

Week 1 NP-complete Problems and Reduction

Definitions : (1) P : set of all (decision) problems solvable in polynomial time
: means can be solved in $O(n^k)$ time where $k > 0$ is some constant independent of n.

(2) Decision Problems : \Rightarrow Yes or No Question

- Given a list of numbers check if in order / if a given number is largest
- Given a weighted graph . what is the MST ?

Ex1 : How to convert sorting into a yes / no problem?

\Rightarrow Given a list of n numbers : is it sorted?

$$O(n) : 1 \leq i \leq n \quad \text{is } A[i] \leq A[i+1] ?$$

$$\text{sort} : O(n \lg n)$$

Ex2 : How to convert shortest path into a yes / no problem?

Input: weighted graph , start vertex : S finish vertex : T

Goal : length of shortest path $S \rightarrow T$.

convert: Add a parameter : length: L

is there a path $S \rightarrow T$ of length L .

Ex3 : How to convert graph coloring into a yes / no problem?

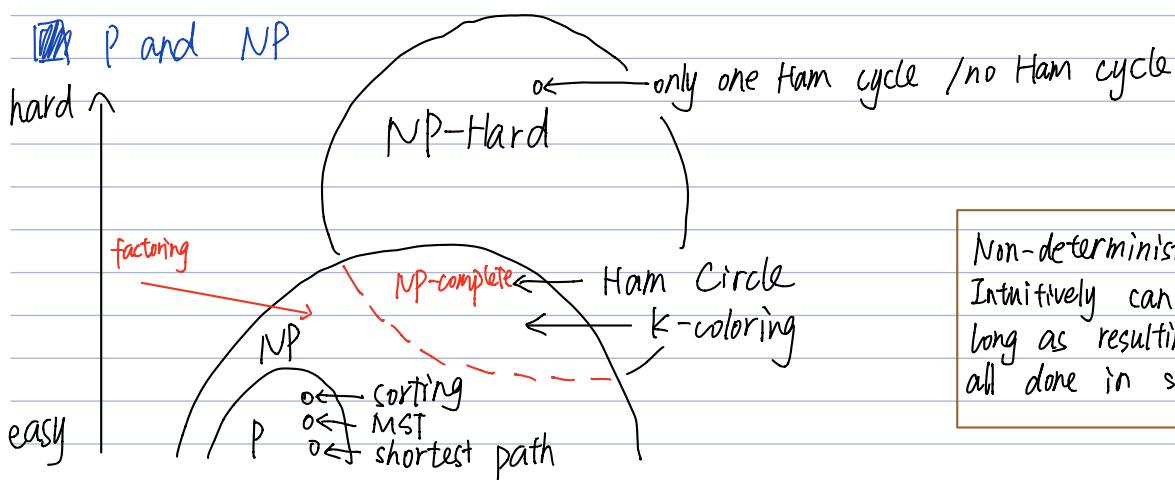
Given a graph, how many colors to color it with a minimum number of colors?

"color" : adjacent vertices must have different colors

More specific : Given a graph , can graph be colored with 3 colors?
Can this problem solve in polynomial time? $4 \leftarrow 5 \rightarrow 1$

Ex4 : How to convert Hamilton cycle into a yes / no problem? $\Rightarrow O(V!)$ (all permutation)

Euler cycle (NP complete) \Rightarrow even degree $\Rightarrow O(n)$



Non-deterministic Algorithm
Intuitively can do things in parallel as long as resulting connected by OR's all done in same time.

Example : $(x \vee y \vee z) \wedge (x \vee \bar{y}) \wedge (u \vee \bar{v} \vee z)$ $n = \# \text{ variables}$
 $M = \# \text{ clauses}$

Problem: Is there a way to assign T/F to n variables to make Formula True? (2^n : Truth table)

2^n processors in parallel : — V — V — V — ... one of them be true \Rightarrow true

② NP: set of all (decision) problems whose solutions can be verified in polynomial time. i.e. a problem A is in NP if there is a polynomial verification ALG V that takes input an instance X of problem A and a proposed solution S and outputs true IFF S is a solution to instance A $V(X, S) = 1$

[e.g. find Ham cycle $\Rightarrow V$: check Ham cycle $O(n)$]

④ NP-complete

① NP-Hard:

X is NP-hard if every decision problem in NP reduces to X.

NP-complete 1st problem: SAT (Cook 1971).

A tedious, difficult, detailed proof

Any problem in NP could be translated into a boolean formula

$ALG \rightarrow$ Boolean Circuit \rightarrow Boolean Formula



Took arbitrary algorithm and showed how to convert to a boolean formula

NP NPC NP-Hard

NPC: X is NP-complete if

① X is in NP

② For every decision problem D in NP $D \leq_p X$.
(X is NP-Hard)

⑤ Reduction

Example: Maximum Bipartite Matching (MBM)

$MBM \leq^* \text{Max Flow}$

$x \uparrow f(x)$

"Reduces to"
Reduction: $O(V^2)$

Max-flow is at least as hard or harder than MBM

$A \leq B$ if B solve in polynomial
then A solve in polynomial
so "reduction" can not
more than $O(B)$

Property: New input $f(x)$ to max-flow should be solvable if and only if the original input to MBM is solvable

\uparrow If f : value of max flow

size of maximum bipartite matching

Reduction: Bipartite Graph: $G = (X \cup Y, E)$



$G' = (V', E')$

① directed edges $x \rightarrow y$

② add new edges $S \rightarrow$ each x to E'

③ add new edges each $y \rightarrow t$ to E'

④ capacity function: $c: V' \times V' \rightarrow \{0, 1\}$

every edge in E' , $c(e) = 1$

⑤ F-F(G')

↳ How we do reduce efficiently? \Rightarrow adj matrix $O(V^2)$

Cook: every problem in NP reduces to SAT, i.e. SAT is NP-complete
(because SAT is in NP)

$\Rightarrow \text{SAT} \leq_p 3\text{SAT}$

SAT: Given a boolean formula in n variables $f(x_1, \dots, x_n)$ is there an assignment of truth values (T/F, 1/0) to the variables so that the formula evaluates to true?

Boolean Formula in CNF. $n = \# \text{variables}$, $m = \# \text{clauses}$. Length of clause $\leq h$

↳ SAT: sample problem with $n=3$

Ex. SAT: $(x) \wedge (y \vee \bar{x} \vee w \vee u) \wedge (\bar{x} \vee y \vee z) \wedge (x \vee \bar{z})$

1. clauses with 1 variable

(x) : use auxiliary variables a, b

$$\Leftrightarrow (x \vee a \vee b) \wedge (x \vee a \vee \bar{b}) \wedge (x \vee \bar{a} \vee b) \wedge (x \vee \bar{a} \vee \bar{b}) \Leftrightarrow x$$

2. clauses with 2 variables

$(x \vee \bar{z})$: use auxiliary variable c

$$(x \vee \bar{z} \vee c) \wedge (x \vee \bar{z} \vee \bar{c}) \Leftrightarrow x \vee \bar{z}$$

3. clauses with 3 variables

$(\bar{x} \vee y \vee z) \Rightarrow$ just copy

4. clauses with more than 3 variables

$(y \vee \bar{x} \vee w \vee u)$: at least one variables is true

use auxiliary variable d .

$$(y \vee \bar{x} \vee d) \wedge (\bar{d} \vee w \vee u)$$

5. what we do if more than 4?

$$x \vee \bar{y} \vee \bar{w} \vee \bar{z} \vee u \vee v \quad \text{use } d_1, d_2, \dots$$

$6 \Rightarrow 4$

$$(x \vee \bar{y} \vee d_1) \wedge (\bar{d}_1 \vee \bar{w} \vee d_2) \wedge (\bar{d}_2 \vee \bar{z} \vee d_3) \wedge (\bar{d}_3 \vee u \vee v)$$

$\text{SAT} \leq_p 3\text{SAT} \Rightarrow 3\text{SAT is NP-complete}$

↓

Reduction: $O(m \cdot (n-2)) = O(m \cdot n)$
polynomial time

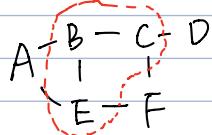
Does 3SAT \leq SAT Yeah! it already is!

Harder Reduction

3SAT \leq_p vertex cover \Rightarrow

$$A \rightarrow B \Leftrightarrow \bar{B} \rightarrow \bar{A}$$

DEF: a vertex cover of $G = (V, E)$ is set $S \subseteq V$ such that if $(u, v) \in E$ then $u \in S$ or $v \in S$



$\{B, C, E\}$ is a vertex cover of size 3

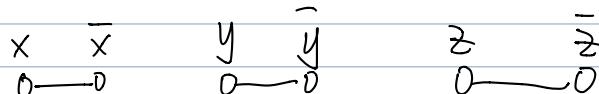
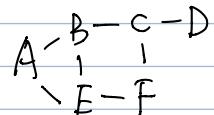
最小顶点覆盖: Minimum vertex cover

\hookrightarrow 3SAT $n = \#$ variables $m = \#$ clauses

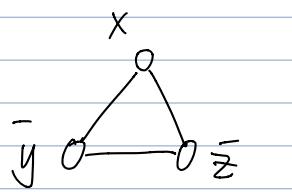
$$(x \vee \bar{y} \vee \bar{z}) \wedge (\bar{x} \vee y \vee z)$$

$$\therefore n=3 \quad m=2 \quad \therefore \Rightarrow x_1 y_1 z_1$$

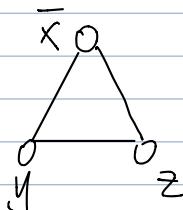
DEF: independent set
a set $S' \subseteq V$ such that
if $(u, v) \in E$, then u, v are
not both in S' .
 $\{A, D, F\}$ is an Independent Set



variable gadget

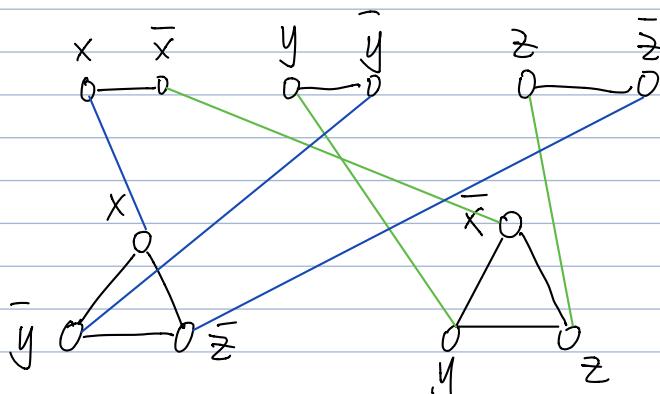


$$(x \vee \bar{y} \vee \bar{z})$$



clauses gadget

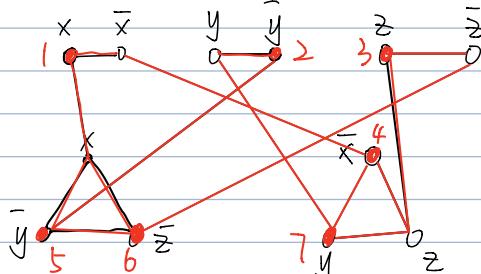
$$(\bar{x} \vee y \vee z)$$



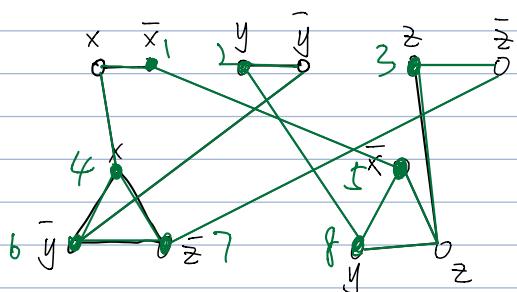
$$\begin{aligned} & \text{变量-系边一个} \quad \text{变量-个三角形约2个} \\ & k = n + 2m \\ & (\min \text{ vertex cover}) \\ & = 3 + 2 \times 2 = 7 \end{aligned}$$

If we can find a min vertex cover with $n+2m$ node in below graph

\Rightarrow There is a solution to original SAT



\Rightarrow we know $x_1 \bar{y}_1 \bar{z}_1 = 1$ is a solution;
start with it, check the min vertex cover
 \Rightarrow see. we finish vertex cover with
 $7 = 2 \times 2 + 3$ vertex



\Rightarrow we also know $\bar{x} \wedge y \wedge z = 1$ is not a solution; start with it
 see what we get
 \Rightarrow we need at least 8 vertex to finish the cover.

\Rightarrow Until Here : We show $3\text{SAT} \leq_p$ vertex cover with k .

But we have not shown \Rightarrow vertex cover $\leq_p 3\text{SAT}$ with k

Why?

3SAT is $\begin{cases} 1. X \text{ in NP} \\ 2. X \text{ is NP-Hard} \end{cases}$
 NPC For all decision problems D in NP , $D \leq_p X$
 so we need to show vertex cover $\leq_p 3\text{SAT}$

Week 2. Approximation For NP-Hard Problems

↳ vertex cover
TSP
set cover

■ Design an algorithm

■ Analyze the algorithm

- Bound Running Time = Big O notation
- Prove Correctness ↑
 poly time

Correct on all input

■ Goal:

- 1. Hard - Problem (NP-complete)
- 2. Fast Algorithm (poly-time) } \Rightarrow insist on
- 3. Exact Solution (correct) \rightsquigarrow relax correctness

Unless $P=NP$, cannot achieve all three at same time

Approximation Bound optimization problems

■ a set of valid outputs (Solutions), the cost / value $c(x)$ of each output x
whether to minimize or maximize cost

Ex. vertex cover : minimum size of vertex

clique : maximum size of vertex

TSP : tour of min weight

■ Optimal vs. Approx ALGs for optimization problems

Optimal ALG (usual notion of "correct")

compute for every input, a best valid output $\xrightarrow{\text{OPT}}$

$$\begin{aligned} C(\text{OPT}) &= \min \{ c(x) : \text{All valid input } x \} \\ &= \max \{ c(x) : \text{All valid input } x \} \end{aligned}$$

Ex. For TSP: $C(\text{OPT}) = \text{tour length} / \text{weight}$ (Not running time)

p: "rho" \Rightarrow approximation ratio

DEF: An ALG is a " ρ -approximation" ALG if on every input it computes a valid output APX whose cost/value is within A $\rho > 1$ factor of OPT.

