# Sufficient Algorithms

## Charlotte Xia

## November 27, 2023

## Contents

# 1  Fibnacci Problem

## 1.1  Problem Description

The well known Fibnacci Problem is defined as follows:

> **Definition 1.1.**
> $$F_n = \begin{cases} F_{n-1}+F_{n-2} & \text{if } n > 1 \\ 1 & \text{if } n = 1 \\ 0 & \text{if } n = 0 \end{cases} \tag{1}$$

## 1.2  Solution

### 1.2.1  Naive Case and Time Complexity

Intuitively, we count the next number by adding the previous two numbers.

---

**Algorithm 1:** fib1(n)

---

**1** if n=0: return 0
**2** if n=1: return 1
**3** return fib1(n-1)+fib1(n-2)

---

However, it has an exponential time complexity.

> **Lemma 1.1.** $T_n > F_n$ ,where $T_n$ stands for the steps needed to calculate $F_n$.

**Proof.** Firstly, we have:

$$T(n) = T(n-1) + T(n-2) + 3 \textbf{ for } n > 1$$

where 3 extra steps are needed for 2 'if' judgements and the addition operation.
In the Base Step, for $n = 0$, we have $T_0 = 1 > 0 = F_0$, and for $n = 1$, we have $T_1 = 2 > 1 = F_1$.
In the Inductive Step, suppose $T_i > F_i$ for $i = 1, 2, \ldots, n$.

$$\begin{aligned} T_{n+1} &= T_n + T_{n-1} + 3 \\ &> F_n + F_{n-1} + 3 \quad \text{(induction hypothesis)} \\ &= F_{n+1} + 3 \\ &> F_{n+1}. \end{aligned}$$

Therefore, by mathematical induction, we have proved that $T_n > F_n \; \forall n \in N$. □

Furthermore, according to, $F_n = 2^{0.694n} = 1.6^n$, which gives out an exponential time complexity.

### 1.2.2  Dynamic Programming and Big Number Addition

We notice that every time we calculate $F_n$, we have to recalculate $F_{n-1}$ and $F_{n-2}$, which is a waste of time. By creating an array to store the past values, we've come with an polynomial solution ,where approximately, every $T_n$ is obtained by the one step addition of $T_{n-1}$ and $T_{n-2}$.

However, when $n$ is large enough, the number of digits of $F_n$ will exceed the maximum number of digits that a computer can store(eg. 64 digits), which brings out the problem of big number addition. At present we assume that adding two n-digit numbers produces complexity of $O(n)$, with $O(\log n)$ methods displayed in Therefore, for $n > 200$, consider $T_{\frac{n}{2}} > 2^{0.694 \times 100}$ . Adding $T_{n-1}$ and $T_n$ takes $O(n)$ time. In total, T(n) takes $O(n^2)$ time.

---

**Algorithm 2:** fib2(n)

---
**1** create an array f[0... n]
**2** f[0] = 0, f[1] = 1
**3 forall** $i \geq 2$ **do**
**4**     $f[i] = f[i-1] + f[i-2]$
**5 end**

---

### 1.2.3   Matrix Multiplication and Fast Power Algorithm

We can rewrite (1) as:
$$\begin{bmatrix} F_n \\ F_{n+1} \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ 1 & 1 \end{bmatrix} \cdot \begin{bmatrix} F_{n-1} \\ F_n \end{bmatrix},$$
which implies
$$\begin{bmatrix} F_n \\ F_{n+1} \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ 1 & 1 \end{bmatrix}^n \cdot \begin{bmatrix} F_0 \\ F_1 \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ 1 & 1 \end{bmatrix}^n \cdot \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$
Therefore, it only remains to compute $T^n$ for
$$T = \begin{bmatrix} 0 & 1 \\ 1 & 1 \end{bmatrix}, T_n = \begin{bmatrix} t_{11} & t_{12} \\ t_{21} & t_{22} \end{bmatrix}$$
. Let $M(n), A(n)$ be the time complexity for computing the multiplication and addition of two n-bit integers respectively. Notice that $M(n) \geq n$, as we at least need to read the inputs. We have seen $A(n) = O(n)$. Therefore, $M(n) = \Omega(A(n))$.

We will design an algorithm to compute $T^n$ and analyze the time complexity in terms of $M(n)$.

Let $k = \lceil \log_2 n \rceil$. We will compute all the $k+1$ matrices $T^{2^0}, T^{2^1}, \ldots, T^{2^k}$. Let $a_k a_{k-1} \cdots a_1 a_0$ be the binary representation of k. It is easy to see that

$$T^n = \prod_{i:a_i=1}^{k} T^{2^i} \tag{2}$$

The algorithm for computing $T^n$ is thus split into two steps:

1. Compute $T^{2^i}$ for $i = 0, 1, \ldots, k$.

2. Compute the product by (2).

**Time Complexity for Step 1.** Notice that it requires only one matrix multiplication to get $T^{2^{(i+1)}}$ from $T^{2^i}$ if we've computed $T^{2^i}$, namely, $T^{2^{i+1}} = (T^{2^i}) \times T^{2^i}$. Now we analyze the time complexity for this multiplication by investigating the length of the for integers in $T^{2^i}$.

**Lemma 1.2.** The length of the integers in $T^{2^i}$ is at most $2^{i+1} - 1$.

**Proof.** For the base step, if $i = 0$, then $T^{2^0} = T^1 = T$. The elements in $T$ are less than $2^1 - 1 = 1$. For the induction step, suppose the lemma holds for $i - 1 (i > 2)$, then $T^{2^i} = (T^{2^{i-1}}) \times T^{2^{i-1}}$. Since th length of product of m-digit number and n-digit number is less than $m + n - 1$, the length of element t in T is less than $(2^i - 1) + (2^i - 1) - 1 = 2^{i+1} - 3 < 2^{i+1} - 1$. $\qquad\square$

Since multiplying two $2 \times 2$ matrices requires 8 integer multiplications and 4 integer additions, and $M(n) = \Omega(A(n))$, computing $T^{2^{i+1}}$ from $T^{2^i}$ has complexity bounded by $c \cdot M(2^{i+1} - 1) \leq c \cdot M(2^{i+1})$ for some constant $c > 0$. Therefore, the overall time complexity for Step 1 is

$$c \cdot (M(2^2) + M(2^3) + \ldots + M(2^{2k+1})) \tag{3}$$

**Lemma 1.3.** For $i = 2, 3, \ldots, k$, $M(2^{i+1}) \geq 2M(2^i)$

**Proof.** To see this intuitively (although the following argument is not rigorous), if otherwise, we will have

$$M(2^i) < 2M(2^{i-1}) < 4M(2^{i-2}) < \ldots < 2^i M(1) = 2^i,$$

which contradicts $M(n) \geq n$ as mentioned earlier. $\qquad\square$

Therefore, we have

$$M(2^{k+1}) \geq 2^1 M(2^k) \geq 2^2 M(2^{k-1}) \cdots \geq 2^i M(2^{k-i+1}) \cdots \geq 2^{k-1} M(2^2)$$

By dividing the inequation by $2^{k-i}$, and pluging k with $\lceil \log_2 n \rceil$, we have

$$M(2^i) \leq \frac{1}{2^{k-i}} M(2^{k+1}) \leq \frac{1}{2^{k-i}} M(2n)$$

the time complexity for (3) is

$$c \cdot (M(2^2) + M(2^3) + \ldots + M(2^{2k+1})) \leq c \cdot M(2n) \left( 1 + \frac{1}{2} + \frac{1}{4} + \ldots + \frac{1}{2^{k-1}} \right) = O(M(2n))$$

, where the last equality is due to the fact that $M(2n) \leq 4M(n)$ (even if we use naive grade school multiplication) and $1 + \frac{1}{2} + \frac{1}{4} + \ldots + \frac{1}{2k-1} < 2$.

**Time Complexity for Step 2.** We will analyze the time complexity for computing

$$\prod_{i=0}^{k} S_k := T^{2^i}$$

We need $k$ matrix multiplications for $S_k$. Suppose we have computed $S_i := T^{2^0} \times T^{2^1} \times \ldots \times T^{2^i}$. Obviously, each entry of $S_i$ is dominated by $T^{2^{i+1}}$. Thus, computing $S_{i+1} = S_i \times T^{2^{i+1}}$ requires at most $c \cdot M(2^{i+2})$ time. The overall time complexity for Step 2 is also at most

$$c \cdot (M(2^2) + M(2^3) + \ldots + M(2^{2k+1})),$$

which, by previous analysis, is $O(M(n))$.

