



JAPNA EXPORTS

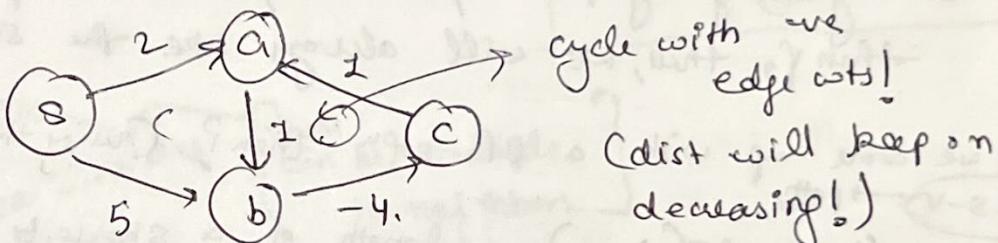
→ Dijkstra's doesn't work with negative edge weights!
 Even with cycles with -ve edge wt, it will be problematic).

When cycles :-

- $\text{dist}(v)$ goes to $-\infty$ for every v on the cycle (a, b, c, a)
- no solⁿ to shortest paths when -ve cycles.

⇒ need to detect -ve cycles!

So, we first detect this, if not found, then we find shortest path!



Properties of shortest paths:-

Suppose the problem has a solⁿ for an input graph.

- Can there be -ve cycle in the graph? No Yes
- Can there be +ve u u u u? Yes No
- Can the shortest paths contain the cycles? No
- Consider a simple shortest path; are its subpaths shortest? Yes

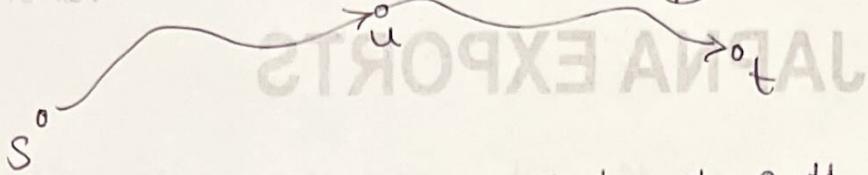
for a simple shortest path, can we have more than $n-1$ edges if n is the no. of vertices from $s \rightarrow v$? No.

Does problem exhibit optimal substructure?

ASPIREOTS-II-18 1st
ASPIREOTS-II-18
ASPIREOTS-II-18 xcf

(P)

(rest)



Let r be simple shortest st. path,
that goes thru. "u".

$$\text{OPT}(s, t) = \text{OPT}(s, u) + \text{OPT}(u, t)$$

[for optimality of full path,
all the sub-paths must be
optimal.]

$$r > \underline{\text{OPT}(s, u)}$$

- ④ any dist. of any path b/w $s \& u$, it will always have shorter
than r , then, we will always take the shortest path.

- ⑤ Can we come up with a DP soln then? (using the fact above
about $\underline{\text{OPT}}$).
from $s \rightarrow v$ path.

Let $\text{OPT}(i, v) =$ length of a shortest $s-v$ path
that was $\leq i$ edges.

{ Subproblem
that expresses
this! }

(try & we more & more edges
to see whether we can find the
shortest path).

$$\text{OPT}(n-1, v)$$

as we can go to at most $(n-1)$ edges
to see what can be the shortest path.

→ We want $\text{OPT}(n-1, v)$, $\forall v \in V$ → all set of vertices.
 $\nwarrow_{\max.}$
 $\underbrace{(n-1)}_{\text{number of edges in complete graph } G}$.

Boundary condition
 $\text{OPT}(0, v) = \begin{cases} 0 & (\text{as no. edges are being considered}), \text{ if } v = s \\ \infty, & \text{if } v \neq s \end{cases}$

$\text{OPT}(1, v) = \begin{cases} 0, & \text{if } v = s \\ w_{sv}, & \text{if } \exists (s, v) \in E \\ \infty, & \text{otherwise.} \end{cases}$

Progressively we
move forward
...



JAPNA EXPORTS

Now, we want to find a recurrence,

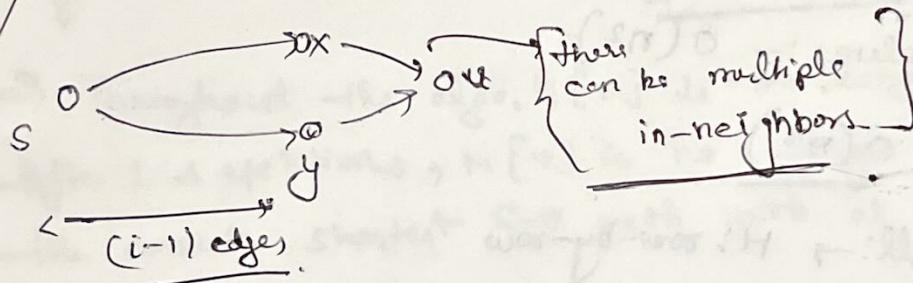
$$OPT(i, v) = \begin{cases} & \end{cases}$$

$$OPT(i-1, v)$$

if shortest s-v path that uses $\leq i$ edges that actually uses $\leq i-1$ edges.

$$\min_{x \in (u, v) \in E} \{ OPT(i-1, x) + w_{uv} \},$$

if shortest s-v path that uses $\leq i$ edges actually uses i edges.



bcz we don't know which one to take, we will consider minimum of those two!

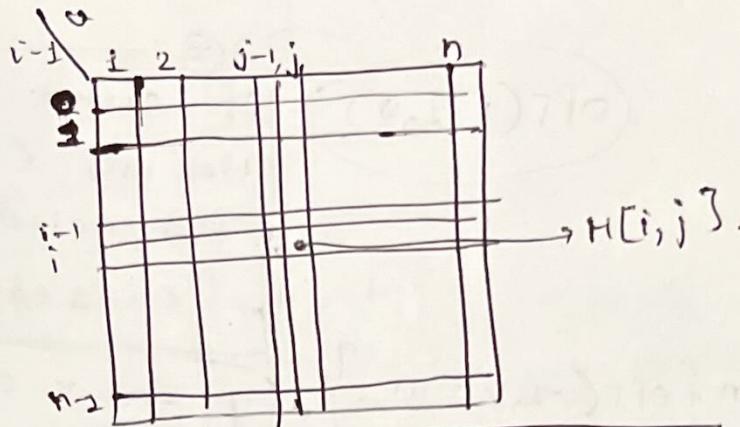
Recurrence

$$OPT(i, v) = \min \left\{ \begin{array}{l} OPT(i-1, v) \\ \min_{x \in (u, v) \in E} \{ OPT(i-1, x) + w_{uv} \} \end{array} \right\}$$

$$\left. \begin{array}{l} 0 \leq i \leq n-1 \\ 1 \leq v \leq n \end{array} \right\} \Rightarrow i * v \Rightarrow \left\{ O(n^2) \text{ subproblems.} \right\}$$

\Rightarrow As max^m no. of in-neighbors of v can possibly be n (no. of vertices)
 \therefore Time for each subproblem $\Rightarrow O(n)$ [$\because OPT(i-1, v)$ is constant time, we already computed that].

* Time to fill in all the sub-problems \Rightarrow $O(n^3)$ time
 (dense graph)
 for filling a matrix $\Rightarrow M[i, v] = OPT[i, v]$.



We want: $M[n-1, v]$, $\forall v \in V$

Time per subproblem: $\sim O(n)$

subproblems: $\sim O(n^2)$

Total time: $\underline{O(n^3)}$

{ Order to fill: $\rightarrow M$: row-by-row
 Total space: $\sim O(n^2)$. (we are only considering 2 rows at a time) $\hookrightarrow O(n)$ time

\therefore we only consider in-degree of j , $\rightarrow O(\deg(v))$

\therefore for $\min_n \{ OPT(i-1, v) + w_{nv} \}$

\therefore time: $\sim O(\deg(j))$

(for tighter analysis)

\therefore for entire row \Rightarrow $O(\sum_i \deg(j))$
 $= O(m)$ all the edges.

Time to fill in a row: $\sim O(m)$

if n rows in M of Total time: \sim

$\boxed{O(nm)}$

Any



JAPNA EXPORTS

→ To reconstruct actual shortest paths, also keep array prev of size n such that,

{ $\text{prev}[v] = \text{predecessor of } v \text{ in current shortest s-v path.}$ }

Improving space reqs:-

Only need two rows of M at all times.

→ Actually, only need one. Thus, we can drop the index i from $M[i, v]$ and only use it as a counter for # repetitions.
(\because 1-D array)

⇒ Throughout the algo, $M[v]$ is the length of some s-v path.
After i repetitions, $M[v]$ is no larger than the length of
the current shortest s-v path with at most i edges.

Early termination condns:- If at some iteration i , no value in M changed, then stop! (This can be used to find if there are any -ve cycles in the graph) \Rightarrow after completed the flow and we have found the shortest path $(n-1)$ iterations. If we run it for 1 more iteration, technically it shouldn't update! If anything is updated, this would show presence of -ve cycles (\because path size can then grow upto $-\infty$), \rightarrow would never end in terms of no. of iterations.

Thus, single loop of Bellmann-Ford should work!

Pseudocode:-

$n \times n$ dynamic prog. table M s.t.

$$M[i, v] = OPT(i, v)$$

$BF(G = (V, E, w), s \in V)$

for $v \in V$ do

$$M[0, v] = \infty$$

end for

$$M[0, s] = 0$$

for $i = 1, \dots, n-1$ do

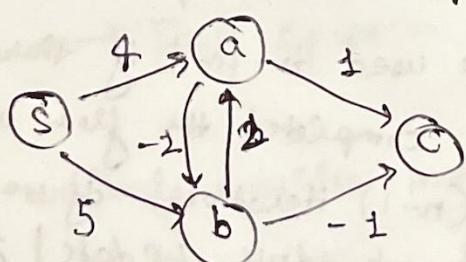
for $v \in V$ (in any order) do

$$M[i, v] = \min \left\{ \begin{array}{l} M[i-1, v] \\ \min_{x: (v, x) \in E} \{ M[i-1, x] + w(x, v) \} \end{array} \right\}$$

end for

end for.

\Rightarrow



	s	a	b	c
s	0	∞	∞	∞
a	0	4	5	∞
b	0	4	3	4
c	0	4	3	2

we can remove
this for space
conditions!

bcoz this
will only
consider
prev. array
formed!

for using i as just
repetition, we
form 1-D array!

	s	a	b	c
i=0	0	∞	∞	∞

	s	a	b	c
i=1	0	4	3	2

(concrete into
this!)

this is the same as 1-D array for
distance



JAPNA EXPORTS

Pseudo code becomes :-

Initialize (G, s)

for $v \in V$ do

$dist[v] = \infty$

$prev[v] = NIL$

end for

$dist[s] = 0$

Update (u, v)

if $dist[v] > dist[u] + w(u, v)$ then:-

$dist[v] = dist[u] + w(u, v)$

$prev[v] = u$

end if

Bellman-ford ($G = (V, E, w), s$)

Initialize (G, s)

for $i = 1, \dots, n-1$ do

for $(u, v) \in E$ do

Update (u, v)

end for

end for.

$s \Rightarrow$ source node

$O(n)$

$O(nm)$

If : Time $\Rightarrow O(nm)$
 : Space $\Rightarrow O(n)$

⇒ Floyd-warshall algos {All pair shortest paths) with -ve edges}

Runs in $O(n^3)$

Ans.

Discrete Random variables 2)

Takes finite no. of values with same probability \Rightarrow

$$E[X] = \sum_j j \cdot P[X=j]$$

Bernoulli trial-

→ flip a biased coin! ←

→ heads with prob. (p)

\rightarrow fails with prob. $(1-p)$

Qnd what is the expected no. of heads?

Let $X \Rightarrow \begin{cases} 1, & \text{if coin flip comes up H} \\ 0, & \text{if } \dots \end{cases}$

$$\Rightarrow E[x] = \sum j \cdot h[x=j] \Rightarrow 0 \cdot h[x=0] + 1 \cdot h[x=1] \\ = 0 \cdot (1-p) + 1 \cdot p = p \quad \underline{\text{Ans.}}$$

Indicator random variable :- a discrete rv. that only takes on values of 1.

Denote occurrence (or not) of an event),

If X is an indicator r.v., then

$$E[X] = \Pr[X=1].$$

filip coin n times.

what is exp. no. of threads?

what is exp. no. of heads?
let X be a rv. that counts the # heads in n coin flips.

Let X be a rv. that counts the # heads.

$$E[X] = \sum_{j=0}^n j \cdot \binom{n}{j} p^j \cdot (1-p)^{n-j}$$

Value of {
that } event

Probability of
that event

$\exists p \geq \text{yes}$, iff. n has a perfect matching.

→ Encoding of a set S of abstract objects:-

Mapping of elements from S to binary string.

Should we be concerned about the particular encoding used?

→ Yes!

Reasonable encodings :-

(not unary encodings).

→ Encoding of an integer is polynomially related to its binary representation (V)
 $1024 \rightarrow$ in binary but not
 $(11 \dots 1)(x) \rightarrow$ not reasonable
102ubits

→ Encoding of a finite set is polynomially related to its "standard" encoding as a list of its elements.
 $\{\dots, \dots, \dots\}$

Example of encodings :-

Input to sorting → $\rightarrow n$ elements

Standard encoding $\langle x \rangle$:- $\{x_1, x_2, \dots, x_n\}$

Length of standard encoding! - $|\langle x \rangle| = O(n)$

poly-time algo runs in:- $O(n^k)$ → (in length of bit string that encodes the I/P).

Input to st connectivity →

Standard encoding → $\langle G, s, t \rangle$

$\{v_1: \{v_{11}, \dots, v_{1K_1}\}, v_2: \{v_{21}, \dots, v_{2K_2}\}, \dots, v_n: \{v_{n1}, \dots, v_{nK_n}\}, s, t\}$

Length of standard encoding of $|\langle G, s, t \rangle|: O(n+m)$

Poly-time algo runs in:- $\underline{\underline{O((n+m)^k)}}$ Any

Input to max flow

Standard encoding $\rightarrow \langle G, s, t, c \rangle \Rightarrow U = \max_{e \in E} c_e$.

$O(nmU)$ & graph rep. as before!

Length $\rightarrow |\langle G, s, t, c \rangle| \geq \underbrace{O(n+m)}_{\text{for } G} + \underbrace{m \log_2 U}_{\text{for capacities}}$

for G for capacities.

Poly-time algo runs in:-

$\delta((nm \log_2 U)^k)$

for FF algo.

Now, coming to the decision problem

we have already designed algo that solves BPM

i. Given:- bipartite graph G (IP to BPM), we construct a flow N/W G' (IP to max flow).

\rightarrow we transform the IP instance n of BPM to an IP instance y of max flow.

we can use this itself to answer the decision, basically if it equals (flow) the number of vertices in X & Y , the $|X| = |Y| = |f|$, then, return "yes", otherwise "no".

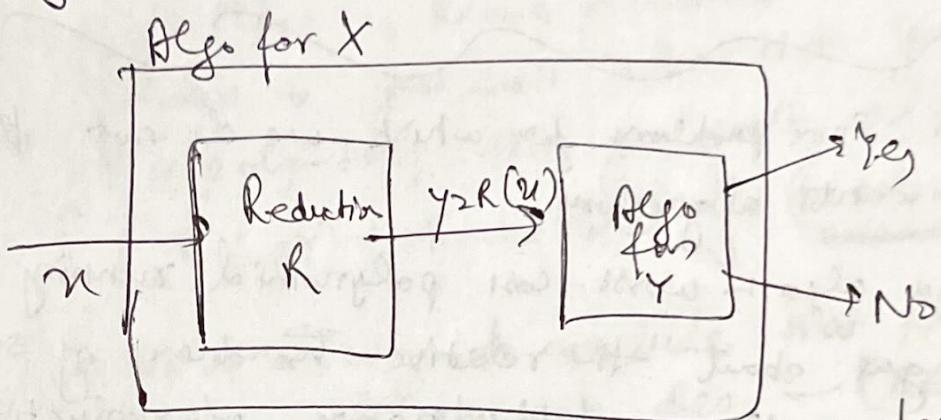
About Reductions)

\rightarrow a transformation for every IP n of X produces an equivalent input $y = \tau(n)$ of Y .

\rightarrow By equivalent \rightarrow we mean that answer to $y = \tau(n)$ considered as an IP for Y is a correct yes/no to n , considered as an IP for X .

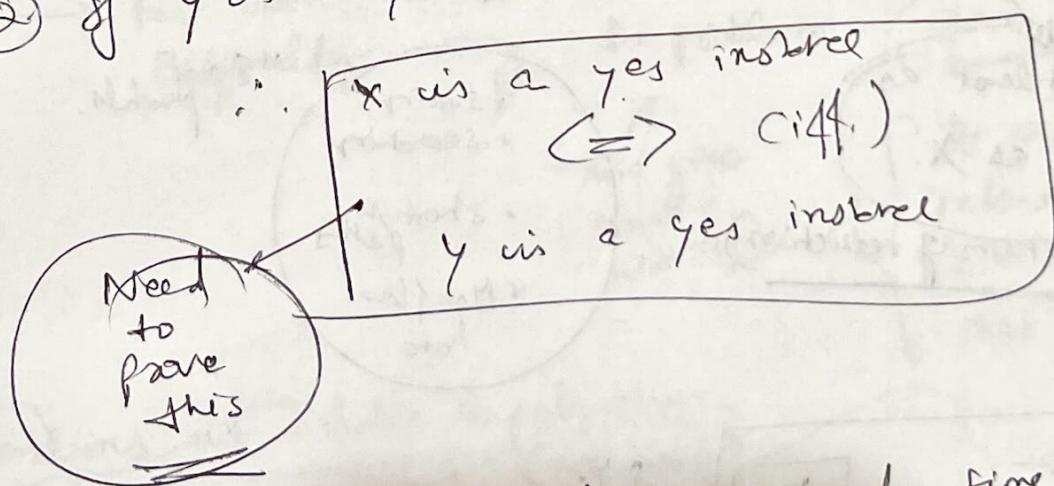
we will require that the reduction is completed in a polynomial time in $|x|$ number of computational steps.

If X & Y are computational problems,
 x, y are I/Ps to X, Y respectively.



To prove equivalence of instances for such a decision problem:-

- ① If x is a yes instance, then y is a yes instance.
- ② If y is a yes instance, then x is a ~~yes~~ instance.



If both Red^n & $\text{Algo for } Y$ are poly-time, then algo for X is considered ~~to be~~ to be an efficient algo-polynomially to Y ;

Thus, we say that X reduces

$$X \leq_p Y$$

If Y is solvable in poly-time, then X is solvable in poly-time ($\because \text{Red}^n$ has to poly-time)

NOTE: To solve BPM we made exactly one call to the black box for max flow. This is mostly same for our cases!

(Reductions provide a powerful technique for designing efficient algos!)

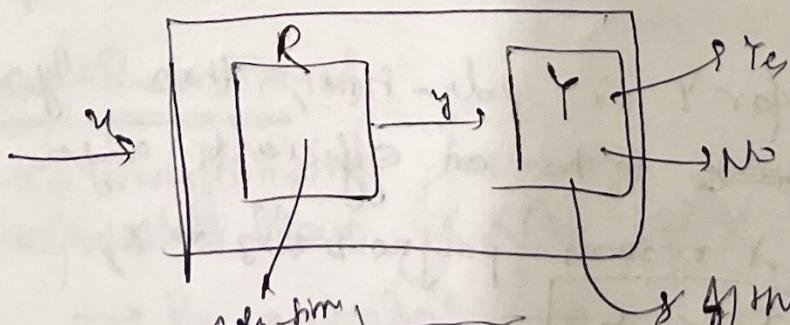
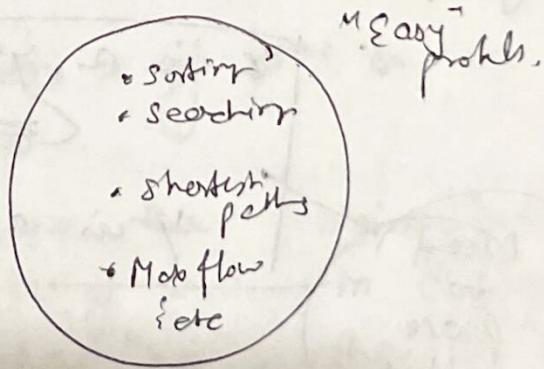
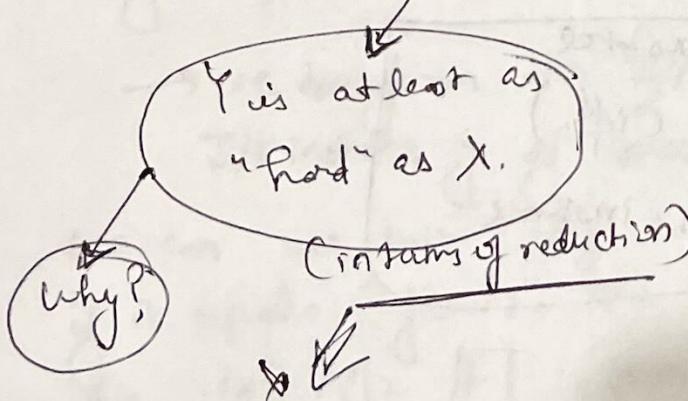
→ There are some problems for which we do not know of any efficient algorithms.

