A'[J] = key$

Suppose I continues to $I - 1$. Condition 1 holds true because: $A[1..I] = A'[1..I]$ and (1) gives $A[1..I - 1] = A'[1..I - 1]$. Condition 2 holds true because: $A[I] = A'[I]$ (from condition 1) and $A[I + 1] = A[I]$ (after an inner loop iteration), therefore $A[I + 1] = A'[I]$, combined with (2) gives $A[I + 1..J] = A'[I..J - 1]$, which is equivalent to $A[(I - 1) + 2..J] = A'[(I - 1) + 1..J - 1]$. Condition 3 holds true because: $A'[I] = A[I] > key$ and $A[I + 1] = A'[I]$ (after an inner loop iteration), therefore $A[I + 1] > key$, combined with (3) gives $A[I + 1..J] > key$, which is equivalent to $A[(I - 1) + 2..J] > key$. \square

Suppose the inner loop terminates when $i = I$. This can happen because of two reasons, either $I = 0$ or $A[I] \leq key$. In the former case, Conditions 2 and 3 of the inner loop invariant gives $A[2..J] = A'[1..J - 1]$ and $A[2..J] > key$. In the latter case, the inner loop invariant gives: $A[1..I] = A'[1..I]$, $A[I + 2..J] = A'[I + 1..J - 1]$ and $A[I + 2..J] > key$.

Claim 3.1.2 (Outer loop invariant). At the start of an outer loop iteration with $j = J$, $A[1..J - 1]$ is the sorted list of elements originally in $A[1..J - 1]$.

Proof. The proof is by induction. The claim is trivial for the base case $J = 2$. Suppose the claim holds for J , and we want to argue for $J + 1$. By the induction hypothesis, at the start of the run with $j = J$, $A[1..J - 1]$ is a sorted arrangement of the original elements in these positions. Call this array A' . After the inner loop terminates in this run of the outer loop, by the paragraph after [Claim 3.1.1](#), there can be two cases:

- $I = 0$. In this case, $key < A[2] \leq \dots \leq A[J]$ where $A[2..J]$ is a right-shift of $A'[1..J - 1]$.
- $A[I] \leq key$. In this case, $A[1] \leq \dots \leq A[I] \leq key$ because these elements are the same as in A' and they were sorted, while $key < A[I + 2] \leq \dots \leq A[J]$ because these elements are a right-shift of $A'[I + 1..J - 1]$ which were sorted.

In either case, line 8 ($A[I + 1] = key$) ensures that $A[1..J]$ is sorted and consists of the elements in $A'[1..J - 1]$ and key . \square

Applying [Claim 3.1.2](#) for $J = A.length + 1$ proves the correctness of [INSERTION-SORT](#).

Notes

1. comparing functions, tricks. $n^{f(n)} = 2^{\lg n^{f(n)}} = 2^{f(n) \lg n} \Rightarrow$ compare powers!
2. $\lg[(\lg n)!] \rightarrow$ translates to addtl. $\lg n \Rightarrow \lg n \cdot \lg \lg n$ because of Stirling's approx.

Analysis of recursive algorithms

① Divide and conquer

↳ divide and conquer is an algorithm design paradigm that involves the following:

(divide) the problem into sub-problems that are smaller instances of the original problem

(conquer) the subproblems by solving them recursively (at some base level)

(combine) the solutions in a way that is correct to get a solution to the original problem

② correctness of recursive algorithms

- ↳ idea is based on induction. A recursive algorithm solves problems by splitting them into smaller problems until some base case is met, then combines them upwards.
- ↳ some recursive algorithms do the work at divide (e.g. quicksort, local maxima) while others do it at combine (e.g. mergesort)
- ↳ to prove correctness, in general, if applicable, using strong induction to show that:

1. (divide) step does indeed reach the base case eventually, and maintains some invariant property (e.g. quicksort - subarrays on left all \leq pivot)

2. (conquer) step maintains some invariant property of base problems, based on divide that got it to reach base step (e.g. mergesort, single element trivially sorted)

3. (combine) step (if applicable) maintains some invariant property, based on property of previous combine (e.g. if subarrays sorted, merge will return another sorted array)

3.1 show that combine is correct for base case (1)

3.2 show that combine is correct for inductive steps by strong MI

The **merge sort** algorithm closely follows the divide-and-conquer paradigm. Intuitively, it operates as follows.

Divide: Divide the n -element sequence to be sorted into two subsequences of $n/2$ elements each.

Conquer: Sort the two subsequences recursively using merge sort.

Combine: Merge the two sorted subsequences to produce the sorted answer.

The recursion “bottoms out” when the sequence to be sorted has length 1, in which case there is no work to be done, since every sequence of length 1 is already in sorted order.

MERGE(A, p, q, r)

```

1    $n_1 = q - p + 1$ 
2    $n_2 = r - q$ 
3   let  $L[1..n_1 + 1]$  and  $R[1..n_2 + 1]$  be new arrays
4   for  $i = 1$  to  $n_1$ 
5      $L[i] = A[p + i - 1]$ 
6   for  $j = 1$  to  $n_2$ 
7      $R[j] = A[q + j]$ 
8    $L[n_1 + 1] = \infty$ 
9    $R[n_2 + 1] = \infty$ 
10   $i = 1$ 
11   $j = 1$ 
12  for  $k = p$  to  $r$ 
13    if  $L[i] \leq R[j]$ 
14       $A[k] = L[i]$ 
15       $i = i + 1$ 
16    else  $A[k] = R[j]$ 
17       $j = j + 1$ 
```

We can now use the MERGE procedure as a subroutine in the merge sort algorithm. The procedure MERGE-SORT(A, p, r) sorts the elements in the subarray $A[p..r]$. If $p \geq r$, the subarray has at most one element and is therefore already sorted. Otherwise, the divide step simply computes an index q that partitions $A[p..r]$ into two subarrays: $A[p..q]$, containing $\lceil n/2 \rceil$ elements, and $A[q + 1..r]$, containing $\lfloor n/2 \rfloor$ elements.⁸

② conquer

① divide

invariant

At the start of each iteration of the **for** loop of lines 12–17, the subarray $A[p..k - 1]$ contains the $k - p$ smallest elements of $L[1..n_1 + 1]$ and $R[1..n_2 + 1]$, in sorted order. Moreover, $L[i]$ and $R[j]$ are the smallest elements of their arrays that have not been copied back into A .

We must show that this loop invariant holds prior to the first iteration of the **for** loop of lines 12–17, that each iteration of the loop maintains the invariant, and that the invariant provides a useful property to show correctness when the loop terminates.

Initialization: Prior to the first iteration of the loop, we have $k = p$, so that the subarray $A[p..k - 1]$ is empty. This empty subarray contains the $k - p = 0$ smallest elements of L and R , and since $i = j = 1$, both $L[i]$ and $R[j]$ are the smallest elements of their arrays that have not been copied back into A .

③ combine

Maintenance: To see that each iteration maintains the loop invariant, let us first suppose that $L[i] \leq R[j]$. Then $L[i]$ is the smallest element not yet copied back into A . Because $A[p..k - 1]$ contains the $k - p$ smallest elements, after line 14 copies $L[i]$ into $A[k]$, the subarray $A[p..k]$ will contain the $k - p + 1$ smallest elements. Incrementing k (in the **for** loop update) and i (in line 15) reestablishes the loop invariant for the next iteration. If instead $L[i] > R[j]$, then lines 16–17 perform the appropriate action to maintain the loop invariant.

Termination: At termination, $k = r + 1$. By the loop invariant, the subarray $A[p..k - 1]$, which is $A[p..r]$, contains the $k - p = r - p + 1$ smallest elements of $L[1..n_1 + 1]$ and $R[1..n_2 + 1]$, in sorted order. The arrays L and R together contain $n_1 + n_2 + 2 = r - p + 3$ elements. All but the two largest have been copied back into A , and these two largest elements are the sentinels.

③ running time of recursive algorithms

(substitution method)

0. guess that $T(n) = \gamma(n)$

1. assume divides one rule. Show that $T(n) = \gamma(n)$ (induction step) for values of c and n_0

$$\text{eg. } T(n) = T(\frac{n}{2}) + T(\frac{n}{2}) + n \Rightarrow \text{assume } T(\frac{n}{2}) = \gamma(\frac{n}{2}), T(\frac{n}{2}) = \gamma(\frac{n}{2}) \Rightarrow \text{prove } T(n) = \gamma(n)$$

2. show that boundary solution - based on n_0 as base step for values of c and n_0

recursion tree

1. draw recursion tree, find no. of operations per layer

$$(k)^h \cdot n = 1 \Leftrightarrow h = \log_k n$$

split eg. $\frac{2}{3}$ leaf

imbalanced tree

2. using weight as bounds, argue for asymptotic bounds with

2.1 max weight : longest simple path from root to leaf

2.2 min weight : shortest simple weight

2. compute summations as $h \cdot g(n)$

substitution method for vigorous

master method

$$T(n) = aT(\frac{n}{b}) + f(n), \quad a \geq 1, \quad b \geq 1,$$

$$\frac{n}{b} = \lfloor \frac{n}{b} \rfloor \text{ or } \lceil \frac{n}{b} \rceil$$

$$1. \text{ Find } \log_b a = \frac{\ln a}{\ln b} = k$$

2. cover

$f(n)$ grows faster than height

$$f(n) \geq c \cdot n^{k+\varepsilon} = \Omega(n^{k+\varepsilon}), \quad \varepsilon > 0$$

