

Recursive Algorithm :

"An algorithm to solve a problem by repeatedly 'calling itself' with subproblems."

Ex) : mergesort, quicksort ...

Methods for Recurrences :

1. Iteration

$$\begin{aligned}T(n) &= n + 3T(n/4) \\&= n + 3(n/4 + 3T(n/16)) = n + 3n/4 + 9T(n/16) \\&= n + 3n/4 + 9(n/16 + 3T(n/64)) = n + 3n/4 + 9n/16 + 27T(n/64) \\&= \dots \\&= n + (3/4)n + (9/16)n + \dots + (3^{(\log_4 n)}T(n/(4^{(\log_4 n)}))) \\&= n + (3/4)n + (9/16)n + \dots + (3^{(\log_4 n)}T(1)) \\&\leq n * (\sum_{i=0}^{\infty} (3/4)^i) + n^{(\log_4 n)} * \Theta(1) \\&= 4n + o(n) \\&= \Theta(n)\end{aligned}$$

Therefore, $T(n) = O(n)$. QED.

2. Guess & Verification

This method uses mathematical induction.

Ex)

$$T(n) = 2T(n/2) + n, T(1) = 1$$

Let's guess that $T(n) = O(n * (\log n))$.

i.e. $T(n) \leq c * n * (\log n)$ for some $c > 0$ and large enough n .

Pf.

Base case :

$$\begin{aligned}T(2) &= 2T(1) + 2 \\&= 4 \\&\leq c * 2 * (\log 2)\end{aligned}$$

c for above inequality exists.

Inductive case :

Assume that for arbitrary $k < n$, there exists $c > 0$ s.t. $T(k) \leq c * k * (\log k)$.

Then,

$$\begin{aligned}T(n) &= 2T(n/2) + n \\&\leq 2 * c * n/2 * (\log n/2) + n \\&= c * n * (\log n) - c * n * (\log 2) + n\end{aligned}$$

Here, if we choose c large enough, $-c * n * (\log 2) + n$ becomes negative.

And that holds for $c \geq 1$.

Therefore, $T(n) \leq c * n * (\log n)$ for $c \geq 1$

We choose $c = 2$, $n_0 = 2$ (n_0 : large enough n)

* We can always verify any claim for a small fixed n . (e.g. $n=2$)

e.g. $T(n) = 10 * T(n/10) + n$, $T(1) = 1$

$$T(n) \geq c * n * (\log n)$$

$$T(10) = 10 * T(1) + 10 = 20 \leq c * 10 * (\log 10) \text{ (there exists } c \text{ that satisfies this.)}$$

Thus, we don't have to explicitly prove the boundary case in 'Guess & Verification' methods.

EX)

$$T(n) = 2 * T(n/2 + 17) + n$$

Let's guess that $T(n) = O(n * (\log n))$.

i.e. $T(n) \leq c * n * (\log n)$ for some $c > 0$ and large enough n .

Pf.

$$\begin{aligned}T(n) &= 2 * T(n/2 + 17) + n \\&\leq 2 * c * (n/2 + 17) * (\log(n/2 + 17)) + n \\&= c(n + 34) * (\log(n/2 + 17)) + n \\&\leq c(n + 34) * (\log(3n/4)) + n \\&\text{(This assumption holds when } n/2 + 17 \leq 3n/4, n > 68) \\&= c * n * (\log n) + c * n * (\log 3n/4) + 34c(\log(3n/4)) + n \\&= c * n * (\log n) + n * ((c * \log(3/4)) + 1) + 34c(\log 3n/4) \\&\leq c * n * (\log n) \\&\text{(Last inequality holds when } n * ((c * \log(3/4)) + 1) \text{ goes to negative.} \\&\text{And that holds for } c \geq (1 / (\log 4/3)).)\end{aligned}$$

Above inequality holds for sufficiently large n . (choose $c = 5$)

CAUTION : Do Not Introduce new constant such as c' .

3. Master Theorem

Given a recurrence relation such as

$$T(1) = 1 \text{ AND } T(n) = a * T(n/b) + f(n) \text{ (for } n > 1, f(n) \text{ is polynomial function of } n.)$$

(Where a, b are positive constants, $f(n) = O(g(n))$ for some polynomial function $g(n)$.)

$$\begin{aligned} T(n) &= f(n) + a * T(n/b) \\ &= f(n) + a * (f(n/b) + a * T(n/(b^2))) \\ &= f(n) + a * (f(n/b) + a * (f(n/(b^2)) + a * T(n/(b^3)))) \\ &= \dots \\ &= (\sum_{i=0}^{k-1} (a^i * f(n/b^i))) + a^k * T(n/b^k) \end{aligned}$$

Assume $n = b^k$

$$\begin{aligned} T(n) &= (\sum_{i=0}^{k-1} (a^i * f(n/b^i))) + n^{(\log_b a)} \\ &= (\text{Particular solution}) + (\text{homogeneous solution}) \end{aligned}$$

Homogeneous solution is a key standard for determining the running time of given algorithm.

This solution is the cost of solving the boundary subproblems of case 1.

Particular solution is the sum of overheads.

Master Theorem :

Given a recurrence relation of

$$T(1) = c \text{ AND } T(n) = a * T(n/b) + f(n) \text{ for } n > 1.$$

(When $a \geq 1, b > 1$ and $c > 0$ are constants,

$f(n) = O(g(n))$ for some polynomial function $g(n)$).