**Overall Time Complexity:** $O(M(n))$

**Remark 1.1.** An Intuitive way to calculate the overall time complexity is to draw analogy between

$$1 + 2 + 2^2 + \cdots + n < 1 + \cdots + 2^{\lceil \log_2 n \rceil} = 2^{\lceil \log_2 n \rceil + 1} - 1 < 2n$$

and

$$M(2^2) + M(2^3) + \cdots + M(2^{k+1}) \leq M(2^2) + \cdots + M(2n)$$

# 2   Divide and Conquer

## 2.1   Multiplication

### 2.1.1   Problem Description

For naive two n-digit multiplication with base 2, suppose n is 2's power.

$$xy = (a \times 2^{n/2} + b)(c \times 2^{n/2} + d) = ac \times 2^n + (ad + bc) \times 2^{n/2} + bd$$

$$1234 \times 5678 = (12 \times 100 + 34) \times (56 \times 100 + 78)$$

We only consider multiplication, ignoring multiplying by 100 (which can be done simply by shifting a few digits) or addition, which are linear compared to Multiplication. Eventually, we need $\log_2 n$ steps to divide into the base step (1-digit times 1-digit).With each step forward, we multiply to number of elements on the level by 4. The multiplication is done on the end base, so in total we have $4^{\log_n 2} = O(n^2)$ complexity.

### 2.1.2   Karatsuba algorithm

Rather than calculating $ad, bc, ac, bd$ separately; we obtain $(ad + bc) = (a + b) \times (c + d) - ac - bd$. Therefore, instead of dividing it into 4 parts, we divide it into 3 groups, through 3 multiplication operations on $\frac{n}{2}$, we get $3^{\log_2 n} = n^{\log_2 3} = n^{1.6}$ Specially, if n is not 2's power, we add zero to max bits, which at most doubles the original length, $(2n)^1.6$. No adding Zero in real coding, since carrying

---

**Algorithm 3:** Karatsuba algorithm

---

1   **Input:**Two n-digit numbers $x, y$.
2   **Output:** their product
3   $x_L, x_R = \text{leftmost} \lceil n/2 \rceil, textrightmost \lfloor n/2 \rfloor bits of x$
4   $y_L, y_R = \text{leftmost} \lceil n/2 \rceil, textrightmost \lfloor n/2 \rfloor bits of y$
5   P1=multiply$(x_L, y_L)$
6   P2=multiply$(x_R, y_R)$
7   P3=multiply$(x_L + x_R, y_L + y_R)$
8   return $P1 \times 2^n + (P3 - P1 - P2) \times 2^{\frac{n}{2}} + P2$

---

occurs, (eg. 64-(32,(32+1)-¿64)). Time complexity is:
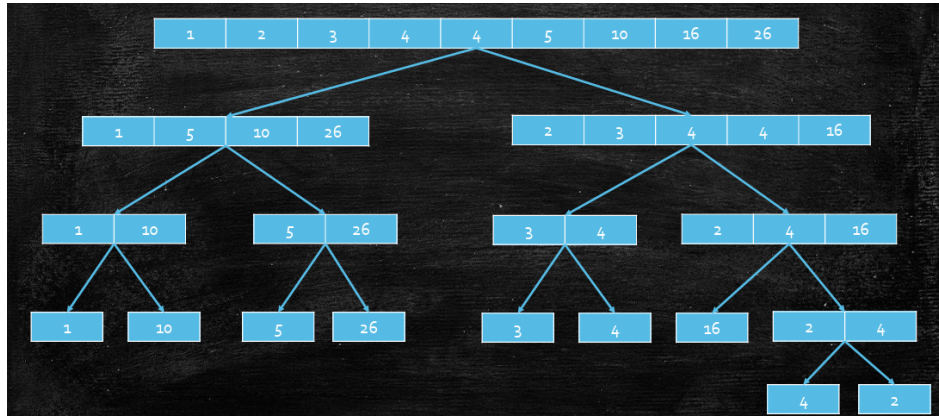
$$T(n) = 3T(\frac{n}{2}) + O(n) = 3T(\frac{n}{2}) + cn$$

Figure 1: Merge Sort Pipeline

Due to carry, $T(\frac{n}{2})$ should be $T(\frac{n}{2}+1) = T(\frac{n}{2}) + O(n)$.

$$T(n) = 3 \cdot (3T(\frac{n}{4}) + \frac{cn}{2}) + cn$$
$$= 3^{\log_2 3} T(1) + cn(1 + \frac{3}{2} + \frac{3^2}{2^2} + \cdots + \frac{3^{\log_2 n}}{2^{\log_2 n}})$$
$$= O(n^{\log_2 3}) + O(n^{1+\log_2 1.5})$$
$$= O(n^{\log_2 3}) = O(n^{1.6})$$

1. **Toom-Cook**:$O(n^{1.465})$
   Breaking into 3 parts $5 \times \frac{n}{3} \times \frac{n}{3}$

2. **Fast-Fourier Transform**

3. **Strassen's magical idea for Matrix Multiplication**
   Squeeze eight multiplications into seven with complex calculation, better than the naive case ($n$ for a row times a column, with $n^2$ numbers, in total $O(n^3)$. By normal divide and conquer, we divide one'n-size multi' into 8 $\frac{n}{2}$ multi, with total operations of $8^{\log_2 n} = O(n^3)$)

---

**Remark 2.1.** **Divide and Conquer** is a general algorithm design paradigm.

1. **Divide** Divide the problem into small size subproblems.

2. **Recursive** Solve small problems recursively

3. **combine** Combine the output of small size subproblems to get the answer for the original problem.

4. **Basic solver** If the problem size is small enough, solve it directly.

## 2.2   Sorting Problem

A very familiar problem, which can be easily picked up by analyzing pipeline<span style="color:red">1</span>.

### 2.2.1   Time complexity

For the 'merge' part, merging two arrays( length m and n) requires $O(m+n)$ time. For each level, we have $O(n)$ merging operations. There are $\log_2 n$ levels, so the total time complexity is $O(n \log_2 n)$

**Remark 2.2.** Never use asymptomatic notations in an induction-based analysis!
An example of a wrong analysis is as follows:

**Proof.** $T(i) = O(i)$ holds for $i = 1, \ldots, n$
**Base Step:** T1 =O(1) holds trivially.
**Inductive Step:** Suppose $T_i = O(i)$ holds for $i = 1, \ldots, n-1$
$T(n) = 2T(\frac{n}{2}) + O(n) = 2O(\frac{n}{2}) + O(n) = O(n)$                                   □

There are two mistakes:
Firstly, the Inductive Step is meaningless, since $O(\cdot)$ takes in a function (with variable i). Given a specific value, $O(i = n - 1) = O(1)$
Secondly, when the recursion happens numerous times, a constant can be a function of n. When we write $O(n)$, what we actually mean is $cn$.
Therefore, the correct proof should be:

**Proof.** $T(i) \leq ci$ holds for $i = 1, \ldots, n$
**Base Step:** $T(1) \leq c$ holds trivially.
**Inductive Step:** Suppose $T(i) \leq ci$ holds for $i = 1, \ldots, n-1$

$$T(n) = 2T(\frac{n}{2}) + O(n) \leq 2T(\frac{n}{2}) + d \cdot n \text{(for some constant d)}$$

$$\leq 2 \cdot c \cdot \frac{n}{2} + d \cdot n (\text{ by induction hypothesis})$$

$$= (c+d) \cdot n$$

Induction fails as constant changes from $c$ to $c + d$.                                   □

A correct proof should be:

**Proof.** If $T(n) \leq 2T(\frac{n}{2}) + c \cdot n$ for constant $c$, $T(n) \leq B \cdot n \log n$ for some constant $B > c$.
**Base Step:** $T(2) \leq B \cdot 2 \log 2$ holds trivially.
**Inductive Step:** Suppose $T(i) \leq Bi \log i$ holds for $i = 2, \ldots, n-1$
For $i = n$, we have

$$T(n) \leq 2T(\frac{n}{2}) + cn \leq 2 \cdot B \cdot \frac{n}{2} \log \frac{n}{2} + cn = Bn \log n - Bn + cn < Bn \log n$$

□

**Remark 2.3.** At least $\Omega(n \log n)$ comparisons are needed. // The time complexity is largely contributed by the running time of comparison. Total possible outputs is $3^{K(n)}$, n! numbers of permutations. Therefore, $3^{K(n)} \geq n!$, $K(n) \geq \log_3 n! \geq \Omega(n \log n)$

## 2.3  Inversions

> **Definition 2.1.** if $a_i > a_j$ when $i < j$, we call $(a_i, a_j)$ an inversion.