$a f(\frac{n}{b}) \leq c f(n)$ for some $c < 1$ and all sufficiently large n

$$\Rightarrow T(n) = \Theta(f(n))$$

regularity condition,
ensures sum of
subproblems $c < f(n)$,
root dominates

$f(n)$ grows slower than height,
height dominates

$$f(n) \leq c \cdot n^{k-\varepsilon} = O(n^{k-\varepsilon}), \quad \varepsilon > 0$$

$$\Rightarrow T(n) = \Theta(n^k)$$

$f(n)$ grows about as fast as height, both height and $f(n)$ matter

$$c_1 n^k \leq f(n) \leq c_2 n^k = \Theta(n^k)$$

$$\Rightarrow T(n) = \Theta(n^k \cdot \lg n)$$

Inseful tricks

$T(\text{base})$ undefined \Rightarrow set $T(n_0) = q > 0$. set $c = \max/\min \{q_i, c\}$

Notes

1. given $T(f(n)) \Rightarrow$ manipulate

$$\text{eg. } T(2^n) = T(2^0) + T(2^1) + \dots + T(2^{n-1}) + n^2$$

$$\Rightarrow T(2^n) - T(2^{n-1}) = T(2^{n-1}) + n^2 - (n-1)$$

manipulate to get
tractable expression

$$\Rightarrow (\text{set } N = 2^n)$$

$$\text{then } T(N) = 2T(2^{n-1}) + 2n-1$$

$$= 2T(N/2) + 2\lg N - 1$$

2. height of T_n is $\lg \lg n$.

$$\sqrt{\sqrt{T_n}} \dots = n^{\frac{1}{2^i}} = 2$$

$$\frac{1}{2^i} = \log_n 2$$

$$= \frac{\lg 2}{\lg n} = \frac{1}{\lg n}$$

$$2^i = \lg n$$

$$i = \lg \lg n$$

3. change base to get familiar terms in \log .

$$q \cdot \frac{n^3}{\lg^2 n} + \frac{n^3}{(\lg n - \lg 3)^2} + \frac{n^3}{(\lg n - \lg 9)^2} \dots \Rightarrow \text{change to base 3}$$

4. function int

$T(n) = \dots$ (some n term) (some sum or product of n)

Show it converges to constant

then $T(n) = \Theta(n \text{ term})$

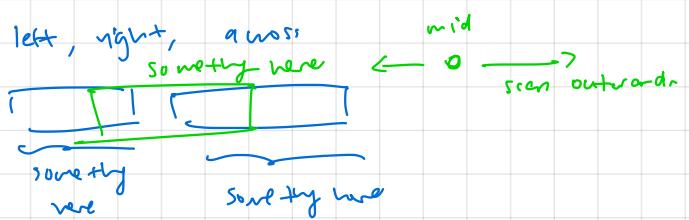
5. merge sort / merge subroutine is a good way to think about divide & conquer algorithms

divide → until ↴

conquer → trivially

combine → correctly by invariant / induction

6. useful divide & conquer paradigm: left, right, aross



Some recursive algorithms

(1D peak finding)

(2D peak finding)

Probabilistic analysis

① tools of analysis

(indicator variables)

provide a convenient way of converting between probability and expectation

$$x_i = \begin{cases} 1 & \text{if event } A \text{ occurs for some } i \\ 0 & \text{otherwise} \end{cases}$$

$P(x_i = 1)$ = same probability of A_i

$$E(x_i) = P(x_i = 1) (1) + P(x_i = 0) (0)$$

$$E(x) = E\left(\sum_{i=1}^n x_i\right) = \sum_{i=1}^n E(x_i)$$

(linearity of expectations)

$$E(X+Y) = E(X) + E(Y)$$

$$E(Y) = \sum_x E(Y|x=x) \cdot P(X=x)$$

$$E(aX) = a E(X)$$

Geometric distribution

$$\hookrightarrow X \sim G(p \text{ success})$$

number of times to get 1 success

$$P(X=k) = \underbrace{(1-p)}_{k-1 \text{ failure}} \underbrace{p}_{1 \text{ success}}$$

$$\begin{aligned} E(X) &= \left(\sum_{k=1}^{\infty} k (1-p)^{k-1} p \right) \cdot k \xrightarrow{\text{reshape}} \\ &= p \left[\sum_{k=1}^{\infty} (1-p)^{k-1} + \sum_{k=L}^{\infty} (1-p)^{k-1} \dots \right] \\ &= p \left(\frac{1}{p} + \frac{(1-p)}{p} + \frac{(1-p)^2}{p} \dots \right) \\ &= \frac{1}{p} // \end{aligned}$$

$$X \sim G(p \text{ success}) \Rightarrow E(X) = \frac{1}{p}$$

Useful integrals

$$\int x \ln x \, dx = \frac{x^2 \ln x}{2} - \frac{x^2}{4}$$

② common analysis problems

(counting occurrences) The hiring problem

suppose at all times, we want to best person for the job. We sequentially interview them, and hire if he is the best so far. With n interviews, how many will we hire?

1. define indicator variable

$$x_i = \begin{cases} 1 & \text{i-th candidate is best so far} \\ 0 & \text{otherwise} \end{cases}$$

probability changes \bar{a} :

$$P(x_i = 1) = \frac{1}{i}, \text{ assuming random \& i.i.d.}$$

2. expectation over sequence

$$\sum_{i=1}^n E(x_i) = \sum_{i=1}^n \frac{1}{i} = \ln n + o(1)$$

(counting paired occurrences) the birthday problem

In a room of n people, how many have the same birthday?

1. Define indicator variable

$$X_{ij} = \begin{cases} 1 & \text{if } i \text{ and } j \text{ have same birthday, } i \neq j \\ 0 & \text{otherwise} \end{cases}$$

$$P(X_{ij} = 1) = \sum_{d=1}^{365} \frac{1}{365^2} = \frac{1}{365}$$

$\frac{1}{365}$ chance on any day
for either person

$$E(X_{ij}) = \frac{1}{365}$$

2. sum to count for room of n people

$$X = \sum_{i=1}^n \sum_{j=i+1}^n X_{ij}$$

$$E(X) = \sum_{i=1}^n \sum_{j=i+1}^n E(X_{ij})$$

$$= [(n-1) + (n-2) \dots 1] \frac{1}{365}$$

$$= \frac{n(n-1)}{2} \frac{1}{365}$$

(combining geometric variables) the coupon collector problem

There are n coupons. How many coupons to collect until we get k unique if each draw is random?

Idea: sum of different geometric RV.

1. i th geometric RV

Let X_i be no. of times to get i th unique ball. Then $P(\text{success}) = \frac{n-i+1}{n}$

$$X_i \sim \text{Geo}(p = \frac{n-i+1}{n})$$

2. how long to get all k ?

$$\sum_{i=1}^k E(X_i) = \sum_{i=1}^k \frac{1}{\frac{n-i+1}{n}}$$

$$= \sum_{i=1}^k \frac{n}{n-i+1} \Leftrightarrow \frac{n}{i}$$

symmetric!

$$n-1+1, n-2+\dots+n-n+1$$

$$1, 2 \dots n$$

$$= n \sum_{i=1}^k \frac{1}{i} \leq n \lg k$$

Randomised algorithms

① randomised algorithms: motivation

1) probabilistic analysis

- ↳ best and worst case analysis give us bounds on our algorithm's performance
- ↳ but with knowledge of the input distribution, we can analyse the average performance of the algorithm, given the probability & performance of the algorithm
- ⇒ problem: how do we know the distribution of inputs?

2) worst case inputs and adversaries

- ↳ suppose the worst case time is very bad, and we have an adversary. e.g. quicksort given sorted sequence.
- ⇒ can we take control away from adversaries?

3) inject randomness into algorithms

- ↳ by enforcing randomised bottlenecks at points in our algorithm, we enforce an a priori distribution and take control away from adversaries.
- ⇒ we sacrifice certainty for controlled uncertainty

4) types of randomised algorithms

(Las Vegas algorithms) These algorithms are always correct. Their running time is a random variable that is bounded in expectation.
e.g. randomised quicksort

(Monte Carlo algorithms) correctness is a random event — they may give an incorrect output with some small probability.

② average case analysis of Las Vegas algorithms

- ↳ idea: we inject controlled randomness into the input - we then analyse the expected running time over these inputs.

(correctness)

1. show that randomisation preserves validity of input
2. prove correctness as usual

running time

let x be input drawn from some distribution D_n . $x \sim D_n$.

$$A(n) = \underset{x \sim D_n}{\mathbb{E}} (\text{run time of } A \text{ on } x)$$

1. consider types of input and their probability of occurrence.

2. For each type of occurrence, analyse running time

3. formulate into expectation $A(n) = \sum_{\pi \sim D_n} A(\pi) p(\pi)$



average case analysis of quick sort

0. suppose we have an input of length n of distinct elements. Then for any input, we observe that the execution of quick sort depends on the rank of a given element, in ascending order.