→ Effi algos worst-case polynomial running time.

→ To argue about the relative hardness of such problems we will use the notion of reductions!

Questions

Suppose $X \leq_p Y$ (prob. X reduces polynomially to prob. Y).



if this is poly-time, then X is also poly-time,
if this is super polynomial
then X is the same

gives us for X, the above fact holds

Suppose $X \leq_p Y$. If Y is solvable in polynomial time,
then X is solvable in polynomial time.

Contrapositive of fact above:

If $X \leq_p Y$ & X is not solvable in poly-time,
then Y is not solvable in poly-time.

Thus, if you have a known problem X that is non-solvable & you can't reduce it to your problem

Y in poly-time, then you can't solve Y as well.

$\{ (X \leq_p Y) \rightarrow \text{should hold.} \} \rightarrow \text{To solve } Y$

* This provides us a way to conclude that problem Y does not have an efficient algo. → [Red helps us to establish for problems discussed], we have no proof that they do not have efficient algo

this
subtly}

→ fact above can be used to establish relative levels of difficulty among the problems (thru. Reductions)

This gives us the answer!
Is there a way to reduce the prob.
 $X \rightarrow Y$ & $Y \rightarrow X$ or not?

Problems not solvable in poly-time:
(not proved)

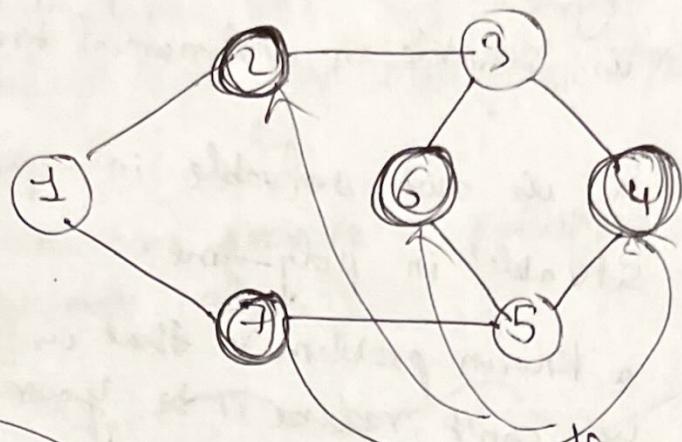
Independent Set:-

→ An independent set in $G = (V, E)$ is a subset $S \subseteq V$ of nodes s.t. there is no edge b/w any pair of nodes in the set.

That is, for all $u, v \in S$, $(u, v) \notin E$.

→ easy to find small IS but difficult to find a larger one!

Eg:-



{ Max^m
Independent
set
problem }

These form the max^m
independent set,

∴ we have to go over each diff. pair
of vertices to check this takes,

2^n time

$\Rightarrow \Omega(2^n)$

(lower bound of 2^n)

(exponential time!).

Another problem:-

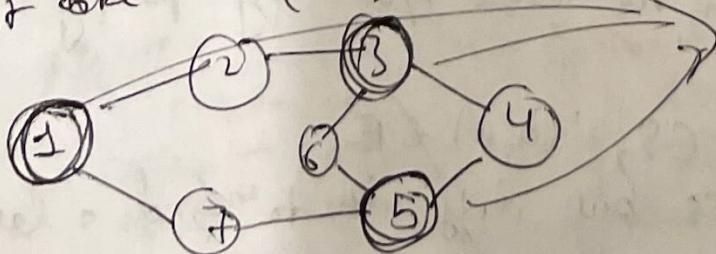
Vertex Cover ↗

A vertex cover in $G = (V, E)$ is a subset $S \subseteq V$ of nodes s.t. every edge $e \in E$ has at least one endpoint in S .

(sets of vertices that covers up all the edges)
→ another way of looking at this.

→ Easy to find a large vertex cover (\because it always equals set of all the vertices \forall). But it is really hard to find a small one.

→ It takes $\Omega(2^n)$ time as well.



This is the min^m
vertex cover for
this graph G_1 .
(hard to find!)

relationship b/w vertex cover & independent set
if S is a vertex cover, then $V-S$ is the independent set.

(complement of each other)

Optimization versions of IS & VC

Max^m IS \Rightarrow Given G , find an IS of max size.

Min^m VC \Rightarrow Given G , find a VC of min^m size

(Requires exponential time :- 2^n candidate subsets
- every vertex may/may not belong to
such a subset)

Now, we can transform these ~~decision versions of~~ optimization problems into decision versions (to try & solve these)

→ An optimization problem may be transformed into a roughly equivalent problem with a yes/no answer, called the "decision version" of the optimization problem by \exists

① supplying a target value for the quantity to be optimized,

② asking whether this value can be attained.

→ denoted by $x(0)$ if prob. was X .

↳ decision version of prob X .

Examples of decision versions of optim. prob's

→ Max flow (D):- Given a flow N/W G and an integer K (target flow), does G have a flow of value at least K ?

- Max Bipartite Matching (D) :- Given a bipartite graph & an integer k , does G have a matching of size at least k ?
- IS(D) :- Given a graph G & an integer k , does G have an independent set of size at least k ?
- VC(D) :- Given a graph G & an integer k , does G have a vertex cover of size at most k ?

So, in a max. prob. the target value is a lower bound, while in a min. prob., it is an upper bound.

- Shortest s-t path (D) :-

Given $(G = (V, E, w), s, t)$ and k , is there an $s-t$ path in G of length $\leq k$?

- Min cost sequence alignment (D) :-

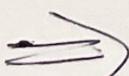
Given $(X, Y, \delta, \alpha) \xrightarrow{\text{costs}} k$, is there an alignment of X, Y that costs $\leq k$?

- Knapsack (D) :-

Given $(n, v_i \leftarrow v_1, v_2, \dots, v_n, w_i \leftarrow w_1, w_2, \dots, w_n, W)$

& K , is there a subset of the item with weight $\leq w$ & value $\geq k$?

(at least k . i.e. this was a max. prob.)



Now, let's try and convert an optimization problem to its decision version:-

Maxⁿ Independent set problem (maxIS) :-

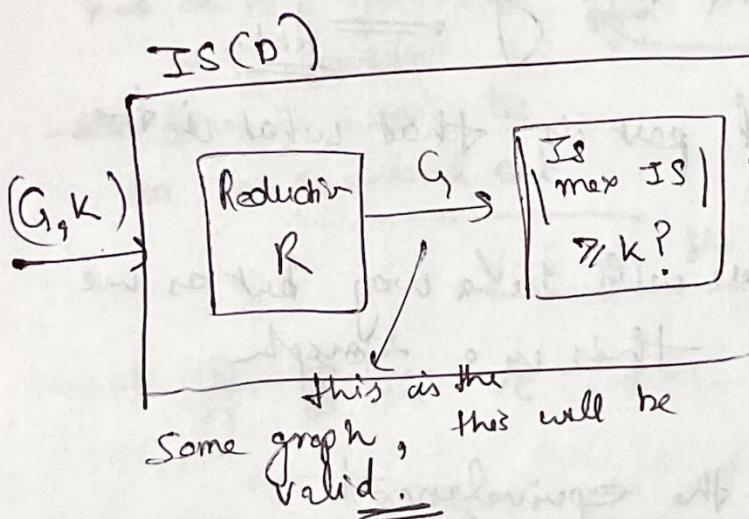
Its decision version denoted by IS(D).

convert this, we first need to prove that these 2 instances are equivalent, i.e., both $\text{IS}(D)$ and max IS follow the fact $X \leq_p Y$ (whether they can be converted in an equivalent manner) \Rightarrow

$$\boxed{\text{IS}(D) \leq_p \text{max IS}}$$

(using reduction, we can convert)

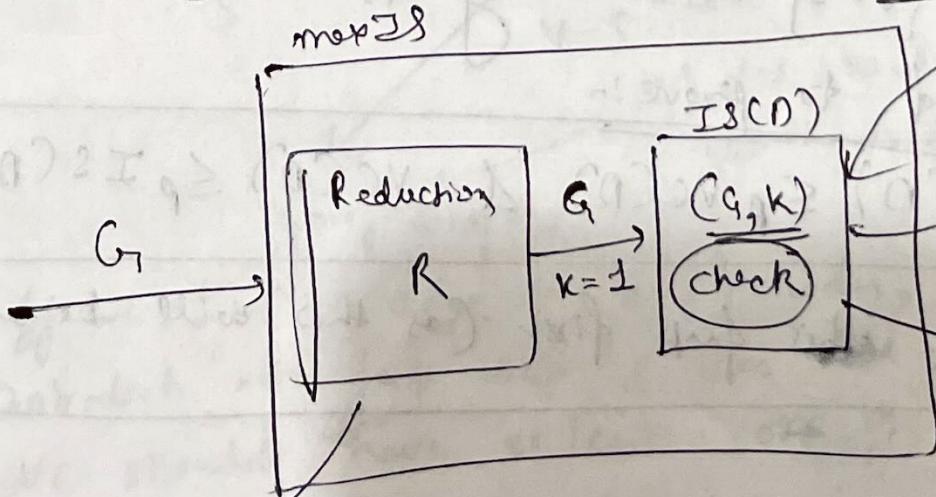
→ Intuitively, it makes sense that max the value of independent set and finding k is at least as hard (actually harder) than when we just need to tell whether an IS of size k is present or not in the I/P graph G .



$$\boxed{\text{max IS} \leq_p \text{IS}(D)}$$

(using reduction)

(In this, calls are made again & again to $\underline{\text{IS}}$)



\therefore using reductions input also for $\text{IS}(D)$
we can make it as hard as max IS problem. (polynomial time)

Yes: $K = K + 1$
(loop again)

No: $K - 1$

(this tells that we could only find IS of max size $K-1$ as we cannot find $\text{max IS} \geq k$.)

~~Notes~~ for max flow problem \rightarrow nU
 reduction will be of
 this size.
 (this is not poly-time).

what we can do
 is search the search-space
 from $k \rightarrow n^k$ using binary search)

$$\log_2(nU) \geq \log n + \log U$$

this becomes polynomial time after this
Ans.

What we are missing in and part is that what is the set, we only return k .

To reconstruct the set, there will be a way but as we don't prove that here, this is a rough approximation!

We have proven the equivalence!

Reduction from Independent Set to VC:-

FACT Let $G = (V, E)$ be a graph. Then $\{S\}$ is an IS of G if and only if $\{V - S\}$ is a VC of G .

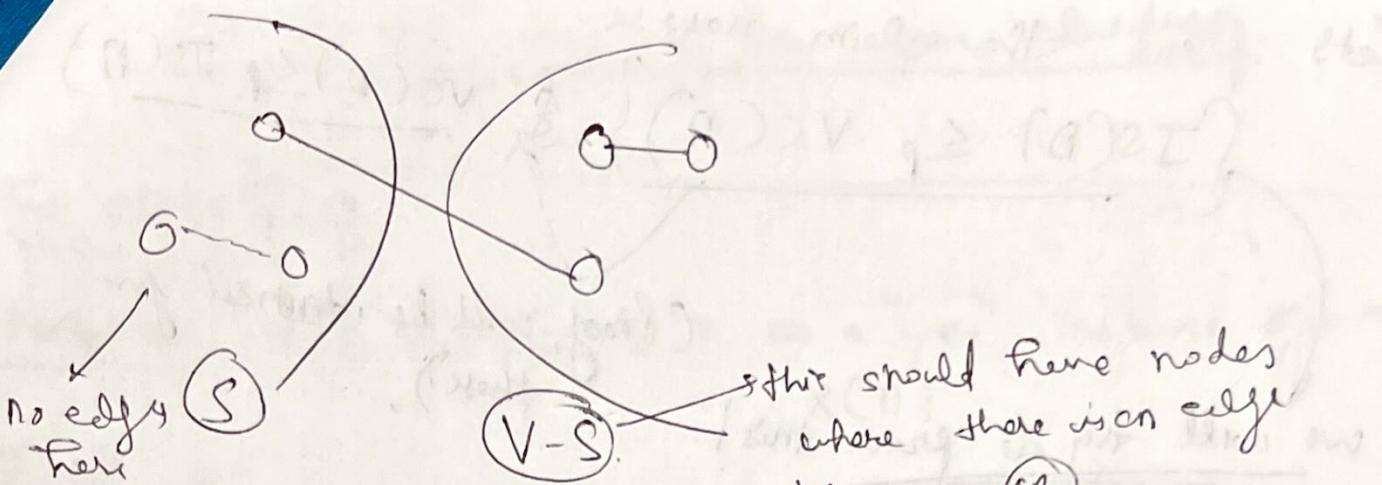
Now, we need to prove :-

$$IS(D) \leq_p VC(D) \text{ & } VC(D) \leq_p IS(D)$$

Let's prove this fact first (as this will be used).

Proof of FACT :-

\rightarrow If S is an indep. set $\Rightarrow V - S$ is a VC.



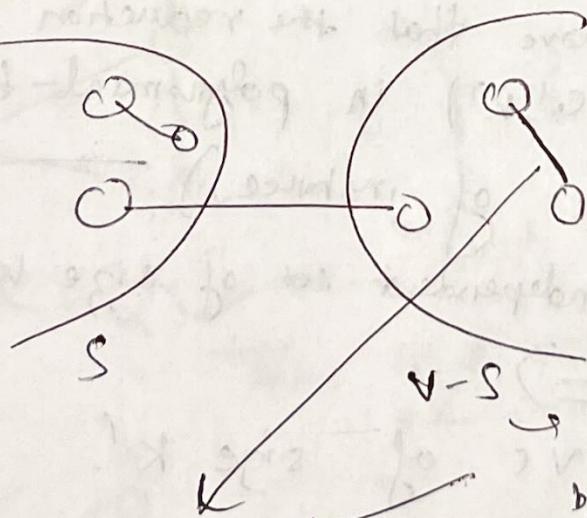
If we had an edge in nodes in S , then
~~then S cannot be an independent set.~~

this should have nodes where there is an edge b/w them OR

one endpt lies in $\{V-S\}$, defn of VC

$\therefore V-S$ is the vertex cover.

\rightarrow If S is in VC , then $V-S$ is an IS.



$\rightarrow V-S$ should have no edges b/w the nodes.

going contrary

If we had an edge in $\{V-S\}$ nodes, then this contradicts the fact that S was the VC bcoz, VC should have at least one of those nodes in the set S .

As this contradicts the premise, we have proven the 2nd dir as well.

Thus, we have proven the fact! \square

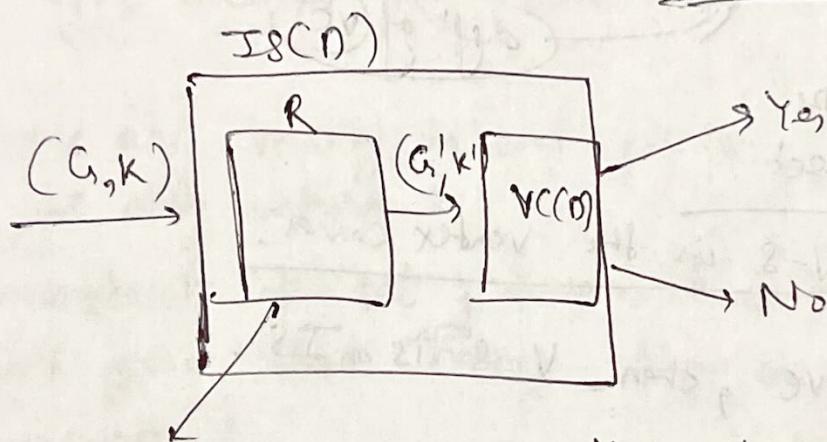
Let's prove the claim now?

$$\{ \text{IS}(\rho) \leq_p \text{VC}(\rho) \} \models \underline{\text{VC}(\alpha) \leq_p \text{IS}(\eta)}$$

(Proof will be identical for these),

we will try to prove this!

for IS problem $\Rightarrow \underline{IS(D)} \leq_p \underline{VC(D)}$



~~(*)~~ Thus, we need to prove that the reduction can convert $(G, k) \rightarrow (G', k')$ in polynomial time.

s.t. (Equivalence of instances)

$\left\{ \begin{array}{l} G \text{ has an independent set of size } k \\ \hookrightarrow \\ G' \text{ has a VC of size } k' \end{array} \right.$

 k
 (meaning)

(Need to prove this!)

we know by the fact proven before that -

$$(G = G') \rightarrow \text{some}$$

$$k' = n - k' \quad (\text{proven before})$$

Thus, the instances are equivalent!

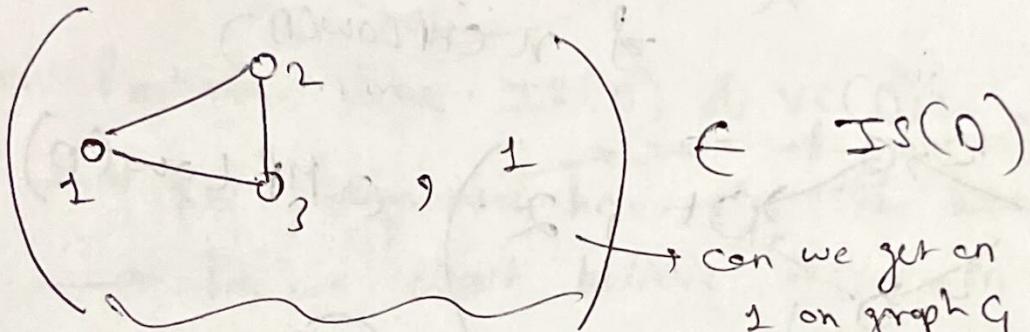
(As we are not changing G)

\therefore IS & VC are of the same hardness.

The class P-d

Notation: $n \in X(D) \iff n$ is a "yes" instance of $X(D)$.

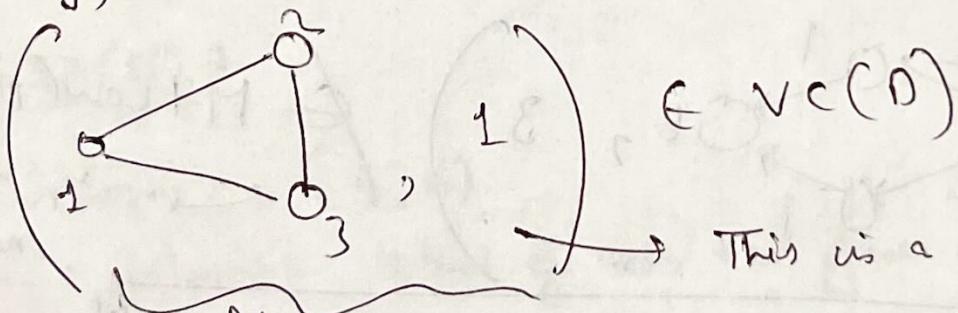
Exs



→ can we get an IS of size 2 on graph G
↳ "Yes"

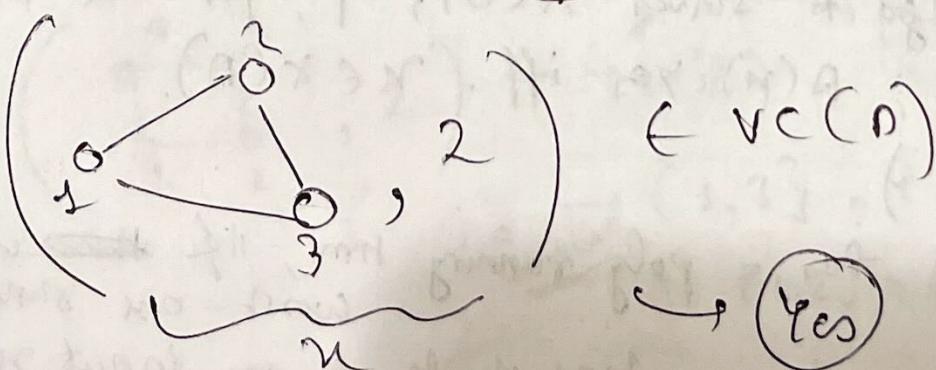
thus $n \in IS(D)$.