Let the homogeneous function(solution) to be $n^{(\log_b a)} = h(n)$

Case 1 :

If $f(n)/h(n) = O(1/(n^\epsilon))$ for some $\epsilon > 0$,

(Homogeneous function overwhelms overhead function.)

then $T(n) = \Theta(h(n))$

Case 2 :

If $f(n)/h(n) = O((\log n)^k)$ for some non-negative integer k ,

then $T(n) = \Theta(h(n) * (\log n)^{(k+1)})$

Case 3 :

If $f(n)/h(n) = \Omega(n^\epsilon)$ for some $\epsilon > 0$

and $a * f(n/b) < f(n)$ for sufficiently large n ,
 (Original problem's overhead dominates next-upper problem's overhead.)
 then $T(n) = \Theta(f(n))$

Ex :

$$T(n) = 3 * T(n/2) + n$$

$$a = 3, b = 2, f(n) = n, h(n) = n^{\log_2 3}$$

$$f(n)/h(n) = n / (n^{\log_2 3}) = n^{(1 - \log_2 3)} = O(n^{(1 - \log_2 3)})$$

$$\text{Therefore, } T(n) = \Theta(h(n)) = \Theta(n^{\log_2 3})$$

$$T(n) = T(n/2) + d$$

$$a = 1, b = 2, f(n) = d, h(n) = n^{\log_2 1} = 1$$

$$f(n)/h(n) = d = \Theta((\log_2 n)^0)$$

$$\text{Therefore, } T(n) = \Theta(\log n)$$

$$T(n) = T(n/2) + (\log n)$$

$$a = 1, b = 2, f(n) = \log n, h(n) =$$

$$f(n)/h(n) = (\log n)/1 = \Theta(\log n)$$

$$\text{Therefore, } T(n) = \Theta((\log n)^2)$$

$$T(n) = T(n/2) + n$$

$$a = 1, b = 2, f(n) = n, h(n) = n^{\log_2 1} = 1$$

$$f(n)/h(n) = n/1 = n = \Sigma(n^1) \text{ (epsilon = 1)}$$

$$1 * f(n/2) = n/2 < n = f(n)$$

$$\text{Therefore, } T(n) = \Theta(f(n)) = \Theta(n)$$

Cf. Master Theorem holds for every natural number n . It does not have to be in form of $n = 2^k$

Without Loss of Generality, assume $n = 2^k, 4^k, \dots$

and $T(n)$ is non-decreasing function.

For $2^k, 2^{(k+1)}$, we can find that $f((2n)^p) = f(2^p * n^p)$ and 2^p will be constant.

So, for enough large n , in asymptotic aspect, we can neglect 2^p and this will give us 2^k and $2^{(k+1)}$ are similar.

Others :

Changing variables -

$$T(n) = 2T(\lfloor \sqrt{n} \rfloor) + (\log n)$$

Now let $n = 2^m$

$$T(2^m) = 2T(2^{m/2}) + m$$

And rename $T(2^m) = S(m)$. Then we obtain

$$S(m) = 2 * (m/2) + m$$

Now we can use Master Theorem.

By Master Theorem, $S(m) = \Theta(m * (\log m))$

Therefore, $T(n) = \Theta((\log n) * (\log \log n))$.

Using a stronger inductive hypothesis -

$$T(n) = 2 * T(n/2) + 1$$

Guess that $T(n) \leq c * n$ for enough large n .

$$T(n) \leq 2 * c * (n/2) + 1 = c * n + 1$$

This is not LEQ than $c * n$.

Instead, guess $T(n) \leq c * n - 2$

$$T(n) \leq 2 * (c * n/2 - 2) + 1 = c * n - 3 \leq c * n - 2.$$

Sorting and Order Statistics

Quicksort :

Pseudo Code :

```
Quicksort(A[], p, q)
    if p < q then
        r <- partition(A, p, q)
        Quicksort(A, p, r-1)
        Quicksort(A, r+1, q)
partition (A[], p, q)
    x <- A[p] // pivot
    last S_1 <- p // set S_1 and S_2 are empty.
    for i <- p+1 to q do // i : index of the 1st unknown
        if (A[i] < x) then
            Last S_1++;
            A[Last S_1] <-> A[i] // swap.
    A[Last S_1] <-> A[p]
    return Last S_1
```

Quicksort divides itself to two smaller part. So if that smaller problems return correct(sorted) results, overall sorting will return correct result.

Partition function needs THETA(n) time.

Cf. The Partition function does not guarantee a balanced partition.

It depends on pivot.

Worst-Case running time :

$$\begin{aligned} 1) T(n) &= \max_{(0 \leq k \leq n-1)} \{T(k) + T(n - k - 1)\} + \text{THETA}(n) \text{ (max function is used due to worst-case analysis).} \\ &\geq T(0) + T(n-1) + \text{THETA}(n) \\ &= T(n-1) + \text{THETA}(n) \\ &= \sum_{k=1}^n \{\text{THETA}(k)\} \\ &= \text{THETA}(\sum_{k=1}^n \{k\}) \\ &= \text{THETA}(n^2) \end{aligned}$$

Therefore, $T(n) = \text{SIGMA}(n^2)$

2) Guess $T(n) = O(n^2)$.

In other words, $T(n) \leq c * (n^2)$ for some $c > 0$.

$$\begin{aligned} T(n) &= \max_{(k=0 \text{ to } n-1)} \{T(k) + T(n - k - 1)\} + \text{THETA}(n) \\ &\leq \max_{(k=0 \text{ to } n-1)} \{c*(k^2) + c(n - k - 1)^2\} + \text{THETA}(n) \\ &= c * (\max_{(k=0 \text{ to } n-1)} \{k^2 + (n - k - 1)^2\}) + \text{THETA}(n) \\ &= c * n^2 - 2cn + c + \text{THETA}(n) \\ &\leq c * n^2 \text{ for some } c \text{ and large enough } n \text{ such that } 2cn \text{ dominates } c + \text{THETA}(n). \end{aligned}$$

Therefore, $T(n) = O(n^2)$.

From above 1 and 2, $T(n) = \text{THETA}(n^2)$

Average-Case running time :

Cf. If always well balanced,

$$T(n) = 2 * T(n/2) + \text{THETA}(n)$$

By Master Theorem, $T(n) = \text{THETA}(n * (\log n))$.

Cf. If always partitioned to 1 : 2?

$$T(n) = T(n/3) + T(2n/3) + \text{THETA}(n)$$

Result is same with Average-Case.