① divide 1. let S_i be the set of permutations of input such that x_i , the i^{th} ranked element is the start (and is thus picked as the pivot)

$$1.1 \text{ observe that } p(S_i) = \sum_{i=1}^n \frac{1}{n!} = \frac{1}{n}$$

② analyse individual case 2. let $G(i)$ be the average running time of quicksort when x_i is the pivot.

$$G(i) = \underbrace{(n-1)}_{\substack{\text{comparisons} \\ \text{to partition}}} + \underbrace{A(i-1) + A(n-i)}_{\substack{\text{recursive calls to either half,} \\ \text{excluding pivot}}}$$

③ formulate

$$\text{weighted sum. } 3. A(n) = \sum_{i=1}^n p(S_i) \cdot G(i)$$

$$= \frac{1}{n} \sum_{i=1}^n G(i) \text{ since uniform distribution, } p(S_i) = \frac{1}{n}$$

$$= \frac{1}{n} \sum_{i=1}^n \underbrace{(n-1)}_{\substack{\text{comparisons} \\ \text{to partition}}} + \underbrace{A(i-1) + A(n-i)}_{\substack{\text{recursive calls to either half,} \\ \text{excluding pivot}}}$$

$$= (n-1) + \frac{1}{n} \sum_{i=1}^n \underbrace{A(i-1) + A(n-i)}_{\substack{\text{observe that} \\ = A(0) + A(1) \dots A(n-1) \text{ by equal} \\ + A(n-1) + A(n-2) \dots A(0)}}$$

$$\begin{aligned}
 &= n-1 + \frac{2}{n} \sum_{i=1}^n A(i-1) \\
 &= n-1 + \frac{2}{n} \sum_{i=1}^{n-1} A(i)
 \end{aligned}
 \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{useful trick to keep for substitution}$$

3.1 We guess that $A(n) = O(n \lg n)$

$$\begin{aligned}
 A(n) &= n-1 + \frac{2}{n} \sum_{i=1}^{n-1} A(i) \\
 &\leq n-1 + \frac{2c}{n} \sum_{i=1}^{n-1} i \lg i \\
 &\leq n-1 + \frac{2c}{n} \int_{i=1}^{n-1} i \lg i \, di \\
 &= n-1 + \frac{2c}{n} \left[\frac{i^2 \ln i}{2} - \frac{i^2}{4} \right]_1^n \\
 &= n-1 + \frac{2c}{n} \left[\frac{n^2 \ln n}{2} - \frac{n^2}{4} + \frac{1}{4} \right] \\
 &= n-1 + c n \lg n - \frac{cn}{2} + \frac{c}{2n} \leq c \cdot n \lg n \quad \forall c \geq 2.
 \end{aligned}$$

Hashing

① hashing

i) hash functions

↪ the idea behind hashing is to drastically reduce the size of the direct access table.

$$h: U \rightarrow \{0, 1, \dots, m-1\}$$

ii) collisions and chaining

↪ suppose $h(x) = h(y)$, $x \neq y$. we call this a collision.

↪ chaining we handle collisions by adding an element to a linked list if an item maps to the same value as another item.

iii) complexity of hashing and chaining

↪ for a hash table of size M and N elements, hash table is of size $O(M+N)$.
 M pointers to linked list, total of N elements.

iv) properties of good hash functions

1. minimize collisions. query(x) and delete(x) need to traverse the list, and so take $\Theta(|h(x)|)$

2. minimise storage space, $M = O(N)$, N is stored items.

3. $h(x)$ should be easy to compute

② universal hashing

i) motivation

↪ if U is large such that $|U| \geq M(N-1)$, then $\forall h: U \rightarrow [m]$, there exists a set of N items to be stored that have the same hash value.

⇒ for a large enough U , an adversary can always force the worst case scenario by picking N items to store that map to same hash value

(proof) 1. pigeonhole principle: suppose not. That there exists a h that maps all $|U|$ items to the hashtable s.t. every slot at most $N-1$ items.

2. Then total $M(N-1) < M(N-1) + 1$ items. contradiction.

v) universal hashing

↳ key idea: randomization. Fool the adversary by not fixing the hash function!

Universal hashing: Suppose \mathcal{H} is a set of hash functions $h_i: U \rightarrow [m]$. Then we say \mathcal{H} is universal if $\forall x \neq y$,

$$\frac{|\{h \in \mathcal{H} : h(x) = h(y)\}|}{|\mathcal{H}|} \leq \frac{1}{m}$$

\Rightarrow for any pair of x, y where $x \neq y$, if h is chosen uniformly randomly from \mathcal{H} , then there is at most $1/m$ probability that $h(x) = h(y)$

values	x_1	x_2	x_3	\dots	for a fixed (x, y) , no. of hash in collisions no. of hash functions
h_1	○	○○	○		
h_2					
:					

$\leq \frac{1}{m}$
size of hashtable

\Leftrightarrow for a given pair x, y s.t. $x \neq y$, the proportion of $h \in \mathcal{H}$ that would make $h(x) = h(y)$ is at most $1/m$

3) properties of universal hashing

1. For $x_1, x_2, \dots, x_n \in U$, $E(\text{number where } x_N \text{ collides with other elements}) < \frac{N}{m}$

\Leftrightarrow if I have inserted $N-1$ elements, the expected no. of collisions when I insert the next element is bounded by $\frac{N}{m}$, since there are at most N distinct x_i .

for $i < N$, let $c_i \begin{cases} 1 & \text{if } x_N \text{ collides with } x_i \\ 0 & \text{otherwise} \end{cases}$

$$E(c_i) \leq \frac{1}{m}. \quad \sum_{i=1}^{N-1} E(c_i) \leq \frac{N-1}{m} < \frac{N}{m}$$

2. $M > N \rightarrow E(\text{cost of } N \text{ operations}) = O(N)$

$$E(\text{cost of 1 operation}) = O(1) + E(\text{no. of collisions when hashing } x_i)$$

if $M > N$, then $M > N > i$, where i is the i^{th} insertion. Then for the i^{th} operation, $E(\text{collisions}) < 1$, and $E(\text{insertion of } i^{\text{th}} \text{ element}) = O(1)$

$$E(\text{cost of } N \text{ operations}) = \sum_{i=1}^N E(\text{cost of 1 operation}) = N \times O(1) = O(N)$$

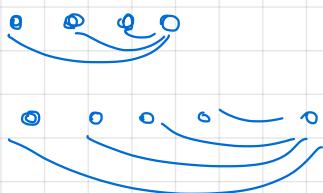
$$3. M > N \rightarrow E(\# \text{ of collisions for } x_i \in N) < 2N$$

\Leftrightarrow if $M > N$, then although at the i^{th} insertion we expect less than 1 collisions, there could still be some collisions over a sequence of N insertions. So the total number of collisions across N insertions is bounded by $2N$.

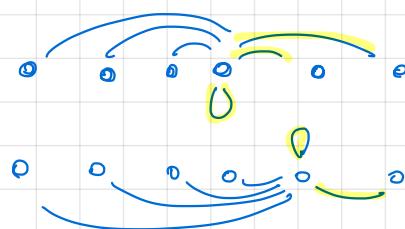
At the i^{th} insertion, we have inserted $i-1$ elements, and that could cause a collision. We would find the total number by finding the bound at the i^{th} insertion and then summing them to N .

An easier alternative is to recognise that in N elements, there are a fixed number of paired permutations that could be drawn that would cover every collision scenario, and use that as a bound instead.

true bound: only backwards



permutation bound: both directions and reflect



Let $A_{ij} = \begin{cases} 1 & \text{if } h(x_i) = h(x_j) \\ 0 & \text{otherwise} \end{cases}$

$$E\left(\sum_{i=1}^N \sum_{j=1}^N A_{ij}\right) = \sum_{i=1}^N E(A_{ii}) + \sum_{i \neq j} E(A_{ij})$$

$$NP_2 = \frac{N!}{(N-2)!} = N(N-1)$$

since $M > N$

$\sum_{i \neq j} E(A_{ij}) \leq N(N-1) \cdot \frac{1}{m} < 2N$

1/m chance of collision

$$4. M > N \rightarrow \text{expected maximum no. of elements hashed to a given slot} < \sqrt{2N}$$

\Leftrightarrow if x_1, x_2, \dots, x_N are hashed, and m/N , we expect that a given slot will have at most $\sqrt{2N}$ elements hashed into it'

In the worst average case, all the expected collisions will be hashed to the same slot — that is our expected maximum load case.

if m/N , we expect that the total no. of collisions $< 2N$. If all to the same slot, and recognising that any pair of elements at the same slot contributes to a collision, the expected no. of elements at the slot is $\sqrt{2N}$.

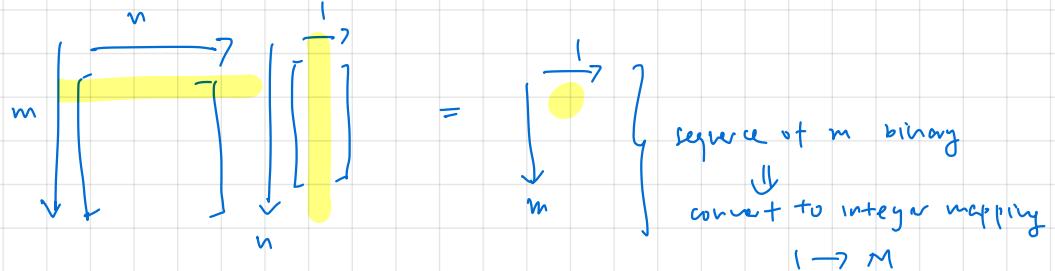
$$E(\max) = E\left(\sqrt{\max^2}\right) \leq \sqrt{E(\max^2)} \leq \sqrt{E(\text{total collisions})} \leq \sqrt{2N}$$

③ Constructing universal hash functions

1. Let $n = \lceil \log_2 m \rceil$ and m be powers of 2, where $\lceil \log_2 m \rceil = 2^n$, $M = 2^m$. Then we can represent any $x \in N$ as a binary sequence of length n .