Similarly,



→ This is a "no" instance.

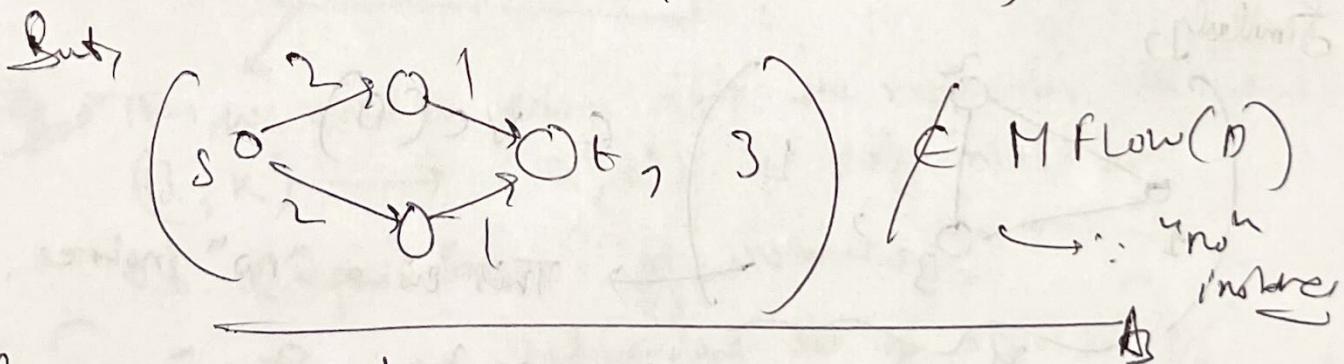
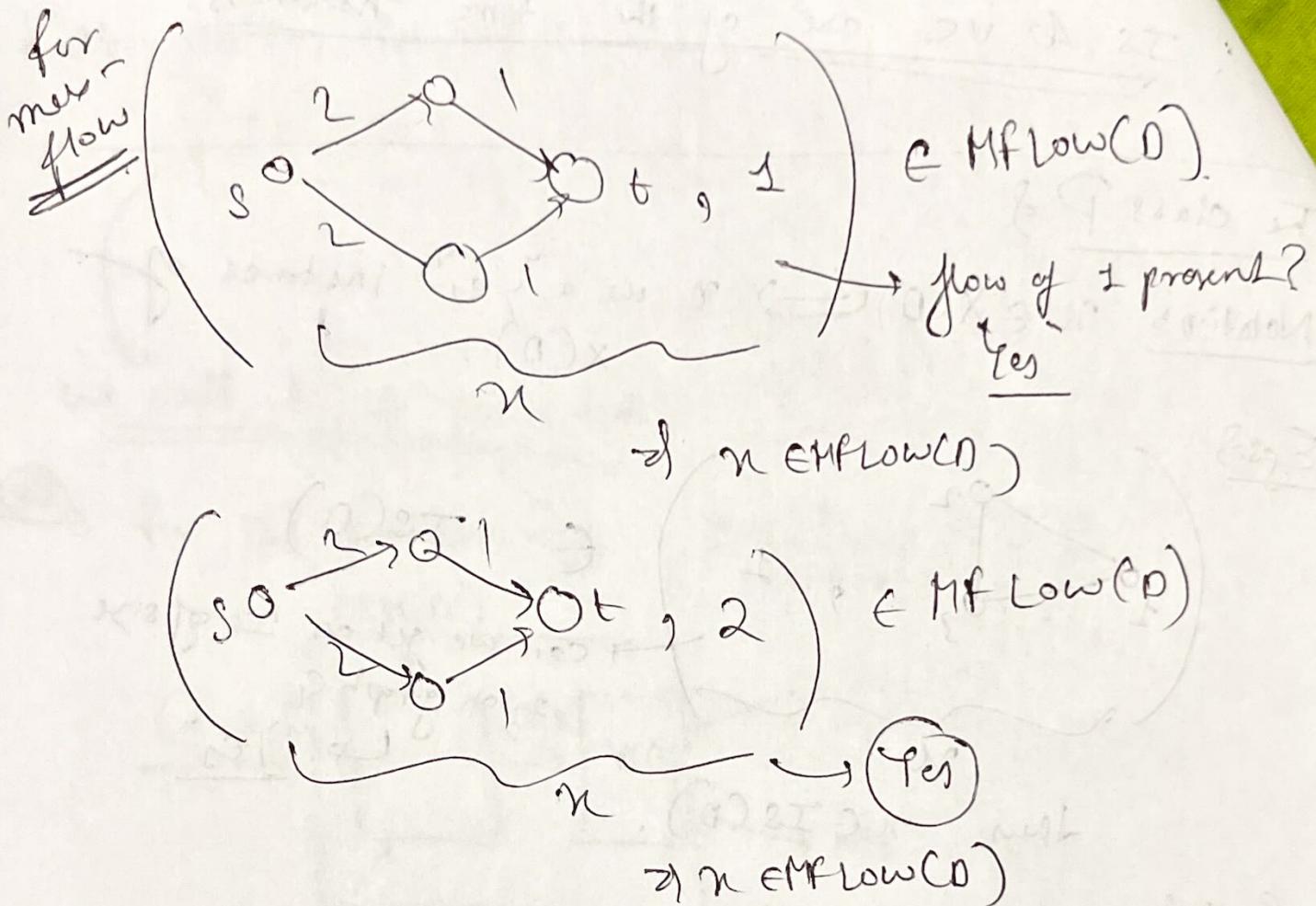
thus $n \notin VC(D)$.



→ Yes

$n \in VC(D)$.

Similarly,



for correctness \rightarrow

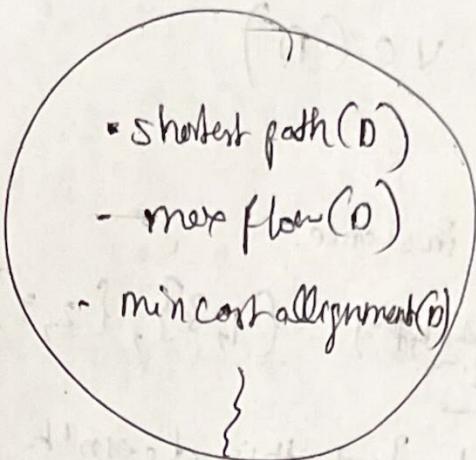
An algo A solves $x(D)$, if, for all input n : $A(n) = \text{yes}$ iff. $n \in x(D)$.

For efficiency \rightarrow

\rightarrow A has a poly running time, if ~~the~~ worst-case running time of algo A on input n is $O(1n)$

Then, we define P to be the set of decision problems that can be solved by poly-time algos.

P. \nrightarrow



clen P

- \Rightarrow for problems like $IS(D)$ & $VC(D)$:
(worst-case instances)
- No polynomial time algo has been found!
 - So, we don't believe they are in P.

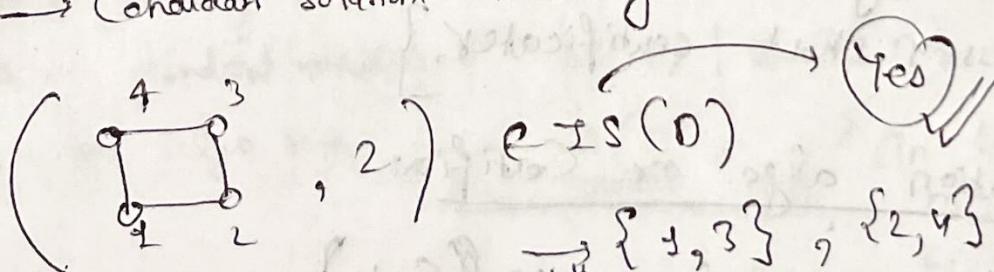
\nrightarrow Is there anything +ve we can say about such problems?

\Rightarrow $IS(D)$ \rightarrow

→ Input: (G, k)

→ Solution: subset of nodes that form an ~~clen~~ of size k .

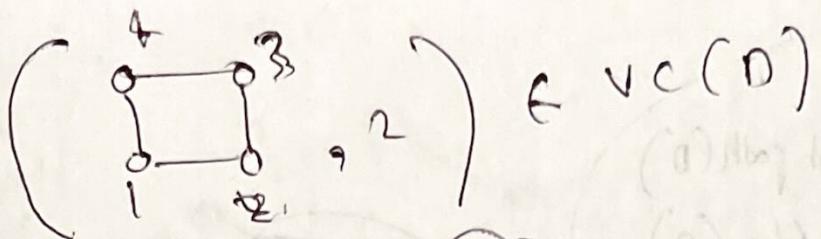
→ Candidate solution: subset of nodes.



a candidate
Is $\{1, 2\}$ a soln? \rightarrow No

instance

Just on basis of a "no" candidate soln, you can't say where it is a "no" instance.



(Yes) instance.

~~Candidate soln~~. $\{1, 3\}, \{2, 4\}$.

~~Candidate $\{1, 2\}$~~ → this doesn't give as "Yes".

for both these problems :-

CERTIFICATE + for input x :

Short soln that CERTIFIES that $x \in X(D)$.

Certificates are PROPERTIES of "Yes" instance

$x \in X(D) \Leftrightarrow \exists$ short t for x.
 \uparrow poly-time

① Can define a complexity class,
NP = {set of decision problem where Yes instances possess short certificates.}

② Verification algo. or Certifier:-

→ Input to IS(D) :- $x = \{a_i, k\}$

→ Input to certifier for IS(D) :- (x, t)

\uparrow
certificate

→ What does the certifier do with its I/P?

- ① $\exists k$ nodes in t
- ② t forms an independent set.

Network flows:-

→ edges carry traffic; nodes → switches where traffic gets diverted!

(source → sink)

flow networks ↗

A flow network $G = (V, E)$ is a directed graph s.t.

1 → every edge has a capacity $c_e \geq 0$

2 → There is a single source $s \in V$.

3 → There is a single sink $t \in V$



Assumptions

A1: - integer capacities.

A2: - no edge enters s .

A3: - no edge leaves t .

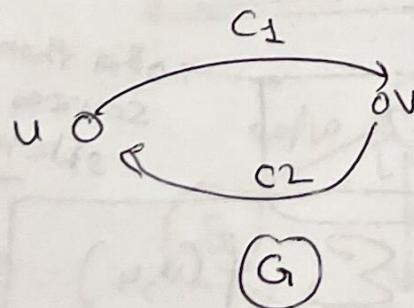


A4: - if $(u, v) \in E$ then $(v, u) \notin E$.

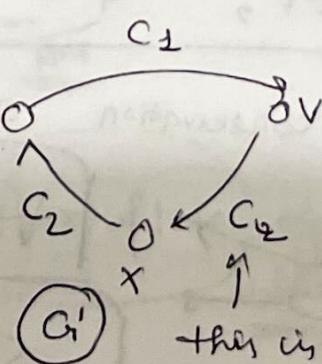
A5: - Every $v \in V - \{s, t\}$ is on some $s-t$ path.

Hence, G has $m \geq n-1$ edges.

If we do have $\boxed{(u, v) \neq (v, u)} \in E$



⇒

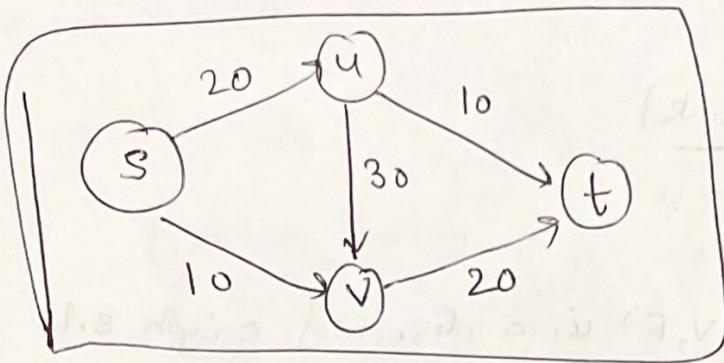


(we would have to convert this).

this is upper bound of the flow thru on edge / path.

This transformation will be polynomial time!

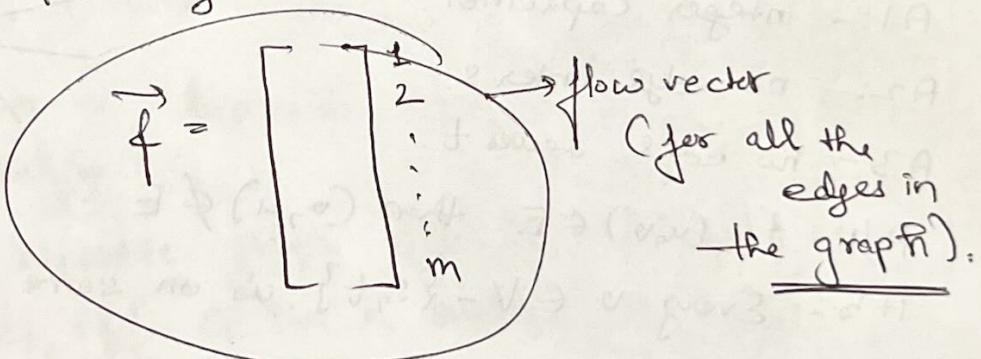
→ An efficient transformation!



flows — given a flow network G , an s - t flow f in G is a funcⁿ

$$f: E \rightarrow \mathbb{R}^+$$

Intuitively, the flow $f(e)$ on edge e is the amount of traffic that edge e carries.

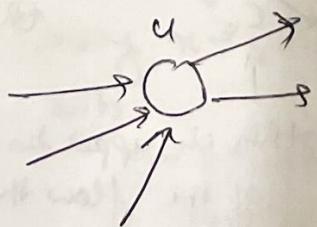


Constraints to satisfy —

① Capacity constraint.

$$\forall e \in E, 0 \leq f(e) \leq c_e$$

② flow conservation constraint



$$\sum_{(u,v) \in E} f(u,v) = \sum_{(v,u) \in E} f(v,u) \quad \text{for } \forall v \in V - \{s, t\}$$

other than sources and sink nodes

$$f^{\text{in}}(u) \triangleq \sum_{e \in \text{in}(u)} f(e); f^{\text{out}}(u) \triangleq \sum_{e \in \text{out}(u)} f(e)$$

Now, flow conservation constraints of

$\forall v \in V - \{s, t\}$,

$$f^{\text{in}}(v) = f^{\text{out}}(v)$$

Ans.

$$|f| \triangleq f^{\text{out}}(s)$$

$$= \sum_{(s, u) \in E} f(s, u)$$

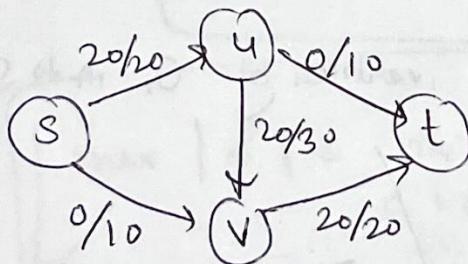
, only out of the source.

$$\text{Value of flow } f \Rightarrow |f| = f^{\text{in}}(t)$$

(flow into the sink t)

(2nd part
of flow-
conservation)

Ans.



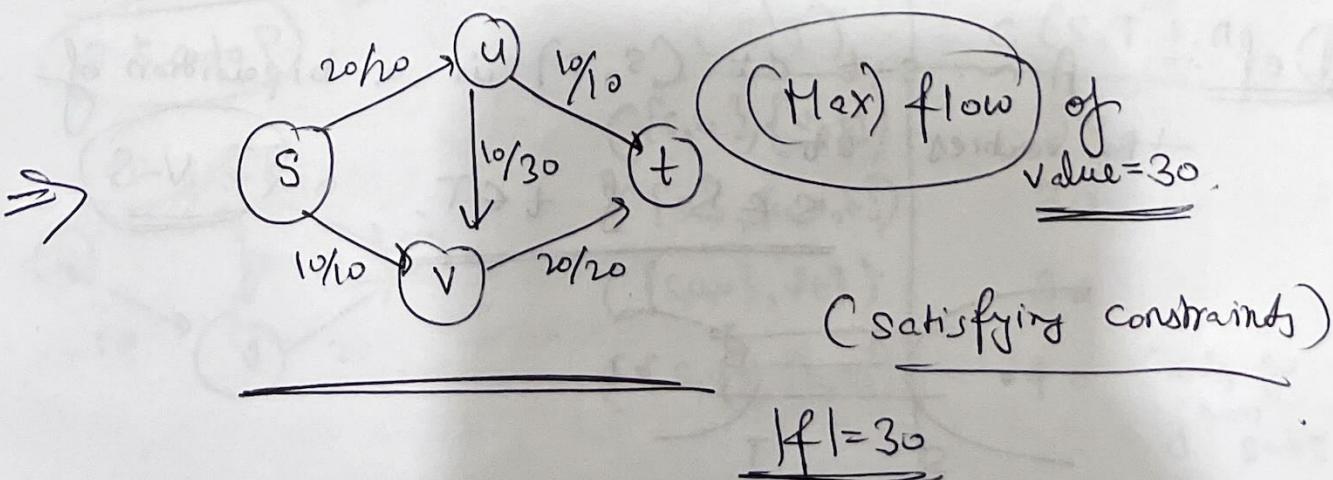
$$\vec{f} = \begin{bmatrix} 20 \\ 0 \\ 20 \\ 0 \\ 20 \end{bmatrix} \begin{matrix} su \\ sv \\ uv \\ ut \\ vt \end{matrix}$$

flow vector

(5 edges)

$$|f| = f^{\text{out}}(s) = f^{\text{in}}(t) = 20$$

Ans.



(Satisfying constraints)

$$|f| = 30$$

Max flow prob

Inputs: (G, s, t, c) s.b.

→ $G = (V, E)$ is a flow N/W.

→ $s, t \in V$ are source & sink.

→ c is the (integer-valued) capacity func'.

OR → a flow of maximum possible value.

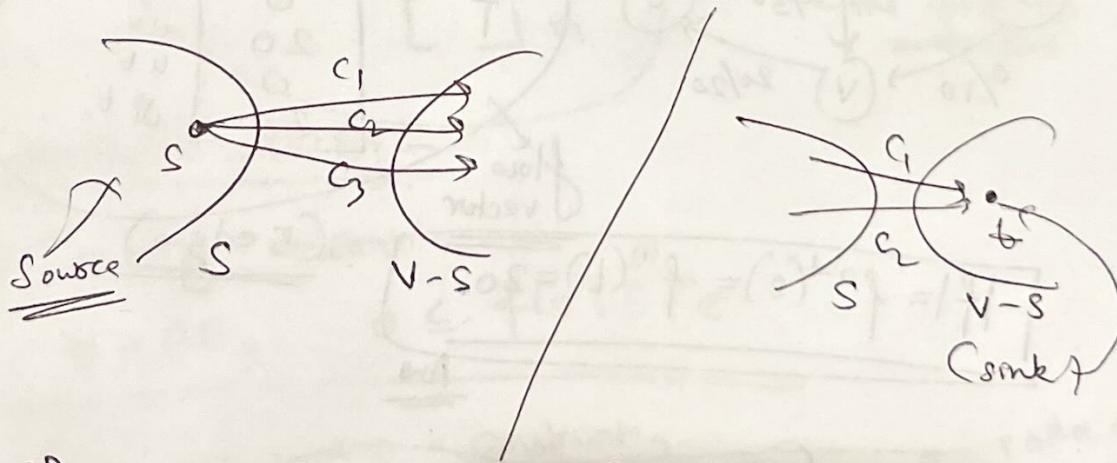
(what prevents flow to be very large value).