For min problem: $C(\text{OPT}) \leq C(\text{APX}) \leq \rho \cdot C(\text{OPT})$

For max problem: $\frac{1}{\rho} C(\text{OPT}) \leq C(\text{APX}) \leq C(\text{OPT})$

ρ can be:

$\rho = O(n^c)$

$\rho = O(\lg n)$

$\rho = O(1)$

$p=1$ means $\Rightarrow C(OPT) = C(APX)$

so: p more closer to 1, more better

But we never know $C(OPT) \Rightarrow$ we need another thing to compare.

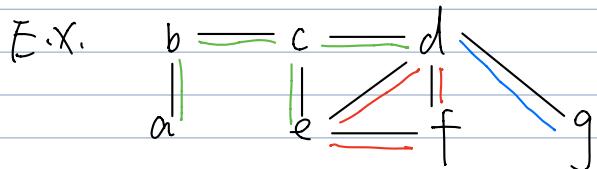
IV Vertex Cover Problem

- Input: undirected graph $G = (V, E)$
- Output: a set $C \subseteq V$ of min size \Rightarrow such that every edge in E has at least one endpoint in C .
- Object: minimize $|C|$.

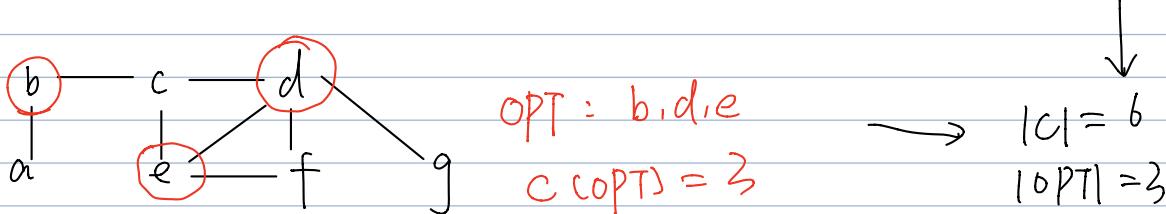
Approx - Vertex - Cover (G) // Adj List format

```

01  $C = \emptyset$ 
02  $E' = E_G$            // mark all edges as uncovered
03 while  $E' \neq \emptyset$     // while there still exists an uncovered edge  $e$ 
04   let  $(u, v)$  be arbitrary edge of  $E'$ 
05    $C = C \cup \{u, v\}$  // add both  $u, v$  to  $E$ 
06   remove from  $E'$  every edge incident to either  $u$  or  $v$ 
07 return  $C$            // Appropriate vertex cover
  
```



Approx Vertex Cover
 $\left. \begin{array}{l} (b,c) \rightarrow \text{add to } C \\ (e,f) \rightarrow \text{add to } C \\ (d,g) \rightarrow \text{add to } C \end{array} \right\} \Rightarrow C = \{b, c, e, f, d, g\}$
 $C(APX) = 6$



$$|C| = 6$$

$$|OPT| = 3$$

Theorem: Appr - Vertex - Cover is a 2-approximation ALG

Proof: output C is a vertex-cover of G . NTS: $p=2$

let $A :=$ set of edges picked on line 4

To cover edges in A , any vertex cover in particular OPT must include at least one endpoint of each edge in A .

A: 不共享节点的边集

No two edges in A share an endpoint, all other edges incident to its endpoint and removed from E' in line 6.

对于不共享节点的边集，在OPT中也至少需要n个点

No two edges in A are covered by the same vertex from OPT vertex cover so $|OPT| \geq |A|$

This is a lower bound on OPT V.C.

对于共享节点的边集，我们加入了2n个点

Further, each execution of line 4 picks and edge None of whose endpoints already in C

$|C| \leq 2|A|$. In fact $|C| = 2|A|$

Combine : $|C| = 2|A| \leq 2|\text{OPT.V.C.}|$.

$$\therefore \frac{|C|}{|\text{OPT}|} \leq 2 = \rho_{\text{ho}}$$

What if we do it greedy? constantly remove vertex with largest degree

Greedy - vertex - cover

- 01 $C := \emptyset$
 - 02 $E_0 = E$
 - 03 while $E_i \neq \emptyset$
 - 04 $i := i + 1$
 - 05 $v_i = \text{vertex in } E_{i-1} \text{ with max degree}$
 - 06 $d_i = \deg(v_i)$
 - 07 $E_i = E_{i-1} - v_i$
 - 08 $C = C \cup \{v_i\}$
 - 09 return C
- $\rho = O(\lg V)$

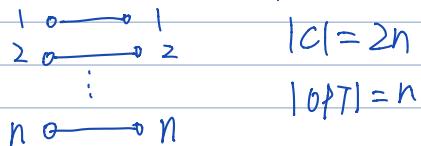
Recall : Matching : a matching M for a graph $G = (V, E)$ is a subset $M \subseteq E$ such that each vertex is incident to at most one edge in M .

Maximal : if not possible to add more edges in M .

$\therefore \text{APPROX-Vertex-Cover} : |A| \text{ is a matching } \star$

Can we decrease ρ $2 \rightarrow 1.9999\ldots$ or smaller? No!

Because we can find an example, that ρ must be 2.



we can do better than 2. But the upper bound must be 2

Currently best known approx for ρ is $2 - O(\frac{1}{\sqrt{\ln n}}) = 2 - O(1)$ // cite: Halperin 2000

TSP (Traveling Salesman Problem)

Input: Complete undirected graph $G = (V, E)$ with weight function: $w: E \rightarrow \mathbb{R}^*$ $|V| = n$

Output: Hamilton cycle H in G

Objective: Minimize $w(H)$

Ham Cycle \leq_p TSP (因 TSP 能解决 Ham Cycle)

Ham Cycle is NP-Hard \Rightarrow TSP is NP-Hard

HamCycle: $G = (V, E)$ Arbitrary (no need to be complete)
Undirected
Unweighted

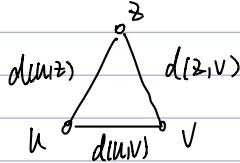
$$\text{let } w(e) = \begin{cases} 1 & \text{if } e \in E \text{ in } G \\ 0/W & \text{o/w} \end{cases}$$

\therefore if there is a Ham Cycle in G . Then there must be a Hamcycle in K_n with weight exactly n . (o/w. it will weight $\geq n+1$)

Extra Condition: (Strong Restriction)

Distance Function $d: E \rightarrow \mathbb{R}^{>0}$ satisfy triangle:

$$\forall u, v, z \in V. \quad d(u, v) \leq d(u, z) + d(z, v)$$



$$d(u, u) = 0 \quad \forall u \in V$$

$$d(u, v) = d(v, u) \quad \forall u, v \in V$$

$$d((x_1, y_1), (x_2, y_2)) = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2} \Rightarrow \text{Euclidean metric}$$

Idea A: Approximate Optimal TSP tour using MST's and graphs obey Δ -inequality

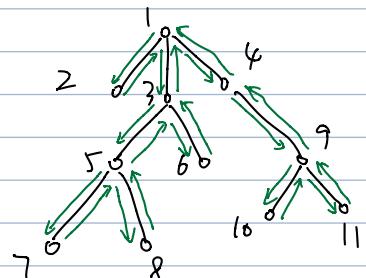
1. weight of MST is lower bound on cost of optimal tour

remove an edge from a TSP tour get a path P is a spanning tree

$$\text{TSP cost} \geq \text{Path } P \text{ cost} \geq \text{MST cost}$$

2. use MST to build a TSP Tour.

consider a DFS traversal of MST

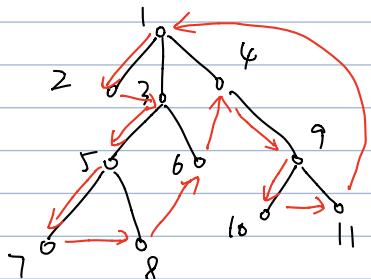


we use each edge twice
 \Rightarrow call it: twice-around tour W
 Tour W weight = $2 \times$ weight(MST)

$$2 \text{TSP cost} \geq 2 \text{MST cost} = \text{Tour } W \text{ cost}$$

3. Remove extra vertices

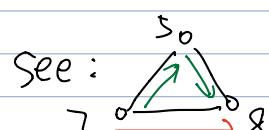
Take a shortest path to next unvisited vertex



$$\begin{array}{ccccccccccccc} 1 & - & 2 & - & 3 & - & 5 & - & 7 & - & 8 & - & 6 & - & 4 & - & 9 & - & 10 & - & 11 & - & 1 \end{array}$$

Tour H

4.



Due to the triangle restriction

\therefore we can say $7 \rightarrow 8$ less than $7 \rightarrow 5 \rightarrow 8$

\therefore Tour W cost \geq final tour H cost (APX)

$$\frac{2 \text{TSP}}{2 \text{OPT}} \geq \text{Tour W} \geq \text{Tour H}$$

$\frac{|\text{APX}|}{|\text{OPT}|}$

$$P = \frac{2}{2} \geq \frac{|\text{APX}|}{|\text{OPT}|}$$

Pseudocode

Approx-TSP-Tour(G, d) \leftarrow Δ -inequality

- 01 $T_1 = \text{MST-Prim}(G, d, S)$ $\parallel O(V^2)$
- 02 $W = \text{DFS}(T_1)$ \parallel tour traversal of T_1 starting at S $\parallel O(V)$
- 03 $H = W$
- 04 for each vertex v in H (in traversal order) $\parallel O(V)$
- 05 if v has already been visited
- 06 then delete v from H
- 07 return H \parallel add up weight: weight of APPROX TOUR
 $\Rightarrow O(V^2)$

Homework: Christofides $\frac{3}{2}$ Appr to TSP

Set Cover Problem

Set system (X, F) $X = \{x_1, x_2, \dots, x_n\}$ Finite Set, $|X| = n$, called universe

$F = \{S_1, S_2, \dots, S_m\}$ Family of subsets of X , such that every element of X belongs to at least one set in F .

Undirected Graph: $G = (V, E)$

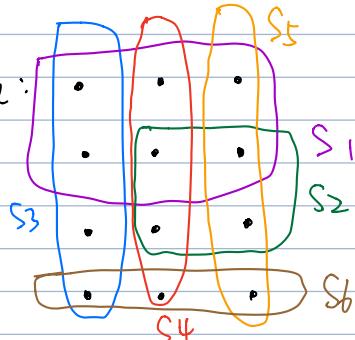
Universe \uparrow
 X \uparrow
 wireless \uparrow
 subset of cardinality 2
 集合中元素的个数

DEF: A cover of F is a subcollection of sets whose union comprises X : $\bigcup_{i \in C} S_i = X$ SieF

Set Cover Problem: Given a set system (X, F) $F = \{S_1, \dots, S_m\}$
 Compute a set C of minimum cardinality such that $\bigcup_{i \in C} S_i = X$

NP-Complete

Example:



X : set of black dots
 F : $\{S_1, S_2, \dots, S_6\}$

Set cover of size 3: S_3, S_4, S_5 [exact cover]

Universe •

Approx : Greedy \Rightarrow

- ① Pick the set that covers the most element
- ② Throw out all elements covered
- ③ Repeat until no element left

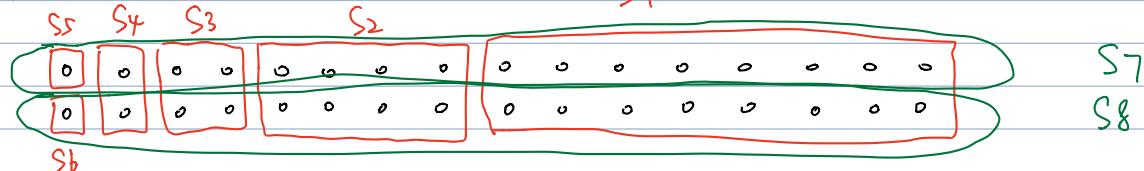
This algorithm has approx $P = O(\ln h)$

Approx-set-cover (X, F) // $F = \{S_1, \dots, S_m\}$ $|X| = n$

- ① $U = X$ // U stores uncovered element
- ② $C = \emptyset$ // C stores set of cover
- ③ while $U \neq \emptyset$:
- ④ select S_i in F that maximizes $|S_i \cap U|$
- ⑤ $U = U - S_i$ // remove element of S_i from U
- ⑥ $C = C \cup \{S_i\}$ // add S_i to C
- ⑦ return C // sub family of F that covers X .

Running Time: $O(n^2 m)$

$|X| = 32 \rightarrow$ Case that APX-ALG will always make mistakes



$$\text{OPT}(S_7, S_8) \quad \text{APX-ALG}(S_1, S_2, S_3, S_4, S_5, S_6) \Rightarrow |\text{OPT}| = 2 \quad |C| = 6$$

Theorem: Suppose $|X| = n$ if OPT uses k sets, greedy APX ALG find a solution with at most $O(k \ln k)$ set.

Proof: Since OPT uses k sets, there must be some set that covers at least $\frac{1}{k}$ of the elements.

Greedy chooses set that covers the most, so it covers at least $\frac{n}{k}$, i.e. after first iteration at most $n(1 - \frac{1}{k})$ left.

Again, OPT uses k sets there must be a set that covers $\frac{1}{k}$ of the remainder. Again Greedy chooses set that covers most elements. after 2-nd iteration there are $n(1 - \frac{1}{k})^2$ left at most.

Generally, after t iterations at most $n \cdot (1 - \frac{1}{k})^t$ elements left.

after $t = k \ln n$ iterations there are at most $n(1 - \frac{1}{k})^{k \ln n} < \frac{1}{k}$ left

$$n(1-\frac{1}{k})^{kn} < n(\frac{1}{e})^{kn} = 1 \text{ element left}$$

which means we are done.

Week 3 Randomize Algorithms

Can flip coins pick a random bit or random element as a basic step.