2. Then for a binary matrix $A \in \{0, 1\}^{m \times n}$, we define n as

$$h_A(\pi) = \text{elementwise-mod } 2(A\pi)$$



3. Then $\mathcal{H} = \{ h_A : A \in \{0, 1\}^{m \times n} \}$, where A is any binary $m \times n$ matrix.
 There are 2^{mn} possible h_A .

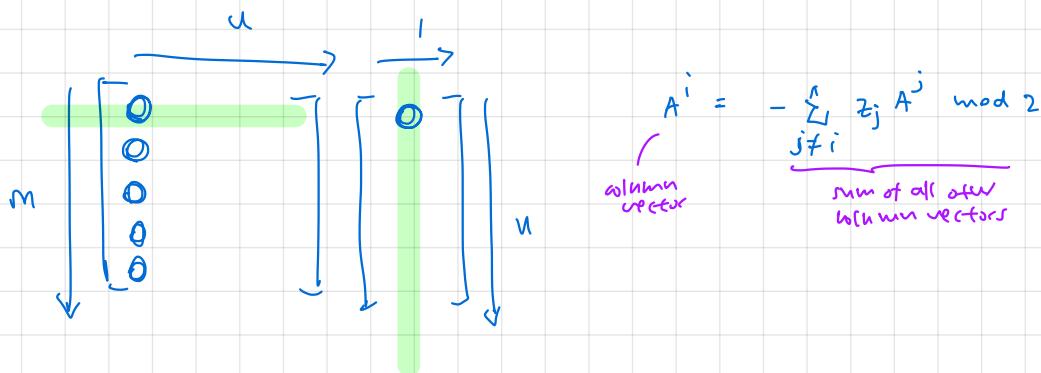
proof that it is universal

i. collision occurs when $A(x) = A(y) \Leftrightarrow A(x-y) = 0$, though $x \neq y$

1. For convenience, let $z = x - y$. so we seek to show $P(Az = 0) \leq \frac{1}{m}$

1.1 $x, y \in \{0, 1\}^n$. i.e. x and y are distinct binary sequences of length n .
 obviously, $\tilde{x} \neq \tilde{y}$.

1.2 observe that to produce $\tilde{0}$, a given entry i of the n -bit input multiplied by a column in A must perfectly offset that of the $\text{sum}(\text{mod } 2)$ of all other columns. (assuming the entry is non-zero)



2. Since A is uniformly chosen in $\mathbb{Z}^{m \times n}$ possibilities, the choice of the (i, j) position is independent of all other positions \rightarrow and the i^{th} column also independent of all other columns

2.1 so the probability that A_i^t is chosen to exactly offset sum of others, given a fixed t , where $\bar{x} = x - y \wedge x \neq y$, is $\frac{1}{2^m} = \frac{1}{m}$ exactly entire column

2.2 So far in universal

complexity of universal hashing

↪ the space required to specify a particular hash function is
 $O(um) = O(\log u \lceil \log m \rceil)$.

↪ time complexity is that of matrix multiplication

④ pairwise independent hashing

(pairwise independent family) If it is a pairwise independent family of hash functions $h: U \rightarrow [m]$ if it has the property that

$$P(h(x) = i_1, h(y) = i_2, x \neq y) \leq \frac{1}{m^2}.$$

⇒ the probability that a given (x, y) pair map to values $i_1, i_2 \leq \frac{1}{m^2}$

for each distinct $x, y \in U$

for each i_1, i_2 they map to

$$\frac{\text{no. of times } (i_1, i_2) \text{ appear}}{\text{no. of hash functions}} \leq \frac{1}{m^2}$$

→ also no. of mappings, if you think about it

pairwise → universal

$$P_{\text{unif}}(h(x) = h(y)) = \sum_{i=1}^m P(h(x) = h(y) = i) \leq m \cdot \frac{1}{m^2} \leq \frac{1}{m} \Rightarrow \text{universal}$$

/
 m possible values to be equal

after 1...N, expected @ given slot

$$\text{hash } x_1, x_2 \dots x_N, E(\text{no. of elements to } j) < \frac{N}{m}$$

↳ i.e. collisions pertaining specifically to j

$$1. c_i = \begin{cases} 1 & \text{if } h(x_i) = j \\ 0 & \text{otherwise} \end{cases}$$

From pairwise independence, we know $P(h(x) = i_1, h(y) = i_2) \leq \frac{1}{m^2}$.

we don't care, so we sum over the m (law of total probability)

$$P(c_i = 1) = \sum_{i=1}^m P(h(x) = j, h(y) = i)$$

$$\leq m \cdot \frac{1}{m^2} = \frac{1}{m}$$

for each insertion, whether will go to j

$$2 - N \text{ elements. } E(\text{hashes to } j) = \sum_{i=1}^N E(c_i) \leq \frac{N}{m}$$

⑤ perfect hashing

↪ idea: suppose we know a priori the set of x_1, \dots, x_N elements we want to hash. Can we construct a hash function so that query (x) can be answered in $O(1)$ worst case time?

i) do perfect hash functions exist?

↪ observation = if we pick an appropriate m , then the $E(\text{total collisions for } x_1, \dots, x_N)$ for family of universal hash functions < 1

↪ previously:

Let $A_{ij} = \begin{cases} 1 & \text{if } h(x_i) = h(x_j) \\ 0 & \text{otherwise} \end{cases}$

$$E\left(\sum_{i=1}^N \sum_{j=1}^N A_{ij}\right) = \sum_{i=1}^N E(A_{ii}) + \sum_{i \neq j} E(A_{ij})$$

$$\sum_{i \neq j} E(A_{ij}) \leq N(N-1) \cdot \frac{1}{m}$$

$$NP_2 = \frac{N!}{(N-2)!} = N(N-1)$$

$\sum_{i \neq j} E(A_{ij}) \leq N + \underbrace{N(N-1)}_{\text{will collide}} \cdot \frac{1}{m}$
1/m chance of collision

$E(\text{collisions}) \leq 1$ if $m = N^2 \Rightarrow$ there must exist at least one $h \in \mathcal{H}$ that
i.e. when $x \neq y$ but $h(x) = h(y)$

$E(\text{collisions}) = 0$, since expectation is less than 1.

ii) double hashing for $O(N)$ space complexity

↪ we know that for any x_1, \dots, x_N , there exists at least one perfect hash function and we need $m = N^2$ to get that.

↪ But $O(N^2)$ space is large. If we set $m = N$, we are going to get collisions.
 \Rightarrow since m grows by a square, why not use perfect hash at second level?

1. choose $h: U \rightarrow [m = N]$ from \mathcal{H} .

2. let L_k be the number of x_i for which $h(x_i) = k$, $k \in [m = N]$
size of each collision chain

3. we know that for each bucket, there must exist a universal function

$h_k: [L_k] \rightarrow [L_k^2]$ such that there are no collisions among the L_k elements mapped there by h initially

4. total space complexity = $O(N) + O(N)$ $\xrightarrow{\text{original table}} \text{sum of all secondary tables}$ pairs! same line of reasoning as T^2N vs. $2N$

4.1 $A_{ij} = \begin{cases} 1 & \text{if } h(x_i) = h(x_j) \\ 0 & \text{otherwise} \end{cases} \quad \sum_k L_k^2 = \sum_{i=1}^N \sum_{j=1}^N A_{ij}$

4.2 if $m > N$, $\sum_k E(L_k^2) \leq N + N(N-1) \cdot \frac{1}{m} \leq 2N$ by choice of $m = N$

⑥ division hashing

(definition)

Given a prime number p , $h_p : \mathbb{Z} \rightarrow \{0, \dots, p-1\}$, $h_p = x \bmod p$
 value $\rightarrow p$ slots

prime uniformly
selected from $1 \dots K$

(choice of p) to avoid adversarial arguments, we choose p uniformly from $1 \dots K$.

CHOOSE-PRIME(K)

- 1 Choose random number p uniformly from $\{1, \dots, K\}$
- 2 if p is prime
- 3 return p
- 4 else CHOOSE-PRIME(K)

- } 1. There are $> \frac{K}{\lg K}$ prime numbers from $1 \dots N$
 2. probability that prime is chosen in one try $> \frac{K/\lg K}{N} = \frac{1}{\lg N}$

2.1 viewing as Bernoulli random variable, expected no. of trials before first success is $< \lg K$.

$\frac{b \lg K}{K}$ collision bound

Suppose we have $0 < x, y < 2^b$. Then if we choose p uniform randomly from $1 \dots K$ collision when $h_p(x) = h_p(y) \Rightarrow (y - x) \bmod p = 0$

so $y - x$ must be divisible by one of its prime factors.

we know that for $z < 2^b$, there are $t < b$ prime factors.

since there are $> \frac{K}{\lg K}$ primes from $1 \dots K$ and p is uniformly chosen,
 $P(p \text{ is one of the } t \text{ prime factors of } z) < t \cdot \frac{\lg K}{K}$

(applying division hash to strings)

↳ the division hash takes in integers.