Capacity of the edges are the upper bounds!

$$\text{Upper bounds} \quad |f| \leq \sum_{e \text{ out of } s} c_e \quad |f| \leq \sum_{e \text{ into } t} c_e$$

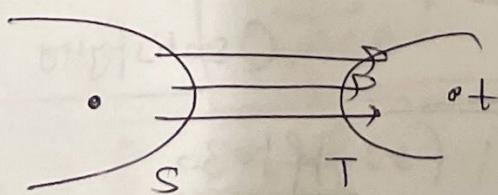
(1 more!)

Considering a bipartite graph of vertices of G , into $S \rightarrow \underline{V-S}$.



Def': A \sim s-t cut (S, T) is a bipartition of the vertices st.

$s \in S \quad \& \quad t \in T$.



$|f| \leq \sum_e c_e$ (for every $S \triangleleft T$)
 e crosses from $S \rightarrow T$
 $\triangleq c(S, T)$

for every (S, T) , value of the flow cannot exceed the capacity of the cut,

$\forall S-T$ cut (S, T)

$|f| \leq c(S, T)$

$|f| \leq \min_{S-T \text{ cut } (S, T)} c(S, T)$

min^m of all capacity of the cuts formed.

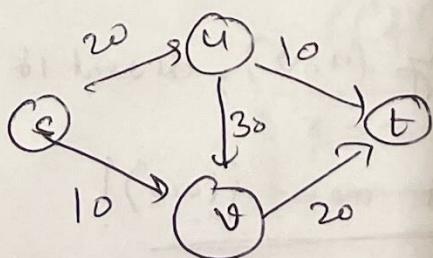
$\max |f| = \min_{S-T \text{ cut } (S, T)} c(S, T)$

max^m flow possible.

min-cut

we need an efficient algo. to find this

S-T cuts?



(S, T)	$c(S, T)$
$(\{S, U\}, \{V, T\})$	50
$(\{S\}, \{U, V, T\})$	30
$(\{S, U, V\}, \{T\})$	30
$(\{S, V\}, \{U, T\})$	40
	Only go from $S \rightarrow T$

we are trying to find upper bound of the flow, thus,
only consider flows from $s \rightarrow t$.

We always have,

$$2^{n-2} \text{ s-t cuts}$$

(for n vertices)

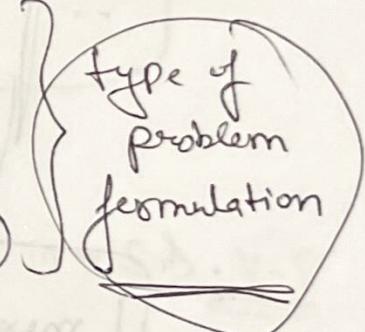
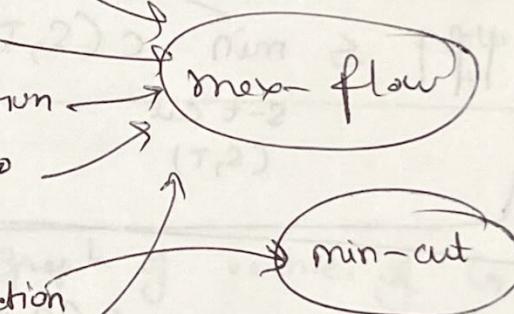
This is a lot

of cuts to go over to find max-flow!

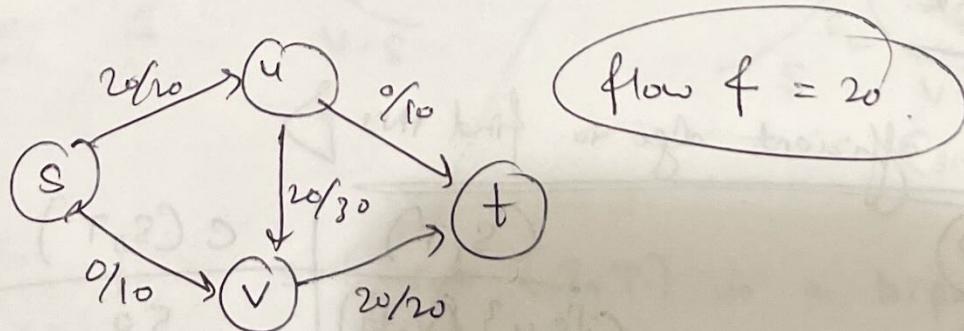
Thus, we need an efficient algo. to compute!

App's of

- Bipartite matching
- Alline sat.
- Baseball elimination
- Dist. of goods to cities
- Image segmentation
- Survey design



"Undeleting" flow →



Goal:— undelete 10 units of flow along (u, v) , divert it along $u - t$.

(to get max-flow)!

To undo that (identify s-t path, s.t. we can flow max amount maybe from diff. dirns).

Identify an s-t path in a s.t.-push flow forward to edges w. leftover capacity.

— push flow backward on edges that already carry some flow. (Undo flow)

Thus, we need diff. kind of Graph(G_f)!

Residual graph

G_f w.r.t a flow network G and a flow f .

→ has the same vertices as G .

→ $\forall e \in G$, s.t. $f(e) < c_e$, introduce $e \in G_f$ with residual capacity $c_f(e) = c_e - f(e)$

IT IS strictly
the
number

leftover
capacity.

(So, that we can
saturate an edge)

FORWARD edges
 (u, v)

→ $\forall e \in G$ s.t. $f(e) > 0$, introduce $e^r = (v, u) \in G_f$ with residual capacity $\underline{c_f(e^r)}$.

$$\underline{c_f(e^r)} = f(e)$$

→ to undo a flow.

$c_f(e^r) > 0$ → strictly +ve.

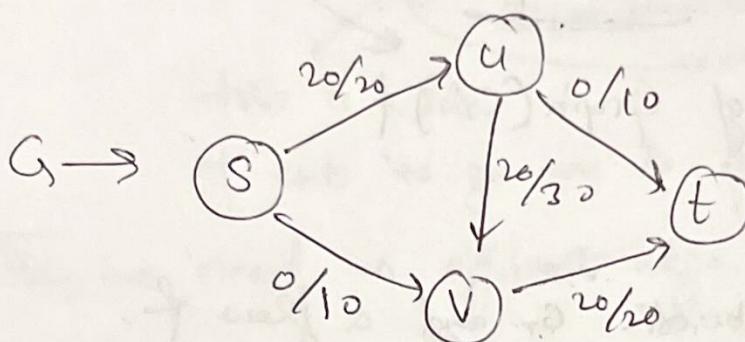
BACKWARD edges.

Thus, we have introduced both kinds of flow.

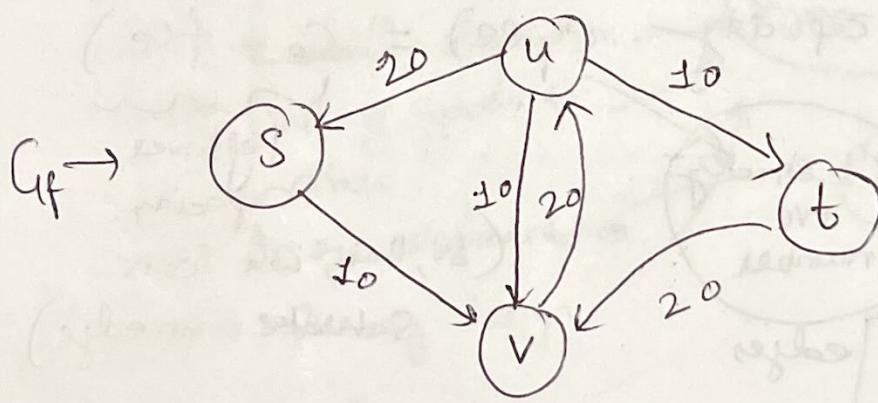
$$\text{max}^n \text{ no. of edges} = \underline{2|E|} \text{ in } G_f$$

→ Thus, we generally don't take anti parallel edges for analysis. (bcz we need backward edges!)

Construction of $G_f \rightarrow$



↓ G_f transform



(not necessarily
a flow N/w)

(there are anti-parallel
edges here!)

(a roadmap to
augment the flow)

→ We will use G_f as a roadmap to augment f in G .

Goal → Identify a simple s-t path in G_f .

— Then, we define the bottleneck capacity,

$$C(P) = \min_{e \in P} C_f(e) \quad \text{in residual graph.}$$

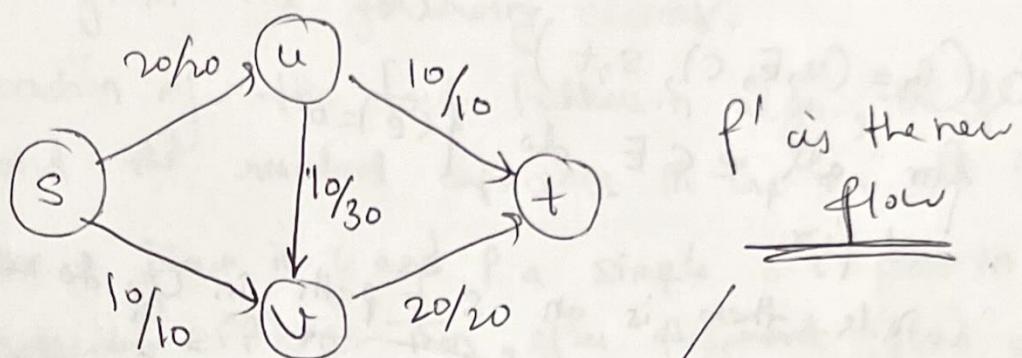
— Then, we update f in G as follows

$\forall e \in P, \rightarrow$ If e is FORWARD, $f'(e) = f(e) + C(P)$

\rightarrow If e is BACKWARD, $f'(v, u) = f(v, u) - C(P)$

so, after updating,

G becomes G' \Rightarrow



\rightarrow If $e \notin P$, $f'(e) = f(e)$.

Proof \Rightarrow

we need to prove that both constraints are valid,

① capacity constraint \rightarrow we don't exceed $c(e)$.

② flow conservation \rightarrow $f^{in}(v) = f^{out}(v)$

Pseudocode of subroutine Augment

I/P \Rightarrow a flow f , and an augmenting path P in G_f

O/P \Rightarrow augmented flow f' .

Augment (f, P) \Rightarrow

for each edge $(u, v) \in P$ do,

if $e = (u, v)$ is a forward edge then

$$f(e) = f(e) + c(P)$$

else

$$f(v, u) = f(v, u) - c(P)$$

endif

end for return $f \rightarrow$ new flow

Ford-Fulkerson Algo

FF ($G = (V, E, c)$, s, t)

for all $e \in E$ do $f(e) = 0$
end for

while there is an $s-t$ path in G_f do

Let $P =$ simple $s-t$ path in G_f

$f' =$ augment (f, P)

Set $f \leftarrow f'$

Set $G_f \leftarrow G_f'$

every vertex appears at-most once!

↑
finishes
in
some
finite
steps!

end while

Return f .

Let $U = \max_{e \in E} c_e$

~~max f~~

$$\therefore |f| = f^{\text{out}}(s) \leq \sum_{e \text{ out of } s} c_e \leq U$$



$$\max |f| \leq nU$$

this seems a finite upper bound?

(If we atleast inc. flow every time, then atleast nU steps will be needed).

Thus, algo ends after finite # steps

increases by an integer amount after Augment (f, P) of

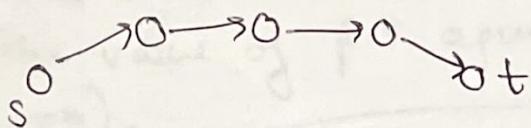
It follows from the following claims.

During execution of the Ford-Fulkerson algo., the flow values $\{f(e)\}$ and the residual capacities in G_f are all integers.

→ Let f be a flow in G and P a simple s-t path in G_f with residual capacity $c(P) > 0$. Then, after Augment (f, P)

$$|f'| = |f| + c(P) \geq |f| + 1 \quad \boxed{3}$$

~~Let f be~~



(s, t) is FORWARD.

$$f'(s, u) = f(s, u) + c(P) \xrightarrow{\text{capacity of the path } "P"}$$

$$|f'| = |f| - f(s, u) + f'(s, u)$$

$$= |f| + c(P) \geq |f| + 1.$$

↳ bcz ~~capacity~~ is strictly +ve.

⑩ Running time of Ford-Fulkerson :-

③ guarantees at most nU iterations.

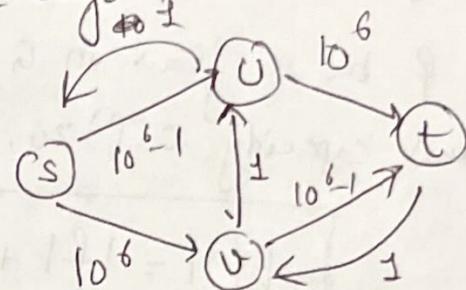
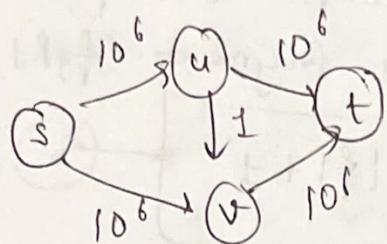
$O(mn)$ to create G_f using adjacency list representation
 $O(m+n)$ to run BFS/DFS to find augmenting path.
 $O(n)$ for Augment (f, P) since P has at most $n-1$ edges.
∴ one iteration requires $O(m)$ time.

∴ Running time of Ford-Fulkerson is $O(mnU)$. Ans

Pseudo-polynomial algo

- * An algo is pseudo-poly algo if it is polynomial in the size of the I/P when the numeric part of the I/P is encoded in unary.

→ FF's pseudo-polynomial time algo.



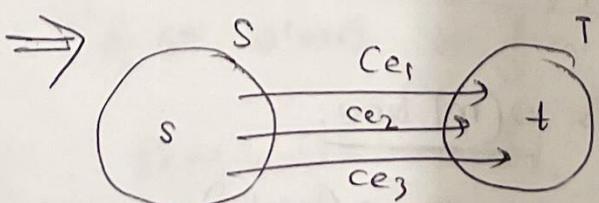
U → more capacity

Here the degree of the capacity affects the running time in an exponential way.

U	Size of U (# digits to describe)	# iterations
1	0	2
10	1	$2 \cdot 10$
100	2	$2 \cdot 10^2$
1000	3	$2 \cdot 10^3$
10000	4	$2 \cdot 10^4$

- * A natural upper bound for max value of a flow?

→ An s-t cut (S, T) in G is a bipartition of the vertices into two sets $S \subseteq T$, s.t. $s \in S$ and $t \in T$.



The capacity $c(S, T)$ of s-t cut (S, T) is

$$c_{\text{out}}(S) + c_e$$

Then intuitively:-

$$\max_f |f| \leq \min_{\substack{(S, T) \text{ cut} \\ \text{in } G}} c(S, T)$$

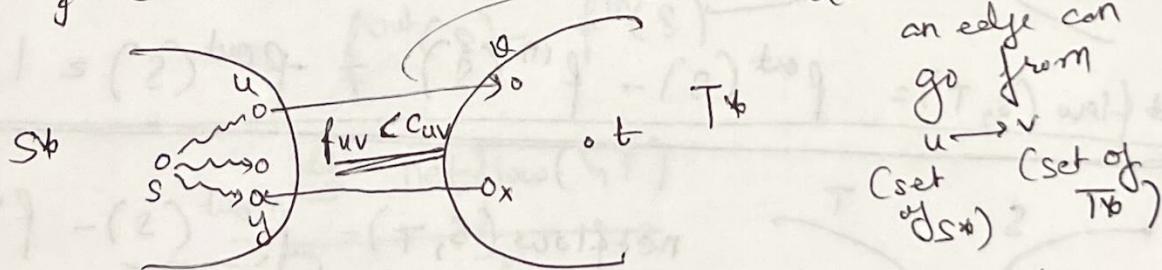
* Roadmap to prove optimality of FF:-

$|f| = f^{\text{out}}(S)$, we will prove the following:-

- (1) for any flow f , the value of the flow $|f|$ cannot exceed the capacity of any cut in G .
- (2) Let f^* be the flow upon termination of the FF algo. We will exhibit a specific cut (S^*, T^*) s.t. the value of f^* equals the capacity of (S^*, T^*) .

$$|f^*| = c(S^*, T^*)$$

So, we conclude that f^* is max flow.
 $\& (S^*, T^*)$ is a cut of minimum capacity.



Let G_{f^*} be the residual graph upon termination of FF
and f^* the flow upon termination of FF.

S^* = set of nodes reachable from S in G_{f^*} .

(S^*, T^*) is an $s-t$ cut in G_{f^*} . \rightarrow both has some set of vertices!

$\Rightarrow \nexists (u, v) \text{ in } G_{f^*} \text{ with } u \in S^*, v \in T^*$.

$\rightarrow \text{If } (u, v) \in G \text{ with } u \in S^*, v \in T^*, \text{ then } f(u, v) = c_{uv}$

$\rightarrow \text{If } (x, y) \in G, \text{ with } x \in T^*, y \in S^*$.

Data:-

$$\text{netflow}(s, t) = f^{\text{out}}(s) + f^{\text{in}}(s)$$

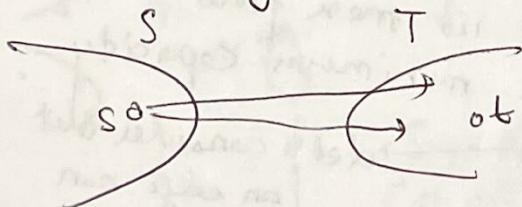
$$f^{\text{out}}(s) = \sum_{e \text{ out of } s} f(e)$$

$$f^{\text{in}}(s) = \sum_{e \text{ into } s} f(e)$$

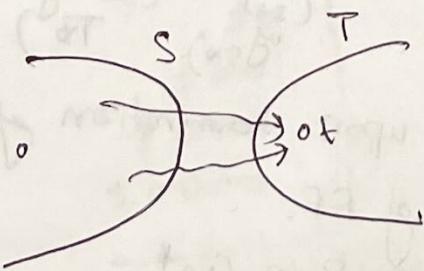
$$\text{netflow}(s^*, t^*) = \sum_{e \text{ out of } s^*} f(e) - \sum_{e \text{ into } s^*} f(e)$$

$$\Rightarrow \sum_{e \text{ out of } s^*} c_e \Rightarrow C(s^*, t^*)$$

flow is netflow(s, t) related to $|f|$,



$$\text{netflow}(s, t) = f^{\text{out}}(s) - f^{\text{in}}(s) = f^{\text{out}}(s) \geq |f|$$



$$\begin{aligned} \text{netflow}(s, t) &= f^{\text{out}}(s) - f^{\text{in}}(s) \\ &= f^{\text{in}}(t) \\ &\geq |f| \end{aligned}$$

LEMMA

$$\text{netflow}(s, t) = |f|$$

\forall s-t cut (S, T), $\forall f$

COROLLARY

$$|f| = \text{netflow}(s, t) = f^{\text{out}}(s) - f^{\text{in}}(s) = \sum_{e \text{ out of } s} c_e = C(s, t)$$

Proof of Zemmer

\forall s-t cut (S, T) $\forall f$,

$|f| = \text{netflow}(S, T)$

Proof \rightarrow

$|f| = f^{\text{out}}(S) - f^{\text{in}}(S)$

To prove

$$= \sum_{v \in S} f^{\text{out}}(v) - f^{\text{in}}(v)$$

$$= \sum_{v \in S} \left(\sum_{(v, u) \in E} f(v, u) - \sum_{(u, v) \in E} f(u, v) \right)$$

$$= \sum_{\substack{(v, u) \in E \\ v \in S, u \in T}} f(v, u) - \sum_{\substack{(u, v) \in E \\ u \in T, v \in S}} f(u, v)$$

$$= \sum_{e \text{ out of } S} f(e) - \sum_{e \text{ into } S} f(e)$$

$$= f^{\text{out}}(S) - f^{\text{in}}(S)$$

$$\triangleq \text{netflow}(S, T)$$

Therefore, proved

* The max-flow min-cut theorem:-

THEOREM:- If ' f ' is an s-t flow such that there is no s-t path in G_f , then there is an s-t cut (S^*, T^*) in G such that $|f| = C(S^*, T^*)$,

$\therefore f$ is a max-flow and (S^*, T^*) is a cut of min capacity.