The idea is that we divide the input into two subsets ( A:$x_1, x_2, \ldots, x_m$; B:$x_{m+1}, x_{m+2}, \ldots, x_{m+n-1}$), and count the inversions of each subset and **across** subsets, which takes $O(mn)$ with each $a_i$ scanning the whole list B.

A magic trick is that we **mix merging and counting together**. We first count the inversions in A and B, then sort them separately. If A,B are sorted, we change a 'comparison' complexity problem into a 'search' complexity problem. In other words, we insert the additional step of counting inversion in the mergesort, by following operations, with time complexity of O(nlogn).

---
**Algorithm 4:** countInversions

**Input:** A list of n integers: S
**Output:** Number of Inversions: count

**1** count=0
**2** Divide the input into two subsets: A:$x_1, x_2, \ldots, x_{n/2}$; B:$x_{n/2+1}, x_{n/2+2}, , x_n$
**3** **if** $n == 1$ **then**
**4**     **return** *0*
**5** **end**
**6** **if** $n == 2$ **then**
**7**     **if** $x_1 < x_2$ **then**
**8**         swap $x_1, x_2$
**9**         **return** *1*
**10**     **end**
**11**     **else**
**12**         **return** *0*
**13**     **end**
**14** **end**
**15** count+=countInversions(A)+countInversions(B)
**16** Maintain 2 pointers $i = 1, j = n/2 + 1$
**17** **while** $i < n/2$ *and* $j < n$ **do**
**18**     **if** $x_i < x_j$ **then**
**19**         $i+=1$
**20**     **end**
**21**     **else**
**22**         $j+=1$
**23**         $count+=j-n/2$
**24**     **end**
**25** **end**
**26** count+=(n/2-i)×(j-n/2)
**27** **return** *count*
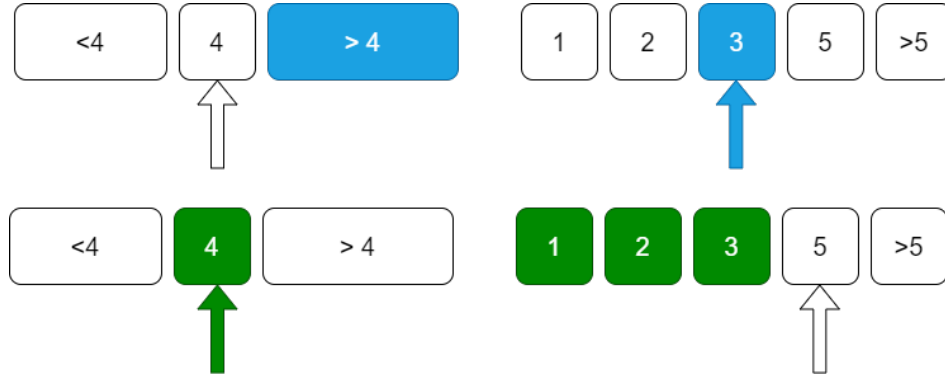
---

## 2.4  Master Theorem
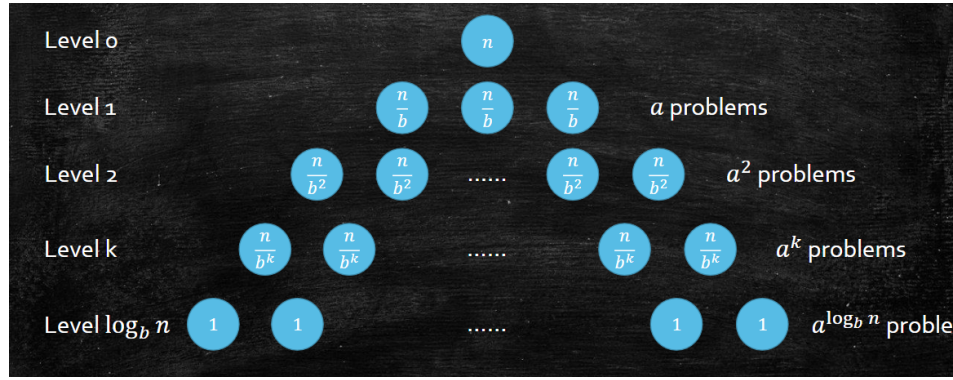
Figure 2: Inversion



Figure 3: Master Theorem

**Theorem 2.1.** If $T(n) = aT(\frac{n}{b}) + O(n^d)$.
a:Divide into a subproblems
b:subproblem size n/b
$O(n^d)$:combining time complexity

$$T(n) = \begin{cases} O(n^d) & \text{if } a < b^d \\ O(n^{\log_b a}) & \text{if } a > b^d \\ O(n^d \log n) & \text{if } a = b^d \end{cases} \tag{4}$$

**Proof.** The running time of solving size-1 problem is $O(a^{\log_b n}) = O(n^{\log_b a})$. The total running time for combining is $c(n^d) + c((\frac{n}{b})^d)\dot{a} + \cdots a^k \dot{c}(\frac{n}{b^k})^d + \cdots a^{\log_b n}\dot{c}$. The size-1 problem complexity is eaten up by the last element of the combining complexity. Eventually, we have $c\dot{n}^d(1 + \frac{a}{b^d}) + \cdots (\frac{a}{b^d})^k + \cdots (\frac{a}{b^d})^{\log_b n}$. When $a < b^d$, intuitively the head is heavier(longer when imaging the operation of every layer as a line). When $a > b^d$, intuitively the tail is heavier. When $a = b^d$, $T(n) = O(n^d)\dot{(}1 \times \log_b n) = O(n^d \log_b n)$.

□

## 2.5 Selection

Input: A set S of n integers $x_1, x_2, \ldots, x_n$ and an integer k.
Output: The k-th smallest integer $x^*$ among $x_1, x_2, \ldots, x_n$.
naive case: Sorting and choosing the k-th element: O(nlogn)
alg 1: quick sort? have the first number to be the pivot.
alg 2: Divide $x_1, x_2, \ldots, x_n$ into three subsets, recursively pick out the subset $x^*$ is in.

---
**Algorithm 5:** Select

---
**1** Choose an **arbitrary** value $v$ among $x_1, x_2, \ldots,$
**2** Divide $x_1, x_2, \ldots, x_n$ into three subsets $L, M, R$
**3 if** $k \leq |L|$ **then**
**4**   | return Select(L,k)
**5 end**
**6 else if** $|L| < k \leq |L| + |M|$ **then**
**7**   | return $v$
**8 end**
**9 else if** $|L| + |M| < k$ **then**
**10**   | return Select(R,$k - |L| - |M|$)
**11 end**

---

**Proof. the scale of the problem shrinks:**
    First, we show that the algorithm always terminates. If it does not terminate in the current recursive iteration, it will either call Select(L,k) or Select(R, $k - L - |M|$). It suffices to show that $|L| < $ n and $|R| <$ n (so that the problem size is strictly decreasing). This is guaranteed by $v \notin L$ and $v \notin R$ (v is in the set, so at least v is excluded).
**the algrithm is correct inductively:**
    Suppose the correct value is output for any array S of length n and any $k \in 1, 2, \ldots, n$. **Base Step:** n = k = 1, which is obviously true.
**Inductive Step:** Suppose the algorithm correctly outputs the k-th smallest value for any array S with any length $l \in 1, 2, \ldots, n-1$ and any $k \in 1, 2, \ldots, n$.
Select(L,k) returns the k-th smallest value in L and Select(R, $k - |L| - |M|$) returns the $(k - |L| - |M|)$-th smallest value in R by induction hypothesis. For any input array of length n and any integer v in it, it is straightforward to check that the k-th smallest value in L is the k-th smallest value in S if $L \geq k$, and the remaining two cases are also correct.
    □

### 2.5.1 Time Complexity

For division,$\Theta(n)$ with each element compared to v. For recursion, $T(|L|)$ or $T(|R|)$ or $O(1)$.
For the worst case, consider $1, 2, \ldots, n$,where k=n. The algorithm is $O(n^2)$.

Figure 4: The white an are doesn't matter terms. We only consider the last step when $a_m, m \in [i,j]$ will be chosen as a pivot. There are (j-i+1) choices, and only 2 choices (namely, $a_i$ or $a_j$) can result in comparison between $a_i$ and $a_j$.

**Remark 2.4.** Randomness in input and randomness in algorithm
Commonly, we choose to analyze the mean time complexity of a random algorithm. Because in real case, randomness of input is a strong condition(block-like features commonly).

By estimation, for every layer, n comparisons. If we are lucky enough to have $\frac{1}{3}$ in the minimum separation, then $n\frac{2}{3}^n > 1$, $k < \log_1 .5n$.