Two Type

① Monte Carlo Algorithm : for every input correct with high probability

② Las Vegas Algorithm : for every input, always correct

String Equality

Multiple copies of huge dataset : n-bit string for some very large n

Problem : Alice (Processor A) : a_1, a_2, \dots, a_n

Bob (Processor B) : b_1, b_2, \dots, b_n

Are string equal? with as few communications as possible

Deterministic Algorithm : Alice transmits n bits to Bob
Bob checked. send Alice answer
Requires $n+1$ bit communication.

Randomized Algorithm : $O(\log n)$ bits of communication.

A Field : a set of elements : additions, multiplication operations, satisfy commutative,
associative, distributive laws. 封闭律

Addition

Multiplication

each element in field has additive and multiplicative identity
additive and multiplicative inverse. EX : \mathbb{R}, \mathbb{Q}

A Finite Field : has a finite # of elements.

Ex. field of integers mod P , P a prime; two integers a and b
are equal iff they have same remainder when divided by P :
 $a \equiv b \pmod{P}$ \mathbb{Z}_p has p elements $\{0, 1, \dots, p-1\}$

Randomized Algorithm: n-bits in server A : a_1, a_2, \dots, a_n
n-bits in server B : b_1, b_2, \dots, b_n .

Consider polynomials $A(x) = \sum_{i=0}^n a_i x^i$
 $B(x) = \sum_{i=0}^n b_i x^i$

Consider a polynomial in \mathbb{Z}_p : $A(x) = \sum_{i=0}^n a_i x^i \pmod{p}$

a root or zero of a polynomial is a value of X for which $A(x)=0$

Theorem: a polynomial $A(x)$ of degree n has at most n zeros (\Leftarrow Using Horner's Rule
秦九韶算法)

Algorithm: ① Alice and Bob agree on a large enough field F
(Modular arithmetic over a large enough prime p)

② Alice picks a random element $r \in F$ evaluates $A(r)$, send r and $A(r)$ to Bob

③ Bob computes $B(r)$ check whether $A(r) = B(r)$
Return { No , $A(r) \neq B(r)$
yes , $A(r) = B(r)$

cost: # bits (digits) communicated; cost of local computation by either Alice and Bob is ignored

Analyze success probability

Analysis: If string same, ALG always return yes

If Return No, ALG will never make mistakes

But the ALG makes a mistake when two strings are not same and it returns yes

Upper bound this error probability:

$A(x) \neq B(x)$ but chosen r that $A(r) = B(r)$

\Rightarrow Happens when r is a root of $(A-B)(x)$: poly of degree at most n

\Rightarrow There are at most n roots

error probability is at most $\frac{n}{|F|} = \frac{n}{p}$

if choose $|F| = 1000n$, error probability is at most 0.001

$|F| = 1000n^2$ error probability is at most $\frac{0.001}{n}$

Repeated Trials: k -trials, independent [further decrease error rate]

let ϵ : Margin of error

T_i : R.V. indicating result of i -th Trial (out of k)

Constraint: $\epsilon \geq \Pr[T_1=0 \text{ and } T_2=0 \dots \text{ and } T_k=0]$

$$\geq \left(\frac{n}{p}\right)^k$$

Solving: $\left(\frac{n}{p}\right)^k \leq \epsilon \Rightarrow k \geq \frac{\lg \epsilon}{\lg(\frac{n}{p})} \because n < p \Rightarrow \lg(\frac{n}{p}) < 0$

Pseudo Code: Return whether $A(x) \equiv B(x) \pmod{p}$ $\forall x \in \{0, 1, \dots, p-1\} = \mathbb{Z}_p$
with correction rate at least $1-\epsilon$

Poly-EQUAL ($a[0..n]$, $b[0..n]$, p, ϵ)

Using Horner's Rule

$$01 \quad K = \left\lceil \frac{\lg \epsilon}{\lg(p/n)} \right\rceil \quad \Theta(1)$$

秦九韶算法

```

02 for i = 1 to k
03   r = random(0, p-1)           Θ(1)
04   A' = A(r)                   ] Θ(n)
05   B' = B(r),
06   if A' - B' ≠ 0 (mod p)     Θ(1)
07   return no
08 return yes
  
```

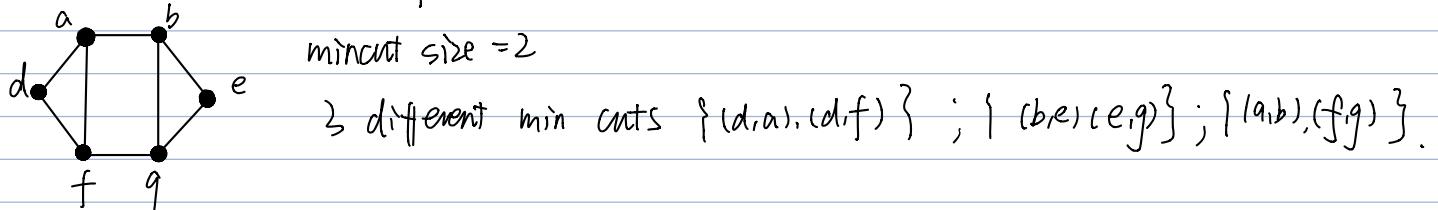
$$T(n) = \Theta(n \cdot \lceil \lg \epsilon / \lg(p/n) \rceil) \quad \Theta(1)$$

Karger's Algorithm

Minimum Cut Problem : $G = (V, E)$ connected, undirected, unweighted graph

① A cut in G is a subset of edges $F \subseteq E$ whose remove disconnected G .

② A minimum cut is a cut of minimum size



③ Properties of min cut:

① Claim 1: The size of a min cut is at most the min degree of any vertex in graph

Proof: if v has min degree. the cut $\{ \{v\}, V - \{v\} \}$ has size $\deg(v)$

since min cut no larger than any cut in the graph

$$\text{min-cut-size} \leq \deg(v)$$

② claim 2: in graph with n vertices if there is a min-cut of size k , G must have at least $\frac{nk}{2}$ edges

Proof: if there is a mincut of size k . every vertex must have degree at least k . by hand shake theorem $|E| = \frac{\sum \deg(v)}{2} \geq \frac{nk}{2}$

③ Min-cut-problem : Find a cut of min size in G

Simple idea : mincut \rightarrow maxflow

A(s) a min $s \rightarrow t$ cut for some $s, t \in V$
compute min $s \rightarrow t$ cut for all pairs $s, t \in V$ and take the min one

n^2 calls of maxflow

Deterministic : $\binom{n}{2}$ maxflow $\Rightarrow \Theta(V^3 E)$ Slow

Randomized ALG: Karger: Random Contraction

\Rightarrow improving Running Time to $\tilde{\Theta}(V^2)$ faster than max flow $\Theta(V^2 \log^3 V) = \tilde{\Theta}(V^3)$

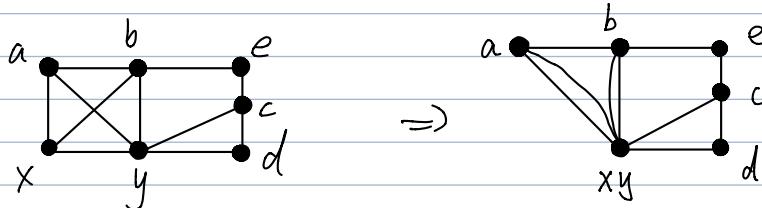
Algorithm: Repeat

pick a random edge e
contract 2 endpoint of e
until only 2 vertices left
By "contracting" an edge, identify the endpoints as a single vertex

Edge Contraction

Given an edge (x, y) in a multigraph G , contract x and y as follows

- Delete all edges between x and y
- Replace x and y with a new "supernode" xy
- replace all edges incident to x or y with edges incident to supernode xy
- Remove loops



Karger's ALG

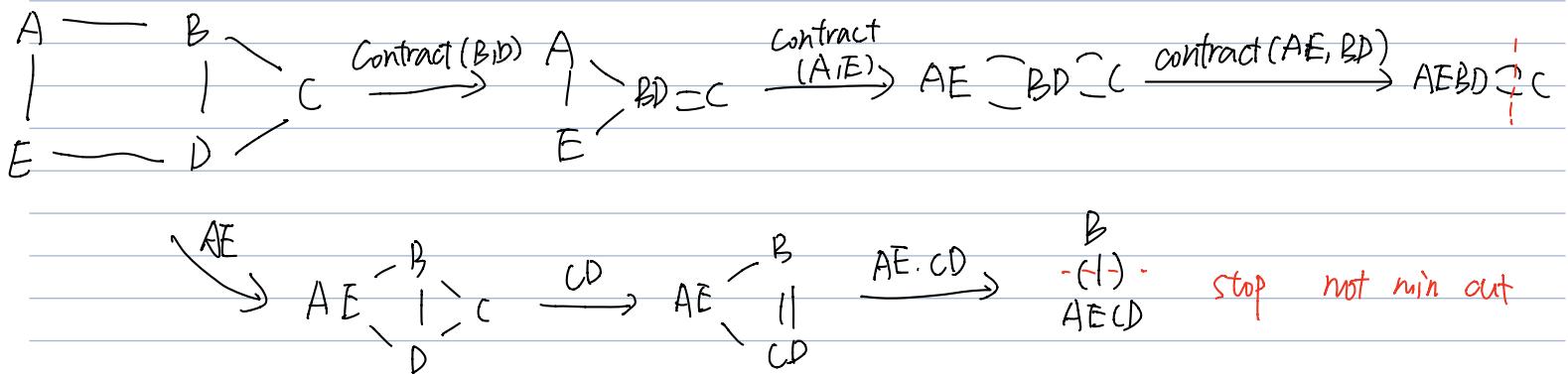
1. pick edge (x, y) at random in G
2. contract edge (x, y)
3. if there are more than 2 vertices \Rightarrow go back to 1.
4. else output edges remaining as the cut

Contract ($G, (u, v)$)

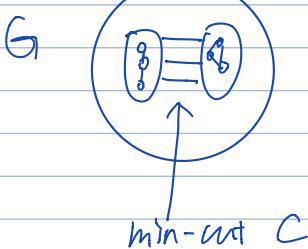
- 01 replace u and v by a new vertex uv
- 02 delete all edges (u, v)
- 03 for every edge (w, u) replace it by (w, uv)
- 04 for every edge (w, v) replace it by (w, uv)

Karger (G)

- 01 for $i=1$ to $n-2$
 - 02 pick an edge (u, v) in G at random
 - 03 $G = \text{contract}(G, (u, v))$
 - 04 return edges between 2 vertices that remain
- $\Theta(1)$ $\Theta(n)$ $\Theta(n^2)$



Correctness Idea



Key idea: ALG should never contract one of the min-cut edges

Don't pick edge in C

Performance Analysis: $\Pr(\text{ALG is correct}) \geq \Pr(\text{C output at end}) = \Pr(\text{no edge of C contracted})$

Theorem: Karger's ALG outputs a min-cut with probability at least $\frac{2}{n(n-1)} = \frac{1}{\binom{n}{2}} \geq \frac{1}{n^2}$
where $n = \# \text{ of vertices}$

Proof: There may be many min-cuts in G. Fix some min-cut C. Let C have k edges.

Lemma 1: So long as none of the edges in cut C are contracted, C remains a cut

Proof: Initially all vertices are in left or right side of C, contracting an edge that does not cross cut can only connect vertices on same side of cut

Lemma 2: Any cut in new graph is also a cut in original graph

Proof: The min degree in all intermediate graphs is at least k.

O/W \rightarrow min size k

Lemma 3: When have $n-i$ vertices remaining (for some $i \geq 0$) the remaining are at least $\frac{(n-i)}{2} k$ edges in current graph.

Proof \Rightarrow : Denote edges contracted by e_1, e_2, \dots, e_{n-2} , C survives if these edges not in cut

① Prob C survive 1-st round ($e_i \notin C$)

$$1 - \frac{k}{\# \text{edges}} \geq 1 - \frac{k}{\frac{n(n-1)}{2}} = \frac{n-2}{2}$$

② Prob C survive after i-th round ($e_1, e_2, \dots, e_i \notin C$)

$$1 - \frac{k}{\# \text{edges}} \geq 1 - \frac{k}{\frac{(n-i)k}{2}} = 1 - \frac{2}{n-i} = \frac{n-i-2}{n-i}$$

$\text{Prob}(c \text{ survives}) = \text{Prob}(e_1 \notin c, e_2 \notin c, \dots, e_{n-2} \notin c)$

$$\begin{aligned}&= \left(1 - \frac{2}{n}\right) \left(1 - \frac{2}{n-1}\right) \left(1 - \frac{2}{n-2}\right) \cdots \left(1 - \frac{2}{3}\right) \\&= \frac{n-2}{n} \times \frac{n-3}{n-1} \times \frac{n-4}{n-2} \times \cdots \times \frac{3}{5} \times \frac{2}{4} \times \frac{1}{\cancel{2}} = \frac{2}{n(n-1)} \\&\geq \frac{1}{n^2}\end{aligned}$$

How can we boost success probability? Repeated Trials

$\text{Prob}(\text{single run of KARGER will return a wrong answer})$:

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

Run Karger t times independently and output smallest cut found:

Prob that all t runs failed: $(1 - \frac{2}{n(n-1)})^t$

Using Useful Inequality $1+x \leq e^x \forall x \in \mathbb{R}$

$$\Rightarrow (1 - \frac{2}{n(n-1)})^t \leq e^{-t \cdot \frac{2}{n(n-1)}}$$

If run Karger $\frac{n(n-1)}{2}$ times, prob of failure is

$$(1 - \frac{2}{n(n-1)})^{\frac{n(n-1)}{2}} \leq \frac{1}{e}$$

if run Karger $\frac{n(n-1)}{2} \cdot \ln h$ times, prob of failure

$$(1 - \frac{2}{n(n-1)})^{\frac{n(n-1)}{2} \cdot \ln h} \leq (\frac{1}{e})^{\ln h} = \frac{1}{h} \quad (\text{small failure})$$

\Rightarrow Run Time $O(n^4 \log n)$

Week 4 Randomized Approximation Algorithm

Max-3-SAT

3-SAT (NP-complete)

Input: a set of m clauses, each with exactly 3 distinct variables, over a set $X = \{x_1, \dots, x_n\}$ of n boolean variables.

Output: Is there a satisfying TVA?