↳ how do we convert a sequence of characters from alphabet $\{0, \dots, l-1\}$ to an integer?

$$X = \langle x_1, x_2, \dots, x_m \rangle \text{ int}(X) = \sum_{i=1}^m x_i \cdot l^{m-i}$$

power of base
 base

conversion to base- l bit representation

$$\text{int}(X) < 2^b, b = m \lg l$$

Notes

universal hashing & $E(\text{no. of pairs that collide})$

$$X_{ij} = \begin{cases} 1 & h(x_i) = h(x_j) \\ 0 & \text{otherwise} \end{cases}$$

$$P(X_{ij} = 1) = \frac{1}{m} = E(X_{ij}), m \text{ is no. of slots, by universality}$$

$$E(X) = {}^n C_2 \cdot E(X_{ij}) = {}^n C_2 \cdot \frac{1}{m} = \frac{n(n-1)}{2} \cdot \frac{1}{m}$$

pattern matching

① Monte Carlo equality checks with hashing

↪ recall that Monte Carlo algorithms are correct in a certain probability.

↪ also observe that uniform and division hashing can both have their probability of collisions controlled by some hyperparameter. In division hashing, controlled by sampling range of $p \in [1 \dots K]$.

i) equality checking

```
if  $h(x) = h(y)$ 
    return true
return false
```

ii) error rates

(hashed word size)

when we use division hashing, we map to $[p] \leq K$ since p is chosen from $1 \dots K$.

Then we need $\lg K$ bits to represent the hash value. Often $\lg K = O(\lg N)$, by choice of K .

0) false negative rate $x \neq y \Rightarrow h(x) \neq h(y)$. obvious.

1) false positive rate

1. Let X and Y be length- m sequences from alphabet $\{0, 1 \dots l-1\}$.

2. Then $x = \text{int}(X)$, $y = \text{int}(Y) < 2^b$, $b = m \lg l$

3. Then we know $P(h_p(x) = h_p(y) | x \neq y) < b \cdot \frac{\lg K}{K}$ length of sequence

4. If we want $b \cdot \frac{\lg K}{K} \leq f$, we set

$$K = \frac{2^m}{f} \cdot \lg l \cdot \lg \left(\frac{2^m}{f} \lg l \right)$$

error rate alphabet size

② Rolling hashes in division hash

↪ suppose we want to repeatedly hash adjacent sequences. Then repeated hashings of $T[s+1 \dots s+m]$ to $T[s+2 \dots s+m+1]$ would take $O(m)$ time every time, even though they only differ by 1 character.

(linearity of int. conversion)

$$\text{int}(T[s+1 \dots s+m]) = \sum_{i=1}^m T[s+i] l^{m-i}$$

$$\Rightarrow \text{int}(T[s+2 \dots s+m+1]) = l \cdot \text{int}(T[s+1 \dots s+m])$$

everything ↑ by power 1

remove first char l^m → $-1^m \cdot T[s+1]$

add new char, who is l^0 → $+T[s+m+1]$

Linearity in modulo operation

$$(a+b) \bmod p = [(a \bmod p) + (b \bmod p)] \bmod p$$

$$c \cdot a \bmod p = c \cdot (a \bmod p) \bmod p$$

$$hp(T[s:t \dots s+m+1]) = hp\left(\ell \cdot hp(T[s+1 \dots s+m]) - hp(\ell^m) \cdot T[s+1] + T[s+m+1]\right)$$

next previous

③ pattern recognition

i) naive pattern recognition

↳ 1D array of length n , looking for sequence of length m .

```
for i in 1 to n-m+1 // O(n)
    isEqual = T[i, ... i+n-m+1] == pattern // O(m) } O(nm)
```

ii) Karp-Rabin algorithm

(idea) we use hashing to make equality check run in $O(1)$ time, and rolling hashes to hash next in $O(1)$ time. Then the algorithm runs in $O(n)$ time.

Karp-Rabin (T, P, f, ℓ)
repage $\xrightarrow{\text{pattern}} \text{hash}$
alphabet

$n = T.\text{length}$

$m = P.\text{length}$

$$f' = f/n \quad // \text{by union bound}$$

$$K = \frac{2m}{f'} \cdot \lg \ell \cdot \lg \left(\frac{2m}{f'} \lg \ell \right)$$

$P = \text{choose-prime}(K)$

$hp(x) = x \bmod p \quad // \text{set hash function}$

hash-pattern = $hp(\text{int}(P)) \quad // \text{hash pattern}$

hash-text = $hp(\text{int}(T[1 \dots m])) \quad // \text{hash first } m$

for $j = 1 \text{ to } n-m$

if hash-pattern == hash-text

return

else

hash-text = $hp(\ell \cdot \text{hash-text} - T[j+1] \cdot hp(\ell^m) + T[j+m+1]) \quad // \text{update}$

FPR for a single comparison is bounded by f .

$$\mathbb{E}_i = \begin{cases} 1 & \text{if false positive} \\ 0 & \text{otherwise} \end{cases}$$

$$\mathbb{P}(\mathbb{E}_i = 1) < f$$

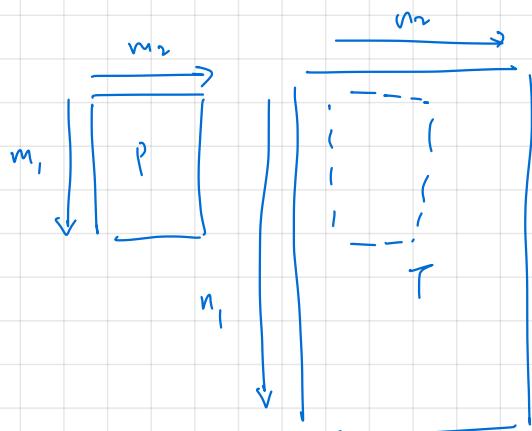
By union-bound, since probability at most as large as sum (where all mutually exclusive)

$n-m$

$$\sum_{j=0}^{n-m} \mathbb{P}(\mathbb{E}_j = 1) < \frac{f}{n} = f'$$

3) Karp-Rabin algorithm in 2D

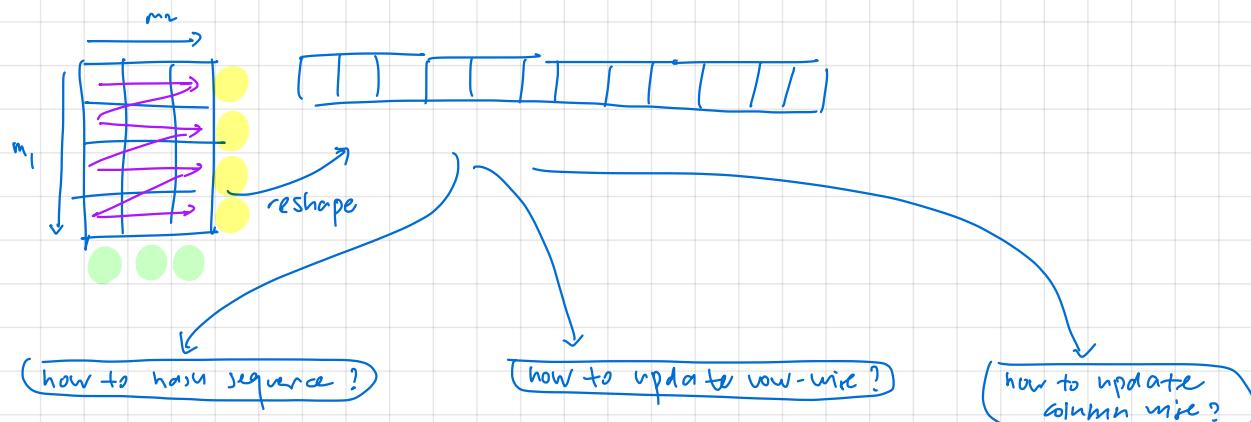
idea we build on the idea in the 1D case by reshaping the blocks into a 1D sequence, and accordingly adjusting our rolling hash



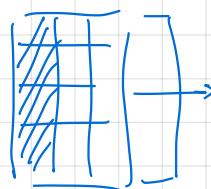
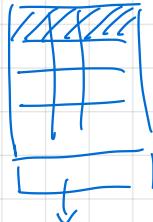
↳ in the 1D case, we convert sequences into base λ by:

$$\dots \Rightarrow \underbrace{\lambda^0 \lambda^{m_1} \lambda^{m_1+1} \dots \lambda^{n_1}}_{l^0}$$

↳ in the 2D case, we can convert blocks by thinking of them as sequences



$$\text{int(region)} = \sum_{i=0}^{m_1-1} \sum_{j=1}^{m_2} x_{ij} \cdot \lambda^{m_1 m_2 - (im_2 + j)}$$



(FP rate) we will run total of $O(n_1 n_2)$ times. $f_1 = \frac{P}{n_1 n_2}$

runtime	$O(m_1 m_2) + O((n_1 - m_1 + 1)m_2) + O((n_2 - m_2 + 1)m_1)$	
initial hashes	row-wise updates	column-wise updates
	$+ O((n_1 - m_1 + 1)(n_2 - m_2 + 1))$	
	equality checker	

Streaming

① streaming model

stream a sequence of insertions and deletions of items from a universe U . We will consider U to be the set $\{1 \dots |U|\}$.