Assume every input permutation occurs with an equal probability.

All possible cases are $n!$.

Assume all keys are distinct.

Therefore, $T(1) = \text{THETA}(1)$.

$$\begin{aligned} T(n) &= (1/n) \text{sigma}_{(i=0 \text{ to } n-1)} \{T(i) + T(n - i - 1)\} + \text{THETA}(n) \text{ (This means pivot is } (i+1)\text{th element.)} \\ &= (2/n) \text{sigma}_{(i=0 \text{ to } n-1)} \{T(i)\} + \text{THETA}(n) \end{aligned}$$

Guess $T(n) \leq c * n(\log n)$ for some $c > 0$.

$$\begin{aligned} \text{Then, } T(n) &= (2/n) \text{sigma}_{(k=0 \text{ to } n-1)} \{T(k)\} + \text{THETA}(n) \\ &\leq (2/n) \text{sigma}_{(k=2 \text{ to } n-1)} \{T(k)\} + \text{THETA}(n) \text{ (THETA}(n) \text{ absorbed } T(0), T(1).) \\ &\leq (2/n) \text{sigma}_{(k=2 \text{ to } n-1)} \{c * k * (\log k)\} + \text{THETA}(n) \\ &= (2c/n) \text{sigma}_{(k=2 \text{ to } n-1)} \{k * (\log k)\} + \text{THETA}(n) \\ &\leq (2c/n) \text{integral}_{(1 \text{ to } n)} \{x * (\log x)\} dx + \text{THETA}(n) \\ &= (2c/n) ([1/2 * x^2 * (\log x)]_{(1 \text{ to } n)} - [x^2 / 4]_{(1 \text{ to } n)}) + \text{THETA}(n) \end{aligned}$$

$$= c * n * (\log n) + (-cn/2 + c/2n + \text{THETA}(n))$$

$$\leq c * n * (\log n) \text{ (We can choose } c > 0 \text{ such that } cn/2 \text{ dominated } \text{THETA}(n).)$$

Therefore, $T(n) = O(n * (\log n))$.

Randomized Quicksort :

It is a quicksort that pivot element is not effected by input list.

Pivot is selected after input list is determined.

Running time is the same as quicksort.

Decision-Tree methods of computation.

Assume all elements for sorting are distinct.

Then, Every sorting algorithm that is based on comparison has time complexity of $\text{SIGMA}(n * (\log n))$.

Def. A decision tree is a binary tree in which the internal nodes are labeled with a comparison of the form $a_i : a_j$ and the edges are labeled with "<" or ">"

On an input sequence(a_1, a_2, \dots, a_n) the computation of a decision tree starts at the root.

Suppose we are at node n with labeled $a_i : a_j$, then a_i and a_j are compared with each other.

If $a_i > a_j$ then move to the right child of n , else move to the left child of n .

The computation stops when we reach a leaf node.

Ex) A decision tree for sorting 3 numbers a, b, c

This decieion tree has 6 possible outputs.(For n numbers : $n!$)

Height of this tree is worst-case running time of decision.

Fact 1 :

A decision tree that sorts n numbers must have $n!$ leaves.

Fact 2 :

A binary tree of height h have at most 2^h leaves.

Theorem :

A decision tree for sorting n numbers must have height at least $(\log_2 (n!))$.

Pf. Straight forward from Fact 1 and Fact 2.

Corollary :

Any comparison-based sorting algorithm has to take time $\text{SIGMA}(n * (\log n))$ in the worst case.

Pf.

$$n! = \sqrt{2 \cdot \pi \cdot n} \cdot (n/e)^n \cdot (1 + \Theta(1/n)) \text{ (From Stirling's approximation.)}$$

$$> (n/e)^n$$

We take log of base 2 to both side.

$$\log(n!) > n \cdot (\log n) - n \cdot (\log e) = \Theta(n \cdot (\log n)).$$

Linear-Time Sortings :

Note.

The lower bound $\Omega(n \cdot (\log n))$ is for sortings based on comparisons.

However, when some conditions are met, lower-than- $O(n \cdot (\log n))$

sortings are possible.

Counting sort :

when each input is an integer in $[1, O(n)]$, $\Theta(n)$ -time sorting is possible.

Pseudo-Code :

```
Counting-Sort(A[], n, k)
// n : Number of elements
// k : A[i] in [1, k], forall i = 1, ..., n
    for i in [1 .. k]
        C[i] <- 0
    for i in [1 .. n]
        C[A[i]]++
    for i in [1 .. k]
        C[i] <- C[i] + C[i-1]
    // Here, C[i] : Number of elements not greater than i.
    for i in [n .. 1]
        B[C[A[i]]] <- A[i]
        C[A[i]]--
```

Running time is $\Theta(n)$ obviously.

Radix sort :

Pseudo-Code :

```
Radix-Sort(A[], d, n) =
// d : Number of digits in A[i] for all i = 1, ..., n
```

for i in [1 .. d]

sort A[1 .. n] on digit i by a stable sort.(This sort should finish in linear time.)

Running Time : $\Theta(k * n) = \Theta(n)$ (k : constant).

Bucket sort :

Pseudo-Code :

Bucket-Sort(A[], n) =

// A[] : array A[1 .. n] with uniform distribution in [0, 1).

for i in [1 .. n]

Insert A[i] into list B[$\lfloor n * A[i] \rfloor$]

for i in [0 .. n-1]

Sort list B[i] using an appropriate sorting algorithm

print the list B[0], ..., B[n-1] in order.

*It does not have to be [0, 1). Transformation can be taken.

Running time(with uniform distribution) :

The expected numbers of elements n_i in each list B[i] is $np=1$ with variance $\sigma^2 = np(1-p) = 1-(1/n)$. ($p = 1/n$)

$\sigma^2 = E(n_i^2) - (E(n_i))^2$

Therefore, $E((n_i)^2) = (1 - (1/n))^2 + 1^2 = 2 - (1/n) = \Theta(1)$

Every insertion sort takes $O((n_i)^2) = \Theta(1)$

Therefore, $O(n)$ in total for bucket sort.