Max cut - min flow theorem:-

In every flow network, the max^m value of an s-t flow equals the min^m capacity of an s-t cut.

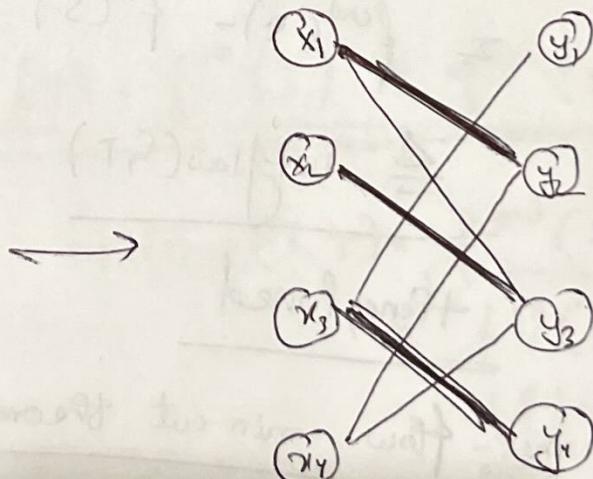
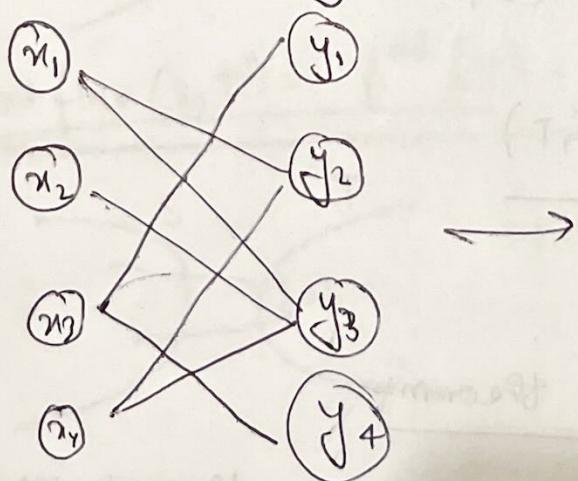
~~and~~ we can efficiently find a min path because we already have the residual graph G_f and can form $S*$ using BFS/DFS.

* Integrality Theorem:-

→ In ff, the flow values & residual capacities in G_f are I.

If all capacities in a flow N/W are ints, then there is a max^m flow for which every flow value $f(e)$ is an integer.

* Bipartite Matching :-



$\left. \begin{array}{l} \text{This is a matching of size 3} \\ M = \{(x_1, y_1), (x_2, y_2), (x_3, y_3)\} \end{array} \right\}$

⇒ A matching M is a subset of edges where every vertex in $X \cup Y$ appears at most once.

We can get a matching of size ≤ 4 here, nothing more since these matchings are one-to-one.

Perfect Matching:- Every vertex in $X \cup Y$ appears exactly once.

Not always possible, e.g.: $|X| = |Y|$.

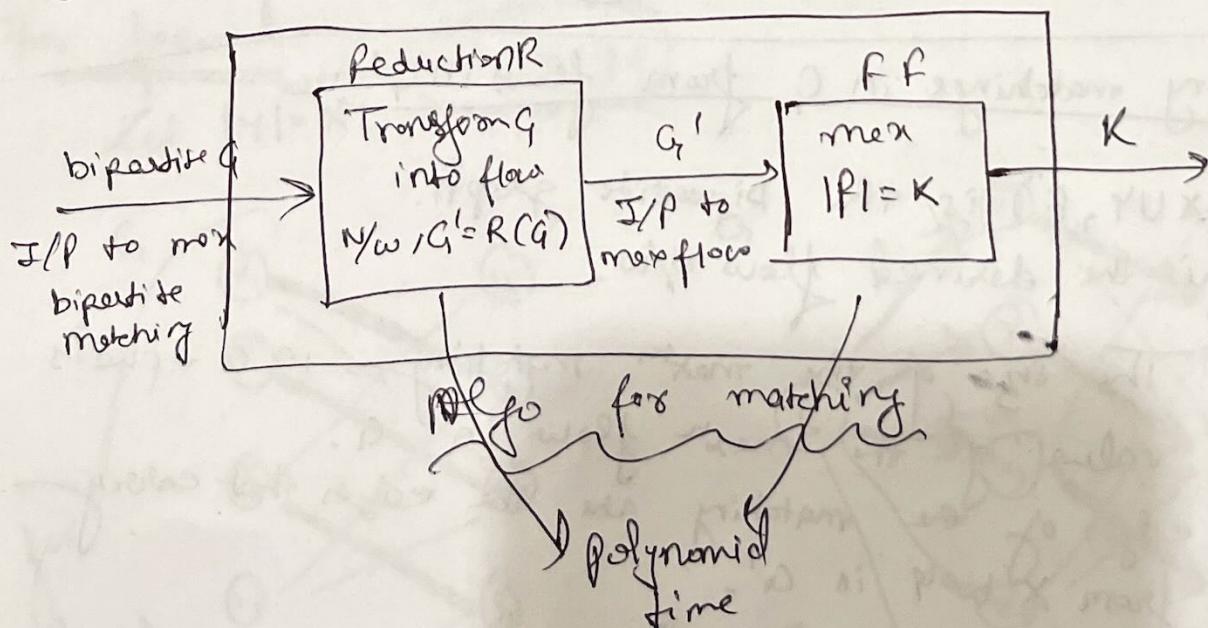
Maximum Matching desirable in applications.

→ If we had an algo to find max^m matching then we could also find a perfect matching, if one exists.

*) finding max matchings in bipartite graphs:-

IDEA:- Use FF algo to find max (or perfect) matching.

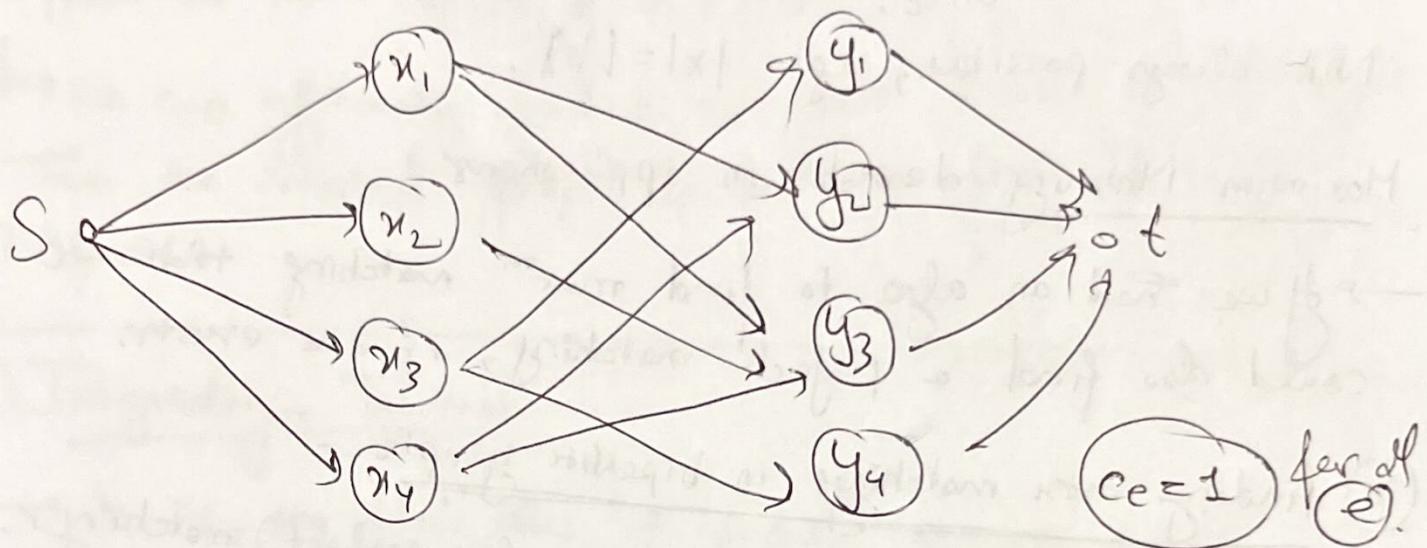
STRATEGY:- reformulate the problem as a max-flow problem which we know how to solve (reduction). To this end, we need to transform our input bipartite graph into a flow N/W.



→ We need to prove the following for correctness:-

- ① Reduction R must be efficient (polynomial in the size of G)

(2) G and G' should be equivalent, in the sense that
 G has a max matching of size k iff the max flow
in G' has value k .



bcoz it is the smallest value ($\text{int} > 0$) to use to

keep it easy. Also, running FF depends on V ,
then, it will be the most efficient.

(V) Computing matchings in G from flows in G' :-

$\rightarrow G = (X \cup Y, E)$ is the bipartite graph.

$\rightarrow G'$ is the derived flow N/W.

CLAIM :- The size of the max^m matching in G equals the value of the maxⁿ flow in G' .

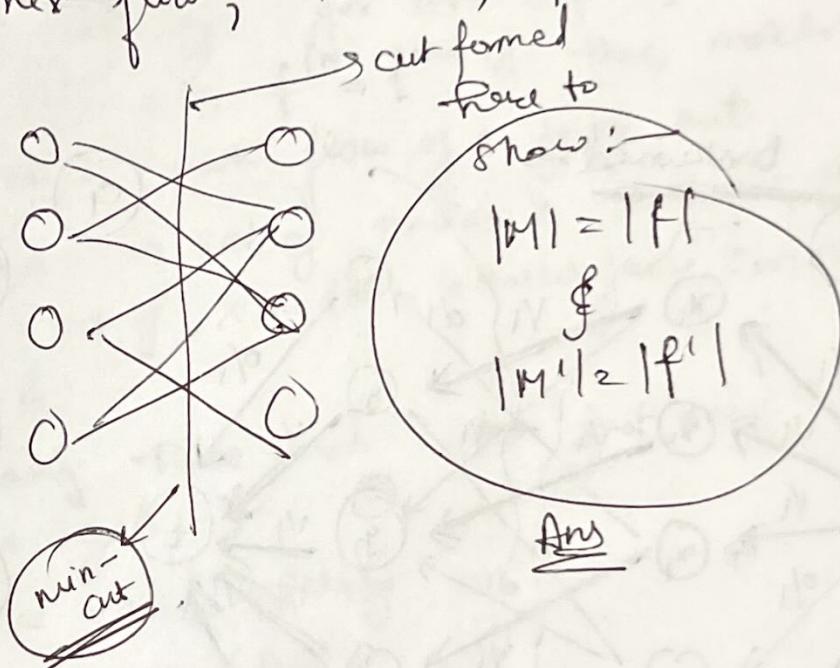
The edges of the matching are the edges that carry flow from X to Y in G' .

(Equivalence of instances)

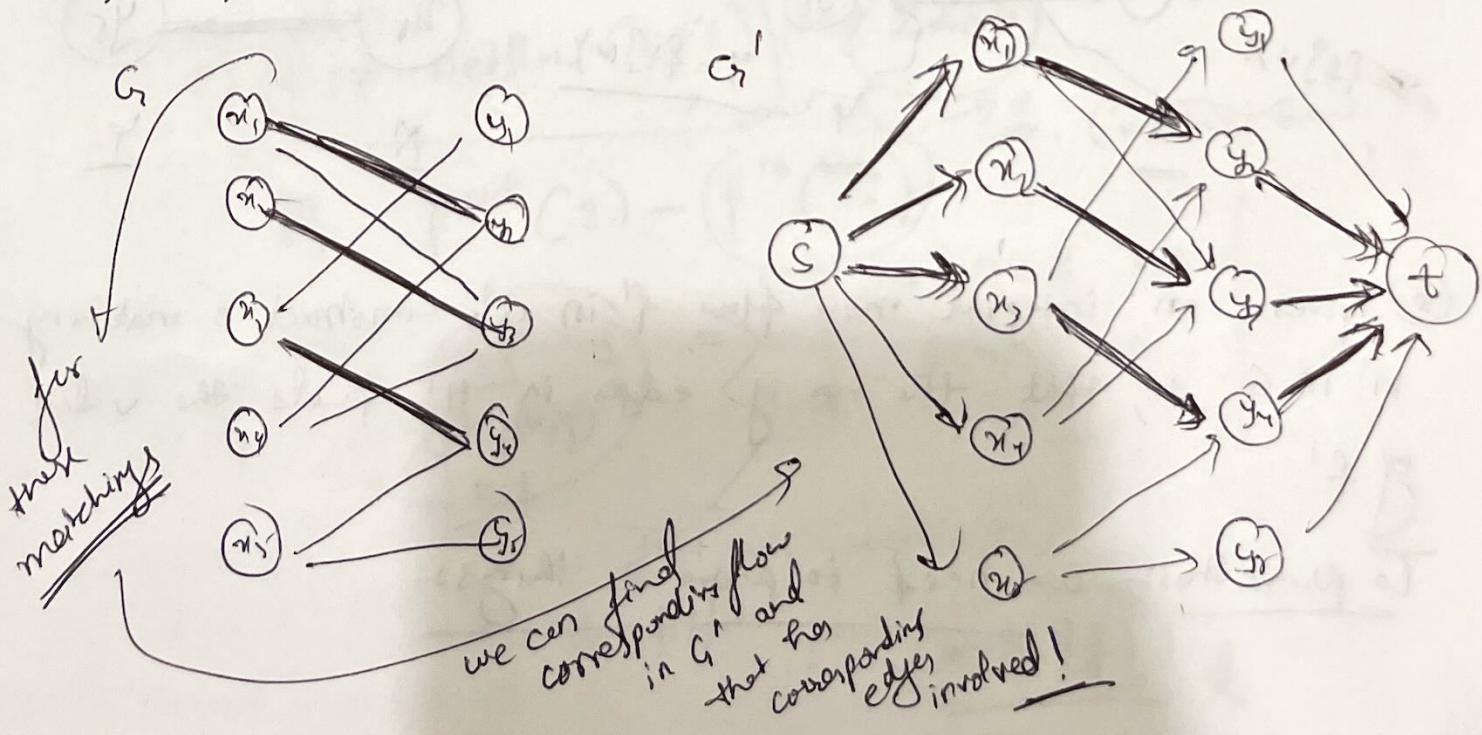
Proof

forward direction \Rightarrow for any matching M in G , we can construct a flow f in G' with value equal to the size of M , that is $|f| = |M|$.

(2) Reverse direction \Rightarrow Given a max flow f' in G' , we can construct a matching M' in G , with size equal to the value of the max flow, that is, $|M'| = |f'|$.



① for forward \rightarrow matching size
Let $|M| = k$,



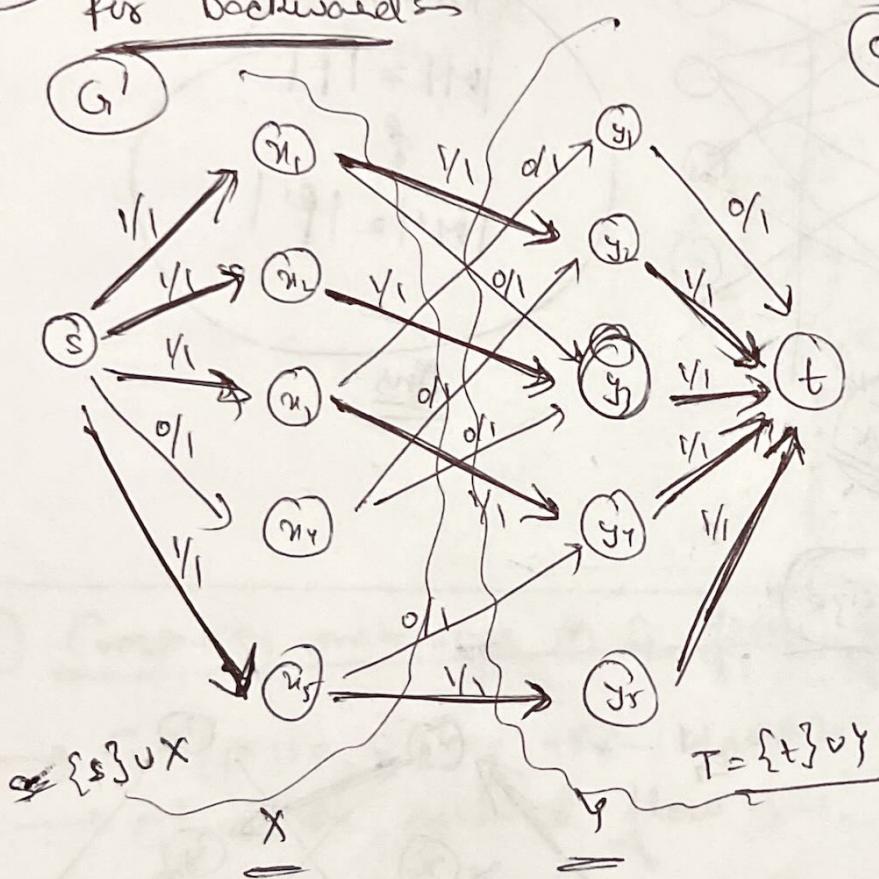
Thus, given a matching M (darkened edges in G), we can construct an integral flow f in G' s.t. the value of f equals the no. of edges in M .

$$|f| = |M|$$

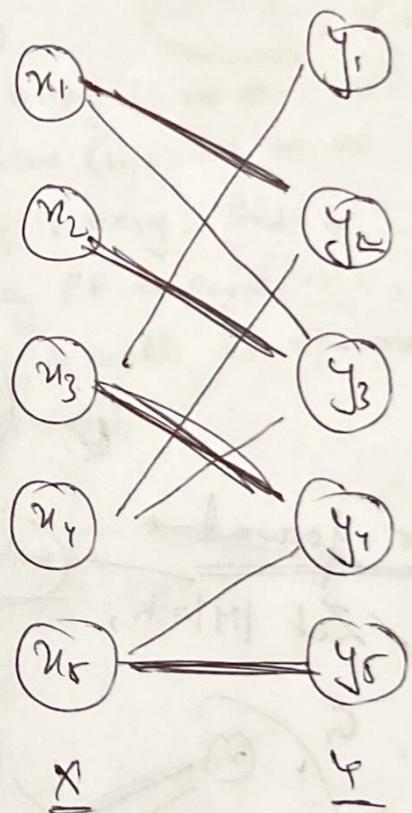
Hence, forward dirⁿ equivalence of instance is proved

②

for backwards



③



- ④ Given an integral max flow f' in G' , construct a matching M' in G so, that the no. of edges in M' equals the value of f' .

To prove this we need to prove 2 things:

first we need to prove that M would be a valid matching that we are getting thru the flow network.

for that, M can be formulated as follows (to get a matching):

$$M = \{ (x, y) \in E \text{ s.t. } f(x, y) = 1 \}$$

for this,

the matchings that we get, due to flow conservation, we can't have 2 matchings for let's say x_1 as a flow of 1 comes into (x_1) , so, $f^{in} = f^{out} = 1$ thus, matching found cannot have 2 edges with flow of 1 going out.

Similarly, any node in set of nodes Y cannot have 2 matchings from set of X nodes otherwise flow conservation will be invalid.

Thus, we can only have each vertex just once in the matchings that we get.

\therefore It is a valid matching.