Stream can be: $\text{Add}(3), \text{Add}(1), \text{Delete}(3) \dots$

separately

② frequency estimation

frequency estimation - f_i : number of times i occurs at end of stream.
 $\text{no. insertions of } i - \text{deletions of } i$

\Rightarrow can we design a smart way to occupy small amt. of space but query efficiently f_i ?

hashing-based frequency estimation

idea: use hashcodes, maintain counter
at each hashcode. increment &
decrement when add or delete

direct access
table. $O(|U|)$.

BST
 $O(m)$, where
 m is no. of distinct

INITIALIZE(U, k)

- Choose hash function $h : [U] \rightarrow [k]$ from a universal family \mathcal{H}
- Create table T of k entries, each initialized to be 0

ADD(i)

$$1 \quad T[h(i)] = T[h(i)] + 1$$

DELETE(i)

$$1 \quad T[h(i)] = T[h(i)] - 1$$

QUERY(i)

1 return $T[h(i)]$

Then we know that will only overestimate in some error.

$$f_i \leq E(\hat{f}_i) \approx f_i + \frac{m}{k}$$

no. of distinct
hash-table size

1. by twostile assumption, for every element, no. insertion \geq no. deletion

$$\text{1.1 so } E(\hat{f}_i) \geq f_i$$

$$\text{2. } X_{ij} = \begin{cases} 1 & \text{if } h(i) = h(j) \\ 0 & \text{otherwise.} \end{cases}$$

$$\text{2.1 } \hat{f}_i = \sum_{j=1}^n f_j X_{ij} - \text{so } i \text{ either collides w all others or itself}$$

$$\text{3. if } i=j, X_{ij}=1, \text{ if } j, P(X_{ij}=1) \leq 1/k \text{ by universality of } h$$

$$\text{3.1 so } E(\hat{f}_i) = E\left(\sum_j f_j X_{ij}\right) = f_i + f_j \sum_{j \neq i} E(X_{ij}) \leq f_i + \left(\sum_{j \neq i} f_j\right) \cdot \frac{1}{k} \leq f_i + \frac{m}{k}$$

f. If we set $k = \frac{1}{\varepsilon}, \varepsilon > 0$:

$$f_i \leq E(\hat{f}_i) \leq f_i + \frac{m}{k} \Leftrightarrow f_i \leq E(\hat{f}_i) \leq f_i + m\varepsilon.$$

Space complexity

(hash table)

$k = \frac{1}{\varepsilon}$ rows, each row with $\lg m$ bits

since there are m distinct, need $\lg m$ bits to index each

(hash function)

A is a $|U| \times |M|$ matrix, so $\lg |M| \cdot \lg |m|$ size.

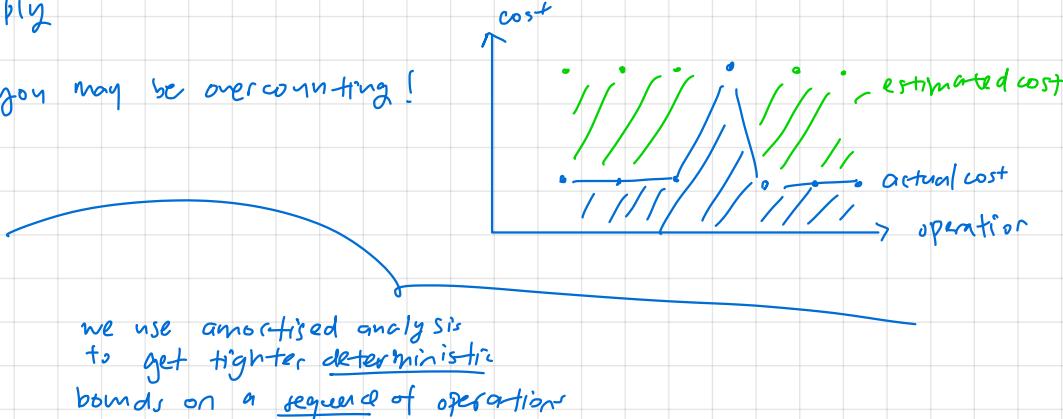
$$O\left(\left(\frac{1}{\varepsilon} + \lg |M|\right) \lg m\right) \text{ bits} \approx O\left(\frac{1}{\varepsilon}\right)$$

Amortized analysis

① amortized analysis

↳ in analysing a sequence of operations on a data structure, an easy way to compute the total worst-case cost is to take the individual worst case costs and multiply

↳ but you may be overcounting!



Amortised analysis

a strategy for analysing a sequence of operations to show that the average cost per operation is small, even though a single operation within the sequence may be expensive

⇒ does not involve probability; guarantees the average performance of each operation (in a known sequence) in the worst case

↗ after worst case sequence

② aggregate method

technique we show that for all n , a sequence of n operations takes worst case time $T(n)$. Then the amortized cost is $\frac{T(n)}{n}$.

③ accounting method

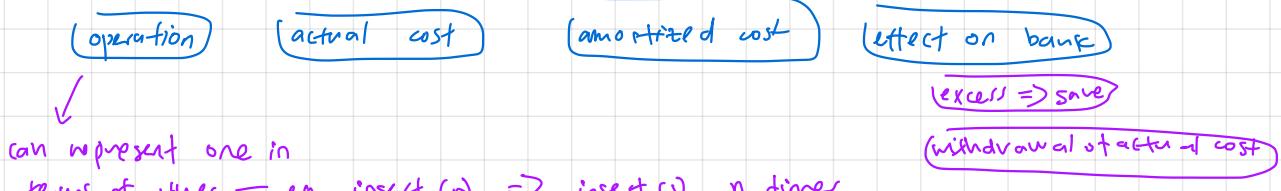
idea impose an extra charge on cheap operations and "save it", use that guaranteed savings (show that never dips below zero) to pay for later expensive operations, during which you charge 0

$$1. \text{ formulate } c(i) = \begin{cases} \dots & \text{for diff. types of operations} \\ \dots & \rightarrow \text{overestimate cheap, underestimate exp., build credit} \end{cases}$$

2. show that $\sum c(i) \geq \sum t(i)$, that is, total amortized cost bounds total actual cost \rightarrow directly, or that "savings" is always

3. $c(i)$ gives amortised cost of diff. operation types

useful visualization using table



technique
 figure out how much you need to overpay per unit up front so that
 later operation (that can happen, since this is happening) is paid for.
 Technique gives avg. worst case cost of each operation type

④ The potential method

- idea use a potential function modelling the state of the data structure to track
 excess / withdrawals of "credit", rather than using explicit operation
 amortised costs to track "credit"
 ↴ i.e. some adjustment to true cost tracking
 excess / shortfalls from previous operations implying
- $c(i) = t(i) + \phi(i) - \phi(-1) = t(i) + \Delta\phi(i)$
 amortized of i^{th} operation is actual cost + change in potential after operation
 - show that by design of representative ϕ , $\phi(n) - \phi(0) \geq 0$
 - then total cost bounded by $\sum c(i) = \sum t(i) + \phi(n) - \phi(0) \geq \sum t(i)$.
 - then amortised cost of each operation is respective $c(i)$

(why does the technique work?)

- we set amortised cost $c(i) = t(i) + \phi(i) - \phi(-1)$
 - then sequence of n operations, total amortised :
$$\begin{aligned} \sum c(i) &= \sum_{i=1}^n t(i) + \phi(i) - \phi(-1) \\ &= t(n) + \phi(n) - \phi(-1) \\ &\quad + t(n-1) + \phi(n-1) - \phi(n-2) \\ &\quad \vdots \\ &\quad + t(2) + \phi(2) - \phi(1) \\ &\quad + t(1) + \phi(1) - \phi(0) \end{aligned} \quad \left. \right\} = \sum_{i=1}^n t(i) + \phi(n) - \phi(0)$$

∴ then so long as $\phi(n) - \phi(0) \geq 0$, $\sum_{i=1}^n t(i)$ bounded by $\sum_{i=1}^n c(i)$

2. Then we can find the amortized cost of each of the i operations

(Technique) Design a potential function such that using state of data structure

1. ϕ (increases) in (small increments) \rightarrow cheap operations
2. (falls) by (large amounts) \rightarrow expensive operations
3. $\phi(0) = 0$ (for convenience)

\Rightarrow then $c(i) = t(i) + \Delta\phi(i)$ gives cost of specific operation

⑤ application to dynamic tables

1) motivation for dynamic tables

- ↳ we do not always know how many objects will be stored
- ↳ when not enough space, we create new larger table \rightarrow copy into it
- ↳ when too much space, we create new smaller, copy into it

\Rightarrow how can we design choices of new table size and when to do it, so that we have reasonable amortized time complexity?

2) dynamic expansion

(idea) everytime table is full, double the size

amortised cost of n insertions

1. set $\phi(i) = 2i - \text{size}(T)$. let us analyse each i of the n insertions.

2. case 1: table is full when trying to insert i

$$2.1 \text{ then } c(i) = \underbrace{i}_{t(i)}, \text{ copy } i-1, \text{ put } i \underbrace{\text{new table}}_{2(i-1) - (i-1)}$$

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

3. case 2: table is not full when trying to insert i

$$3.1 \text{ then } c(i) = 1 + (2i - 2T) - (2(i-1) - 2T)$$
$$= 3$$

4. In both cases, $c(i) = 3 = O(1)$

(OP) accounting method : $c(i) = \begin{cases} (i-1) + 1 & \text{if full and doubling} \\ 1 & \text{otherwise.} \end{cases}$

then $\sum_{i=1}^n c(i) = n + \sum_{j=0}^{\log n - 1} 2^j \leftarrow \text{doubling}$
 $\leq 3n$

3) dynamic expansion and contraction

(idea) table full, double; @ $\frac{1}{4}$ capacity, halve ↗ piecewise potential function!