Order Statistics :

Given a list of n elements,

Selecting minimum : needs n-1 comparisons.

Selecting maximum : needs n-1 comparisons.

Selecting i(th) smallest element :

A trivial upper bound is $O(n * (\log n))$ (by sorting)

A trivial lower bound is $\Omega(n)$

Average-Case $O(n)$ algorithm :

Select(A[], p, r, i)

```

// Return i(th) smallest from A[p .. r]
if i > r - p + 1 or i <= 0 then
    return "error!"
if p = r then
    return A[p] // i = 1
q <- partition(A, p, r)
k <- q - p + 1
if i = k then return A[q]
else if (i < k) then select(A, p, q-1, i)
else select(A, q+1, r, i-k)

```

Worst-Case :

$$\begin{aligned}
 T(n) &= \text{THETA}(n) + \max_{1 \leq k \leq n-1} \{T(k)\} \\
 &\geq \text{THETA}(n) + T(n-1) \\
 &\geq \text{THETA}(n) + \text{THETA}(n-1) + T(n-2) \\
 &\geq \text{THETA}(\sum_{i=2}^n i) + T(1) \\
 &= \text{THETA}(n^2)
 \end{aligned}$$

Therefore, $T(n) = \text{SIGMA}(n^2)$

Average-Case :

$$\begin{aligned}
 T(n) &\leq \text{THETA}(n) + (1/n) * (\sum_{k=1}^n \{T(\max(k-1, nk))\}) \\
 &\leq \text{THETA}(n) + (2/n) * (\sum_{k=\lceil n/2 \rceil}^n \{T(k-1)\}) \\
 &= \text{THETA}(n) + (2/n) * (\sum_{k=\lceil n/2 \rceil-1}^{n-1} \{T(k)\})
 \end{aligned}$$

Now, guess $T(n) \leq cn$ for some $c > 0$ and large-enough n .

$$\begin{aligned}
 T(n) &\leq \text{THETA}(n) + (2/n) * (\sum_{k=\lceil n/2 \rceil-1}^{n-1} \{ck\}) \\
 &= \text{THETA}(n) + (2c/n) * \{(\sum_{k=1}^{n-1} \{k\}) - (\sum_{k=1}^{\lceil n/2 \rceil-1} \{k\})\} \\
 &\leq \text{THETA}(n) + (2c/n) * ((n^2 - n)/2 - ((n/2)^2 - \dots)/2) \\
 &= cn - c - (cn/4) + \text{THETA}(n) \\
 &\leq cn
 \end{aligned}$$

We can choose $c > 0$ s.t. $cn/4$ dominates $\text{THETA}(n) - c$.

Therefore, $T(n) = O(n)$.

($T(n)$ is definitely $\text{SIGMA}(n)$, because we have to at least see n elements 1 time.)

Worst-Case $O(n)$ algorithm :

Idea :

- 1) Devide the input into $n/5$ groups each having 5 elements.
- 2) Find the medians of each group.
Let the medians be $m_1, m_2, \dots, m_{\lceil n/5 \rceil}$.
- 3) Let median of $m_1, m_2, \dots, m_{\lceil n/5 \rceil}$ be M .
- 4) Use M as the pivot element and partition the whole elements as in Quicksort.
- 5) Select the proper side of the partitions and recursively conduct steps 1~5.

(Number of elements that are equal to M or less) $\geq 3(n/10 - 1) = (3n/10) - 3$

(Number of elements that are equal to M or greater) $\geq 3(n/10 - 1) = (3n/10) - 3$

Linear-Select($A[]$, p , r , i) =

if ($p = r$) then

return $A[p]$

Find the pivot M as described in Steps 1~5.

Exchange $A[p]$ and M

$q \leftarrow \text{partition}(A, p, r)$

$k \leftarrow q - 1 +$

if ($i = k$) then

return $A[q]$

else if ($i < k$) then

Linear-Select(A , p , $q-1$, i)

else

Linear-Select(A , $q+1$, r , $i-k$)

Worst-Case running time :

$T(n) \leq \text{THETA}(n) + \text{THETA}(n) + T(\lceil n/2 \rceil) + \text{THETA}(n) + T((7n/10) + 3)$

// $T((7n/10) + 3)$ is the number of remaining elements from equation $n - ((3n/10) - 3)$

// First $\text{THETA}(n)$: step 1

```

// Second THETA(n) : step 2
// Third T(ceiling(n/2)) : step 3
// Fourth THETA(n) : step 4
// Fifth T((7n/10) + 3) : step 5 (due to (7n/10) + 3 > (3n/10) - 3).
<= THETA(n) + T((n/5) + 1) + T((7n/10) + 3)

```

Guess $T(n) \leq cn$ for some $c > 0$, $n_0 > 0$ and $n \geq n_0$.

Therefore,

$$\begin{aligned}
 T(n) &\leq \text{THETA}(n) + c((n/5) + 1) + c((7n/10) + 3) \\
 &= c(9n/10) + 4c + \text{THETA}(n) \\
 &= cn - (cn/10) + 4c + \text{THETA}(n) \\
 &\leq cn
 \end{aligned}$$

We choose c to make $(cn/10)$ dominates $4c + \text{THETA}(n)$.

Therefore, $T(n) = O(n)$.

Dynamic Programming

Elements of dynamic programming :

Optional substructure :

an optimal solution contains within it optimal solutions to the problems.

Overlapping Subproblems :

A recursive algorithm revisits the same subproblems over and over again.

Matrix-chain Multiplication :

Given a matrix $\langle A_1, A_2, \dots, A_n \rangle$, we want to compute the product $A_1 * A_2 * \dots * A_n$.

Note that Matrix multiplication is associative.

If A is $p \times q$ matrix and B is $q \times r$ matrix,

the cost of computing the product AB is pqr .

Object :

We want to find a sequence of multiplications with the minimum cost.