**Definition 2.2.** $E[x] = \sum_{\omega \in \Omega}^n \Pr(\omega) \cdot X(\omega)$

**Theorem 2.2. Linearity of expectation**

$$E|\sum_{i=1}^n X_i| = \sum_{i=1}^n E[X_i]$$

holds when $X_i, \ldots X_n$ are dependent.

eg. For dices, $E[x] = \frac{1}{6} \cdot (1 + 2 + \cdots + 6)$. When rowing two dices, by the Linearity of expectation, $E[x] = 3.5 + 3.5 = 7$. If two dices output the same number magically, we still have $E(x) = (2 + 4 + 6 + ... + 12)/6 = 7$.

For our case, since comparison takes up most of the running time, we calculalte the average times of comparison. Firstly, if $a_i, a_j$ are compared, one must be pivot and the other goes to the next layer of the division tree, which will only be compared once; When we consider the ordered sequence, $a_i, a_j$ are compared iff no element betweeen them are chosen as pivots, else they will be in different subsets, as is shown in 4.

$$E[X_{ij}] = \Pr(X_{ij} = 1) \cdot 1 + \Pr(X_{ij} = 0) \cdot 0 = Pr(X_{ij} = 1) = \frac{2}{j - i + 1}$$

$$i = 1 : \frac{2}{2} + \frac{2}{3} + \cdots + \frac{2}{n} = \sum_{i=1}^n \frac{2}{i} = 2 \cdot \sum_{i=1}^n \frac{1}{i} = 2 \cdot \log n$$

$$i = 2 : \frac{2}{2} + \frac{2}{3} + \cdots + \frac{2}{n-1} \cdots$$

To sum up, worst-case $O(n^2)$, expected average running time $\Theta(n \log n)$

### 2.5.2   Median of the medians(1973)

How to pick a **good pivot**? In other words, how to make sure $|L|$ and $|R|$ are approximately equal?
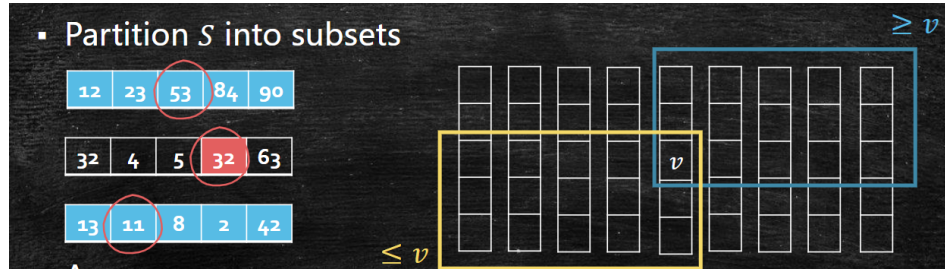
Figure 5: Median of the medians

Partition $S$ into subsets with size 5. O(n) Find the medians of each subset: $O(n/5)$ (since operations in a fixed-sized set is O(1)).
Find the median of the medians: $T(n/5)$.
Elements remaining: According to 5, there is $n - \frac{n}{5} \cdot \frac{1}{2} \cdot 3 = \frac{7n}{10}$ elements left.

$$T(n) = T(\frac{n}{5}) + T(\frac{7n}{10}) + O(n)$$

**Proof.** $T(n) = O(n)$
   Assume that $T(n) \leq cn$
$T(\frac{n}{5}) < \frac{cn}{5}$, $T(\frac{7n}{10}) < \frac{7cn}{10}$, $T(n) = T(\frac{n}{5}) + T(\frac{7n}{10}) + O(n) < 0.9cn + bn \leq cn$. which holds when we set c to be bigger then b.

$\square$

If the size of subset is i=2m+1, $m \in N$. To obtain the median of medians, we have $O(\frac{n}{2m+1})$. There are $n - \frac{n}{2m+1} \cdot \frac{1}{2} \cdot (m+1) = \frac{(3m+1)n}{2(2m+1)}$ elements left.

$$T(n) = T(\frac{n}{2m+1}) + T(\frac{(3m+1)n}{2(2m+1)}) + O(n) = T(\frac{(3m+3)n}{2(2m+1)})$$

However, in real case, the const is much bigger than 'quick sort' approach with a random pivot.

## 2.6   Clostest Pair

We can improve the naive algorithm by **sorting**. Imagine in 1-dim case, sorting $(O(n \log n))$ plus comparison between neighbouring points $(O(n))$ results in an $O(n \log n)$ algorithm.
   Consider 2-dim case, apparently using only x sequence or y sequence cannot lead us to the right answer. Some tricks are implemented:

1. **First sort by x, then by y**
   The only use of sorting by x is to divide the original dots into x-neighbouring subsets. After the division, any changes inside the subsets are allowed, so we merge the sort by y-axis step into the conquer step, which eliminates the $O(n \log n)$ sorting operation in each layer.

2. **Bound the search plane**
   If we've found the minimum distance inside A and B separately, we only need to consider the distance between points across A and B. Additionally, $|x_i - x_{mid}| < dist$ ,which forms the $2 - \delta strip$ shown in fig. Similarly, during the conquer step, for each point $p_i$, we only consider $p_j$ that satisfies $|y_i - y_j| < dist$. By proposition 2.1, we magically found that only the distance between **four** neighbouring points are considered for each point (as is shown in fig 7)

> **Proposition 2.1.** At most 4 points can appear in a $\delta \times \delta$ square, so that the distance between any two points are no smaller than $\delta$.

> **Proof.** Divide the square into four smaller squares, as in shown in fig 6. Since two points are at most $\frac{\delta}{\sqrt{2}} < \delta$ apart, at most one point can exist in each small square. So in total, at most four points appear in the square. $\qquad\square$
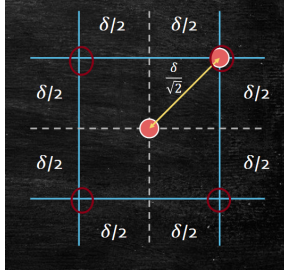


Figure 6: At most 4 points appear in the square



Figure 7: Only four neighbouring points need to be considered for each point a

This brings out an $O(n \log n)$ alogrithm. The sorting by x-coordinate takes $O(n \log n)$, and the 'divide and conquer' step amounts to adding constant steps( since each point only has four neighbouring points to count) to MergeSort.

---

**Algorithm 6:** FindClosestPair

**Input:** A list of pairs P: $p_1 = (a_1, b_1), \ldots p_i = (a_i, b_i), \ldots, p_n = (a_n, b_n)$
**Output:** distance of closest two pairs in P: dist

**1** **if** $|P| == 1$ **then**
**2** $\quad$ return
**3** **end**
**4** Sort P by x-coordinate.
**5** Divide P into two equally sized subsets A,B, the point in the middle $p_{mid} = (a_{mid}, b_{mid})$.
**6** dist=min(**FindClosestPair(A)**,**FindClosestPair(B)**)
**7** Merge A and B by y-cordinate to form C.
**8** **foreach** $c_i = (a_i, b_i)$ *in C* **do**
**9** $\quad$ **if** $|a_i - a_{mid}| < dist$ **then**
**10** $\quad\quad$ add $c_i$ to S
**11** $\quad$ **end**
**12** **end**
**13** **foreach** $s_i = (a_i, b_i)$ *in S* **do**
**14** $\quad$ j=i+1
**15** $\quad$ **while** $|b_i - b_j| < dist$ **do**
**16** $\quad\quad$ dist=min(dist, $\|b_i - b_{i+1}\|$)
**17** $\quad\quad$ j=j+1
**18** $\quad$ **end**
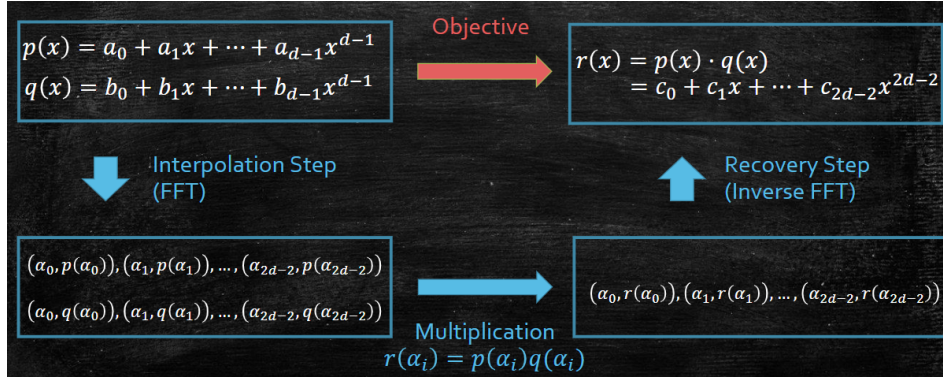**19** **end**
**20** return dist

---

Figure 8: FFT

## 2.7   Fast Fourier Transform and Polynomial Multiplication

We use array $(a_1, a_2, \ldots, a_{d-1})$ to describe $p(x) = \sum_{i=0}^{i=d-1} a_i x^i$

$$r(x) = \sum_{i=0}^{2d-2} c_i x^i, \text{where } c_i = \sum_{k=0}^{i} a_k b_{i-k}$$

By Karasuba Algorithm, we can have $O(n \log_{1.5} n)$ algorithm. The key of interpolation is to use the minimum number of points to describe the original subject. Any polynomial is a node in the n-dim space, with its coordinates $(a_1, a_2, \ldots, a_{d-1})$. We need $2d - 1$ interpolation points to fully store the answer polynomial by theorem .