\therefore Now, we need to prove that $|M| = |f|$

$$|M| = |f|$$

$$|f| = \text{netflow}(S, T) \rightarrow \text{Lemma}$$

$$\Rightarrow f^{out}(S) - f^{in}(S)$$

$$\Rightarrow \sum_{(x,y) \in E} f(x,y) = 1 \Rightarrow \begin{cases} (x,y) \in E \text{ s.t.} \\ f(x,y) = 1 \end{cases}$$

This was actually M

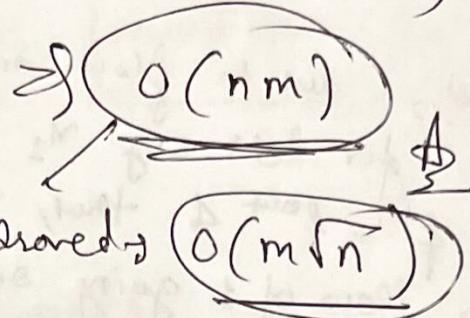
$$\Rightarrow |M| = |f| \therefore |M| = |f|$$

∴ Hence, we have proved the equivalence!

Thus is a poly-time algo.

Reduction $\rightarrow O(n+m)$

FR $\rightarrow O(nm)$



\Rightarrow If we convert this into a decision problem!

BPM: Does a bipartite graph have a perfect matching?

I/P \rightarrow a bipartite graph $g = (X \cup Y, E)$.

O/P \rightarrow yes, iff g has a matching of size $|X| = |Y|$.

How to do this?

Terminology \rightarrow

\rightarrow An instance of BPM is a specific I/P graph g .

\rightarrow we may think of an I/P instance g as a binary string x encoding the graph g , with length $|x|$ bits.

\rightarrow we assume reasonable encodings: e.g. a binary number b requires a logarithmic no. of bits to be encoded.

\rightarrow reasonable encodings are related polynomially.

\rightarrow An algo. that solves BPM admits:-

\rightarrow I/P:— a binary string x encoding a bipartite graph

Hashing, Bloom filters

→ A data structure maintaining a dynamic subset S of a huge universe U .

$$|S| \ll |U|$$

support:

→ efficient insertion

→ efficient deletion

→ efficient search

holds all those → { dictionary }

Operations supported

① Create: Initialize a dict with $S = \emptyset$.

② Insert: add x to S , if $x \notin S$

→ Additional info about x might be stored in the dict as part of a record for x .

③ Delete(x): delete x from S , if $x \in S$.

④ Lookup(x): determine if $x \in S$.

for IP addresses

$$|U| = 2^{32}$$

$$|S| = 250$$

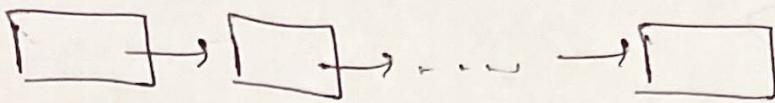
→ sample size of IP addresses.

* If we have an array for storing the bits as booleans, then time complexity is very good, but unrealistic to store it. (Space considerations are huge!).

lookup $\rightarrow O(1)$ ← Insert
Delete.

Space $\rightarrow O(|U|)(X)$

* If Linked List,



lsl nodes

Space $\rightarrow O(lsl)$

Lookup $\rightarrow O(lsl)$

Insert

Delete

So, space complexity gets better, but time goes higher!

Thus, we are trying to both optimize time & space.

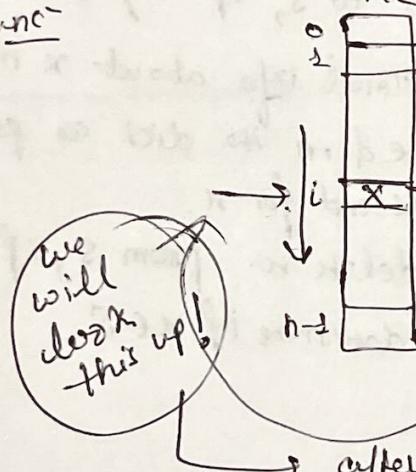
$$lsl \ll lul$$

$$lsl = n$$

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

hash func

$H \in \text{hash table}$



$x \in S$

$$h(x) = i$$

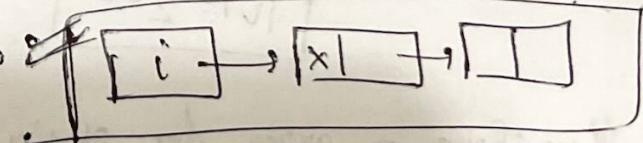
ptr to a linked list

$$h(y) = i$$

for dealing with collisions!

after hashing it!

Hashing as chaining



Thus, what would be the time required?

$\Rightarrow \text{Lookup}(x) :=$

(1) Compute $h(x) \rightarrow O(1)$ time (might be diff!)

(2) Traverse the linked list at $h(x)$ to find x .

Time :- $O(\# \text{elements in linked list at } h(x))$

it should spread out elements as much as possible

insert & delete are constrained by Lookup!

⇒ Simple hash fn might not work!

I) (2 simple hash func' that hash an IP address n from

$$\{0, \dots, 2^{32}-1\} \text{ to } \{0, \dots, 255\}.$$

→ assign the last 8 bits of n as $f(n)$.

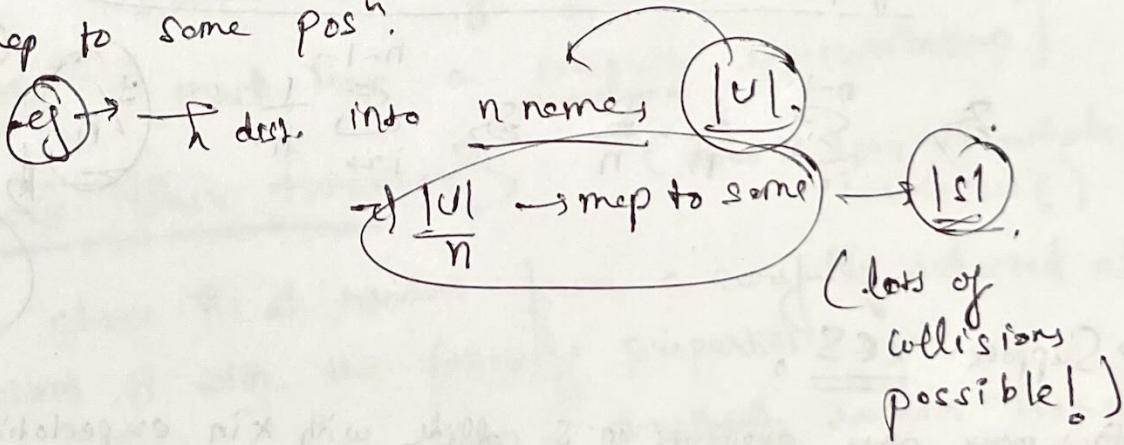
→ assign the 1st 8 bits of n as $f(n)$.

Nothing is inherently wrong ⇒ the problem is that our 2³² IP address might not be drawn uniformly at random from among 2³² possibilities.

ii) No single hash func' can work well on all data sets.

FACT:- for any fixed (deterministic) $h: U \rightarrow \{0, 1, \dots, n-1\}$

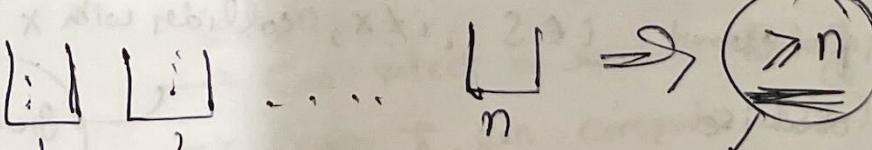
where $|U| > n^2$, there exists some set S of n elements that all map to same posⁿ.



2) There are n^2 balls into n bins in any way!

→ what is lower bound on the max load of any bin?

Proof
(n
collisions
always!)



to optimize max
across all the bins!

for any bin, to optimize,
can at least give n balls to all the bins.

→ we are going to introduce randomness here to decrease
collisions!

Completely random hash fn:

$\forall x \in U, \forall 0 \leq j \leq n-1$
independently.

$$\boxed{\Pr[h(x)=j] = \frac{1}{n}}$$

for $x, y \in U$,

$$\Pr[h(x)=h(y)] \Rightarrow \Pr[h(x)=i \wedge h(y)=i]$$

$$\sum_{i=0}^{n-1} \Pr[h(x)=i \wedge h(y)=i]$$

$$\Rightarrow \boxed{\sum_{i=0}^{n-1} \Pr[h(x)=i] \cdot \Pr[h(y)=i]}$$

$$\Rightarrow \sum_{i=0}^{n-1} \frac{1}{n} \cdot \frac{1}{n} \Rightarrow \sum_{i=0}^{n-1} \frac{1}{n^2} = \frac{1}{n}$$

(when collision would happen)

probability for that?

\Rightarrow Suppose $x \in S$,

How many other elements in S collide with x in expectation?

Let $X = \#$ other elements from S that collide with x .

$$\underline{\mathbb{E}[X]} = ?$$

$X_i = \begin{cases} 1, & \text{if element } i \in S, i \neq x, \text{ collides with } x \\ 0, & \text{otherwise.} \end{cases}$

Indicator func

$$\boxed{X = \sum_{i=1}^{n-1} X_i}$$

$$E[X] = E\left[\sum_{i=1}^{n-1} X_i\right] = \sum_{i=1}^{n-1} E[X_i] = \sum_{i=1}^{n-1} \Pr[X_i=1]$$

$\Rightarrow \frac{n-1}{n} < 1$

\therefore Size of linked list is const. (w.r.t the collisions)

(\because it works in const time) \rightarrow time is const. in expectation!
~~As space also is.~~

* What completes the problem?

\hookrightarrow Hash assigns a random value! (it's not a func as we don't know where to store it).

{ Randomness is valuable but is not feasible! }

\rightarrow So, we will randomize the choice of hash function from a suitable class of funcs into $[0, n-1]$.

~~func~~ $\{ h \text{ must have a compact representation} \}$.

\Rightarrow Universal hash function \rightarrow (does it span the whole universe?)

Idea: choose h at random from a carefully selected class of functions H with the following properties:-

i. h behaves almost like a completely random hash func.
 \rightarrow for $x, y \in U$. The probability that a randomly chosen $h \in H$ satisfies $h(x) = h(y)$ is at most $\frac{1}{n}$.

\rightarrow Can select a random h efficiently.

\rightarrow Given h , can compute $h(x)$ efficiently.

\Rightarrow Such hash funcs are called universal; their design relies on number theoretical facts.

(8)

$$U \Rightarrow \{a, b, c, d, e, f\}$$

$n = 2$ (range of hash fn).

$k=2$ (no. of hash function).

$H = \{\tilde{h}_1, \tilde{h}_2\} \rightarrow$ universal class of hash funcs?

Is this?

$$\tilde{h}_1: U \rightarrow \{0, 1\}$$

$$\tilde{h}_2: U \rightarrow \{0, 1\}$$

	a	b	c	d	e	f
\tilde{h}_1 :	0	0	0	1	1	1
\tilde{h}_2 :	0	1	0	1	0	1

$$\Rightarrow \text{Prob. of } \tilde{h}(x) = \tilde{h}(y) \Rightarrow \Pr[\underline{x}] \Rightarrow \frac{1}{n} \quad \boxed{\Pr_{\tilde{h}}[\tilde{h}(x) = \tilde{h}(y)] \leq \frac{1}{n}}$$

fix $a, b,$

$$\Pr_{\tilde{h}}[\tilde{h}(a) = \tilde{h}(b)] \Rightarrow$$

$$\Pr[\tilde{h}_1(a) = \tilde{h}_1(b) | \tilde{h}_1] \cdot \Pr[\tilde{h}_1]$$

$$+ \Pr[\tilde{h}_2(a) = \tilde{h}_2(b) | \tilde{h}_2] \cdot \Pr[\tilde{h}_2]$$

$$\Rightarrow 1 \cdot \frac{1}{2} + 0 \cdot \frac{1}{2} \Rightarrow \boxed{\frac{1}{2}}$$

$$\Pr_{\tilde{h}}[\tilde{h}(a) = \tilde{h}(b)] \Rightarrow \frac{1}{2} \leq \frac{1}{n} \quad \text{when } n=2,$$

$k \leq k$ → this holds!

→ This should hold for all pair of values!

$x \in \{a, c\}$

$$\Pr_{f_1} [f_1(a) = f_1(c)] \rightarrow$$

$$\Pr [f_1(a) = f_1(c)]. \Pr [f_1]$$

$$+ \Pr [f_2(a) = f_2(c)]. \Pr [f_2]$$

$$\Rightarrow \frac{1}{2} \cdot \frac{1}{2} + \frac{1}{2} \cdot \frac{1}{2} = \textcircled{1} \leq \frac{1}{n} \textcircled{X}.$$

(doesn't hold)

→ This is not universal class of hash func.

Hence proved

Ans.

~~April~~ we will have to assume a completely random hash fn exists.

(~~start bcos~~ this provide us a good rough idea of how hashing schemes perform in practice)

(in reality, there are no completely random hash fns).

Remark → $h(x)$ is fixed for every x : it just takes one of the n possible values with equal probability.

$(x, h(x))$

listened as pair!

→ Hashing modeled as a balls and bin problem

Q: How many elements can we insert in the hash table before it is more likely than not that there is a collision?

Ans. At least 1 entry in the table that has more than 1 ball (in any of the bins).

$$\Rightarrow \boxed{\Omega(\sqrt{n})}$$

(just an occupancy problem),

hashing

- ① balls correspond to elements from V .
- ② bins are slots in the hash table.
- ③ each ball falls into one of the n bins independently & with prob. $\frac{1}{n}$. (uniform distribution).

→ Analyzing time / space of hash table →

Q What is the expected time for $\text{lookup}(x)$?

$$\Rightarrow \underbrace{O(1)}_{\substack{\text{calculate} \\ \text{hash func}}} + \underbrace{O(1)}_{\substack{\text{scan the} \\ \text{linkedlist} \\ \text{at that} \\ \text{position.}}} \Rightarrow \boxed{O(1)}$$

Calculate hash func scan the linkedlist (bcz we will take completely random hash fn).

Q What is the expected wasted space in the hashed table?

(how many entries are empty) $\Rightarrow \boxed{n/3}$ ($\geq \frac{n}{3}$)
 (expected no. of empty bins).

(hashing with chaining).

Q worst case time for $\text{Lookup}(x)$?

$$\Rightarrow \boxed{O(n)}$$

(all items might be mapped at the same locⁿ).

worst-case with highest prob. case, (w.h.p.)

$$\geq 1 - \frac{1}{n}$$

(independent

$$\Rightarrow O\left(\frac{\ln n}{\ln(\ln n)}\right) \Rightarrow O(\ln n)$$

uniformly
at
random!)

capacity of
highest loaded
bin!

Ave.

cons of hashing to save space \Rightarrow

\Rightarrow Password checker:-

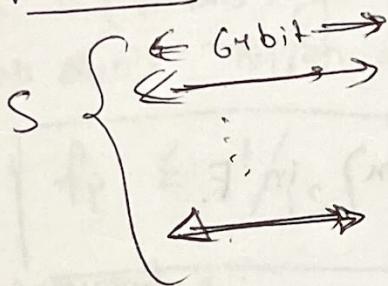
\Rightarrow maintain a dictionary for a set S of 2^{16} bad passwords, so that when a user tries to set one, we check as quickly if it belongs to S and reject it.

\rightarrow Assume:- p/w consists of 8 ASCII chars.
(rep. 8 bytes (64 bits) to represent password)

Q how to store the dictionary?

$\Rightarrow S$ at 2^{16} bad passwords

$\Rightarrow \boxed{NES}$



how to efficiently search?

① Sort S .

② Binary search for X in S .

③ Time:- $O(\log |S|)$

④ Space:- $|S|.64$ bits.

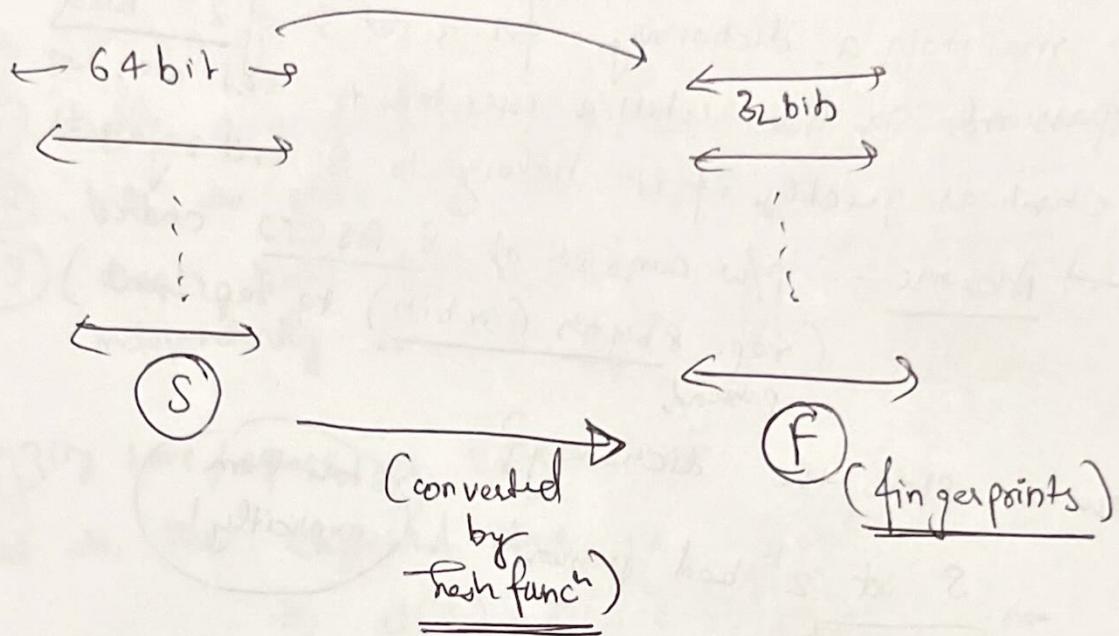
\Rightarrow Now, we want to use less space!

(We could potentially allow for some errors).
 \rightarrow we might reject some good passwords, but never accept a bad password!

Okay to have FPs; certain about PNS.

IDEA \Rightarrow Map passwords to fingerprint.

$$h: \{0, \dots, 2^{64}-1\} \rightarrow \{0, \dots, 2^{32}-1\}.$$



- ① Sort F
- ② Compute $h(x)$
- ③ Binary search for $h(x)$ in F .
- ④ Let's say $h(x) \in F$

How to answer is x in S ?

\Rightarrow $x \in S$ (we will find it in F and reject ~~or~~ accordingly)

$h(x) \notin F$

Answer is $x \notin S$

(correct), so we do not reject!

\rightarrow Time to answer the query?

Same as before!

$O(\log_2 |S|)$

\rightarrow Space := $O(|S|, 32)$ \rightarrow (32 bits fingerprint)

Correctness:

If $x \notin S$ as we do not find that,

* FNs are not possible! (as we are still storing all hashes of n).

② FPs might happen. (as smaller range of b, T strings stored, collisions might be more).
(as info. is less)

$$\therefore P(FPs) \neq 0$$



Then, we have reduced space without FNs and $P(FPs) \neq 0$.
(we want $P(FPs)$ as less as possible as well,
so, we set no. of bits in f accordingly). Ans

Q. what is b (# bits) for each fingerprint s.t.
FNs don't happen & prob. of FPs $\neq 0$ but less.

$$fp \leq 1/|S|$$

fix $x \rightarrow$ good password x .

$$fp = \Pr[\text{h}(x) \text{ collides with } \text{h}(y) \text{ for some "bad" } y]$$

\rightarrow fix a "bad" y ,

$$\Pr[\text{h}(x) = \text{h}(y)] = \frac{1}{n} = \frac{1}{2^b}$$

(for completely random hash function)

$$\Pr[\text{h}(x) \neq \text{h}(y)] \Rightarrow 1 - \frac{1}{2^b}$$

(differs from all
those fingerprints.)

$$\Pr[\text{at least } \text{h}(x) \neq \text{h}(y)] \Rightarrow \left(1 - \frac{1}{2^b}\right)^{|S|}$$

$f_p \Rightarrow \Pr [\exists \text{ body, s.t. } h(x) = h(y)]$

 $\Rightarrow 1 - \left(1 - \frac{1}{2^b}\right)^{|S|}$

So,

$$1 - \left(1 - \frac{1}{2^b}\right)^{|S|} \leq \frac{1}{|S|}$$

Solving for b gives,

$$b = 2 \log_2 |S| \quad \underline{\text{Ans}}$$

Space $\Rightarrow O(|S| \log |S|)$ bit Ans

2) Improved $f_p \Rightarrow$

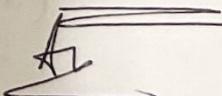
$$b = 2 \log_2 m$$

Derivation original $f_p \Rightarrow$

$$\Rightarrow 1 - \left(1 - \frac{1}{2^b}\right)^m$$

$$\Rightarrow 1 - \left(1 - \frac{1}{m^2}\right)^m \leq 1 - \left(1 - \frac{1}{m}\right) = \frac{1}{m}$$

Thus, if our dict has $|S| = m = 2^{16}$ bad words, using a hash func' that maps each of the m perm to 32 bits yields a false pos. prob of about $1/2^{16}$



* we would use bloom filter (for checking bad URL)

Q:- a large set S , and queries of the form "Is $x \in S$?"
we want a data structure that answers the queries:-

→ fast (faster than searching in S)

($O(1)$)

→ is small (smaller than the explicit reprⁿ provided by hash table).