This is finding an optimal parenthesizing.

Ex)

A : 10×100 matrix

B : 100×5 matrix

C : 5×50 matrix

(AB)C \rightarrow 7500 multiplications.

A(BC) \rightarrow 75000 multiplications.

Let matrix A_k be matrix of dimension $p_{(k-1)} \times p_k$

We consider n matrices, A_1 to A_n .

Then we have $n-1$ choices, when parenthesizing into 2 parts.

General form : $(A_1 * A_2 * \dots * A_k)(A_{(k+1)} * \dots * A_n)$

Let C_{ij} : Minimal computation cost of $A_i * \dots * A_j$

Then our goal is to find C_{1n}

For general form described above, cost is $(p_0 \times p_k \times p_n) + C_{1k} + C_{(k+1)n}$.

So, $C_{1n} = \min_{(k=1 \text{ to } n-1)} \{(p_0 * p_k * p_n) + C_{1k} + C_{(k+1)n}\}.$

Therefore, $C_{ij} = \min_{(k=i \text{ to } j-1)} \{(p_i * p_k * p_j) + C_{ik} + C_{(k+1)j}\}.$

$C_{ij} = 0$ (when $i=j$)

$\min_{(k=i \text{ to } j-1)} \{(p_i * p_k * p_j) + C_{ik} + C_{(k+1)j}\}$ (when $i < j$)

Algorithm :

MC(i, j) =

// Return the minimum cost for Matrix multiplication $A_i * \dots * A_j$

if $i=j$ then

return 0

min <- inf

for $k=i$ to $j-1$ do

q <- $MC(i, k) + MC(k+1, j) + p_{(i-1)} * p_k * p_j + c_{ik} + c_{(k+1)j}$

if $q < \text{min}$ then

min <- q

return min

Running time of MC(1, n) by the number of calculations $p_{(i-1)} * p_k * p_j$:

$T(n) = 0$ (when $n = 1$)

$\sigma_{(k=1 \text{ to } n-1)} \{T(k) + T(n-k) + 1\}$ (when $n > 1$)

= 0 (when $n = 1$)

$2 * \sigma_{(k=1 \text{ to } n-1)} \{T(k) + (n-1)\}$ (when $n > 1$)

$T(n)$ grows exponentially for n.

Guess $T(n) \geq c * 2^n$

$T(n) = 2 * \sigma_{(k=1 \text{ to } n-1)} \{T(k) + (n-1)\}$

$\geq 2 * \sigma_{(k=1 \text{ to } n-1)} \{c * 2^k + (n-1)\}$

= $2c(2^n - 2) + (n - 1)$

= $2c * 2^n - 4c + n - 1$

> $c * 2^n$

for large-enough n.

Therefore, $T(n) = \text{SIGMA}(2^n)$

Observation :

The recursive algorithm $MC(i, j)$ calls the same subproblems repeatedly.

Memorizing the minimum costs of subproblems can dramatically save running time.

Dynamic programming solution :

Construct a two-dimensional array $m[i][j]$ where $m[i][j]$ has the minimum cost for $A_i * \dots * A_j$.

A sample bottom-up implementation :

```
for i=1 to n do
    m[i][i] = 0
for l=1 to n-1 do // l : length
    for i=1 to n-l do // i, j : <A_i * ... * A_j>
        j <- i + l
        m[i][j] <- min_(k=i to j-1){m[i][k] + m[k+1][j] + p_(i-1)*p_k*p_j}
return m[1][n]
```

Running time : $\Theta(n^3)$

Top-down implementation is also possible :

Instead of calling the same subproblems, look at the table entries and get the values.

Called "Memoization"

$MC(i, j)$ = // Assume all $m[i][j]$'s initialized to -1.

```
if m[i][j] <> -1 then
    return m[i][j]
else if i=j then
    m[i][j] <- 0
    return m[i][j]
else
    the same as before. (in  $\Sigma(2^n)$  algorithm)
    m[i][j] <- min
    return m[i][j]
```

Running time : $\Theta(n^3)$

Longest Common Subsequence (LCS) :

Given two sequences $X = \langle x_1, x_2, \dots, x_n \rangle$ and $Y = \langle y_1, y_2, \dots, y_m \rangle$
we want to find an LCS of X and Y.

Note 1)

$\langle BCDB \rangle$ is a subsequence of $\langle ABCBDAB \rangle$

$X = \langle ABCBDAB \rangle$, $Y = \langle BDCABA \rangle$

$\langle BCA \rangle$ is a common subsequence of X and Y.

$\langle BCBA \rangle$, $\langle BDAB \rangle$ are LCSs of X and Y.

Optimal substructure of LCS :

$X_m = \langle x_1, \dots, x_m \rangle$, $Y_m = \langle y_1, \dots, y_m \rangle$

Let $Z_k = \langle z_1, \dots, z_k \rangle$

1. If $x_m = y_n$ then $z_k = x_m = y_n$ and $z_{(k-1)}$ is an LCS of $X_{(m-1)}$ and $Y_{(n-1)}$

2. If $x_m \neq y_n$ then $z_k = \max\{\text{LCS}(X_m, Y_{(n-1)}), \text{LCS}(X_{(m-1)}, Y_n)\}$

c_{ij} = the length of $\text{LCS}(x_i, y_j)$

$c_{ij} = 0$ (if $i = 0$ or $j = 0$)

$c_{(i-1)(j-1)} + 1$ (if $i, j > 0$ and $x_i = y_j$)

$\max\{c_{(i-1)j}, c_{i(j-1)}\}$ (if $i, j > 0$ and $x_i \neq y_j$)

Obviously, exponential-time implementation is possible by a recursive algorithm.

But, because of the existence of $\Theta(mn)$ subproblems, DP is possible.