1. set $\phi(i) = \begin{cases} 2 \cdot T.\text{num} - T.\text{size} & \text{if load } \geq \frac{1}{2} \\ T.\text{size}/2 - T.\text{num} & \text{if load } < 2 \end{cases}$

2. analysis of table insert
 3. analysis of table delete } see CLRS 12. 470

Notes

1. idea: cost of examine & operate

↑
needs to iterate!

↓

Instead, we maintain a pointer \Rightarrow then we can go to pointer & operate, reset pointer if necessary
 \Rightarrow costs 2 each time

dynamic programming

① recursion and optimal substructure

(optimal substructure) we say that a problem exhibits optimal substructure if an optimal solution to the problem contains within it optimal solutions to subproblems

⇒ i.e. fancy way of saying that you can write a recursive expression for the solution that is correct

(notion of optimality) the choice of "optimal" is a misnomer, in the sense that it refers more to correctness, just often in the context of an optimization problem

⇒ e.g. $\text{Fib}(n) = \text{Fib}(n-1) + \text{Fib}(n-2)$ exhibits optimal substructure, since the "optimal" solution is the correct solution.

(proving optimal substructure) "cut sub-optimal subproblem solution, paste optimal"

1. formulate a recurrence relation of a problem e.g. $T(n) = 1 + T(n-1)$

1.1 solution to the problem consists of making a choice (e.g. where to split), leaving one or more subproblems to be solved, then combine.

divide

conquer

2. (suppose not) suppose $T(n)$ is an optimal solution, and suppose $T(n-1)$ is not.

2.1 (show contradiction) then we could improve $T(n)$ by $T(n) = 1 + T^*(n-1)$

2.2 but $T(n)$ is optimal, and cannot be improved - contradiction.

we formalize the notion of "recursion is correct" with "exhibits optimal substructure"

have optimal substructure \iff can be solved w/ divide and conquer

② Recursive problems with overlapping subproblems

↳ in recursive algorithms, we can represent the computation graph as a recursion tree. But in some recursive formulation, some nodes are identical — that is, there exists repeated computation

(overlapping subproblems)

when a recursive algorithm revisits the same independent problem repeatedly, we say that the problem has overlapping subproblems.

↳ note that independent = do not share resources = can be executed in parallel
= a unique node

③ dynamic programming

(tackling repetition in divide & conquer)

- divide and conquer algorithms partition the problem into subproblems, solve them recursively, then combine their solutions. They exploit optimal substructure to get a solution better than brute force.
 - when a divide & conquer algorithm is inefficient (not incorrect) is where there exist overlapping subproblems, because of repeated computation
- ⇒ that's where DP techniques come in

[idea] brute force, but carefully to avoid repeated computation

[memoization]

[bottom up]

implicit memoization

[idea] memoize the natural, but inefficient recursive algorithm.

if value not in table, compute & add it. Else just use

[running time] $O(n \text{ independent problems} \times \text{complexity of solving/combining})$

[correctness] by induction

[idea] compute subproblems from the base case, and proceed in order of dependency

essentially converting recursion tree to graph where each node is unique subproblem, solving in topological order

[technique]

- find recurrence, prove optimal substructure
- initialize base cases
- for loops starting from base cases to n , each loop combining by using recurrence

[running time] running time of iterative algorithm

[correctness]

induction

loop invariant

↳ value always computed before

(indirectly)

④ summary

- 1- formulate recurrence
 - 1.1 break into subproblems
 - 1.2 think how to choose among subproblems (e.g. max, min, 1+ $T(n-1)$)
2. prove optimal substructure
- 3- point out overlapping subproblems
 - 3.1 (count unique, if possible)
- 4- bottom up
 - 4.1 initialize table w base cases
 - 4.2 loop, starting in late night after base
 - 4.3 terminate by return
- 5 - prove correctness by
 - 5.1 loop invariant of optimal solution to subproblem always being computed
 - 5.2 combination gives optimal

knapsack

Given $(v_1, w_1), (v_2, w_2) \dots (v_n, w_n)$, w .

Find a subset of $S = \{1 \dots n\}$ such that $\sum_{i \in S} w_i \leq w$

1. 0. enumeration of 2^n subsets is costly.

1. we can formulate the recurrence in the regular knapsack problem.

$$T(n, w) = \begin{cases} 0 & \text{if } n=0 \text{ or } w=0 \\ \max \left(T(n-1, w), T(n-1, w-w_n) + v_n \right) & \text{if } w_n \leq w \\ T(n-1, w) & \text{otherwise no choice but to skip} \end{cases}$$

mark idea of subset:

if can include → complete includes, not

if cannot → don't, skip

base case

Greedy algorithms

① Greedy choice property & implications

(greedy choice property) A locally optimal choice is a globally optimal choice

) proving that a problem exhibits it

[ideq] prove that greedy choice is always in an optimal solution

\Leftrightarrow greedy choice never leaves you worse off

1. Formulate a recurrence that greedily chooses a subproblem to solve using some measure (e.g. latest activity)

2. Let $S(n)$ be an optimal solution.

2.1 suppose $s(k)$, where k is greedy choice is not in any optimal solution

2.2 show that swapping $s(k)$ with choice leaves solution equal or better

2.3 Then it must be in at least one optimal solution

we use idea to show that while we may not get a specific solution, our solution in this choice is always optimal

) exploiting the greedy choice property

\hookrightarrow observe that DP algorithms improve on divide & conquer problems with overlapping subproblems by memorizing (explicitly or implicitly)

\hookrightarrow but efficiency can be improved further if we reduce the no. of subproblems we need to solve

how do we choose them?

\Rightarrow if a problem exhibits the greedy choice property, we can exploit it by always making the greedy choice and reduce the no. of subproblems to solve

② Greedy algorithms

\hookrightarrow in a greedy algorithm, we choose the stair subproblem(s) that seem best at the moment, and solve only those that are chosen

\hookrightarrow the choice(s) made may depend on choices so far, but it cannot depend on any future choices or solutions to subproblems

(greedy vs. DP algorithms)

- ↪ In DP, we solve subproblems, and then combine them. often by making some choice among the solutions.
- ↪ in greedy, we choose first which subproblems to solve, and then go ahead with it, pruning the computation graph

idea brute force, but even more carefully, by choosing nicely less subproblems

(formulating a greedy algorithm)

1. formulate recurrence, one for optimal, one for greedy
2. prove optimal substructure \Rightarrow can because
3. prove greedy choice property \Rightarrow can reuse this smart way
4. recursive algorithm + time complexity

\rightarrow ie. can reuse,
overlapping

(optimal substructure) optimal solution exactly uses optimal to sub this way

(greedy choice) there exists an optimal solution where direct choice is inside

1. Let S be an optimal solution (sequence).
2. suppose greedy choice not inside
 - 2.1 always some x_j it can swap with.
 - 2.2 swapping always no worse off

Reductions

① computation problems

(problems) a computational problem is defined by a set of (x, s) tuples where x is a valid input (instance) and s is a set of valid solutions to the input x

(decision problems) a problem where for any valid input x , the solution is either a YES or NO

(search problems) general problems where the solution goes beyond YES and NO

② input size & running time

(encodings) An encoding of a set S of abstract objects is a mapping from S to a set of binary strings

↳ either encoding is explicitly specified, or we follow conventions

integer \rightarrow binary

string \rightarrow concatenation of constant-bit length strings

graph \rightarrow adjacency matrix

(running times) we say that an algorithm solves a concrete problem (i.e. encoded explicitly) in $O(T(n))$ if, when it is provided a problem instance i of length $n = |i|$, the algorithm produces a solution in $O(T(n))$ time.

↳ take note that unit size of encoding \rightarrow not value, not abstract representation

③ reductions

Reductions between decision problems

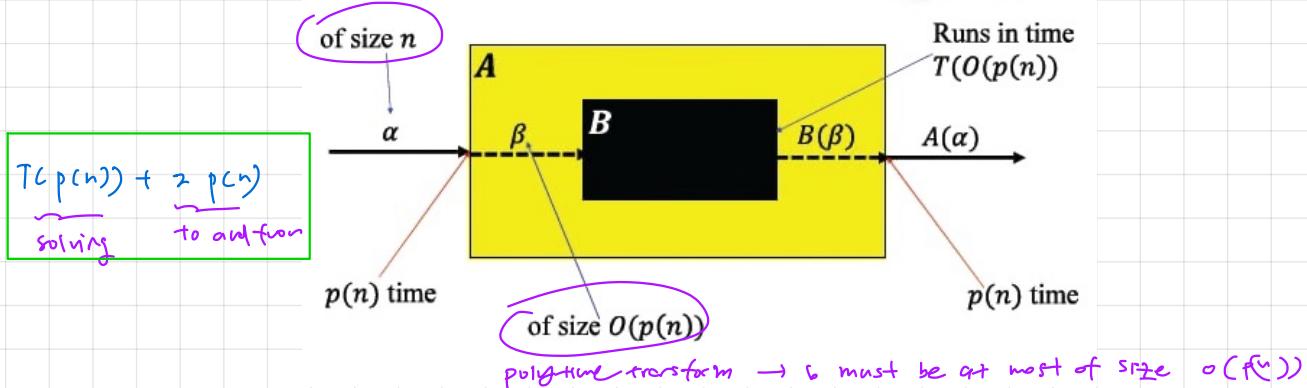
Given two decision problems X and Y , we say that X reduces to Y if there exists an algorithm that takes an input x of problem X and outputs an instance y of problem Y , such that x is a YES instance of X iff. y is a YES instance of Y .