**Theorem 2.3.**

$$
\begin{bmatrix} y_0 \\ y_1 \\ \vdots \\ y_{d-1} \end{bmatrix} =
\begin{bmatrix}
1 & x_0 & x_0^2 & \cdots & x_0^{d-1} \\
1 & x_1 & x_1^2 & \cdots & x_1^{d-1} \\
\vdots & \vdots & \vdots & \ddots & \vdots \\
1 & x_{d-1} & x_{d-1}^2 & \cdots & x_{d-1}^{d-1}
\end{bmatrix} \cdot
\begin{bmatrix} a_0 \\ a_1 \\ \vdots \\ a_{d-1} \end{bmatrix}
$$

Since the determinant of a Vandermonde matrix is $\Pi_{0 \le i < j \le d-1}(x_j - x_i)$. Denote the middle matrix as A, A is an invertible matrix.

interpolation step: For counting $p(a_i)$, we need d addition for (2d-1) interpolations, which brings out $O(n^2)$. However, we can try to compute $p(a_{i+1})$ from $a_i$ by carefully chosen interpolation points. If we choose $a = \alpha$ and $a = -\alpha$, we can calculate $a_0 + a_2 \alpha^2 + \ldots$ and $a_1 + a_3^3 + \ldots$ with d times, and reduce (2d) to obtain both $p(\alpha)$ and $p(-\alpha)$ with d times.

$T(D) = 2T(\frac{D}{2}) + O(n)$ is wrong, because the algorithm cannot be done recursively. In the second layer, $\alpha_0^2 = -\alpha_2^2$ cannot be satisfied. In the complex plane, however, for every element, there exists two roots.

multiplication of two complex numbers $r(a_i) - p(a_i)q(a_i)$

recovery: How to find $A^{-1}$ and calculate matrix multiplication?

$$) ($$

If columns are orthonormal, so are the rows. $AA^* = E, A^{-1} = A^*$. Intuitively, the inverse of a orthonormal matrix can be obtainly easily.

> **Proposition 2.2.** $\frac{1}{\sqrt{D}}A(w)$ is orthonormal for $w = e^{\frac{2\pi i}{D}}$

# 3   Graph

## Terminologies

1. vertex, vertice, edge, edges, index, index

2. Adjacency Matrix, Adjacency List

3. reachability, connected component

4. Tree: a connected and acyclic graph, containing backedge and Tree edge.

5. DAG: Directed Acyclic Graph

6. Sink(Source): A collection of vertice that doesn't have outgoing(ingoing) edges.

7. Partition: $C_1, C_2, \ldots, C_n$ is a partition of $C$ iff their union is C and they don't intersect.

8. cut: Given a partion A,B of P, value defined by $c(A, B) = |E(A, B)|$(number of edges across two sets).

## 3.1   Depth First Search

### 3.1.1   Existence of cycles

Use dfs to traverse all the vertice in a graph. First reach a vertex, then visit its connected component. Most effective when stored in an adjacency list. Dfs can form a tree(forest).

Specially, we analyze the case of directed graph. For directed graph, there exists 4 kinds of edges:

1. **Tree Edges** Edges in the DFS search: (u, v) where marked(v) = false

2. **Forward Edges** Edges that point to non-child descendants

3. **Back Edges** Edges that point to ancestors

4. **Cross Edges** All the other edges: edges point to vertices on other tree doesn't exist for undirected graph, since if $(c, f)$ occurs it will only be tree edge.
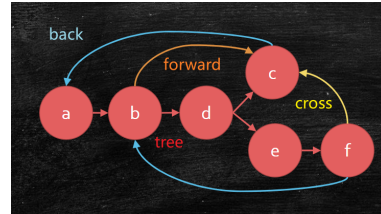


Figure 9: Directed Graph



Figure 10: colored edges

**Correctness of Algorithm:**
If there exists a back edge in directed graph, a cycle exists.

If a cycle exists, then a backedge must be found. Since a tree is acyclic.

By algorithm 7 plus the condition of the existence of a backedge ($pre[v] < pre[u]$ and $post[u] < post[v]$), depicted in 9

---

**Algorithm 7:** DFS on Directed Graph

---

**1** clk=0
**2** **foreach** $v \in V$ **do**
**3**  | **if** *marked[v]=False* **then**
**4**  |  | explore($v$)
**5**  |  | clk++ # used to calculate the number of connected components
**6**  | **end**
**7** **end**
**8** **Function** explore($u$):
**9**  | pre[u]=clk
**10** | clk++
**11** | visited[u]=false
**12** | **foreach** $(u, v) \in E$ **do**
**13** |  | **if** *marked[v]=false* **then**
**14** |  |  | explore ($v$)
**15** |  | **end**
**16** | **end**
**17** | post[u]=clk
**18** | clk++

---

### 3.1.2   Application: Topological Ordering

A DAG must have topological ordering, since:

1. Every DAG have a sink.
   Intuitively, if we start at $v$, since there exists an outgoing edge for every vertex. Since DAG doesn't have cycles, in finite steps, we'll visit every vertex,

2. Find a topological order.
   When deleting sink vertice and edges pointing to them, we obtain a new DAG. By repeting the operation of finding sink (visiting sink at last), and deleting sink to construct new graph, and finding sink in the new graph... We successfully find a topological ordering by following the inverse sequence of the sinks.

**Time Complexity:**
   Finding a sink takes $O(|V|)$, removing a sink to update graph with $|V|$ rounds.
   A trick is that the vertice with the earliest post ordering is a sink. If not, the outgoing point points to another vertex that finishes earlier. (consider the fact that no edge (u,v) exists in DAG that satisfied *post[v] > post[u]*. See 9, backedges satisfy, but doesn't exist in DAG) Therefore, the topological ordering is the inverse order of finishing points, which brings the time complexity of DFS down to $|E| + |V|$.
**Correctness:**
   If (y,x) exists, y is not a sink since x is in the graph, contradiction!

### 3.1.3   Strongly Connected Components for Directed Graph

All the SCC forms a partition of C. i.e. for every $v$, $\exists C_i$ so that $v \in C_i$. For completeness, by construction, every single $v$ can be formed as a SCC. For non-intersection, if $v \in C_i, C_j$, $C_i$ and $C_j$ together forms a bigger C.C., as every vertex from $C_i$ can reach points from $C_j$ through $v$.

- Note that we can't simply employ DFS to find C.C., as although you find reachable points from a certain vertice $v$, the vertice you explore may not be reachable from each other.

- Note that unlike topological ordering, the node with the earliest post time doesn't necessarily occur in sink SCC, as the vertice with outer edges are previsited where our algorithm ends in somewhere else.i.e.,for graph 11 consider DFS in sequence 5,6,7,8, apparently, the 8 has the smallest post value, but is not in sink.

- However, the vertice with biggest post time must be in source SCC. Assume $u$ has the largest finish time, which must be a root of one DFS tree. Suppose $u$ isn't in source SCC, i.e. $\exists v$ s.t. $(v, u)$ exists. $v$ isn't in the tree of $u$, as $(u, v)$ doesn't exist. ,$v$ cannot start earlier than $u$, $v$ isn't in another DFS tree.

- Consider the meta-graph, where connected components form a big node, it must be DAG, else contradict with biggest SCC, therefore, there must be a sink. Also, a way to find whether $v$ is reachable from $u$ is to check the connectivity of bigger SCCs, $U$ and $V$.

---

**Algorithm 8:** Strongly Connected Component

**Input:** DAG:G
**Output:** num of SCC
**1** Deduce G's inverse $G^R$ from G
**2** Run DFS on $G^R$ to find source SSC, which is sink SSC of $G$, record it in the descending finish time as F.
**3** DFS on $G$ in the order of F.

---

**Time Complexity:** Naive Case: Find a sink, delete it and repeat again...which takes $O(|V||E|)$. Better Case: Doing DFS twice, shown in algorithm 8, taking $O(|V| + |E|)$.