$\text{LCS-length}(X, Y) =$

//m, n : the length of sequences X and Y

for $i = 0$ to m do

$c[i, 0] = 0$;

for $j = 0$ to n do

$c[0, j] = 0$;

for $i=1$ to m do

for $j=1$ to n do

if $x_i = y_j$ then

```

        c[i, j] = c[i-1, j-1] + 1;
    else
        c[i, j] = max{c[i-1, j], c[i, j-1]}
    return c[m, n]

```

Optimal binary search tree (static) :

1. $S = \{x_1, x_2, \dots, x_n\}$ where $x_1 < x_2 < \dots < x_n$
2. P_i = the probability that member(S, x_i) is called, $i = 1, \dots, n$
3. q_i : the probability that member(S, x) is called for $x_i < x < x_{(i+1)}$, $i = 0, 1, \dots, n$ (unsuccessful search)
(let $x_0 = -\infty$, $x_{(n+1)} = \infty$)

We want to find a BST that has the minimum expected number of comparisons.

The cost of a BST with $x_1 < x_2 < \dots < x_n$ is

$$\text{SIGMA}_{\{i=1 \text{ to } n\}}(p_i * (\text{depth}(x_i) + D)) + \text{SIGMA}_{\{i=0 \text{ to } n\}}(q_i * \text{depth}(e_i))$$

c_{ij} = the optimal cost for a tree with $x_i, x_{(i+1)}, \dots, x_j$

$$w_{ij} = \text{SIGMA}_{\{l=i \text{ to } j\}}(p_l) + \text{SIGMA}_{\{l = i-1 \text{ to } j\}}(q_l)$$

Optimal substructure :

$$c_{ij} = 0 \text{ (if } i > j)$$

$$\min_{(k=i \text{ to } j)}\{c_{i(k-1)}, c_{(k+1)j}\} + w_{ij} \text{ (if } i \leq j)$$

Running time for $c_{1n} = \text{THETA}(n^3)$ (due to $\text{SIGMA}_{(i=1 \text{ to } n)}\{\text{SIGMA}_{(j=i \text{ to } n)}\{(j-1+1)\}\} = \text{THETA}(n^3)$)

Static Data = Does not change once given.

Dynamic Data = Changes over time.

We want to support the following operations :

- 1) Member(S, x)
- 2) Insert(S, x)
- 3) Delete(S, x)
- 4) Min(S), Max(S)
- 5) Union(A, B, C) ($C = A \cup B$)
- 6) Find(x)

Data structures that supports :

Member, Insert, Delete : Dictionary
Min, Insert, Delete : Priority Queue
Union, Find : UnionFind

Hash Table : Can be used for dictionary

BST : Can be used for both dictionary and priority queue.

Heap : Can be used for priority queue.

Hash Table :

Idea :

Hash function h , universe U , hash table size of m .

$h : U \rightarrow \{0, 1, 2, \dots, m-1\}$

cf. Usually $m \ll |U|$

Collision :

If $x \neq y$ and $h(x) = h(y)$ then we have a collision.

Hash Table is the most ideal data structure for searching data. It gives data in constant time.

Hash functions :

Assume that all keys are natural numbers.

(Non-natural number can be easily transformed to natural numbers.)

Division method :

$$h(x) = x \bmod m$$

Multiplication method :

$$h(x) = \lfloor m(xA \bmod 1) \rfloor, 0 < A < 1$$

Collision Resolution :

1) Chaining : Use linked list when same hash function value occurs.

Let m be the number of keys in a hash table.

Define the load factor $\alpha = n/m$

Assumption :

h is computable in $O(1)$ time.

h distributes keys uniformly in the table.

All elements of U occur with equal probability.

Fact :

Under the above assumption, a search takes $\Theta(\max(1, \alpha))$ on average.

2) Open Addressing :

$h_0, h_1, \dots, h_{(n-1)}$

Insert(S, n) =

$h_0(x)$

if $h_0(x)$ is occupied, $h_1(x)$

if $h_1(x)$ is occupied, $h_2(x)$

...

if $h_i(x)$ is occupied, $h_{(i+1)}(x)$.

Ideally want $h_0(x), h_1(x), \dots, h_{(m-1)}(x)$ be a permutation of $0, 1, \dots, m-1$.

cf. Linear probing : $h_i(x) = (h_0(x) + (c \cdot i)) \bmod m$.

cf. Quadratic probing : $h_i(x) = (h_0(x) + (c_1 \cdot i) + (c_2 \cdot i^2)) \bmod m$

cf. double hashing :

h, h' : hash functions.

$$h_i(x) = (h(x) + i \cdot h'(x)) \bmod m.$$

Double hashing prevents secondary clustering problem of Linear or Quadratic hashing by adding one more hash function h' .

Uniform hashing :

For any key x , the probe sequence $h_0(x), h_1(x), \dots, h_{(m-1)}(x)$ is

a random permutation of $\{0, 1, \dots, m-1\}$

This is most ideal hashing.

Theorem :

Assume uniform hashing, the expected number of probes in an unsuccessful search or an insertion is at most $1/(1-\alpha)$, where α is the load factor. (At most times, $1/(1-\alpha) \leq 2$.)

Proof :

cf. Insertion and unsuccessful search is same, since insertion occurs right after unsuccessful search.

Let $p_i = \Pr(\text{exactly } i \text{ probes access occupied slots})$

Let $q_i = \Pr(\text{at least } i \text{ probes access occupied slots})$

Expected number of probes $= 1 + \sum_{i \geq 0} i \cdot p_i$ (note that $p_0 = 0$)

$$= 1 + \sum_{i \geq 1} (q_i - q_{i+1})$$

$$= 1 + \sum_{i \geq 1} q_i$$

$$\leq 1 + \sum_{i \geq 1} \alpha^i$$

Here, $q_1 = n/m$, $q_2 = (n/m) \cdot (n-1/m-1)$, ..., $q_i = (n/m) \cdot (n-1/m-1) \cdot \dots \cdot (n-i+1/m-i+1)$.

$n/m = \alpha$, $n-1/m-1 > \alpha$, and so on.