No clause contains a variable and its negation

Max-3-SAT

Input: Same e.g. $C_i = (x_{i1} \vee \bar{x}_{i2} \vee x_{i3})$

Output: a satisfying TVA that maximizes number of clauses satisfied

\Rightarrow Also NP-Hard

Randomized Approximation ALG

IDEA: Try the all-true and all-false assignment

\Rightarrow every clause is satisfied by at least one of them

\therefore The best of the 2 assignment satisfies at least half the clauses and gives a 2-approximation

Randomizing - Pseudocode

01 for each variable x_i

02 set its value to 0/1 by flipping a coin

suppose coin is fair

\Rightarrow what is the expected # clauses satisfied by such a random assignment?

\Rightarrow for any clause C_i , the probability that it is not satisfied is exactly $(\frac{1}{2})^3$. $\therefore C_i$ is satisfied with prob $= 1 - \frac{1}{8} = \frac{7}{8}$

Let Z_i be indicator R.V. such that $Z_i = \begin{cases} 1 & \text{if } C_i \text{ satisfied} \\ 0 & \text{o/w} \end{cases}$

$Z :=$ total # of satisfied clause

$$Z = \sum_{i=1}^n Z_i$$

$$\mathbb{E}(Z) = \mathbb{E}\left(\sum_{i=1}^m Z_i\right) = \sum_{i=1}^m \mathbb{E}(Z_i) = \frac{7}{8}m \quad (\text{Theorem. we don't care about independency})$$

key Ideas: $E(Z) = \frac{7}{8}m$ represent a kind of average i.e. not all the Z_i 's can be less than $E(Z) = \frac{7}{8}$

\therefore There has to be some assignment to the variables that achieves it.

Theorem: For any instance of 2-SAT, there exists a truth assignment that satisfy at least a $\frac{7}{8}$ fraction of all the clauses

Proof: $E(Z) = \frac{7}{8}m$. and the R.V. $Z \geq E(Z)$ some of the time.

ERDŐS: probabilistic method

Clearly m is an upper bound on # satisfied clauses

\therefore approx ratio is at most $m/\frac{7}{8}m = \frac{8}{7}$

Approximation factor $\frac{7}{8}$ of optimal

■ Balls and Bins (CLRS 5.4.2)

n balls, m bins

Each ball throw to a random (uniformly) independently

■ Expected # balls in a bin

let B_{ij} be independent R.V.

$$B_{ij} = \begin{cases} 1 & \text{if ball } j \text{ in bin } i \\ 0 & \text{o/w} \end{cases}$$

$$\begin{aligned} \text{then } E[\# \text{ balls in bin } i] &= E\left[\sum_{j=1}^n B_{ij}\right] = \sum_{j=1}^n E(B_{ij}) \\ &= \sum_{j=1}^n \Pr(\text{ball } j \text{ in bin } i) \\ &= \sum_{j=1}^n \frac{1}{m} = \frac{n}{m} \end{aligned}$$

In particular, while $m=n$ expected # = 1

Do we expect that most bins have about one ball most of the time?

- 5 balls, 5 bins

- 0 small prize: if 4 balls in row
- 0 or 4 balls in column

- - - 0 0 0 0 0 2 prize: if 5 balls in a row
1 2 3 4 5 or 5 balls in a column

Binomial

expected # empty bins

let y_i be indicator R.V.

$$y_i = \begin{cases} 1 & \text{if bin } i \text{ is empty} \\ 0 & \text{otherwise} \end{cases}$$

Then: $E(y_i) = \Pr(\text{bin } i \text{ is empty}) = (1 - \frac{1}{m})^n \leftarrow \text{must be missed by all } n \text{ balls}$
 $\approx e^{-\frac{n}{m}} \quad (\because 1-x \leq e^{-x}$
so. $E[\#\text{empty bins}] = E\left[\sum_{i=1}^m y_i\right] = m \cdot e^{-\frac{n}{m}}$ and $1-x \approx e^{-x}$ for small x)

when $n=m$. Expect to see an $\frac{1}{e}$ fraction of the bins are empty

⇒ famous (coupon collectors problem)

Balls and bins (CLRS 5.4.2)



Maximum Load \Rightarrow Hashing, Addressing, collision

Maximum Load

Birthday Paradox:

for what n do expect to see 2 balls in a bin? \Rightarrow a "collision" occurs

Birthday Paradox where $m=365$ (ignore feb 29)

Prob of no collision in 1-st n balls:

$$1 \cdot (1 - \frac{1}{m}) \cdot (1 - \frac{2}{m}) \cdots (1 - \frac{n-1}{m}) \leq e^{-\frac{1}{m}} \cdot e^{-\frac{2}{m}} \cdots e^{-\frac{n-1}{m}} = e^{-\frac{n(n-1)}{2m}} \approx e^{-\frac{n^2}{2m}}$$

This prob is less than $\frac{1}{2}$ when $n = \sqrt{2m \ln 2}$

for $m=365$ when $n \geq 22.49$ the prob that maxload is at least 2 is at least $\frac{1}{2}$

Summarize: expect to see a collision with $n = \Theta(\sqrt{m})$

$\Rightarrow n^2 \approx m$. Some collisions occur.

Hashing

Sequence of insert, search, delete operations

Goal: Do all operations in $O(1)$ time

Formal Setup: keys come from some large universe U

set $S \subseteq U$ of actual keys stored : dynamic

$$|S|=n \quad n \ll |U|$$

Perform insert, search, delete operations

Table T of size m and hash function: $h: U \rightarrow \{0, 1, \dots, m-1\}$

Given key $x \in U$, store it table position $T[h(x)]$

Two key questions

How compute h ? h is Pseudo-Random: $h(x)$ equally likely to be any of $0, 1, \dots, m-1$ independent of n for other keys.

How resolve collision?

collision: $h(x) = h(y)$ for 2 different keys $x \neq y$.

Chaining: each entry in table T is a linked list.

open-addression: No linked lists: all keys stored in table T .

■ Chaining

operation: search, insert, delete

\Rightarrow compute index $i = h(x)$

search: walk down linked list at $T[i]$

insert: insert x at head of $T[i]$

delete: perform deletion on list $T[i]$

1. keys spread out

2. $m = O(n)$

3. function n should be fast to compute. Assume $O(1)$

Average length of lists:

Let n_i = length of list $T[i]$ ($\sum_{i=0}^{m-1} n_i = n$)

Running Time is $O(\max_i n_i)$ (worst case)

length of list should be $O(n)$ (bad news)

switch to expected cost of operation let $\alpha = \frac{n}{m}$ = "load factor"
 $=$ average length of list

Running Times: $O(1+\lambda) = O(1)$ if $m = O(n)$

How to compute a hash function?

① Division method $\underline{h(x) = X \text{ mod } m}$ $m = \text{size of table}$

deterministic hash function. only Randomness assuming

Fact: put data is **not** random

| EX. if m prime for all integers a and x

$$h(x+am) = h(x)$$

All such integer multiples of m collide with certainty

→ Conclusion: Hash function must be random

② Universal Hashing: incorporates randomness into construction of hash function

DEF: a randomized ALG H for constructing hash functions $h: U \rightarrow \{0, 1, \dots, m-1\}$ is **universal** if for **all** $x \neq y$ in U

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

Set H of such hash functions a **universal hash family**
if procedure "choose h from H " is universal.

Theorem: if H is universal, then for any set $S \subseteq U$, $|S| = n$.

for any $x \in U$ if $h \in H$ the **expected** # collisions
between x and all other items in S is at most $\frac{n}{m}$

Proof: Each $y \in S$ ($y \neq x$) has at most a $\frac{1}{m}$ chance of colliding with
 x by def of universal. So

$$\text{let: } C_{xy} = \begin{cases} 1 & \text{if } x, y \text{ collide} \\ 0 & \text{o/w} \end{cases}$$

$$\text{let } C_x := \text{total # collision for } x. \text{ so } C_x = \sum_{\substack{y \in S \\ y \neq x}} C_{xy}$$

$$\mathbb{E}(C_{xy}) = \Pr(x, y \text{ collide}) \leq \frac{1}{m}$$

$$\mathbb{E}(C_x) \leq \frac{n}{m}$$

Corollary: expected time to search / insert X with universal hash function

= expected time to look up a given key X .

= $O(1 + \text{expected # keys in } X's \text{ linked list})$

= $O(1 + \text{expected # keys that collide with } X)$

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

$$= O(1+\alpha) \quad \alpha = \frac{n}{m} = \text{load factor}$$

$$= O(1) \quad \text{if } m = \Theta(n)$$

Theorem: Universal family of hash functions exist. (CLRS Theorem 11.5)

Example: $h_{a,b}(x) = ((ax+b) \bmod p) \bmod m$

for any integers $a \in \mathbb{Z}_p^+ = \{1, 2, \dots, p-1\}$ $b \in \mathbb{Z}_p = \{0, 1, \dots, p-1\}$

a, b Random chosen at set up, $\begin{cases} 0 \leq a, b < p & (a \neq 0) \\ p: \text{large prime} \end{cases}$

$$\Pr_{a,b} (h(x) = h(y)) \leq \frac{1}{m} \quad \text{for all } x \neq y.$$

Theorem: $h_{a,b}$ is universal

Proof: fix 4 integers $r, s, x, y \in \mathbb{Z}$. such that $x \neq y, r \neq s$.

$$\begin{aligned} \text{Linear system: } ax + b &\equiv r \pmod{p} \\ ay + b &\equiv s \pmod{p} \end{aligned}$$

has a unique solution $a, b \in \mathbb{Z}_p$ with $a \neq 0$, v.z.

$$\begin{aligned} a &= (r-s)(x-y)^{-1} \pmod{p} && \text{Mult inverse of } (x-y) \\ b &= (sx-ry)(x-y)^{-1} \pmod{p} \end{aligned}$$

Open Addressing For Collision Resolution

No linked lists : All keys stored in table T

Each $T[i]$ is either NIL (empty) or a key

	T	$\sum_{i=0}^{n-1} T[i] \leq m \quad \alpha \leq 1 \quad (\alpha = \frac{n}{m})$
0	N	
1	N	
	k_2	Table-doubly it's necessary
	N	
	k_3	Typically $\alpha \leq \frac{1}{2}$
	k_1	
	k_4	
$M-1$	N	(python resizes at $\alpha = \frac{2}{3}$)

Hash function

Replace $h(X)$ with $h(x, i)$ for $i = 0, \dots, m-1$

given a probe sequence $h(x, 0), h(x, 1), \dots, h(x, m-1)$

for key x . probe sequence is a permutation of $0, 1, \dots, m-1$ specifying a sequence of slot positions to be probed in search

Search (T, X)

```
01 for i=0 to m-1
02   j = h(x, i)
03   if T[j] == x
04     return j
05   else if T[j] == nil
06     return nil // not in table
else (T[j] != x & T[j] != x)
  pass
```

Time: $O(\alpha)$

Insert (T, X)

```
01 for i=0 to m-1
02   j = h(x, i)
03   if T[j] == nil
04     T[j] = x
05   return j // return # of slot that get key x.
06 return full
```

Deletion Tricky:

- find key
- replaced by "deleted"
- skip over "deleted" in search
- but use "deleted" in insert

Set $T[h(k, i)] = \text{nil}$

may cause search fail

Probe Strategies

• Linear probing:

$$h(x, i) = (h(x, 0) + i) \bmod m$$

• Double Hashing

$$h(x, i) = (h_1(x) + i \cdot h_2(x)) \bmod m$$

↑ ↑
2 different hash functions

Bloom Filters

Setting: huge universe U of possible elements (e.g. password)

Maintain a subset of $S \subseteq U$ where $|S| \leq n$

using an array $B[0, 1, \dots, m-1]$ where $m = |B|$, $m = cn$ for $c \geq 1$

Operations: $\text{insert}(x)$: add x to S

$\text{query}(x)$: is $x \in S$?

Want faster queries
small space

deletion not considered

What if willing to trade an occasional unsuccessful search in exchange for size and speed, nearly all cases?

Idea: use probabilistic data structure: Bloom filter

Bloom filter $\{$: bit vector of length m

\cdot a set of K hash functions that returns value between 0 and 1

$K = \#$ of hash functions

hash functions $h_1, h_2, \dots, h_K : U \rightarrow \{0, 1, 2, \dots, m-1\}$

\hookrightarrow need to be random and mutually independent

① Initial Bloom Filter zeroed out
initialize B to all 0

② Insert(X)

for $i = 1$ to K
 $f = h[i]$
 $B[f(x)] = 1$

Ex: 0000 0000 000

Insert password,
 $\left\{ \begin{array}{l} \text{hash}_1(\text{PW}_1) = 3 \\ \text{hash}_2(\text{PW}_2) = 10 \\ \text{hash}_3(\text{PW}_3) = 8 \end{array} \right.$

0001 0000 101

Insert password 2

$h_1(\text{PW}_2) = 5$

$h_2(\text{PW}_2) = 8$

$h_3(\text{PW}_2) = 1$

0101 0100 101

▲ contains PW1 and PW2

so we can't delete key but in password scenario we don't need delete

⑩ Query (x)

```
01 for i=1 to k  
02   f=h(i)  
03   if B[f(x)] == 0  
04     return false  
05 // If reached here all B[f1(x), f2(x), ..., fk(x)] = 1  
06 return True
```

Note : If $x \in S$. Always return "Yes"

But if $x \notin S$ Also might return "Yes" [all the hints could be covered by other elements in the set]

How to minimize False Positive?

What is prob of False Positive?

$$\Pr(\text{all } nk \text{ balls miss bin } b) = \left(1 - \frac{1}{m}\right)^{nk} \approx e^{-\frac{nk}{m}}$$

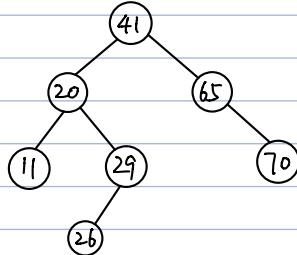
$$\Pr(\text{False Positive for } x \notin S) = \left(1 - e^{-\frac{nk}{m}}\right)^k$$

Week 5 Balanced Search Trees (B-tree)

(nothing but Binary Search Tree)

Quick Review On BST.