($O(1)$)

{ 1) FPS:— answer yes when $x \in S$. } \rightarrow const fp probability.
2) FNS:— answer no when $x \notin S$.

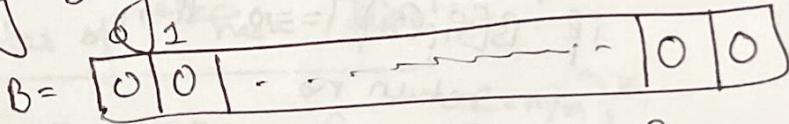
(we need
this!)

O/P:— small prob. of FPs, no FNs.

→ Bloom filters \Rightarrow (used by Google for finding bad URLs).

consists of :-

1) an array B of n bits, initially all set to 0.



2) K independent random hash funcⁿ f_1, \dots, f_K with range $\{0, 1, \dots, n-1\}$

A basic bloom filter supports:-

→ Insert(x)

→ Lookup(x)

\Rightarrow Pseudocode:-

Setup Bloom filter (S, f_1, \dots, f_K):

Initialize arr B of size n to all zeros.

for $i=1$ to m do:-

 Insert(x_i)

end for

Insert(x):

$$|S|=m$$

for $i=1$ to K :

compute $\text{f}_i(x)$

set $B[\text{f}_i(x)] = 1$

end for

~~Ampl~~ An entry of B may be set multiple times; only the first change has an effect.

$$\text{f}_i: U \rightarrow \{0, 1, \dots, n\}$$

Lookup

To check membership of an element x in S do:

lookup(x):-

for $i=1$ to K do:

compute $\text{f}_i(x)$

if $B[\text{f}_i(x)] = 0$ then

return no

return yes

→ If $B[\text{f}_i(x)] \neq 1$, for some i , then $x \notin S$.

→ Otherwise, answer " $x \in S$ ".

↳ might be false since we could get the same values in bloom filter mapped due to other elements which were actually not part of the same element x_i (could be comb of multiple n_i 's).

↳ false negative not possible since we never go back to 0, so, we won't lose any info. & we have all the info. for the complete set S . (which covers all the info. for finding an x) → thus, a kind of a filter!

Probability of false positive:-

→ After all elements from S have been hashed into the bloom filter, the probability that a specific bit is still 0 is:-

$$\left(1 - \frac{1}{n}\right)^k \approx e^{-km/n} = p$$

→ To simplify the analysis, assume that the prob. that a specific bit is still 0 is exactly p .

→ The prob. of a false positive is the probability that all k hashes evaluate to 1:-

$$f = (1-p)^k$$

Optimal number of hash funcⁿ:-

$$f = (1-p)^k = (1 - e^{-km/n})^k$$

~~more hash~~
~~more fp less~~
~~K↑, fp↓~~ → Trade-off b/w k & p : using more hash functions → gives us more chances to find a 0 when $x \notin S$, but we need to reduce sparsity of the word B. (bloom filter).
 → Compute optimal no. K^* of hash functions by minimizing f as a funcⁿ of K :

$$\begin{aligned} n &\rightarrow |B| \\ m &\rightarrow |S| \end{aligned}$$

$$K^* = (n/m) \cdot \ln 2$$

Then the false positive probability is given by :-

$$f = (1-p)^{K^*} \approx (0.6185)^{n/m}$$

Ans

Big savings in space

n/m bits

for ex - set $n = 8m$. Then $k^* = 6$ and $f \approx 0.02$

(thus, constant space gives
fp quite less)

6 no FNS(✓)

→ small constant false pos. prob. by using only 8 bits
(1 byte) / element of S ; independently of the size of S !

Ans. space / input element \rightarrow 8 bits (thus, const. space)

$$\mathbb{E}[X] = n H_n = n \ln n + n \Theta(1)$$

$$= \boxed{\Theta(n \ln n)}$$

Avg. (super-linear no.)

{ # balls needs
to be thrown! }

Randomized algo. \Rightarrow worst-case (deterministic I/p)

- \Rightarrow On the same input, the algo. produces the same o/p but diff. executions may req. diff. running times.
- \Rightarrow The latter depends on — the random choices of the algorithm.
- \Rightarrow Random samples are assumed independent of each other.

\rightarrow worst-case I/p.

\rightarrow expec. run. time analysis: analysis of the running time of the rand. algo. on a worst case input.

① Las Vegas \Rightarrow

- deterministically
- probabilistically

correct $\xrightarrow{?}$ always returns
fast. $\xrightarrow{?}$ correct

② Monte Carlo

- deterministically fast
- probabilistically correct

Avg.

some executions
may fail to
return the
correct
answer.

→ Rand. algos may be:-

- ① simpler
- ② req less memory of the past
- ③ useful for symmetry-breaking
- ④ generally equally efficient as deterministic algos.

they
why
use

Deterministic algos are a special case of randomized algos.

Randomized Quicksort :-

→ Can we use randomization so that quicksort works with an "avg" I/P even when it receives a worst-case I/P?

⑤ Explicitly permute the I/P (there's no def. way to do that, tough to analyze!)

② Use random sampling to choose pivot: instead of using $A[\text{right}]$ as pivot, select pivot randomly.

we won't often pick the largest/smallest item as pivot (prob very less). Thus, most often the partitioning will be "balanced".

→ Random Partition

{ be random ($\text{left}, \text{right}$) }
swap ($A[b], A[\text{right}]$)
, then we call "Partition"
 $\boxed{\text{random}(i, j)}$

this is
a random
variable!

which is an $O(1)$ op,

(n) \rightarrow expected running time of Randomized Quicksort. (based on Quicksort f Partition).

Partition $\rightarrow O(n)$

Can be called n times in worst-case! (# calls to partition) $\rightarrow \underline{\underline{n}}$

* Total time spent in partition calls outside the for loop $\Rightarrow O(n)$

* Total time spent in partition calls outside the for loop

rv \rightarrow x_i comparisons in the i^{th} partition call
 $\Rightarrow O(x_i)$

$x \Rightarrow$ total # comparisons in all partition calls.

$\Rightarrow O(x)$ time

Running time of Rand. QS \Rightarrow depends on random choices made in the algorithm!

$$\Rightarrow O(n) + O(x)$$

$T(n) \Rightarrow$ exp. runtime of Rand-QS.

$$T(n) = E [O(n) + O(x)] \\ = O(n) + O(E[x])$$

* fix x_i & $x_j \Rightarrow$ each of the two can be compared at most once \rightarrow (either once/never)

\therefore compared ≤ 1 throughout the algo.

Pairs of x_i & $x_j \Rightarrow \leq \binom{n}{2} \cdot 1 \Rightarrow \underline{\underline{O(n^2)}}$
comparisons

\Rightarrow for better behavior, we need to know expected # Comparisons before \Rightarrow

— Assume all I/P items are distinct.

— Relabel the I/P as $\{z_1 < z_2 < \dots < z_n\}$

→ # Comparisons (total)

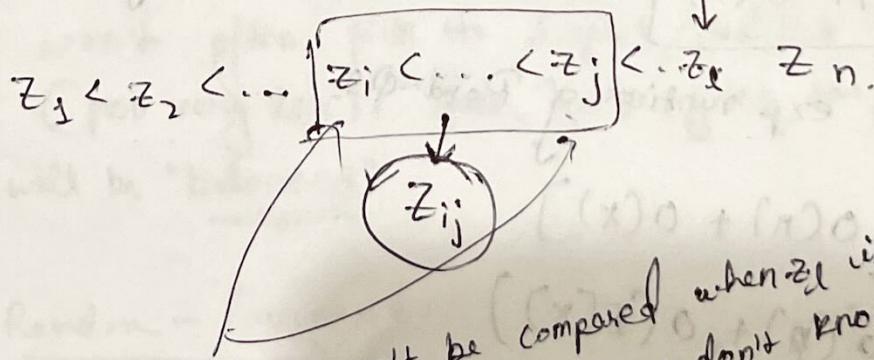
$$E[X] \Rightarrow$$

$\left\{ x_{ij} \right\} \begin{cases} 1, & \text{if } z_i \text{ and } z_j \text{ are never compared} \\ 0, & \text{otherwise.} \end{cases}$

$$X = \sum_{1 \leq i < j \leq n} x_{ij} \quad \xrightarrow{\text{dist. of many r.v's.}}$$

$$E[X] = \sum_{i=1}^n \sum_{j=i+1}^n \Pr[x_{ij} = 1]$$

$$\Pr[x_{ij} = 1] = ?$$



fixed pivot

$z_i < z_2 < \dots \{ z_i < \dots < z_j \} < \dots < z_n$.

z_i and z_j won't be compared when z_l is being considered. Otherwise, we don't know whether they would be!

So, how can we know when z_i and z_j will be compared?

→ This will happen!

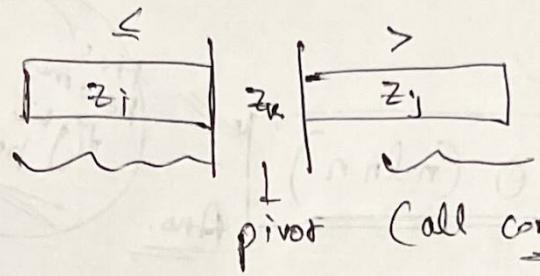
→ Suppose an item from z_{ij} is chosen as pivot.

$\rightarrow z_i = \text{pivot}$

$\rightarrow z_j > \text{pivot}$

Then, they would be compared. ✓

$\Rightarrow z_i < z_k < z_j, z_k = \text{pivot}$ } No, when pivot is being considered.
X (so, never compared, here as well!).



\Rightarrow Comparisons between z_i & z_j ,
 → the first item picked as pivot, item from z_{ij} is z_i/z_j .

So, $P_2[X_{ij} = 1] \Rightarrow P_2[z_i/z_j \text{ as the 1st item picked as pivot from } z_{ij}]$.

NOTE $P_2[\varepsilon_1 \cup \varepsilon_2] = P_2[\varepsilon_1] + P_2[\varepsilon_2] \text{ if } E_1 \cap E_2 = \emptyset$

$$P_2[X_{ij} = 1] = \frac{1}{|z_{ij}|} + \frac{1}{|z_{jj}|} - \frac{2}{|z_{ij}|} = \frac{2}{j-i+1} = P_2[X_{ij} = 1]$$

$$E(x) = \sum_{i=1}^n \sum_{j \geq i}^n \frac{2}{j-i+1}$$

$$\Rightarrow 2 \sum_{i=1}^n \sum_{j \geq i}^n \frac{1}{j-i+1}$$

Let $K = j-i+1$

for $j > i+1 \Rightarrow k = i+1 - j + 1 = 2$.

for $j = 1 \Rightarrow k = n - i + 1$.

$$\therefore \geq 2 \sum_{i=1}^n \sum_{k=2}^{n-i+1} \frac{1}{k}$$

$$\leq 2 \sum_{i=1}^n \sum_{k=1}^n \frac{1}{k}$$

$$\leq 2(n H_n)$$

$$\Theta(n \ln n)$$

H_n

this becomes tight bound?

Ans.

$$X_i = \begin{cases} 1, & \text{if } i\text{-th coin flip is H} \\ 0, & \text{otherwise} \end{cases}$$

$$X = \sum_{i=1}^n X_i \quad \begin{matrix} \text{show many coins} \\ \text{came up as Head.} \end{matrix}$$

$$\Rightarrow E[X] = E\left[\sum_{i=1}^n X_i\right] = \sum_{i=1}^n E[X_i] = \sum_{i=1}^n \Pr[X_i = 1]$$

Linearity of expectation.

\Downarrow
= np

Occupancy problems:-

dist. of balls into bins when m balls are thrown independently & uniformly at random into n bins.

Q How many balls can we throw before it is more likely than not that some bin contains at least two balls?

find K s.t.,

$$\Pr[\exists \text{ bin with } \geq 2 \text{ balls after } K \text{ balls thrown}] > \frac{1}{2}$$

find K s.t.

→ complementary event

$$\Pr[\text{every bin has } \leq 1 \text{ ball after } K \text{ balls thrown}] \leq \frac{1}{2}.$$

A bin

Let \mathcal{E}_k = event that the k -th ball lands in a diff. bin than the first $k-1$ balls.

$$\Rightarrow \mathcal{E}_k \cap \mathcal{E}_{k-1} \cap \mathcal{E}_{k-2} \cap \dots \cap \mathcal{E}_1$$

$$\Pr\left[\bigcap_{i=1}^k \mathcal{E}_i\right] = \Pr\left[\mathcal{E}_k \mid \bigcap_{i=1}^{k-1} \mathcal{E}_i\right] \cdot \Pr\left[\bigcap_{i=1}^{k-1} \mathcal{E}_i\right] =$$

$$\Pr\left[\mathcal{E}_k \mid \bigcap_{i=1}^{k-1} \mathcal{E}_i\right] \cdot \Pr\left[\mathcal{E}_{k-1} \mid \bigcap_{i=1}^{k-2} \mathcal{E}_i\right] \cdot \Pr\left[\bigcap_{i=1}^{k-2} \mathcal{E}_i\right]$$

$$= P_k \left[\varepsilon_k \mid \bigcap_{i=1}^{k-1} \varepsilon_i \right] \cdot P_k \left[\varepsilon_{k-1} \mid \bigcap_{i=1}^{k-2} \varepsilon_i \right] \cdots P_k \left[\varepsilon_1 \mid \bigcap_{i=1}^0 \varepsilon_i \right]$$

$$\approx P_k \left[\varepsilon_i \mid \bigcap_{j=1}^{i-1} \varepsilon_j \right]$$

If $(i-1)$ bins had already i balls, then prob. of failure if total # bins = n ,

$$\therefore \frac{i-1}{n}$$

\therefore Prob. of success (of putting a ball in no bin that already has a ball).

$$\Rightarrow 1 - \frac{i-1}{n}$$

$$\Rightarrow \left(1 - \frac{k-1}{n}\right) \left(1 - \frac{k-2}{n}\right) \cdots \left(1 - \frac{1}{n}\right) \cdots 1$$

$$\text{S } \prod_{i=1}^{k-1} \left(1 - \frac{i}{n}\right)$$

After i bins have been filled, # bins total = n .

By Taylor expansion,

$$1+x \leq e^x, \forall x$$

Expanding this,

$$\leq \prod_{i=1}^{k-1} e^{-\frac{i}{n}} = e^{-\sum_{i=1}^{k-1} \frac{i}{n}} = e^{-\frac{1}{n} \cdot \frac{k(k-1)}{2}}$$

going back to the prob. eqn,

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

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

$$\ln 2 \leq \frac{1}{n} \cdot \frac{k(k-1)}{2}$$

$$2n \ln 2 \leq k(k-1)$$

{Birthday paradox}

$$2n \ln 2 \leq k^2$$

$$\cancel{2n \ln 2} \quad \boxed{k \geq \sqrt{2n \ln 2}}$$
$$\Rightarrow \boxed{k = \Omega(\sqrt{n})}$$

Ans.

Q2 what is expected load of a bin after m balls are thrown?

⇒ Let $x = \# \text{ balls in a bin.}$

we want $E[x],$

let $X_i = \begin{cases} 1, & \text{if ball } i \text{ falls in this bin.} \\ 0, & \text{otherwise.} \end{cases}$

$1 \leq i \leq m.$

for a bin

$$x = \sum_{i=1}^m X_i$$

$$\Rightarrow E[x] = \sum_{i=1}^m E[X_i] = \sum_{i=1}^m 1 \cdot P[X_i=1] = m \left[\frac{1}{n} \right]$$

$\frac{m}{n}$

Q3. What is expected # empty bins after m balls are thrown?

$\hat{m}_n \equiv \det Y = \# \text{ empty bins.}$

we want $E[Y]$

Let $y_j \in \begin{cases} 1, & \text{if bin } j \text{ is empty} \\ 0, & \text{otherwise.} \end{cases}$

$$\Rightarrow E[Y] = E\left[\sum_{j=1}^n Y_j\right] = \sum_{j=1}^n E[Y_j] = \sum_{j=1}^n 1.P_i[Y_j=1]$$

\Rightarrow we want $R[\gamma_j = 1]$ after m balls are thrown!

→ that means

$\{ \text{ball 1 doesn't land in bin } j, \text{ ball 2 doesn't land in bin } j, \dots, \{ \text{ball } m \text{ doesn't land in bin } j \} \}$

$$\therefore \rightarrow (\text{since events are independent}) \rightarrow \underline{\text{prod. of all events}}$$

$$P_e[Y_{j=2}] = P_e[\text{Ball 1 didn't lie in } b_j] \cdot P_e[\text{Ball 2 didn't lie in } b_j]$$

β_3 [unless we didn't lie in b_j].

$\Rightarrow \Pr[\text{ball } i \text{ does not lie in } b_j]$.

$$2) \left(1 - \frac{1}{n}\right)^m$$

$$\text{for } E[Y] \Rightarrow n \left(1 - \frac{1}{n}\right)^m$$

$n \rightarrow$ no. of birds
 $n \rightarrow$ bells thrown

$$\underline{E[Y] \leq n e^{-m/n}} \quad \text{(by exp. series expansion).}$$

If $m=n$,

$$E[Y] \approx \frac{n}{e} > \frac{n}{3}$$

~~high~~

$$E[Y] > \frac{n}{3} \rightarrow \text{more than } \frac{1}{3}\text{rd of the entries are expected to be empty}$$

An

Q4. What is the load of the full bin with high probability?
 [when n balls thrown ind. & uniformly into n bins]
 random

as $n \rightarrow \infty$, prob. happens to be 1.
 that's high prob,
 w.h.p means

$$\Pr \geq 1 - \frac{1}{n}$$

(with even stronger prob. for higher n).

$$\Pr \geq 1 - \frac{1}{e^n}$$

(for weaker)

$$\Pr \geq 1 - \frac{1}{\sqrt{n}}$$

we will consider

this one

here! find k^* such that —
 $\Pr [\text{every bin has } \leq k^* \text{ balls}] \geq 1 - \frac{1}{n}$

Complementary event :-

find k^* st.,

by inclusion-exclusion property.

$$Pr[\exists \text{ bin with } \geq k^* \text{ balls}] \leq \frac{1}{n}$$

a partition

(2) $Pr[\text{bin has } \geq k \text{ balls}] \Rightarrow \therefore \text{max. no. of balls thrown} N$.

$$\Rightarrow \sum_{j=k}^n Pr[\text{bin has } j \text{ balls}]$$

$$\Rightarrow \sum_{j=k}^n \binom{n}{j} \left(\frac{1}{n}\right)^j \left(1 - \frac{1}{n}\right)^{n-j}$$

Ans.

$$\leq \sum_{j=k}^n \binom{n}{j} \left(\frac{1}{n}\right)^j$$

[\because this is less than 1,
so, if we remove this,
we can upper bound
the left-out part].

$$\binom{n}{j} = \frac{n!}{j!(n-j)!} = \frac{(n-j+1)\dots n}{j!} \\ \leq \frac{n^j}{j!}$$

$$j! \approx \sqrt{2\pi j} \cdot \left(\frac{1}{e}\right)^j \left(1 + O\left(\frac{1}{j}\right)\right)$$

$$j! \geq \left(\frac{1}{e}\right)^j \Rightarrow \frac{1}{j!} \leq \left(\frac{e}{j}\right)^j$$

$$\leq \sum_{j=1}^n \binom{n}{j} \left(\frac{1}{n}\right)^j = \sum_{j=1}^n \frac{1}{j!}$$

$$\leq \sum_{j=k}^n \left(\frac{e}{j}\right)^j \leq \sum_{j=k}^n \left(\frac{e}{k}\right)^j$$

$$\leq \sum_{j=k}^{\infty} \left(\frac{e}{k}\right)^j = \frac{(e/k)^k}{1 - e/k}$$

$$\Pr[\text{bin has } \geq k \text{ balls}] \leq \frac{(e/k)^k}{[1 - e/k]}$$

2) $\Pr[\exists \text{ bin with } \geq k \text{ balls}] \rightarrow \dots$ there exist, there events will occur independently

$$\Pr\left[\bigcup_{i=1}^n E_i\right] \leq \sum_{i=1}^n \Pr[E_i]$$

$\Pr[\exists \text{ bin with } \geq k \text{ balls}]$

$$\leq \sum_{i=1}^n \Pr[\text{bin has } \geq k \text{ balls}]$$

$$= \boxed{\frac{n (e/k)^k}{1 - e/k}}$$

Any

$$2) \frac{n (\epsilon/k)^k}{1 - \epsilon/k} \leq \frac{1}{n}$$

Now, we solve this!