So above inequality holds.

$$= 1/(1-\alpha)$$

Theorem :

Under the same assumption the expected number of probes in a successful search is at most $(1/\alpha) \ln(1/(1-\alpha))$.

Proof :

Successful search exactly steps the same footprint of insert operation.

The load factor right after the i th key has been inserted was i/m .

If x was the $(i+1)$ th key inserted, then the expected number of probes in a successful search for x is by the previous theorem, at most $1/(1-(i/m)) = m/(m-i)$

Average overall keys.

$$(1/n) \sum_{i=0}^{n-1} \{m/(m-i)\} = (m/n) \sum_{i=0}^{n-1} \{1/(m-i)\}$$

$$\leq (1/\alpha) \int_0^1 \{1/(1-x)\} dx$$

$$= (1/\alpha) \ln(1/(1-\alpha))$$

Binary Search Tree(BST) :

All keys should have different values.

All values of node that is on left subtree is smaller than given node's value, vice versa.

There are $n!$ sequence when size of n input is given.

We cannot guarantee the balancedness of BST when we do not know the root of BST.

When median is chosen for root of every subtree and the whole BST, BST is guaranteed to be balanced.

Theorem :

A sequence of n inputs into an empty BST takes time $O(n \log n)$ on average.

(Assume that every permutation of the input sequence is equally likely.)

Proof :

Assume the input sequence is a_1, a_2, \dots, a_n .

Let b_1, b_2, \dots, b_n be the input sequence sorted in order.

a_1 will become the root of the resulting BST.

Assume $a_1 = b_j$ for some j ($1 \leq j \leq n$).

Then, $b_1, \dots, b_{(j-1)}$ will go to left subtree of a_1 (root).

And likely, $b_{(j+1)}, \dots, b_n$ will go to right subtree of a_1 (root).

Let $T(n)$ be the expected number of comparisons for inserting n keys into an empty BST.

$T(n) = (1/n) \text{SIGMA}_{(j=1 \text{ to } n)} \{ (j-1) + T(j-1) + (n-j) + T(n-j) \}$ // $j-1, n-j$ corresponds to comparison with root.

$= (1/n) \text{SIGMA}_{(j=1 \text{ to } n)} \{ T(j-1) + T(n-j) \} + \text{THETA}(n)$

$= (2/n) \text{SIGMA}_{(j=1 \text{ to } n-1)} \{ T(j) \} + \text{THETA}(n)$ // This form of formula has seen with quicksort.

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

Note that worst case has $\text{THETA}(n^2)$ time complexity.

Balanced Binary Search Tree :

(Red-Black tree as an example)

Definition : https://en.wikipedia.org/wiki/Red%E2%80%93black_tree

(From PPT of Lecture) :

Every node in the search tree has a color : red or black.

It has to satisfy the following properties(red-black properties) :

- 1) Root is black
- 2) Every leaf is black

3) If a node is red, its children should be black

4) In any path from the root to a leaf, the number of black node on the path is the same. (which is called "Black height")

cf. Here, leaf does not mean ordinary leaf node. We assume that every NIL pointer points to "NIL" leaf node.

cf. In worst case, the longest path length cannot exceed twice length of shortest path.

Theorem :

In a red-black tree of n nodes, the worst-case depth is $O(\log n)$

Def. Black height :

The number of black nodes from, but not including, the root to a (any) leaf

Lemma :

A red-black tree with n internal nodes has height at most $2 * \log(n+1)$.

Proof :

For any node r in a red-black tree, the subtree rooted at r contains at least $2^{BH(T)} - 1$ internal nodes. --- (1)

Let h be height of the tree.

Then by property 3,

$$h \leq 2 * BH(T)$$

$$h/2 \leq BH(T) \text{ --- (2)}$$

$$n \geq 2^{BH(T)} - 1 \text{ --- from (1)}$$

$$\geq 2^{(h/2)} - 1 \text{ --- from (2)}$$

Therefore, $h \leq 2 * \log(n+1)$

Insertion in Red-Black tree : Refer to PPT.

Deletion in Red-Black tree : Refer to PPT.

Disjoint Sets :

Want to handle a collection of disjoint sets.

(Example usage : Kruskal algorithm)

Operations :

1) Make-Set(x) - Create a singleton $\{x\}$.

2) Union(A, B) - Create the set $A \cup B$. Destroy A, B .

3) Find(x) - Return the name(the representative) of the set containing x .

Linked-list representation :

Represent each set by a linked list.

EX : {2, 4, 7} with representative as 2. 4 → 2, 7 → 2 (2 → 2. 2 → next : 4, 4 → next : 7)

Make-Set(x) - THETA(1)

Find(x) - THETA(1)

Union(A, B) - $O(|A| + |B|)$ (set new representative, link representative's set's last element's next to non-representative's head.)

Weighted-Union Heuristic

- When Union(), attach the smaller set to the larger one.

Theorem :

Using linked-list representation and weighted-union heuristic,

a sequence of m make-Set, Union, and Find operations. n of which are Make-Set operations (i.e. n elements)

takes $O(m + n \log n)$ time.

Proof.

<Amortized Analysis>

For element x, when changing link in x occurs n times,

set containing x has size at least $2^{(k-1)}$ (self linking counts.)

But $2^{(k-1)} \leq n$ since there are n elements.

Therefore, $k \leq \log n$. and for all n elements, less than or equal to $n \log n$

So, running time is $O(m + n \log n)$.

Note. The dominating cost for union is spent for updating the pointers to the representative.

For any element x, each time x's pointer is updated (in union) it belongs to the

smaller set (if it belongs to the larger set, no update needed.)

The set size containing x : $1 \rightarrow \geq 2 \rightarrow \geq 4 \rightarrow \dots \rightarrow \geq 2^k$

If the number of elements is n, there can be at most $\log n$ updates for any x.

Therefore, the total time for update is $O(n \log n)$

Since each Make-Set and Find takes $O(1)$ time, the total cost for the entire sequence is $O(m + n \log n) = O(\max\{m, n \log n\})$

The amortized cost per union is $O(\log n)$

Tree representation :

Each set is represented by a tree.