- Rooted Binary Tree
- each node x :
 - key [x]
 - left (x)
 - right (x)
- $p[x]$ (parent)



- any node x : for all nodes y in left subtree of x $\text{key}[y] \leq \text{key}[x]$
for all nodes y in right subtree of x $\text{key}[y] \geq \text{key}[x]$
- operations: search (T, k)
insert (T, k)
delete (T, k)
minimum (X)
maximum (X)
successor (X) \Rightarrow smallest key $> \text{key}[X]$
predecessor (X) \Rightarrow largest key $< \text{key}[X]$

- Analysis: operation all $O(h)$ where $h = \text{height of tree}$

Minimum Running Time:

Problem: worst case : $\Theta(n)$ No better than Linked List

Solution: Guarantee Small Height (balanced)

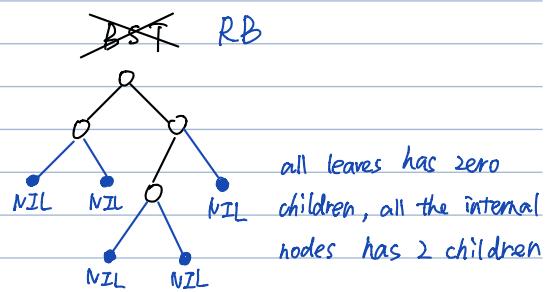
↓

complete binary tree : $O(\lg n)$

Method: Restructure Tree if necessary: insert / delete change the structure of tree

Red-Black Tree (CLRS 13.1)

- Add a color field (one bit) to each node of BST
- RB Tree Properties
 - ① every node is red / black
 - ② root and leaves (nils) are black
 - ③ every red node has a black parent
 - ④ all simple paths from node X to a descendant leaf of X have the same # of black nodes = Black height of X



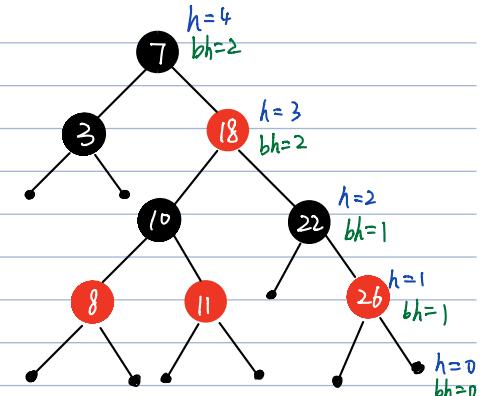
Black Height: $bh(x) = \# \text{ black nodes on path}$
to a leaf, not containing X

(internal nodes)

Theorem: a RB tree with n keys has height $h \leq 2 \lg(n+1)$

Lemma 1: any node with height h has black height $\geq \frac{h}{2}$

proof: height is the length of a longest path to a leaf
by property ③, $\leq \frac{h}{2}$ red nodes on the path
 $\therefore \geq \frac{h}{2}$ are black



Lemma 2: for any node X , the subtree rooted at X contains $\geq 2^{bh(x)-1}$ internal nodes.

proof: By induction on height on X .

basis: $h=0 \Rightarrow X \text{ leaf} \Rightarrow bh(x)=0 \Rightarrow 2^{bh(x)-1}=1-1=0 \checkmark$

subtree at X contains 0 internal nodes \checkmark

Inductive Step: let height of X be h and $bh(x)=b$

Some child of X has height $h-1$ and black height b or $b-1$ (if child is black)

By IH, each child has $\geq 2^{bh(x)-1}-1$ internal nodes

\therefore subtree rooted at X contains $\geq 2 \cdot (2^{bh(x)-1}-1) + 1 = 2^{bh(x)} - 1$
for itself

Proof for theorem: let h and b be height and black-height, desp., of the root

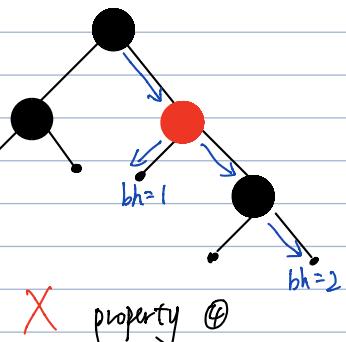
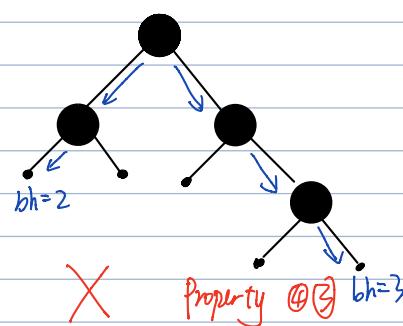
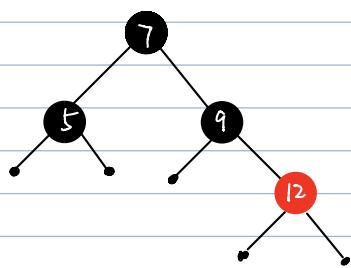
$$n \geq 2^b - 1 \geq 2^{\frac{h}{2}} - 1$$

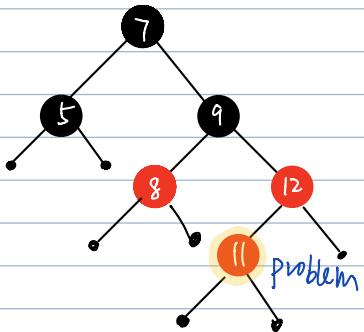
△ △
lemma 2 lemma 1

$$n+1 \geq 2^{\frac{h}{2}}$$

$$\lg(n+1) \geq \frac{h}{2} \quad h \leq 2 \lg(n+1)$$

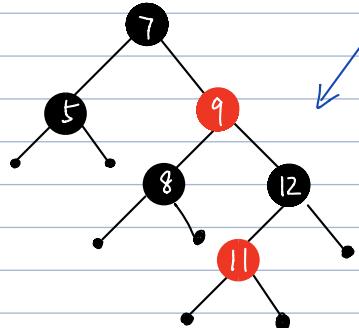
Problem of insertion into RB trees





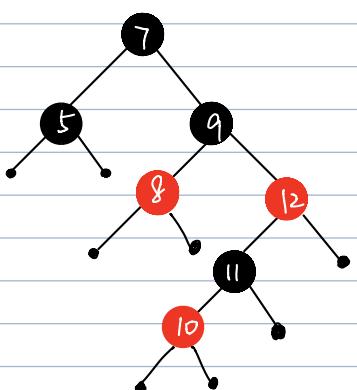
- insert(8) \Rightarrow left for 9: color red
- insert(11) \Rightarrow left for 12 \Rightarrow color red $\Rightarrow \text{X } \textcircled{3}$
- color black $\Rightarrow \text{X } \textcircled{4}$

recolor the tree



insert(10) : left child of 11

\leftarrow recolor? not enough



$$\text{why } bh(7) = \begin{matrix} 2 \\ 3 \\ 3 \end{matrix} \quad (7 \rightarrow 5 \rightarrow \text{nil})$$

$$(7-9-12-11-10-\text{nil})$$

X property $\textcircled{4}$

Goal: Restructure Tree in $O(\lg n)$ time



■ Insert Operation

IDEA: insert x using BST-TREE-INSERT(x)

尽量把这种冲突推到树顶

- color x red
- parent might be red \Rightarrow violate $\textcircled{3}$
- move violation up tree by recoloring \leftarrow 高上层再处理
- because tree imbalance can't satisfy $\textcircled{4}$ (black-height) without violation $\textcircled{3}$ (recoloring)
- must change tree structure in $O(\lg n)$ time

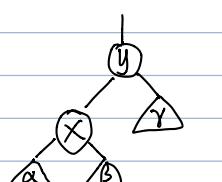
Goal

⇒ ■ Rotation: (CLRS 13.2)

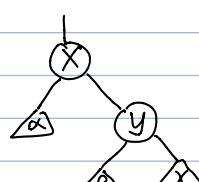
- basic operation for changing tree structure
- key idea: Rotation does not change BST properties



- change local pointer structure



Right-rotate(y)
left-rotate(x)



- Precenve BST Property
 $x \leq a \leq b \leq y \leq c$

- pseudo code

Right - Rotation (T, y)

```

01  x = y.left           // set x
02  y.left = x.right
03  if x.right != nil
04    x.right.parent = y
05  x.parent = y.parent
06  if y.parent == nil
07    root[T] = x
08  else if y == y.parent.right
09    y.parent.right = x
10  else
11    y.parent.left = x
12  x.left = y
13  y.parent = x

```

LEFT-ROTATE(T, x)

```

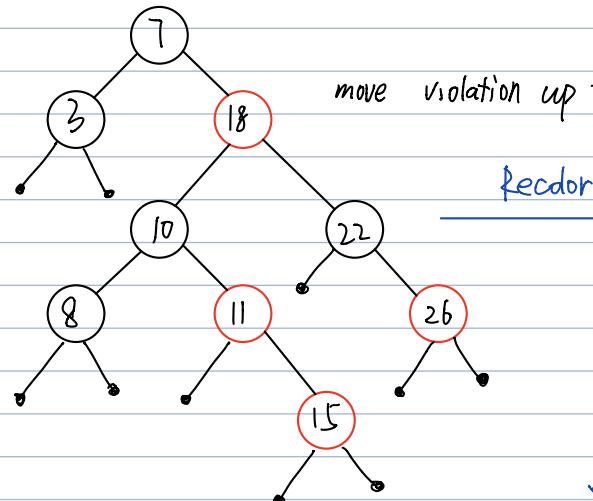
1  y = x.right
2  x.right = y.left
3  if y.left != T.nil
4    y.left.p = x
5  y.p = x.p
6  if x.p == T.nil
7    T.root = y
8  elseif x == x.p.left
9    x.p.left = y
10 else x.p.right = y
11 y.left = x
12 x.p = y

```

$\//$ put x on y 's left

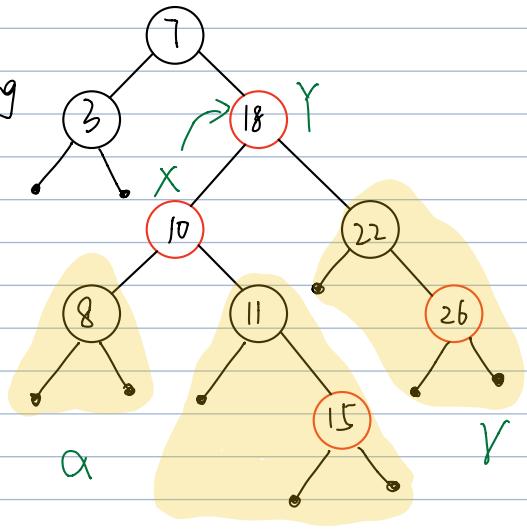
$T = O(1)$

Example Insert (15)



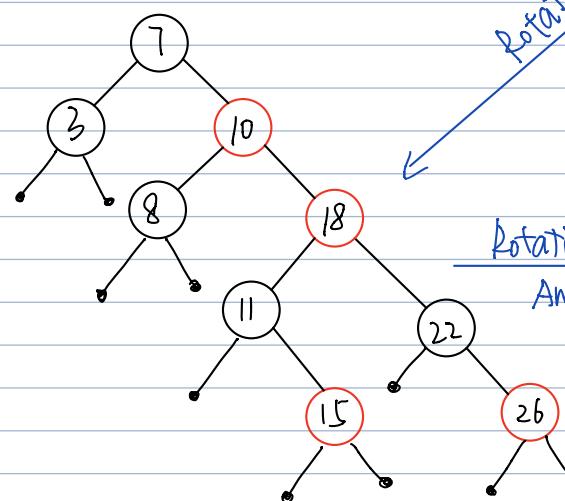
move violation up tree by recoloring

recolor



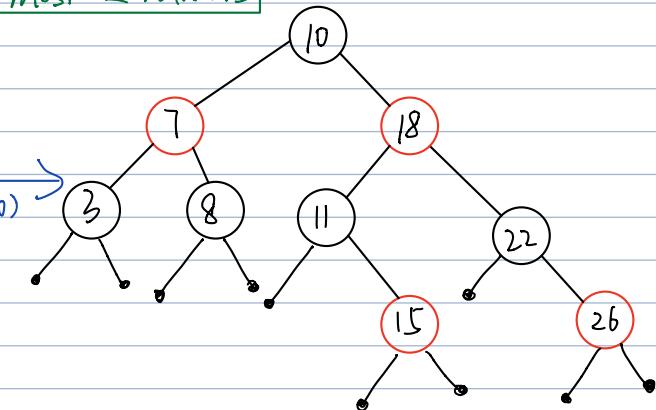
Rotation Right (8)

We only need
at most 2 rotations



Rotation left (7)

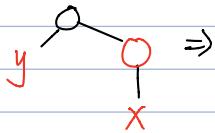
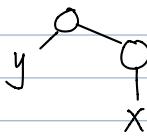
And Recolor(7,10)



• Pseudo code

```

01 BST-insert (T, X)
02 color X red
03 while (X ≠ T.root) and (X.parent.color = red):
04     if X.p == X.p.p.r // X is parent & X is grandparent to T2
        y = X.p.p.l // set y
        if y.color == red:
            X.color = black
            y.color = black
            X.p.p.color = red
            X = X.p.p
        else if X == X.p.left
            X = X.p
            then right-rotate (T, X)
            X.p.color = black
            X.p.p.color = red
            left-rotate (T, X.p.p)
        else
            same as "then" clause, but swap "left" and "right"
    
```



Local Search

- what is local search?
 - a simple heuristic method for solving hard optimization problems
 - solving optimization problem:
 - find a solution minimizing/maximizing an **objective function** in solution space

Local solution method

- Start with an arbitrary solution S for a given problem
- Let $N(S)$ be set of solutions in "neighbour" of S obtained from local moves/changes
- If there is a "better" solution $S' \in N(S)$, switch to S'
- Else stop and output S

"Better": for minimization, S' is better than S if $\text{val}(S') < \text{val}(S)$
for maximization, S' is better than S if $\text{val}(S') > \text{val}(S)$

Example: Local search for 3-SAT

- Given: A 3-SAT instance (n variable, m clause, variables in each clause distinct)
- Let $x \in X = \text{set of candidate solution to problem}$
- Starting Solution:
 - Example: any random assignment $x \in X = \{0, 1\}^n$
- For each $x \in X$ specify $y \in X$ are its neighbour
 - Ex: x, y are neighbouring variable assignments \Leftrightarrow differ in value of single variable
 - better solution in neighbourhood has fewer unsatisfied clauses

Key Questions For Local Search Algorithm

- how long will it run?
 - neighborhood of each solution small
 - need to find local optima in polynomial # of steps.
- how good is local optimal solution?
 - local optima is "nearly as good" as global optima

Local Search Outcomes

- Provably optimal solution

Ex: Network Flow (Ford - Fulkerson)

- start with zero flow: $f=0$
- "local neighbourhood": set of all flows obtained by augmenting current flow along augmented path
- "better": higher flow value
- provably optimal: $O(E^2 V)$ Time
EDMONDS-KARP

local search terminates with local optimum that is a good approximation

- Bound ratio between optimal solutions and solutions that local search returns

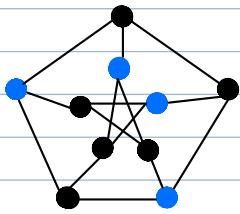
- Example : Maximum Cut

- input : undirected connected graph $G = (V, E)$

- goal : a cut $(S, V-S)$ such that total # edges crossing cut is maximum

$$\text{Maximize } |E'| \quad E' = \{(u, v) \in E : u \in S, v \in V-S\}$$

example



$$|V|=10$$

$$|\text{max-cut}|=12$$

$$\text{max-cut: } S = \{\text{blue}\}$$

$$V-S = \{\text{black}\}$$

(with weight \Rightarrow even harder) \Rightarrow KARP's 21 NP-complete problems.