## 3.2   Breadth First Search

### 3.2.1   Unweighed Shortest Path in Undirected Graph

---

**Algorithm 9:** alg:bfs

**1** Queue Q;
**2** **Function** `bfs(,Q,u)`:
**3**     **for each** $v \in V$ $marked[v] \leftarrow false$
**4**     Q.push(s) $marked[s] \leftarrow true$ **while** *Q is not empty* **do**
**5**         u=Q.top();Q.pop();
**6**         **foreach** $(u, v) \in E$ **do**
**7**             **if** *marked[v]=false* **then**
**8**                 marked[v]=True
**9**                 Q.push(v)
**10**            **end**
**11**        **end**
**12**    **end**

---

Consider doing topological order in BFS. Firstly it cannot tell the difference between crossedge and backedge, and it has no forward edge(as the only chance of visiting unvisited nodes forms a tree edge , i.e. nodes in layer 2 can only reach unvisited nodes and index them with 3). Additionally, DFS cannot detect cycles(=detect backedge), or find SCCs(based on cycles).
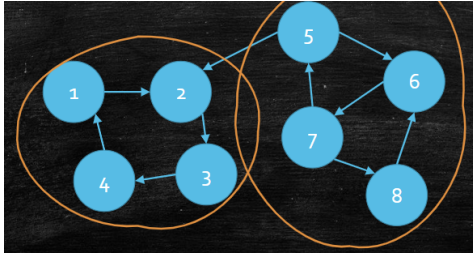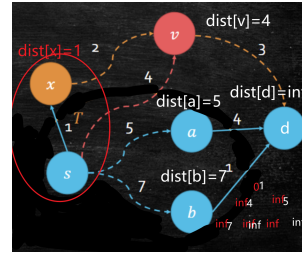
Figure 11: SCC



Figure 12: Correctness of Dijkstra's algorithm

### 3.2.2 Dijkstra: SP. with positive weight

> **Remark 3.1. difference between quantity and numerical value**
> Pay attention to the length of a numerical value. i.e., input length is $n$, the actual number complexity. Storing $n$ quantities takes n storage space, but storing a n-length numerical value takes $\log n$ bits. Therefore, an increase in bits takes exponencial time complexity. $O(|V| + |E|)$ , $|\text{'}|$ is a quantity $O(w_m ax)$ is a numerical value.

- Works by keep dragging the closest vertex to the current SPT and updating the minimum path length.

- Dijkstra algorithm is based on the assumption that evrey step out of SPT(shortest path tree) at least increases the path length, so it only solves graphs without negative cycles.

- A trick is to use heap (binary-heap, d-heap, Fibonachi-heap) to store the shortest distance from other vertice to $s$, so that we can quickly obtain the closest node from SPT, without the need to traverse all outer edges again. Creating the heap takes $O(|V| \log |V|)$, picking the top vertex takes $O(\log |V|)$, and updating all the nodes takes $O((|V| + |E|) \log |V|)$. i.e, for figure **??**, consider the first step. We update x,a,b,v and the closest vertex x natually pops to the top.

**Correctness of algorithm:** If there exists $x \notin SPT, v \in SPT$, s.t. $dist(s \to x \to v) < dist(v)$, $dist(x) < dist(v)$, $x$ should be in $SPT$,contradiction!

### 3.2.3 Bellman-Ford: shortest path with negative weight

Dijkstra's algorithm works because the dist value it maintains is either overestimates or exactly correct. The key idea of Bellman-Ford is that it updates all the edges (kind of 'edgewise' operation) each time (different from Dijkstra's algorithm who updates nodes one step connected to SPT), which takes $O(|V||E|)$.
   **Correctness of algorithm:**

1. After k rounds, dist(v) is the shortest path of all k-edge-paths.($dist[u_k] \le d(u_1, u_2, \ldots, u_k)$) For base step, dist[s] is the shortest dist of all 0-edge-paths.
   For induction step, suppose it is true for $k-1$ rounds, $dist[u_{k-1}] \le d(u_1, \ldots, u_{k-1})$. By Bellman-ford update condition, $dist[u_k] \le d(u_1, \ldots, u_{k-1}) + w(u_{k-1}, v)$, the righthand side means all the other k-edge-paths.

2. For graphs without negative cycle, all the shortest paths will have at most $|V| - 1$ edges. If there are more than $|V|$ paths, a vertex must be visited twice, but by erasing the whole cycle results in a shorter path, contradiction!

3. Negative cycle exists iff some dist is updated in the last($|V|^{th}$) round.

19

---

**Algorithm 10:** Di(j)kstra

---

**Input:** G=(V,E),s

**Output:** dist(u) is the shortest distance from s to u

**1** # initialize **foreach** $u \in V$ **do**

**2**      dist(u)=inf

**3**      pre(u)=Nan

**4** **end**

**5** dist(s)=0

**6** H=makequeue(V)

**7** **while** *H is not empty* **do**

**8**      u=H.top();H.pop()

**9**      **foreach** $(u,v) \in E$ **do**

**10**          **if** $dist(v) > dist(u) + l(u,v)$ **then**

**11**              $dist(v) = dist(u) + l(u,v)$ # update the shortest dist

**12**              pre(v)=u# record previous node

**13**          **end**

**14**      **end**

**15** **end**

---

---

**Input:** G=(V,E),s

**1** dist[s]=0, dist[x]=inf for vertex other than s

**2** pre(v)=Nan

**3** **while** *some dist[x] updates* **do**

**4**      # update at most —V— rounds. **foreach** $(u,v) \in E$ **do**

**5**          **if** *dist[v]¿dist[u]+w(u,v)* **then**

**6**              dist[v]=dist[u]+w(u,v)

**7**              pre(v)=u

**8**          **end**

**9**      **end**

**10** **end**

---

4. If no vertice updates in $V^{th}$ round, then no vertice will be updated then. But if a single vertice is not updated.

5. If a vertex is upgraded after $V$ rounds, then its path must include a negative cycle. A negative cycle can be found by tracing the recorded shortest path until two repeated nodes occurs. Tricks include running till the $2|V|^{th}$ round and see the vertice upgraded in $|V|^{th}$ and $2|V|^t h$.

# 4   Greedy Algorithm

## Homework++ and Partition

**Definition:** Given some homework with start time, end time and duration of time. Give out whether we have a feasible solution to finish all the homework. If greedy algorithm cannot promise a feasible solution, then no solution occurs. different end time? do ddl sooner first. different start time? do start time sooner first/ or arrange different release time/ deadline/size?Judge whether a solution exists(no need to find it)

Indeed, by Complexity Theory (), the problem is NP-hard. Consider a deliberately designed example: In **??**, we only have two types of work. For type one, However, either by choosing bigger or smaller size doesn't guarantee a feasible solution. We claim that Partition problem is NP: Partition the set into two parts, whose sum are equal. Given an input $a_1, \ldots, a_n$,let $w = \sum_{i=1}^{n} a_i$. n normal homework starts at 0 and ends at $w + 1$. Special homework has release time $w/2$ , deadline $w/2 + 1$ and size 1. Therefore, a special case for Homework++ is a Partition Problem.

NP: difficult to find answer, but easy to verify.

## 4.1   Prim and Kruskal: Minimum spanning tree

**Definition:** A Minimum spanning tree (MSP) is the spanning tree( a tree that connects to all the edges) with the minumum total weight.

**Property:** Shift edge operation on a cycle doesn't affect connectivity.

### 4.1.1   Prim's Algorithm

**Prim's Algorithm:**

Adding the edge out of the previous Tree $T_i$ with the minimum weight to form $T_{i+1}$.

**Correctness of algorithm:**

The key idea for Prim is to guarantee that the tree with $n$ edges $T_n$ is a subset of **global** MST $T$. In inductive step, if $T_i$ is part of $T$, then we can construct $T_{i+1}$, s.t. $T_{i+1} \in T$. Consider the first edge $e$ out of $T_i$, if $e = (u, v) \notin T$, then adding $e$ to $T$ constructs a cycle $C$. Since $e$ rather than $e' \in C, e' \notin T_i, e' \in T, w(e') \geq w(e)$. Since $T$ is the smallest spanning tree, $T - w(e') + w(e) \geq T$, $w(e) \geq w(e')$, therefore, $w(e) = w(e')$. Therefore, by shifting the two edges guarantees $T_{i+1}$ is a MST component of $T$.

**Time complexity:** $O(|E| \log(|E|))$.

### 4.1.2   Kruskal's Algorithm

**Kruskal's algorithm:**

Similar to Prim's algorithm, it adds the smallest weighted edge that doesn't form a cycle with $T_{i-1}$ to form $T_i$.