Make-Set(x) : x's link points to x.

Union(A, B) : B's root's link points to A's root.

Make the root of one tree be a child of the root of the other tree.

Find(x) : follow the parent pointers to the root. - $O(n)$

cf. This case, Find(x) operation is the time-dominant operation.

Two heuristics to improve running time :

1) Union by rank (almost same as weighted union.)

Each element x has an integer value rank[x], to maintain an upper bound on the height of x's subtree).

After Make-Set(x), rank[x] = 0

After that, rank[x] changes only when rank[x] = rank[y] and y's tree is merged into x. (increased by 1)

The tree with smaller height becomes a subtree of the other tree.

Fact. With Union by rank, a sequence of n Make-Set, Union, Find takes time $O(n \log n)$

2) Path Compression

Each time we do a Find(x), make every node on the path from x to the root a child of the root.

Fact, With path compression(only), a sequence of n Make-Set, Union, Find takes time $O(n \log n)$.

Theorem :

With union-by-rank and path compression, a sequence of n Make-Set, Union, Find takes $O(n * (\log^* n))$,

where $\log^* n = \min\{k(\log \log \dots \log n \leq 1, \log \text{ for } k \text{ times.})\}$

Greedy Algorithms

Greedy Algorithm :

Choose a sequence of locally best choices.

Usually doesn't guarantee the global optimum.

cf. Examples not guaranteeing optimum : coin change, knapsack

(knapsack : filling knapsack with different value items, each may have different volumes.)

(real-number knapsack problem's optimal solution can be found with greedy algorithm, where integer knapsack doesn't)

cf. Examples guarantees optimum : Seminar-room assignment

Seminar-room problem :

Get reservations from people who wants to use seminar room. People gives start time and finish time of their reservation.

Goal is to satisfy the most reservations.

We assign seminar-rooms in order of fastest finish time avoiding collision of previously reserved times.

Theorem :

Above algorithm guarantees optimal.

Proof.

Claim : There exists at least an optimal schedule that includes the meeting with the earliest finishing time.

Proof for claim :

Assume there exists an optimal schedule S not including the meeting, say m_1 , with the earliest finishing time.

Let m_i be the meeting with the earliest finishing time in S .

Then, $S - \{m_i\} \cup \{m_1\}$ is also an optimal schedule.

Therefore, selecting m_1 is safe in the sense that it does keep a path to an optimal schedule.

Thus, repeat the process with the set of meetings excluding m_1 and those conflicting to m_1 .

Metroid :

It is related with independent set.

1. heredity property

2. exchange property

Def.

A metroid $M = (S, I)$ of a finite set S and a non-empty set I consisting of "independent" subsets of S satisfying :

1. If B in I and A is subset of B , then A is in I . (heredity) (also empty set is in I)
2. If A in I , B in I and $|A| < |B|$, then there is some element x in $B-A$ such that $A \cup \{x\}$ is in I .(exchange)

An example :

Graphic metroid $M_G = (E, F)$ where E is the set of edges, F is the set of acyclic subset of edges(which is same as set of trees = Forest).

Claim : M_G is metroid i.e. the set of forests is an independent set.

Proof for claim :

- (1) A subset of forest is also a forest.(heredity)
- (2) Consider two forests A and B s.t. $|A| < |B|$. Since A is a forest, A consists of zero or more (disconnected) trees.
For each tree of A of size k , B does not contain more than k edges connecting the vertices in the tree.
Thus, there is at least an edge e that connects two different trees of A .
 $A \cup \{e\}$ is also a forest. (e is extension of A)

Def. In a metroid $M = (S, I)$ and A in I , if $A \cup \{x\}$ is in I for x not in A , then we call x is an extension of A .

Def. In a metroid $M = (S, I)$, A in I is called maximal if A does not have any extension.

Theorem :

In a metroid $M = (S, I)$, every maximal subset of I has the same size.

Proof.

Assume for the contradiction there are two maximal subsets A, B in I of different sizes.

Let $|A| < |B|$ without loss of generality.

By the property of exchange, there is an element x in $B-A$ such that $A \cup \{x\}$ in I .

Contradiction to the fact that A is maximal.

Greedy algorithm on a weighted metroid :

On a metroid $M = (S, I)$ where every element has a positive weight, we want to find a subset in I with the maximum weight sum.

cf. A subset with the maximum weight sum is always maximal.(Due to every weight is positive.)

Algorithm_Greedy(M, W) // M : metroid (S, I) , W : weights

```
{
  A <- empty_set
  Sort S in decreasing order by W.
  for each x in S in decreasing order by W
    if (A ∪ {x} in I) then
```

```

        A <- A U {x}
    return A
}

```

An example : kruskal algorithm for minimum spanning trees.

Theorem :

The algorithm greedy guarantees a maximum-weight independent set.

Proof.

Assume to the contrary that there is an independent subset X with a greater weight sum than that of the subset A , the output of greedy.
 since A and X are both maximal subset, $|A| = |X|$.

Let the weights of A be $a_1 \geq a_2 \geq \dots \geq a_n$ and the weights of X be $x_1 \geq x_2 \geq \dots \geq x_n$.

(And let a_1, a_2, \dots, a_n be the sequence added to A .)

Since $W(X) > W(A)$, there is at least an x_i such that $x_i > a_i$.

Consider $A_{(i-1)} = \{a_1, a_2, \dots, a_{(i-1)}\}$ and $X_i = \{x_1, x_2, \dots, x_i\}$

By the property of exchange, there is at least an x_k in $\{X_i - A_{(i-1)}\}$

such that $A_{(i-1)} \cup \{x_i\}$ is in I .

Since all x_k 's are greater than a_i , the algorithm Greedy chooses the greatest x_k among them.

Therefore, if there exists a better subset than A , the algorithm Greedy never returns A .

this means that if the algorithm Greedy returns A , then there is no better solution.