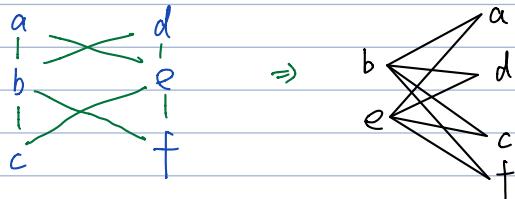
contact with: min-cut \Rightarrow karger: $\tilde{O}(V^2)$ Ford-Fu: $\tilde{O}(E^2V)$

Max-cut : NP-Hard

(reduction from 3-SAT, or. max-independent set)

\Rightarrow computationally tractable case: Bipartite Graphs

Ex:



$$|V|=6 \quad S = \{b, e\}$$

$$|E|=8 \quad V-S = \{a, d, c, f\} \quad |\text{max-cut}| = 8$$

Solve : run BFS in $O(V+E)$ \Rightarrow run BFS on arbitrary vertex

odd layers : $V-S$

even layers : S



Local Search Algorithm For Max Cut

START WITH ANY NONEMPTY SUBSET S OF V :

THIS GIVES A CUT $(S, V-S)$

WHAT IS A SMALL LOCAL CHANGE TO INCREASE $\#$ OF EDGES CROSSING CUT?

\rightarrow WHILE THERE IS A VERTEX V SUCH THAT MOVING V TO OTHER SIDE

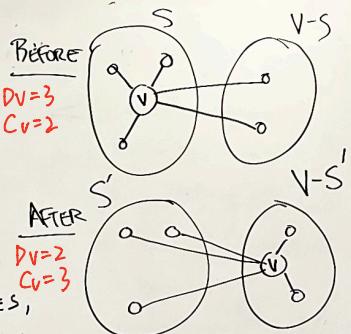
OF CUT INCREASES # CROSSING EDGES,
MOVE V TO OTHER SIDE

REPEAT STEP UNTIL NO SUCH V EXISTS

WHEN : WHEN V HAS MORE INCIDENT EDGES WITHIN S

THAN CROSSING EDGES, i.e. $|\{(u, v) \in E : u \in S, v \in S\}| > |\{(u, v) \in E : u \in S, v \in V-S\}|$

(OR VICE VERSA), MOVE V



真詮不動
都變動
F T

Notation : For a cut $(S, V-S)$ and vertex V . define

$C_V(S, V-S) = \# \text{ edges incident on } V \text{ that cross } (S, V-S)$

$D_V(S, V-S) = \# \text{ edges incident on } V \text{ that don't cross } (S, V-S)$

pseudo code:

Local - Search - Maxcut

01 let $(S, V-S)$ be arbitrary cut of G

02 while there is a vertex V with $D_V(S, V-S) > C_V(S, V-S)$

03 move v to other side of cut // increase # of crossing edges by $D_V - C_V > 0$

04 return find cut $(S, V-S)$

Running Time: THM: local-search-maxcut terminates in $O(n^2)$. $n = \#$ of vertices

Proof: each iteration \uparrow # edges crossing cut at least by 1

since # edges crossing cut $\leq |E| = O(n^2)$

the ALGO terminates after $O(n^2)$

How good?

$$\frac{\#(\text{APX})}{\#(\text{OPT})} \geq \frac{1}{2}$$

\Rightarrow proof:

Theorem: Local-search-maxcut always outputs a cut in which # crossing edges is at least half of the maximum possible

Proof: cut at local opt $(S, V-S)$. arbitrary vertex v .

at local opt, # edges crossing cut must be greater than # edges not crossing cut.

for every vertex v , # v 's edges not crossing cut $<$ # v 's edges crossing
summing over all v , $\sum_{v \in V} \# v$'s edges not crossing $< \sum_{v \in V} \# v$'s edges crossing
counts each not crossing edge twice
 $+ 2 \cdot \# \text{edges crossing}$

$2|E| \leq 4 \cdot \# \text{edges crossing}$
 $|E| \leq 2 \cdot \# \text{edges crossing}$
 $\# \text{edges crossing} - |\text{APX}| > \frac{|E|}{2}$

2-approximation alg

Hash map m_1 : (key: v . value: $D_v - C_v$)

Hash map m_2 : (key: $D_v - C_v$ value: HashSet(v))

random assignment: $V \leftarrow$ hash set

for all v :

$m_1[v] = D_v - C_v$

if $D_v - C_v > 0$:

$m_2[D_v - C_v].add(v)$

while (m_2 is not empty):

key = any key in m_2

$v =$ any element in $m_2[\text{key}]$

for $u \in v.\text{adj}$:

if $u.v$ in same part:

$\text{val} = m_1[u]$

$m_1[u] += 1$

if $\text{val} > 0$

$m_2[\text{val}].remove(u)$

$m_2[\text{val}+1].add(u)$

else if $\text{val} == 0$:

$m_2[1].add(u)$

rand assignment $V \leftarrow$ hash set

heap $\leftarrow v: D_v - C_v$

while $\text{heap}.\text{max} > 0$

$v.c = \text{heap.pop}$

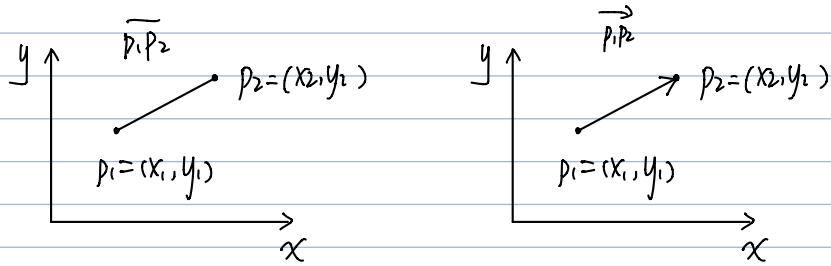
for $v \in v.\text{adj}$

if $u.v$ are in same side

```
| else (u,v not in same part)
| val = m1[u]
| m1[u] -= 1
| if val > 1
|   m2[val] - remove (u)
|   m2[val-1].remove(u)
|
| elif (val == 1) :
|   m2[val].remove(u)
|
val2 = m1[v]
m1[v] -= val2
m2[val2].remove(v)
```

Week 7 Computational Geometry

- Algorithm and data structures for geometric problems
- input / output / intermediate data: Geometric objects represented by real numbers
- work in 2 dimensions i.e. in plane
- Basic Definitions
 - points are (x, y) with $x, y \in \mathbb{R}$
 - line segment is a portion of a straight line between 2 points $P_1(x_1, y_1), P_2(x_2, y_2)$

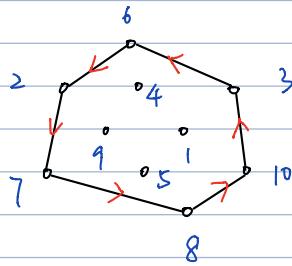


Convex hull

problem: give a set P of n points in the plane

Goal: find the smallest convex polygon containing $P \Rightarrow$ (given as array with 2 fields x-coordinate, y-coordinate)

Example:



output : indices in some order

[7, 8, 10, 3, 6, 2]

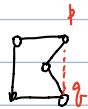
6 vertices on boundary in CCW

Definition :

① convex hull of P : smallest convex polygon containing all points of P

② polygon: a circular sequence of points (vertices) and line segments (edges)

③ convex: X is convex : $p \in X, q \in X$ then $\overline{pq} \in X$ 反例



④ smallest: (inclusion minimal): for every convex set X if $P \subseteq X$ then $\text{conv}(P) \subseteq X$

smallest: any convex proper subset of convex hull excludes at least one point of P

任何凸多邊形的子集如果不是整個都漏掉了一個點，

Representation : represent CH (convex hull)

As a circular doubly linked list of vertices in CCW order, if i -th point is a vertex of CH, next [i] is index of next vertex CCW and pred [i] index of next vertex CW, otherwise next [i] = pred [i] = 0

Q: How tell if a point is on one side or another of a line?

Convex Hull :

$n=1$ what is CH?

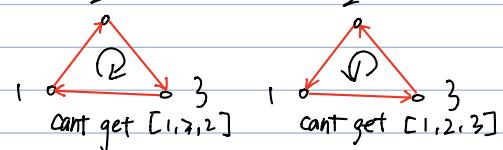
$n=2$

$n=3$

point
line segment
triangle

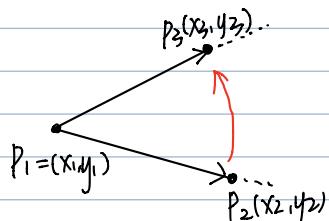
output : []

output : [1,2] or [2,1]



output: 6개면积

Given 3 points, what order are they in?
which line has bigger slope?



Goal: answer in $O(1)$ time

- use only $+$, $-$, \times and comparison
- avoid division and trig functions

$$\frac{y_2 - y_1}{x_2 - x_1} < \frac{y_3 - y_1}{x_3 - x_1} \Rightarrow \underbrace{(x_2 - x_1)(y_3 - y_1) - (x_3 - x_1)(y_2 - y_1)}_{\leq 0} > 0 \text{ CCW}$$

$$< 0 \text{ CW}$$

$$= 0 \text{ collinear}$$

orientation test
or

CCW test

$$\begin{vmatrix} x_2 - x_1 & y_2 - y_1 \\ x_3 - x_1 & y_3 - y_1 \end{vmatrix} > 0 \text{ CCW}$$

or

$$\begin{vmatrix} 1 & x_1 & y_1 \\ 1 & x_2 & y_2 \\ 1 & x_3 & y_3 \end{vmatrix} > 0 \text{ CCW}$$



orientation test is 2-Dim analogue to comparisons

$$a > b ? \quad \begin{vmatrix} 1 & a \\ 1 & b \end{vmatrix} < 0 ?$$

Jarvis March

① find left most point $P_L \Rightarrow$ all other

② Repeat

find next hull vertex

③ until next point is P_L

01 $L = 1$

02 for $i = 2$ to n

03 if $x[i] < x[L]$

04 $L = i$

05 $p = L$

06 repeat

07 $q = p + 1$ // make sure $p \neq q$

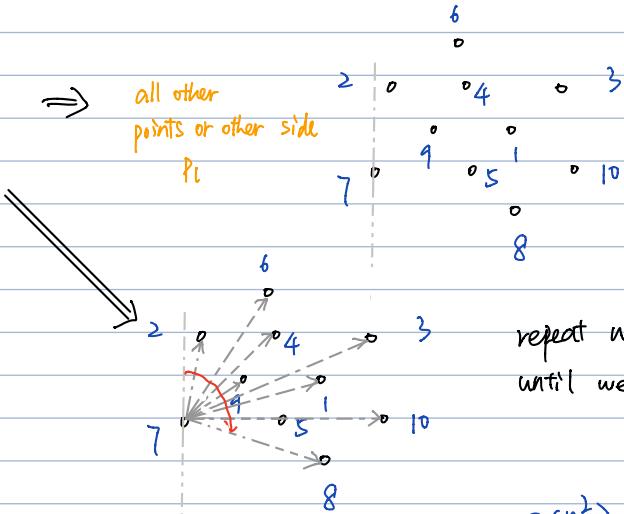
08 for $i = 2$ to n // search for a point q st. $CCW(p, i, q)$ is CCW for all points i

09 if $CCW(p, i, q)$ // if i is more CCW than current q . update q .

10 $q = i$

11 $next[p] = q$: $prev[q] = p$

12 until $p = L$



$O(n^2)$

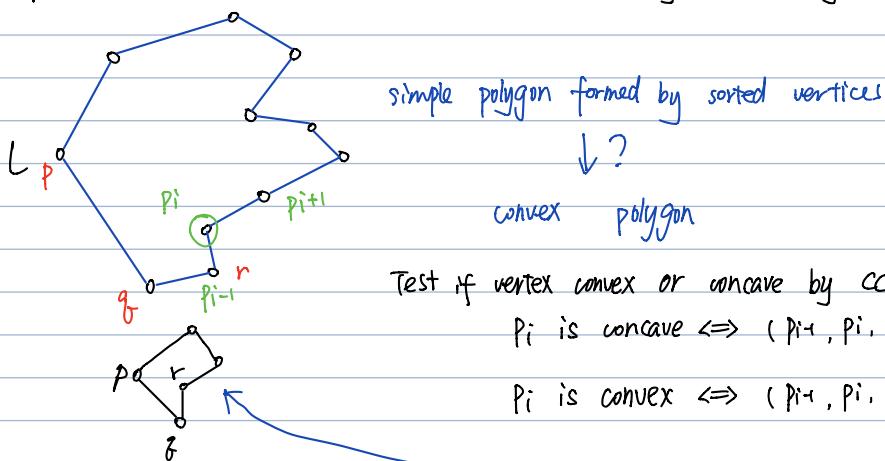
找到一个 q 让其他所有点在 p 后面且逆时针

for $i = 2$ to n // search for a point q st. $CCW(p, i, q)$ is CCW for all points i

if $CCW(p, i, q)$ // if i is more CCW than current q . update q .