Reductions in polynomial time

If transformation algorithm τ reduces X to Y in polynomial time, we say X reduces to Y in polynomial time - $X \leq_p Y$

Running time composition

If there is a $p(n)$ time reduction from A to B, and a $T(n)$ time algorithm to solve problem B on instances of size n , then there is a $T(O(p(n)) + O(p(n)))$ time algorithm to solve problem A on instances of size n .



(proof that A reduces to B in $p(n)$ time)

1. transformation algorithm of $a \rightarrow b$
1.1 proof of run in $p(n)$ time. If obvious, start obvious.
2. YES $a \rightarrow$ YES b
3. YES $b \rightarrow$ YES a

④ decision reduces to optimization

Given an optimization problem, we can convert it into a decision problem

\Rightarrow is there a solution w value $\leq k$? \Rightarrow get solution \rightarrow compare to k

⑤ pseudo-polynomial time reduction

Given an algorithm that runs in polynomial time in the numerical value of the input (but exponential in size of input encoding) is pseudo-polynomial

e.g. $O(n) \rightarrow O(2^{\lg n}) = O(2^{\text{encoding size}})$ for integers, w size $\lg n$

$O(n) \rightarrow O(\frac{n \lg n}{\lg n}) = O(\frac{1}{\lg n} \text{ encoding size})$ for arrays, w size $n \cdot \lg n$

⑥ Reductions & relative hardness

1. $A \leq_p B$.
2. If there is a $p(n)$ time algorithm for B , there is one for A .
 - 2.1 proof by construction. Reduce → solve, transform.
3. If there is no $p(n)$ time algorithm for A , then there is none for B .
 - 3.1 proof by contradiction. None for A , hence for B . Then can convert to B and solve in $p(n)$ time. contradiction.

⑦ general reduction

(Karp-Lemkin reduction) Direct reduction to solve.

1. $X \xrightarrow{R_1} Y$, $Y \xrightarrow{R_2} X$, both run in $p(n)$ time.
2. Then if A is algorithm to solve Y , $R_2(A(R_1(X)))$ is a solution to X .

(Cook reduction) A more general type of reduction where we invoke an algorithm A_Y for Y as a subroutine a polynomial number of times in order to solve X .

P and NP

① P and NP

1) definitions

(P complexity class) The set of problems that can be solved in polynomial time.

(NP complexity class) The set of problems for which polynomial time verifiable certificates of YES instances exist -

↓
(efficient poly) verifier a polynomial time algorithm that outputs YES iff. input is a solution to problem X.

↳ i.e. correct solution \rightarrow YES
wrong solution \rightarrow NO \Leftrightarrow YES \rightarrow correct solution

2) key results

(P \subseteq NP) proof by construction : just return the solution and see if same.

(show X is in NP)

1. propose algorithm that takes the solution (of size $p(n)$) and outputs YES / NO
 - 1.1 in $p(n)$ time
 - 1.2 correctly : YES \Leftrightarrow solution is correct

② NP hard and NP complete

(NP hard) A problem is said to be NP hard if for every problem B in NP, B can be reduced to it. $\forall B \in NP, B \leq_p X$.

(showing NP hard)

1. propose a reduction of NP hard problem Y to X.
2. show it runs in $p(n)$ time.
3. show YES $x \rightarrow$ YES y
4. show YES $y \rightarrow$ YES x
5. By transitivity of reductions, since Y is NP hard, X is NP hard.

(NP complete) A problem that is in NP (verifiable in $p(n)$) and also NP hard.

③ Some NP complete problems

Independent set Given a graph G and integer k , is there a subset of k vertices such that no two are connected by an edge?
ie. max no. of vertices such that no two connected

Vertex cover Given a graph G and integer k , is there a subset of k vertices such that each edge is incident to at least one vertex
ie. smallest no. of vertices such that all edges included

Set cover Given integers k and n , and a set \mathcal{S} of subsets of $\{1 \dots n\}$, are there $\leq k$ of these subsets whose union equals $\{1 \dots n\}$?
ie. minimal joining of subsets to get original set

3SAT Given a CNF formula where each clause contains exactly 3, not necessarily distinct literals, does it have a satisfying truth assignment?

Partition Given a set of positive integers S , can the set be partitioned into 2 sets of equal sums?

Hamiltonian path Given a graph G , is there a path in a directed graph from some beginning vertex to an ending vertex, such that every node is visited exactly once?

Clique A set of vertices U of a graph G is a clique if every pair in U has an edge in G . Given a graph G and integer k , is there a fully connected subgraph with $\geq k$ nodes?

(4) reductions covered in class

Independent set \leq_p clique

1. Given graph G . Take the complement graph of G^1 as an input to the clique problem. set $k^1 = k$.
2. Obvious that this transformation runs in polynomial time, $O(V^2)$.
3. YES independent set \rightarrow YES clique
 - 3.1 Let S be independent set solution.
 - 3.2 Then corr. k nodes will be disconnected in complement graph and form a clique of size $|S| = k$.
4. YES clique \rightarrow YES independent set
 - 4.1 symmetrical argument to 3.

Independent set \leq_p vertex cover

can modify G , or
modify k !



1. Given graph G and k to IS. we check that vertex cover of size $n-k$ (circled) of graph G^1 .
2. obvious that this runs in polynomial time. $O(V^2)$
3. YES independent set \rightarrow YES vertex cover
 - 3.1 Let solution to IS be S . Then $G-S$ is a solution to vertex cover.
 - 3.2 proof: let (u,v) be every edge in G . then since $|S| = n-k$ cannot both be in S .
4. YES vertex cover \rightarrow YES independent set

4.1 Let solution $n-k$ nodes to VC be S . ie. there is a solution \Rightarrow YES instance

4.2 when $G-S$ is solution to IS problem. proof by contradiction: if edge (u,v) both u,v in $G-S$, then edge not covered in S . contradiction.

useful way to think about vertex cover

(vertex cover \leq_p set cover)

defining a suitable transformation such that solution to one is exactly solution to other

1. For a given graph $G(V, E)$ and k to the VC problem, we generate the following as input to the SC problem

1.1 $n = |E|$

1.2 $S_k = \{ \text{edges incident on vertex } k \}$

1.3 Then $f = \{S_k\}_{k \in V}$

2. It is obvious that the generation runs in polynomial time. $O(n^2)$.

3. YES VC \rightarrow YES SC

3.1 Let S be the solution to VC, with k vertices.

3.2 Then we know all edges are covered in S . By definition of S_k , the union of S_k , $k \in S$ will be f .

4. YES SC \rightarrow YES VC

4.1 By definition, the subsets that form solution S have corr. nodes, which cover all edges in G .

(3SAT \leq_p independent set)

1. Input to 3SAT as $x_1 \dots x_n$ in k clauses of 3 literals each. Then we generate an input to IS with each x_i as a vertex, and connect an edge between all literals in the same clause, all literals to negations of itself.

2. Obvious that this runs in polynomial time.

3. YES IS \rightarrow YES 3SAT

idea: use edge to constrain simultaneous presence

3.1 Let S be independent set of size k .

3.2 Then S cannot contain vertices labelled by conflicting literals, since they are connected by an edge.

idea: use edge to isolate one from each clause

3.3 S must contain exactly one vertex from each clause, since exactly one can be chosen from a given triangle (so that no connected vertices)

3.4 Thus it is possible to choose a set of literals that are present in every clause and have no contradictions, and setting these literals to true will make the 3-CNF statement true!

4. YES \rightarrow SAT \rightarrow YES IS

4.1 For every triangle, pick exactly one vertex that is true in the assignment satisfying the 3SAT.

4.2 K clauses. Valid assignment means at least 1 true in each clause, with no contradictions

4.3 Then 4.1 will create an independent set, since x_i and $\sim x_i$ cannot simultaneously be in solution (ie. no edge for $x_i - \sim x_i$) and choosing one per triangle also means no edges.

⑤ Notes about reductions

1. 3SAT \rightarrow boolean expressions that are logically equivalent
 \rightarrow as a graph

2. IS \rightarrow edges as constraints (cannot both be present in solution)
 $\rightarrow (u, v) \in E$, u & v cannot both be in solution, since connected

3. VC $\rightarrow (u, v) \in E$, at least one of u, v must be in solution

4. Set cover \rightarrow by definition, $S_k = \{ \text{some property } y \}$, $\text{union} \{ S_k \} = \{\text{original}\}$

5. Unique \rightarrow complement to IS

6. Partition / subset sum \rightarrow pad w/ $x_i = 0$
 \rightarrow double the set w/ repeat

7. longest simple path / cycles \rightarrow pad w/ dummy vertices that are longer than what could originally exist

3 SAT \leq_p hamiltonian cycle

3 SAT \leq_p subset sum