**Correctness of Kruskal:**

Similar to correctness of Prim, but the shift edge operation happens between components in a forest.

**Time complxity:**

$O(|E|\log(|E|))$ for sorting edges, and $|E|$ iterations of finding the representative element. A faster data structure is by **unified set**.

### 4.1.3   unified set

store groups by their representative element (the children are represented by the root vertice of their tree).

- Find: return the representative element in a group.$(2|E|)$

- Union: Merge two groups.$(|V| - 1)$

---

**Algorithm 11:** Kruskal's Algorithm

**1** Sort the edge set E to ascending order.
**2** X={}
**3** For each $u \in V$, makeset$(u)$.# $|V|$ create group
**4** For each $(u, v) \in E$ in ascending order
**5** **if** $find(u) \neq find(v)$ **then**
**6** $\quad$ Add edge $(u, v)$ to X
**7** **end**
**8** **Function** union$(u,v)$:
**9** **Function** find$(x)$:
**10** $\quad$ **while** $x \neq pre(x)$ **do**
**11** $\quad$ **end**
**12** **Function** union$(v)$:

---

1. **Path Compression:** Set the parent of the nodes checked to representative element.

2. **time compexty analysis by Amortized Cost:**

   **check cycle:**

## 4.2   Set Coverage

**definition:**
$\quad$ **set Cover:** Find a sub-collection $S \subseteq T$ with minimum $|S|$ such that $\bigcup_{A_i \in S} A_i = U$.
$\quad$ **max k coverage:** Given $k \in N$, find a sub-collection $S \supseteq T$ with $|S| \leq k$ that maximizes element number it covers. Denote $f(S) = |\bigcup_{A_i \in S} A_i|$.
$\quad$ Also viewed in bipartite graph, with the left set be set index, and the right set be the node index. Edges across two sets means that $node \in set$.

> **Remark 4.1. Approximation Algorithm:**
> $\quad$ For maximize problem, algorithm $\mathcal{A}$ is called $\alpha$-approximation algorithm if $\frac{\mathcal{A}(I)}{OPTIM(I)} \geq \alpha$. Normally, time complexity should be polynomial.

$\quad$ We deliberately construct the worst case: $|W_1| \geq \frac{1}{5}$biggest of remaining sets. $W_1 \cup W_2 = 1/5 + 1/5(1 - 1/5) = 1 - (1 - 1/5)^2$. $f(S) \geq (1 - (1 - 1/k)^k)$
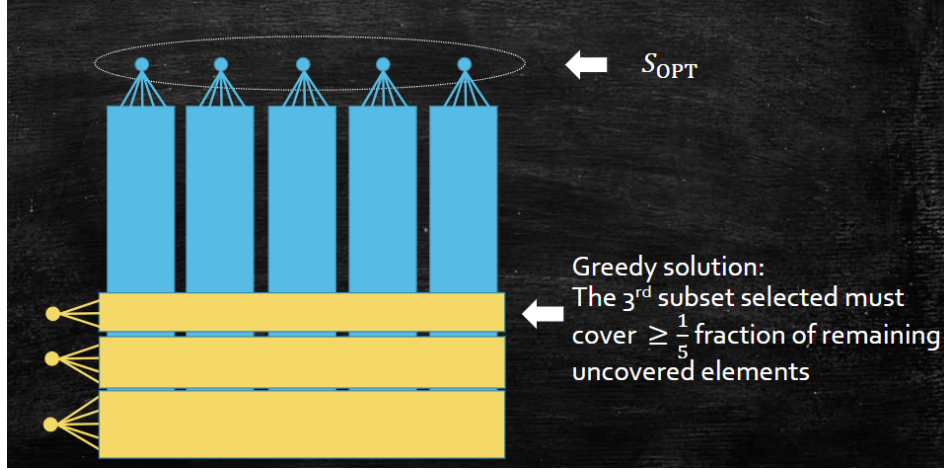
Figure 13: Directed Graph

**Lemma 4.1.** After choosing $l$ sets, $f(S) \geq (1 - (1 - 1/k)^l)f(S^*)$, where $f(S^*)$ is the optim. sol.

**Proof.** Denote $S = \{A_1, \dots A_m\}$, $S_t = \{A_1, \dots A_t\}$, optim solution $S^* = \{O_1, \dots, O_n\}$
For base step, $l = 1$, by greedy nature, $f(S1 = A1) \geq f(O_i)$ for all $O_i$. Thus $f(S_1) \geq \frac{1}{k}\sum_{O_i \in S^*} f(O_i) = (1 - (1 - \frac{1}{k})^1)f(S^*)$ This implies when calculating the sum, union elements are calculated more than once. Also, the size of the first chosen set is larger than the average size of the optimal solution. Denote $\Delta(O_t|S_t) = f(S_t \cup O_i) - f(S_t)$.
$\Delta(A_{t+1}|S_t) \geq 1/k\Delta(S^*|S_t) \geq \frac{1}{k}\Delta(S^*|S_t)$. For inductive step, suppose $f(S_t) \geq (1 - (1 - 1/k)^t)f(S^*)$ $f(S_{t+1}) - f(S_t) \geq \frac{1}{k}(f(S^* \cup S_t) - f(S_t)) \geq \frac{1}{k}(f(S^*) - f(S_t))$. $f(S_{t+1}) \geq \frac{1}{k}f(S^*) + (1 - \frac{1}{k})f(S_t) \geq \frac{1}{k}f(S^*) + (1 - \frac{1}{k})(1 - (1 - 1/k)^t)f(S^*) = ((1 - (1 - 1/k)^{t+1})f(S^*))$.          □

For max-k-coverage, we at most choose $k$ sets. $\lim_{k \to \infty}(1 - (1 - 1/k)^l) = 1 - 1/e$. For set cover problem, it is a $lnn$ approximation. $f(S) \geq 1 - (1 - 1/k)^{k \cdot lnn} > (1 - 1/e^{lnn}) = n - 1$ Therefore, $f(S) = n$.

## 4.3   Huffman Coding

1. prefix-free coding: each encoding is not a prefix of other encodings, i.e.,no vertice is an ancestor of others.

2. optim prefix code(with minimum avg_len):$avg\_len = \sum w(\sigma) \cdot len(\sigma)$, where $w$ is the frequency and $len$ is the length of the encoding.

3. Huffman Tree is constructed by replacing two vertice with minimum $w$ with a new parent vertice whose weight is their weight sum repeatedly.

4. The correctness is guaranteed. An intuitive way is the more frequent an encoding appears, the shorter the encoding is. Denote $T$ as the Huffman Tree, and $T'$ the tree that switched $u, v$, where $len(u) < len(v)$, $w(u) > w(v)$. avg_len(T)-avg_len(T')=$w(u) \cdot len(u) + w(v) \cdot len(v) - w(v)len(u) - w(u)len(v) = (w(u) - w(v))(len(u) - len(v)) < 0$.

   Next, we'll guarantee that every next step is optim, i.e. $T_{i-1}$ is part of $T$. For base step,... For inductive step: If we merge $A$ and $B$ in $T$. If $A$ is paired with $C$ in $T^*$, then $w(C) > w(A)$ ,$w(B) \leq w(C)$. If $w(B) = w(C)$, changing B and C preserves avg_len. If $w(B) \leq w(C)$, if

$len(B) < len(C)$, contradiction! If $len(B) > len(C)$, swapping B and C also preserves avg_len. If $len(B) > len(C)$, D is not processed, $w(D) \geq max(w(A), w(B))$

Another intepretation is by sorting inequality, namely, the sum of $\sum ordered\_sequence \leq \sum random\_sequence \leq \sum inversed\_sequence$.

5. **time complexity analysis** By using binary-heap, every time we conduct pop twice to obtain the top two min vertice and a push to add the new parent vertice in the tree.

### 4.3.1   K-Centers

Find a set of k centers, $S = \{s_1, \ldots, s_k\} \subseteq V$, $|S| = k$ that minimizes $f(S) = \max_{v \in V} \min_{s \in S}(s, v)$ which means the maximim distance of any vertex $v \in V$ to its closest center $s \in S$. Note: k-center(2), k-means(9), k-median(5) are all np-hard, but have relatively good approximation algorithms.

> **Remark 4.2.** (V,d) is a matrix space if:
>
> 1. d(v1,v2)=0 iff v1=v2
>
> 2. d(v1,v2)=d(v2,v1)
>
> 3. $d(v1, v3) + d(v3, v2) \geq d(v1, v2)$
>
> e.g., G=(V,E,$w > 0$)

**Algorithm:** Iteratively pick the center farthest to the existing centers. Pick the first by random. by Dijkstra for weighted graph and BFS for unweighted.
Note: Intuitively, we want to choose the vertex with largest degree. However, this guarantees that more vertice are closer to center, but isn't align with our objective of minimizing the dist(furthest vertice).
**analysis of approximation algorithm** The algorithm is 2-approximation. Denote $OPT = \max \min d(o_1, v)$.