Graham's Scan

1. Find left most point $P_1 = L$ (by x-coordinate)
2. Sort $P_2 \dots P_n$ around P_1 (using CCW test instead of comparison)
3. once points sorted, connect in CCW order starting and ending at $P_1=L$



Test if vertex convex or concave by CCW test:

P_i is concave $\Leftrightarrow (P_{i+1}, P_i, P_{i-1})$ is CW

P_i is convex $\Leftrightarrow (P_{i+1}, P_i, P_{i-1})$ is CCW

Three-penny-algorithm

Repeat: if p, q, r is CCW (q is convex)
move penny p forward to successor r
if p, q, r is CW (q is concave)
remove vertex q add edge Pr
move middle penny back to predecessor

Pseudo Code

```

01  $P_1 = \text{left most point of } P$ 
02  $[P_2 \dots P_n] = \text{remaining points of } P_1 \text{ sorted in CCW order around } P_1$ 
03 if  $n < 2$ 
04 return CH is empty
05 else  $S = \text{empty stack}$ 
06 push  $(P_1, S)$ 
07 push  $(P_2, S)$ 
08 push  $(P_3, S)$ 
09 for  $i = 4 \text{ to } n$ 
10   while CW(next-to-top(S), top(S),  $P_i$ )
11     pop(S)
12   push  $(P_i, S)$ 
13 return S
    
```

$O(n \lg n)$

How many times iteration?

2 cases

1. How many iterations when move penny forward $n-2$
2. How many times move middle penny back?
How many times delete edge and add new? $\Rightarrow n-h$

Total: $2n-h-2$

h : size of (CH vertices)

Line segment intersection

Input: set of n line segments, each specified by x-and-y-coordinate
 $\Rightarrow 4n$ real numbers

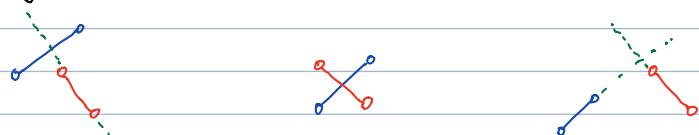
Goal: want to know if any ≥ 2 intersect

Assumptions: General position; i.e. No degenerate cases



这三种情况

Two segments $n=2$ How to tell whether two line segments intersect?



Simple condition: 2 segments \overrightarrow{ab} and \overrightarrow{cd} intersect IFF the end points a and b are on opposite side of line \overleftrightarrow{cd} and end points c and d are on opposite side of line \overleftrightarrow{ab}

Two segment intersect (a, b, c, d)

- 01 if $CCW(a, b, c) = CCW(a, b, d)$
- 02 return False
- 03 if $CCW(a, c, d) = CCW(b, c, d)$
- 04 return False
- 05 return True

Sweep-line Algorithm

1-st: give line segment a unique label

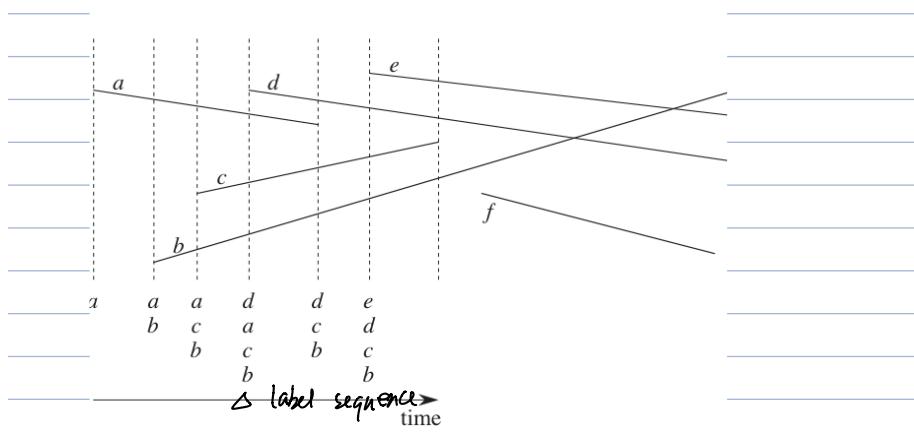
next: image a vertical line sweeping across segment from left to right

at each position of sweep line, look at sorted sequence (by y-coordinate) of labels of segments that line hits

only times this sequence can change is when sweep line passes through endpoints or passes intersection points

ANY-SEGMENTS-INTERSECT(S)

- 1 $T = \emptyset$
- 2 sort the endpoints of the segments in S from left to right,
breaking ties by putting left endpoints before right endpoints
and breaking further ties by putting points with lower
 y -coordinates first
- 3 for each point p in the sorted list of endpoints
 - 4 if p is the left endpoint of a segment s
 - 5 INSERT(T, s)
 - 6 if (ABOVE(T, s) exists and intersects s)
or (BELOW(T, s) exists and intersects s)
 - 7 return TRUE
 - 8 if p is the right endpoint of a segment s
 - 9 if both ABOVE(T, s) and BELOW(T, s) exist
and ABOVE(T, s) intersects BELOW(T, s)
 - 10 return TRUE
 - 11 DELETE(T, s)
 - 12 return FALSE



Any Intersections

```
Any Intersections (S[1...n] // x,y coordinates)
01 sort the end points of S from left to right
02 create an empty label sequence
03 for i = 1 to 2n
04     l = label[i]
05     if isLeft[i]
06         insert(l)
07         if intersect(S[l], S[successor(l)])
08             return true
09         if intersect(S[l], S[predecessor(l)])
10             return true
11     else (isRight[i])
12         if intersect(S[successor(l)], S[predecessor(l)])
13             return true
14         delete(label[i])
15 return False
```

Week 8 Advanced Flow

standard flow network

Input: Directed Graph $G = (V, E)$
 Capacity $c(e) > 0$ for each node $e \in E$ two vertices $s, t \in V$

Define: Feasible Flow

$$f: E \rightarrow \mathbb{R} \text{ st. } 0 \leq f(e) \leq c(e) \quad \forall e \in E \quad (\text{capacity})$$

$$\sum_{u \in V} f(u, v) = \sum_{w \in V} f(v, w) \quad \forall v \in V \quad v \neq s, t \quad (\text{conservation})$$

Goal: Maximize total flow out of source vertex s .

$$\text{Maximize : } |f| = \sum_w f(s, w) - \sum_u f(u, s)$$

How to modify?

Balance $b(v)$ for every vertex $v \neq s, t$

Define: Feasible Flow

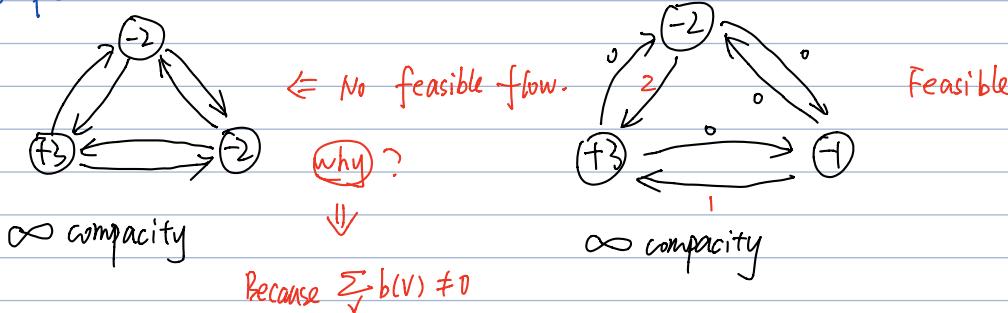
$$f: E \rightarrow \mathbb{R} \text{ st. } 0 \leq f(e) \leq c(e) \quad \forall e \in E \quad (\text{capacity})$$

$$\sum_{w \in V} f(v, w) - \sum_{u \in V} f(u, v) = b(v) \quad \forall v \in V \quad (\text{conservation})$$

Goal: Find a feasible flow if one exists

(从什么问题修改为流量平衡问题)

Example:

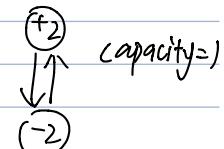
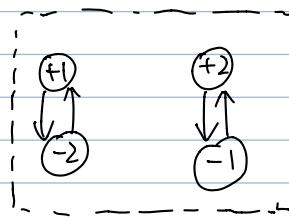


$$f \text{ is feasible} \Rightarrow \sum_v b(v) = \sum_j \left(\sum_u f(u, v) - \sum_w f(v, w) \right)$$

flow along all edges into V flow away all edges out of V

$$= \sum_{u \in V} f(u, v) - \sum_{v \in W} f(v, w) = 0$$

$\sum_v b(v) = 0$: necessary but not sufficient Example:



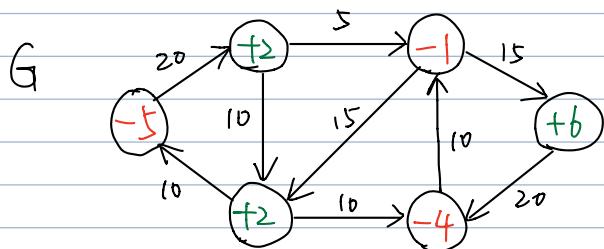
Intuition:

$b(v) > 0 \Rightarrow$ demand, consumption

$b(v) < 0 \Rightarrow$ supply, production

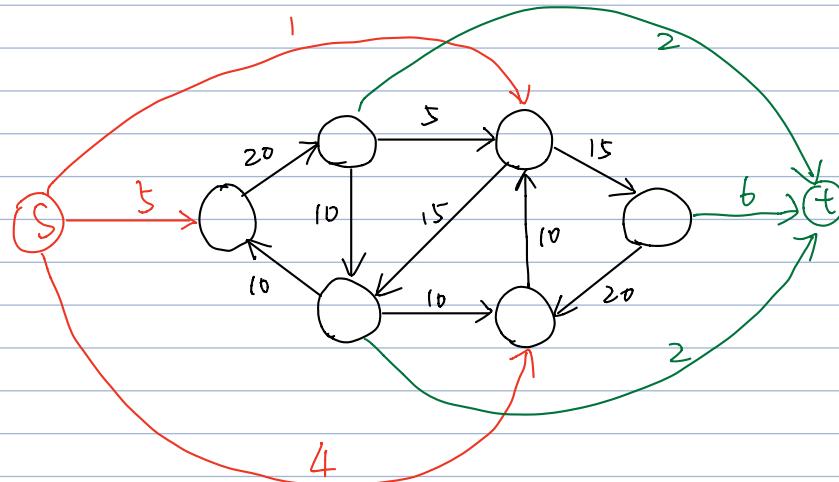
Solve by reducing to standard max flow

Ex.



Given: $G = (V, E)$, $c(e)$ for all $e \in E$, $b(v)$ for all $v \in V$

G'



construct $G' = (V', E')$
 $V' = V \cup \{s, t\}$

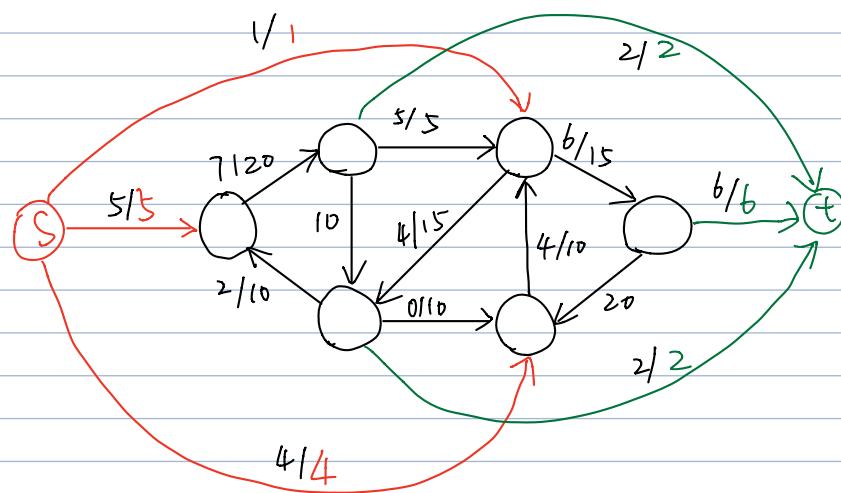
$$E' = E \cup \{(s, v) : b(v) < 0\} \\ \cup \{(v, t) : b(v) > 0\}$$

$c(e)$ for all $e \in E$.

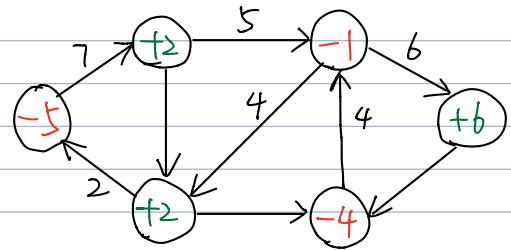
$$c(s, v) = -b(v) \quad | \text{ 取绝对值} \\ c(v, t) = b(v)$$

If there is a saturating flow f' in $G' \rightarrow$ feasible flow f in G .

what is: saturating flow? saturates all edges leaving s in G' (and entering t)
 滿流: means: flow = capacity



(dismiss the capacity)



We found a saturating flow

f' is saturating $\Rightarrow f'|_E = f$ Edmonds-Karp $O(E^2V)$

Week 8 Linear Programming (LP)

- solve optimization problems
- Linear objective function
- Linear constraints

Given:
n variables over \mathbb{R}
m linear equations / inequalities
Linear objective function

Goal: Assign real values to variables satisfy linear inequalities

maximum / minimize linear objective function

Example: Profit Maximization

2 produce A, B

How much of each produce to maximize profit?

- each unit of A: profit \$1
- B: \$b

constraint: demand:

≤ 200 units of A per day

≤ 300 units of B per day

Total supply ≤ 400 units per day

Associated LP: $x_1 = \#$ of units of A

$x_2 = \#$ of units of B

objective function: $\max x_1 + b x_2$

constraints: $x_1 \leq 200$

$x_2 \leq 300$

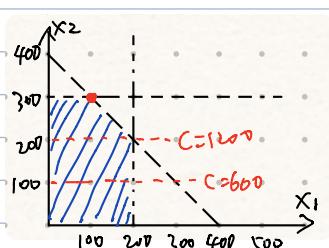
$x_1 + x_2 \leq 400$

$x_1, x_2 \geq 0$

Geometry of LP

five constraints = linear inequality half-space in 2-d plane (2 variable x_1, x_2)

Set of all feasible solutions of LP = points (x_1, x_2) that satisfy all constraints.



feasible region

Goal: find max C where $x_1 + 6x_2 = C$ intersects feasible region
or point (x_1, x_2) in feasible region with max C is optimal solution

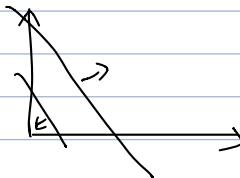
at point $(x_1, x_2) = (100, 300)$ profit = 1900 = optimum C
satisfy all constraints

When is there no optimum solution?

① LP infeasible = feasible region empty

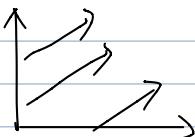
(constraints too tight)

$$\begin{aligned} \max \quad & x_1 - x_2 \\ \text{subject to} \quad & 2x_1 + x_2 \leq 1 \\ & x_1 + x_2 \geq 2 \\ & x_1, x_2 \geq 0 \end{aligned}$$



② unbounded : $\max x_1 + x_2$

$$x_1, x_2 \geq 0$$



Algorithm for solving LPs

● simplex algorithm (1947)

- worst case running time is exponential
- typically linear in size of problem

● polynomial time algorithm

- ellipsoid method, interior-point method

(1972)

(1967)

Simplex Algorithm: (= local search)

- start at a vertex of feasible polygon in other case (0,0)
- search for an adjacent vertex (connected by an edge) of better objective value
- repeat until reach a vertex with no better neighbour
- return solutions as optimum $(0,0) \rightarrow (200,0) \rightarrow (200,200) \rightarrow (100,300)$

■ LP in n dimensions

- constraints = half-space
- feasible region = intersection of n half spaces
= convex polyhedron in n dimensions

- vertices = corners of feasible regions
- = intersections of sides of feasible region
- = points satisfying n inequalities with equality satisfy remaining $m - n$ inequalities
- $\binom{m+n}{n}$ vertices

neighbouring vertices share $n-1$ of their equality constraints.

How do we verify simplex "solution is optimal?"

- one way: check by checking all neighbouring vertices
- alternative: LP duality

$$x_1 = \# \text{ of units of A}$$

$$x_2 = \# \text{ of units of B}$$

LP duality

Notice: 1st constraint + 6 · 2nd constraint
 \Rightarrow upper bound

objective function: $\max x_1 + 6x_2$

$$x_1 + 6x_2 \leq 2000$$

$$\begin{aligned} \text{constraints: } x_1 &\leq 200 \\ x_2 &\leq 300 \\ x_1 + x_2 &\leq 400 \\ x_1, x_2 &\geq 0 \end{aligned}$$

Can we find a tight upper bound systematically

\Rightarrow LP duality: let $y = (y_1, y_2, y_3)$ with $y_1, y_2, y_3 \geq 0$

Multipl/ter	constraint
y_1	$x_1 \leq 200$
y_2	$x_2 \leq 300$
y_3	$x_1 + x_2 \leq 400$

$$\text{Then } \sum y_i \cdot \text{constraint}_i = (y_1 + y_2)x_1 + (y_2 + y_3)x_2 \leq \underbrace{200y_1 + 300y_2 + 400y_3}_{\text{upper bound}}$$

LHS: $x_1 + 6x_2$
OBJ function

RHS: upper bound on OPT solution
minimize to get tight bound

$$\Rightarrow y_1 + y_3 \geq 1$$

$$y_2 + y_3 \geq b$$

original LP:

Dual LP: $\min 200y_1 + 300y_2 + 400y_3$

constraint $y_1 + y_3 \geq 1$

$y_2 + y_3 \geq b$

$y_1 + y_2 + y_3 \geq 0$

objective function: $\max x_1 + 6x_2$

constraints: $x_1 \leq 200$

$x_2 \leq 300$

$x_1 + x_2 \leq 400$

$x_1, x_2 \geq 0$

Dual LP: $(y_1, y_2, y_3) = (0, 5, 1) \Rightarrow 1900$

■ Standard Form for LPs

Variables: X_1, X_2, \dots, X_n

OBJ function: $\max \sum_{j=1}^n c_j x_j$

constraints: $\sum_{j=1}^n a_{ij} x_j \leq b_j$ for each $i=1, \dots, m$

$x_j \geq 0$ for each $j=1 \dots n$

$$\text{LA } x = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}$$

$$c = \begin{pmatrix} c_1 \\ \vdots \\ c_n \end{pmatrix}$$

mxn constraint matrix $A \in \mathbb{R}^{m \times n}$

$$A = \begin{pmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{m1} & \cdots & a_{mn} \end{pmatrix}$$

$$b = \begin{pmatrix} b_1 \\ \vdots \\ b_m \end{pmatrix}$$

original LP

Dual LP

$$\max c^T x$$

$$Ax \leq b$$

$$x \geq 0$$

$$\min b^T y$$

$$A^T y \geq c$$

$$y \geq 0$$

■ Ex: maximum flow as LP

capacity:

input: directed graph $G=(V, E)$, two specified vertices $s, t \in V$, every edge has $c(e) > 0$

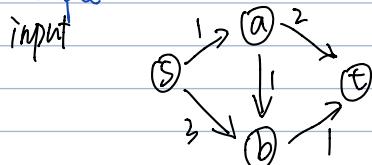
output: flow f specified by $f(e)$ for every edge $e \in E$

$$\text{Maximizes: size}(f) = |f| = \sum_V f(s, v)$$

constraints: for all $e \in E$ $0 \leq f(e) \leq c(e)$

$$\text{for every } V \subseteq V - \{s, t\}: \sum_{u: (u, v) \in E} f(u, v) = \sum_{z: (v, z) \in E} f(v, z)$$

Example:



LP: define flow variables $f(i, j)$ on each edge $(i, j) \in E$

Obj function: $\max f(s, a) + f(s, b)$

constraint: $0 \leq f(s, a) \leq 1$

$$0 \leq f(s, b) \leq 3$$

$$0 \leq f(a, b) \leq 1$$

$$f(s, a) = f(a, b) + f(a, t)$$

$$f(s, b) + f(a, b) = f(b, t)$$

variable: $f(e)$ for each edge

obj function: $\max \sum_{(s, v) \in E} f(s, v)$

constraints: for all $e \in E$, $0 \leq f(e) \leq c(e)$

for every $V \subseteq V - \{s, t\}$:

$$\sum_{u: (u, v) \in E} f(u, v) = \sum_{z: (v, z) \in E} f(v, z)$$

$$0 \leq f(a|t) \leq 2$$

$$0 \leq f(b|t) \leq 1$$

Week 9 String Matching Problem (CLRS CH32)

Given 2 strings A text $T[1..n]$ and a pattern $P[1..m]$, is P a substring of T ? (and if so, where?)

Ex: P : what type in Google
 T : entire web

T is big, typically, and P is small. Treat both T and P as arrays.

$T[1..n]$ - Text $m \leq n$
 $P[1..m]$ - Pattern

Elements \in Finite Alphabet Σ
elements: character

Obvious Way: Check each substring of T of length M
And see if matches P



switch one over and so on

More formally, P occurs with shift s in T .
if $0 \leq s \leq n-m$ and $T[s+1, \dots, s+m] = P[1..m]$

① Simple-string-match ($T[1..n]$, $P[1..m]$) // find small shifts or report no shift

- 01 for $s=0$ to $n-m$:
- 02 : equal = true
- 03 : $i=1$
- 04 : while equal == True and $i \leq m$
- 05 : : if $T[s+i] \neq P[i]$
- 06 : : equal = false
- 07 : : else
- 08 : : $i=i+1$
- 09 : : if equal == true:
- 10 : : return s
- 11 return None

Running Time: $O((n-m) \cdot m) = O(mn)$

#iteration at most n
compares

Worst case actually achieved: search for pattern: aaa...aab in text aa...a } very slow

Need at least $\alpha(nm)$, could be $O(n^k)$, can do better? yes

- IDEA: have "sliding window" of size m
 characters $T[S_1 \dots S_m]$ is "visible"
 what to know: is it the same as $P[1 \dots m]$?

- IDEA:
- compute hash (P)
 - for each length- m window of T , compute hash ($T[S_1 \dots S_m]$) and compute hash ($P[1 \dots m]$)
 - if $=$: check to see if strings really match
 if \neq : move on to next S

compute hash values in $O(1)$ times v.s. $O(m)$ comparisons of characters

How many iteration we need to compute hash? $\Rightarrow n-m \Rightarrow$ #iterations stay the same

- Now we want the hash function such that we can go from one window to another in $O(1)$ time

\Rightarrow RABIN-KARP

- ① idea one: treat string as integers

Given: $|\Sigma| = d$ treat each string of characters in Σ as an integers in Radix- d notation.

Ex1. $\Sigma = \{0, 1, \dots, 9\}$ $d=10$. string "274" has numeric value: $2 \cdot d^2 + 7 \cdot d^1 + 4 \cdot d^0 = 274$

Ex2. $\Sigma = \{a, b\}$ $d=2$ mapping $a \rightarrow 0$ $b \rightarrow 1$. Then string "bab" = $1 \cdot d^2 + 0 \cdot d^1 + 1 \cdot d^0 = 5$

Ex3. DNA: $\Sigma = \{A, C, G, T\}$, $A \rightarrow 0$, $B \rightarrow 1$, $C \rightarrow 2$, $D \rightarrow 3$ $\Rightarrow TAT = 51$

C 1

G 2

T 3

RABIN-KARP ALGO

- need to compute hash values for $P[1 \dots m]$, $T[S_1 \dots S_m]$ for $1 \dots n-m$.
- Need an $O(1)$ update hash value when slide window one position to right.

- $P[1 \dots m] \Rightarrow \Theta(m)$ time (use Horner's Rule)

$$x = \sum_{i=1}^m a_i d^{m-i}$$

$$\downarrow$$

$$P = p[m] + d(p[m-1] + d(p[m-2] + \dots + d(p[2] + d.p[1])))$$

- Given substring $T[1 \dots m]$ let t denote its decimal value, to computed in $\Theta(m)$

- How to do shift:

- ② second key idea:

i-1 window: $[i \dots i+m-1]$

i window: $[i+1 \dots i+m]$

$$t_i = (t_{i-1} - T[i] \cdot d^{m-1}) \cdot d + T[i+m]$$

Rabin-Karp - version 1 ($T[1 \dots n]$, $P[1 \dots m]$, $|\Sigma|=d$)

01 $h = d^{m-1}$

02 $p = 0$ // pattern value

03 $t = 0$ // substring value

04 for $i = 1$ to m :

05 $p = 10 \cdot p + P[i]$

06 $t = 10 \cdot t + T[i]$

07 for $s = 0$ to $n-m$:

08 if $p = t$:

} Horner's Rule

```

09     return S
10     t = 10(t - h · T[s+1]) + T[s+m+1]
11     return None

```

Running Time: $O(m) + O(m) + O(n-m) = O(n+m) = O(n)$

\uparrow \uparrow \uparrow $m \leq n$
 line 1 line 4-6 line 7-10

$O(n)$ worst case

Bounds # of arithmetic operations, since m may be large cannot assume $O(1)$

• Rabin-Karp : Finger Printing

One simple change to make efficient $d = |\Sigma|$
 perform all arithmetic operations some prime q .

Choose q so that $d \cdot q$ fits into a standard integer variable

- $P \bmod q$ fingerprints of $P[1 \dots m]$
- $T_s \bmod q$ fingerprints of $T[s+1 \dots s+m]$
- $T_{s+1} \bmod q = (d(T_s - T[s+1] \cdot h) + T[s+1+m]) \bmod q$

↓

New Problem: collision: $T_i = P \quad T[i+1 \dots i+m] \neq P[1 \dots m]$

- ↳ ① if $(P \bmod q) \neq (T_s \bmod q)$ $P \neq T_s$
 ② if $(P \bmod q) = (T_s \bmod q)$
 $\Rightarrow P \neq T_s$ differ by $k \cdot q$, $k \in \mathbb{Z}$
 \Rightarrow if $P \neq T_s \Rightarrow$ False Match

e.g. $T = \underline{\underline{3 \ 1 \ 4 \ 1 \ 5}} \ 2 \ \underline{\underline{6 \ 7 \ 3 \ 9 \ 9 \ 2}}$
 $31415 \bmod 13 = 17$ $67399 \bmod 13 = 17$

$$q = 13 \quad d = 10$$

So we need brute force comparison $O(m)$

Rabin-Karp - version 2 ($T[1 \dots n], P[1 \dots m], |\Sigma|=d$)

01 $h = d^{m-1} \bmod q$
 02 $p = 0$ // pattern value
 03 $t = 0$ // substring value
 04 for $i=1$ to $n-m$:
 05 $p = (10 \cdot p + P[i]) \bmod q$
 06 $t = (10t + T[i]) \bmod q$
 07 for $s=0$ to $n-m$:
 08 if $p = t$: \rightarrow if $P[1 \dots m] = T[s+1 \dots s+m] \neq O(m)$
 09 return S
 10 $t = (10(t - h \cdot T[s+1]) + T[s+m+1]) \bmod q$
 11 return None

Time: $O(n + F \cdot m)$

△ F : #False Matches

if mod q values are "random": $E(F) = \frac{n-m}{q} < \frac{n}{q}$

\Rightarrow Time = $O(n + F \cdot m) \leq O(n + \frac{n}{q} \cdot m) \Rightarrow$ want: $q > m$