---

**Input:** metric space (V,d)
**Output:** A set of k centers, $S$ that minimizes $f(S) = \max_{v \in V} \min_{s \in S}(s, v)$

---

$A = \{a_1, \ldots a_k\}$ for algorithm result, $X_i = \{v_1, \ldots v_i\}$ as the set of vertice closest to $o_i$, optim solution $O = \{o_1, \ldots, o_k\}$

Case 1:Suppose $A \cap X_i \neq \varnothing$ Then each $X_i$ contains a center. As the biggest length between two vertice in $X_i$ is smaller than 2optim.

Case 2:Suppose $\exists A \cap X_i = \varnothing$. Suppose we choose $a_j$ and $a_r \in X_i$. Let $a_r$ be the second center chosen in $X_i$. By greedy choice, $ALG = d(A, a_r) \leq d(a_j, a_r)$ before $a_r$ was chosen. Additionally, $d(a_j, a_r) \leq d(a_j, o_i) + d(a_r, o_i) \leq 2OPT$.

Also, it is the best algorithm. **dominated set** a dominating set is a subset of vertices $S$ such that, for any $v \in V S$, there is a vertex $u \in S$ that is adjacent to $v$. A dominated set problem(output whether a dominated set exists) is an NP problem. If $\exists$ a polynomial time $(2 - \epsilon)$-approximation algorithm $\mathcal{A}$ for K-center. Denote the algorithm with slight modification running on $(G, k)$ to output dominated set result as If $\mathcal{A}'$. If $\mathcal{A}'(G, k) = Yes$, $OPT \leq 1$. If $\mathcal{A}'(G, k) = No$, $ANS \geq 2, OPT \geq 2/(2 - \epsilon)$. Therefore, $G$ has a dominated set of size k iff K-center solution of $(G, k)$ output by $A$ has cost 1. Also called gap(amplifying)-reduction.

   **local search on k-center**
   Two ways to solve NPH problem, either by hyristic algorithm (with no approximation lower bound or guarantee of optimal solution) or approximation, restrictions on the input,

#### 4.3.2 local search: max-cut

1. **Time Complexity** checking whether to update takes $O(|E|)$, every update you earns at least an edge, so it takes $O(|E|)$ rounds. The total algorithm is $O(|E|^2)$, polynomial!

2. **Approximation** When the algorithm terminates, $\forall u$, at least $0.5\deg(u)$ edges are in the cut, which leads to a 0.5-approximation.

$$c(A, B) \geq \frac{1}{2} \sum_{u \in V} \frac{1}{2} deg(u) = \frac{1}{2}|E|$$

3. better approximation algorithm: 0.878-approximation(semi-definite programming)

palindrome: a string that reads the same backward as forward, e.g., madam, racecar, abba, etc.

## 5 Dynamic Programming

From smaller subproblems to larger subproblems.

- **DAGs:** guarantees topological order, i.e. each sub-problem only depends on previous sub-problem

- **Induction:** solve sub-problems in topological order, implemented by a for loop.

- **Memoization:** (not memorization)"backward" recursion with "storing"

Different from Divide and Conquer, where the scale decreases exponentially by a **'magic' merge operation**, DP reduces the scale a little bit.

Difference between Induction and Memoization: Theoretically, time-complexity is the same. Memoization has the advantage of not needing to compute **topological order**, and some **not-computed sub-problems**. But in coding, 'throwing recursive call' is time costing.

### 5.1 Longest Increasing Subsequence

Formation of a DAG: establish a node $i(i \in [1, n])$ for each element $a_i$, if $i < j$ and $a_i < a_j$, then $(a_1, a_j) \in E$. Then the problem is to find the longest path in the DAG. Construct L(i) by:

$$L(j) = 1 + \max L(i) : (i, j) \in E$$

**Time Complexity** is $O(|V||E|) \leq O(n^2)$.
**Correctness of Algorithm**

We'll prove that L(i) is the length of LIS ending at $a_i$.

For base step, L(0)=0 is the length of LIS for an empty sequence.

For induction step, assume that L(i) is the length of LIS ending at $a_i$, Let S be the LIS ending at $a_i$ and $T = S\backslash\{a_i\}$, the last number of T is $a_t$. Prove by $LIS[i] \geq |S|$, and $LIS[i] \leq |S|$. Since $T$ is a LIS ending at $a_t$, $|T| = L(t)$, so $|S| = L(t) + 1$. $L(i) = 1 + \max_{a_j < a_i, j < i}\{L(j)\} \geq L(t) + 1 = |S|$.

On the other hand, $L(i) = 1 + \max_{a_j < a_i, j < i} L(j) + 1 \leq 1 + |T| = |S|$, so $L[i] \leq |S|$.

Therefore, $L(i) = |S|$.

### 5.2 Edit Distance

Edit distance is the minimum number of editsinsertions, deletions, and substitutions of charactersneeded to transform the first string into the second. We denote $E(i, j)$ as the editing distance between $x[1..i]$ and $y[1..j]$.

$$E(i, j) = \min(1 + E(i - 1, j) + E(i, j - 1) + E(i - 1, j - 1))$$

---

**Algorithm 12:** Edit Distance

---

    **Input:** Two strings $x[1..m]$ and $y[1..n]$
    **Output:** Edit distance $E(m, n)$
**1 for** $i = 1$ *to* $m$ **do**
**2**    $E(i, 0) = i$;
**3 end**
**4 for** $j = 1$ *to* $n$ **do**
**5**    $E(0, j) = j$;
**6 end**
**7 for** $i = 1$ *to* $m$ **do**
**8**    **for** $j = 1$ *to* $n$ **do**
**9**       $E(i, j) = \min(1 + E(i-1, j), 1 + E(i, j-1), diff(i, j) + E(i-1, j-1))$;
**10**    **end**
**11 end**
**12 return** $E(m, n)$;

---

diff(i,j)=0 iff x[i]=y[j]. Forming a $m \times n$ table, and each cell is computed only once, which results in **Time Complexity** of $O(mn)$.

    Likewize, in the palindrome problem, $H(i, j)$ to denote the biggest palindrome length.

$$H(i, j) = \begin{cases} H(i+1, j-1) + 2, & \text{if } A[i] = A[j], \\ \max\{H(i, j-1), H(i+1, j)\}, & \text{o.w.} \end{cases} \tag{5}$$

## 5.3   Nnapsack Problem

Denote $s[i] = v[i]/c[i]$, $s$ means the ratio of value to cost.

1. indivisibility
   If it is divisible, if you choose the object with the largest xingjiabi each time, it is optimal, else you can switch two objects and get a better one, which results in a better solution

   If it is indivisible, it is essentially an NP-Hard problem. As a problem where r of all objects are equal and v=totalSum/2 is a Partition Problem.

    For **time-complexity**, it is not a polynomial problem, with $O(nW)$, where W is a numerical value rather than size.

    We can control the precision of $W$ by rounding, but the new optim solution you obtain may be infeasible in the original problem.

    Howvever, if you round value, the feasibility of solutions doesn't change. Let $V = \max v_i$, $OPT \leq nV$. Denote $A[i, v]$ as the minimum cost cost(S) if we select till the $i^{th}$ item with exactly value $v$.

$$A(i+1, v) = \begin{cases} \min A[i, v], c_{i+1} + A[i, v - v_{i+1}] \, if \, v[i+1] < v \\ A(i, v) & \text{o.w.} \end{cases} \tag{6}$$

    For **space-complexity**, it can be optimized from $O(nW)$ to $O(W)$. A[1,v]=0 if v=0 c1 if v=v1 inf o.w.

$$f[i][v] = \max f[i-1][v], f[i-1][v - c[i]] + w[i]$$

Where $f[i][v]$ means putting the first $i$ objects into a $v$-capacity Knapsack. For every object $i$, you either put it in the bag, where the left capacity is $v - c[i]$, or not put it in, which is equal to $f[i-1][v]$.

**Remark 5.1.** Fully Polynomial Time Approximation Scheme(FPTAS) gives a $(1-\epsilon)-$approximation with running time polynomial in term of $n$ and $\frac{1}{\epsilon}$

Algorithms that can scale down the problem to constant length and use naive explore algorithm are commonly FPTAS.

Polynomial Time Approximation Scheme(PTAS) gives a $(1-\epsilon)-$approximation with running time polynomial in term of $n$ and constant $\epsilon$.

Some have fixed constant ratio