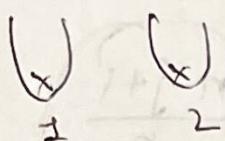
$$\text{if } k^* = \sqrt{\frac{\ln n}{\ln(\ln n)}} \rightarrow \underline{\underline{\mathcal{O}(\ln n)}} \quad \underline{\text{Ans.}}$$

\Rightarrow with hash table, in worst-case you can have $\log n$ time (going by this).

(Applicable to randomized algos.)

what is the expected no. of balls until every bin has at least 1 ball (Coupon Collector's Problem)?
 ↗ non-empty

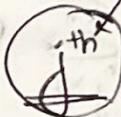
Ans. Lower bound 2) n bins = (n).



Initially, prob will be higher as there are lots of empty bins & then it goes smaller step by step.

Success:— ball that lands in an empty bin.

Epoch:— Sequence of balls starting after a success and ending with a success.



$\rightarrow X_j \Rightarrow \# \text{ balls in epoch } j$.

Let $X \Rightarrow \# \text{ balls until every bin has } \geq 1 \text{ ball}$.

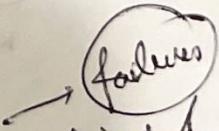
$$X = \sum_{j=1}^n X_j$$

$$E[X] = E\left[\sum_{j=1}^n X_j\right] = \sum_{j=1}^n E[X_j]$$

$X_1 = 1$ (in 1st epoch).

$X_2 = k \rightarrow$ (no. of balls need to throw until we hit an empty bin).
 ↗ don't know!

$$\Pr[X_2 = k] \Rightarrow \left(\frac{1}{n}\right)^{k-1} \cdot \left(1 - \frac{1}{n}\right) \xrightarrow{\substack{k-1 \text{ thrown in bin 1} \\ k \text{th time is a success.} \\ (\text{hit a diff. bin})}}$$



$$P[x_3=k] \Rightarrow \left(\frac{2}{n}\right)^{k-1} \cdot \left(1 - \frac{2}{n}\right)$$

⋮

$$P[x_j=k] \Rightarrow \left(\frac{j-1}{n}\right)^{k-1} \cdot \left(1 - \frac{j-1}{n}\right)$$

⋮

$$P[x_n=k] \Rightarrow \left(1 - \frac{1}{n}\right)^{k-1} \cdot \frac{1}{n}$$

$$\frac{n-j+1}{n}$$

$$E[X] = \sum_{j=1}^n E[X_j]$$

X_j is geometrically distributed with success prob $\frac{n-j+1}{n}$

$$P[X=k] = (1-p)^{k-1} \cdot p$$

$$E[X] \Rightarrow \frac{1}{p}$$

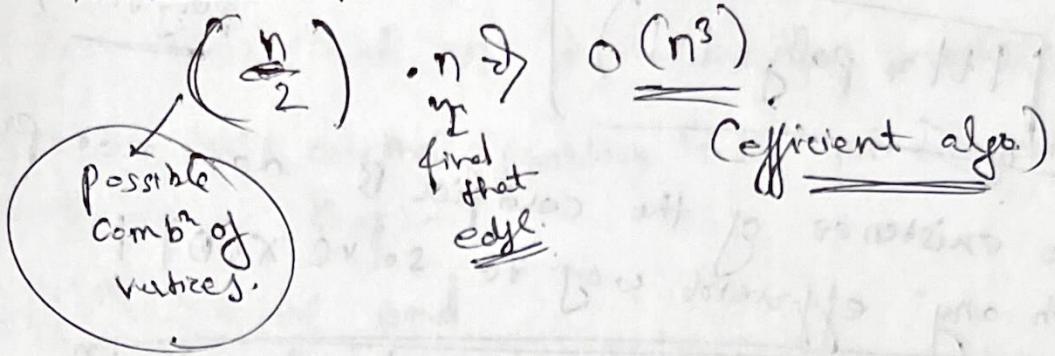
$$\Rightarrow \sum_{j=1}^n \frac{n}{n-j+1} \Rightarrow n \sum_{j=1}^n \frac{1}{n-j+1}$$

$$\Rightarrow n \sum_{k=1}^n \frac{1}{k}$$

$$\left\{ \sum_{j=1}^n \frac{1}{j} = H_n \quad (\text{nth harmonic no.}) \right.$$

$$H_n \approx \ln n + \Theta(1)$$

2) Run-time to find 2).



2) Verification algo. (or Certifier).

→ Input to $VC(D)$: - $x_2(G, k)$

→ Input to certifier for $VC(D)$: - (x, t)

→ What does the certifier do with its input?

① $t \leq k$ nodes in t .

② t forms a VC.

} → Can be solved in linear time to find the vertices set t and for the possible sets of $\binom{n}{2} \cdot n$ $\rightarrow O(n^3)$

thus, efficient algo.

∴ $X(D) \in NP$ if ∃ efficient certifier (verification algo) B

non-deterministic
gives a certificate!

that takes two inputs, the instance x & the certificate t , and it is s.t.

$x \in X(D) \Leftrightarrow$ there is short t s.t. $B(x, t) = \text{yes}$.

* eff:— worst-case poly-time. (bcz we want whole algo
to work in poly time)

④ short:— $|t| = \text{poly}(|x|)$

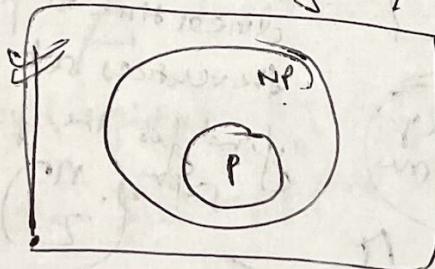
→ Note the existence of the certifier B does not provide us with any efficient way to solve $X(D)$!

(Ans) (Amp) bcz certifiers should be efficient to solve it in polynomial time).

⑤ $P \subseteq NP$

Is $P = NP$?

→ Not! finding a^* should be harder than checking one, especially for hard problems.



→ hardest prob. in NP are least likely to be in P !

NP-complete problems

A problem $X(D)$ is NP-complete if

① $X(D) \in NP$ and

② for all $Y \in NP$, $Y \leq_p X(D)$.

$X(D)$ is at least as hard as \underline{Y} .

(FACT)

→ Suppose X is NP-complete. Then X is solvable in polynomial time (i.e., $X \in P$) if & only if $P = NP$.

Soln to Euclidean TSP, which is

Now do we show that a problem is NP-complete?

Suppose we find an NP-complete problem X .

To show that another problem Y is NP-complete, we only need to show that:-

- (1) $Y \in \text{NP}$ and
- (2) $X \leq_p Y$

→ this provides an easier to prove more NP-complete problems.

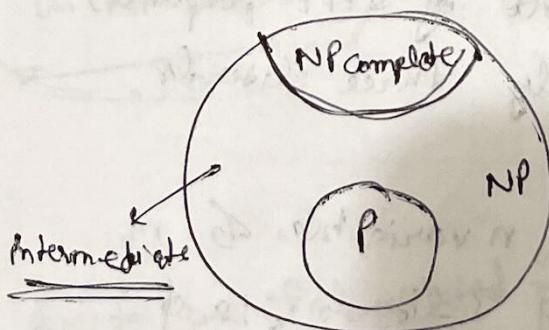
FACT → Transitivity of reductions.

If $X \leq_p Y$ and $Y \leq_p Z$, then $X \leq_p Z$.

→ we know that for all $Z \in \text{NP}$, $Z \leq_p X$. By fact $Z \leq_p Y$. Hence, Y is NP-complete.

* The first NP-Complete problem.

→ Circuit SAT is NP-complete!



CNF:

formula ϕ in CNF if it consists of conjunction of clauses each of which is a disjunction of literals ϕ with m clauses is in CNF if -

$$\phi = C_1 \wedge C_2 \wedge \dots \wedge C_m$$

(each clause C_i is disjunction)

$$C_i = (l_1 \vee l_2 \vee \dots \vee l_n).$$

\therefore we will work in CNF now!

- ⇒ A truth assignment satisfies a clause if it causes the clause to evaluate to 1.
- ⇒ A truth assignment satisfies a formula in CNF if it satisfies every clause in the formula.
- ⇒ A formula ϕ is satisfiable if it has a satisfying truth assignment.

Satisfiability (SAT) & 3SAT

Given a formula ϕ in CNF with n vars & m clauses, is ϕ satisfiable?

A convenient (if not easier) variant of SAT requires that every clause consists of exactly three literals.

3SAT

Given a formula ϕ in CNF with n variables & m clauses s.t. each clause has exactly 3 literals, is ϕ satisfiable?

- ⇒ These problems are $\text{NP}^{\text{complete}}$ problems. (Hardest problem in NP!).

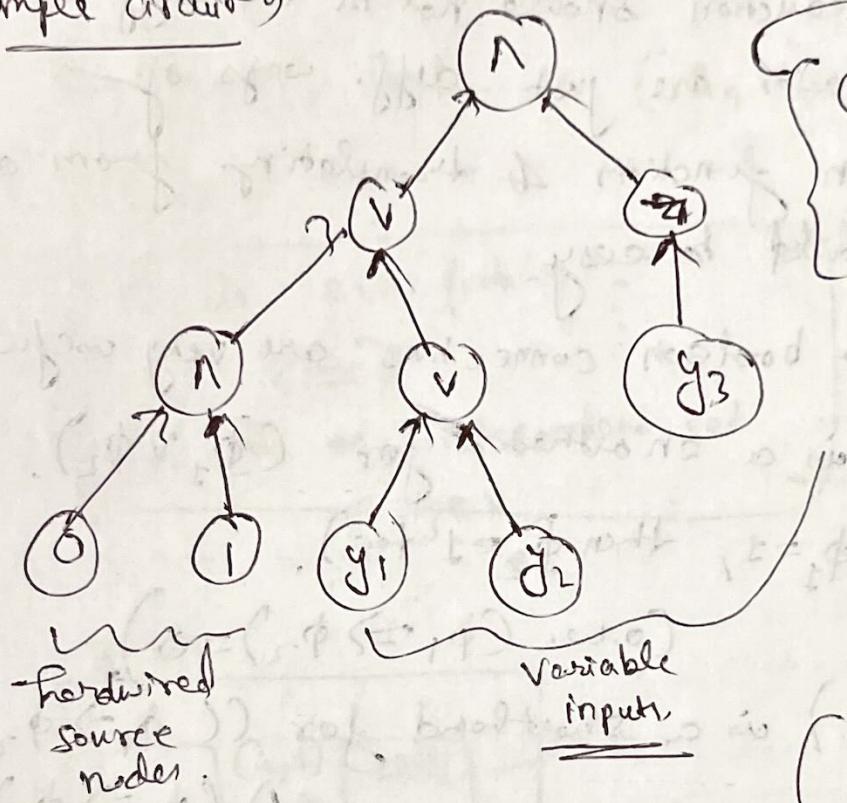
$$B: (\phi, t)$$

Circuit-SAT

Given a circuit C , is there an assignment of truth values to its inputs that causes the o/p to evaluate to 1.

Circuit-SAT is NP. (also NP-complete).

Example circuit)



Circuit-SAT \leq_p SAT
 $C_1 \rightarrow \phi$

Given a circuit C , is there an assignment of truth values to its I/P: that causes o/p to evaluate to 1?

(Inorder traversal of the tree)

$$\text{ENF: } ((0 \wedge 1) \vee (y_1 \wedge \neg y_2)) \wedge (\neg y_3)$$

Could be reduced to CNF if not using laws.

So, far we have stated \Rightarrow

\rightarrow Circuit-SAT is NP-complete

\rightarrow Circuit-SAT \leq_p SAT

\rightarrow SAT \leq_p 3-SAT

\Rightarrow SAT & 3-SAT are NP-complete.

Is ISCM as "hard" as SAT?

Before doing this,

Zemme \rightarrow Circuit-SAT \leq_p SAT

Intuitively, this reduction should not be too difficult.
formulas and circuits are just diff. ways of
representing boolean functions & translating from one
to the other should be easy.

The following two boolean connectives are very useful.

① $(\phi_1 \Rightarrow \phi_2)$ is a shorthand for $(\bar{\phi}_1 \vee \phi_2)$.

Intuition: if $\phi_1 = 1$, then $\phi_2 = 1$ too.

(o.w. $(\phi_1 \Rightarrow \phi_2) = 0$)

② $(\bar{\phi}_1 \Leftrightarrow \bar{\phi}_2)$ is a shorthand for $((\phi_1 \Rightarrow \phi_2) \wedge (\phi_2 \Rightarrow \phi_1))$,

which may be expanded to $(\bar{\phi}_1 \vee \phi_2) \wedge (\phi_1 \vee \bar{\phi}_2)$.

This clause evaluates to 1 iff. $\phi_1 = \phi_2$.

Transforming C into ϕ requires polynomial time \Rightarrow

This completes our construction of the clauses of ϕ .

Construction is polynomial in the size of the IP circuit.

Moreover, every clause consists of at most 3 literals
once ϕ is in CNF.

For independent sets

$\text{IS}(D)$ or hardy

$$\text{3SAT} \leq_p \text{IS}(D)$$

To prove $\text{IS}(D) \in \text{NP}$.

$$\text{3SAT} \leq_p \text{IS}(D)$$

$$\Phi(n, m) \rightarrow (a_i, k)$$

s.t.

Φ is satisfiable

$$\Leftrightarrow$$

or there is an independent set
of size k .

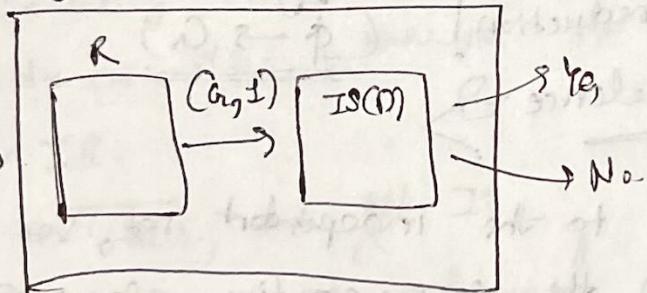
evidence
of
instances

they furnish
gives a
"yes"
instance!



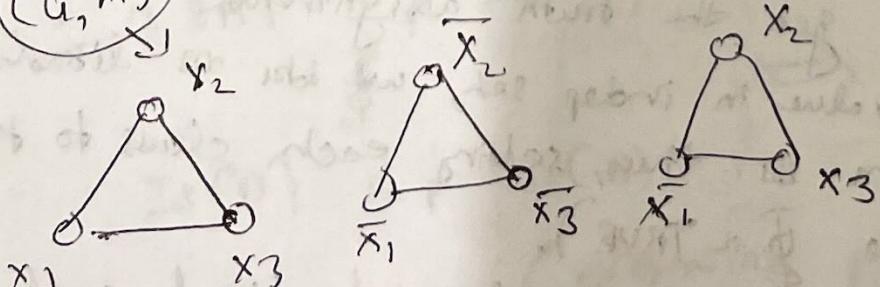
3-SAT

$$\Phi \rightarrow$$



$$\Phi = (x_1 \vee x_2 \vee x_3) \wedge (\bar{x}_1 \vee \bar{x}_2 \vee \bar{x}_3) \wedge (\bar{x}_1 \vee x_2 \vee \bar{x}_3).$$

(a_i, m)



$$x_1 = 1, \bar{x}_2 = 1, \bar{x}_3 = 1$$

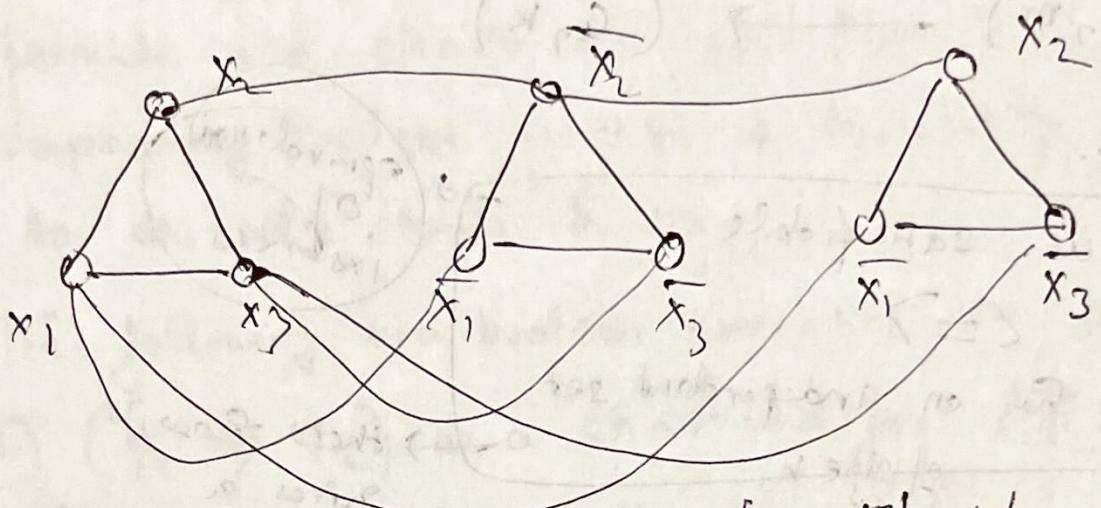
{ we can get
an independent
set of size 3. }

{ max independent
set of size
3. }

→ There are so many ways to satisfy all clauses.
(but each time at most we can get m literals as part of the independent set).

∴ we shouldn't have complemented literals in the same pair of independent set.

for this, we will have to add edges b/w them?



[one \exists s] could still be in
(no. of clauses).

⇒ This is a polynomial-time reduction!

Now, we need to prove equivalence ⇒

Idea → If a node belongs to the independent set, set the corresponding literal in the corresponding clause equal.

Proof

If A has an independent set of size m , then Φ is satisfiable.

We idea to get the truth assignments
(as we take m values in indep sets, we take m literals from each of the m Is then, setting each clause to true which would make $\Phi \rightarrow \text{TRUE}$).

Also, the \exists s is valid as complemented literal, can't both belong to an Is bcz of the edges b/w them in the graph.

... we have proved this in backward direction.

(FORWARD)

If Φ is satisfiable, there exist ind set of size ' m ' in G .

C. e. every clause should evaluate to 1, we need to ~~set~~ for
Set at least one of the listed values in each of the clauses
as 1).

∴ we will select those m nodes to be a part of an
independent set. Now, we need to prove that the set
formed is an ind. set.

∴ we have complemented the literals and formed
edges b/w them, we can't take them & also we are
selecting m literals from each clause, thus, we
take care that there is no edge b/w nodes considered
for IS.

Thus, we get IS.

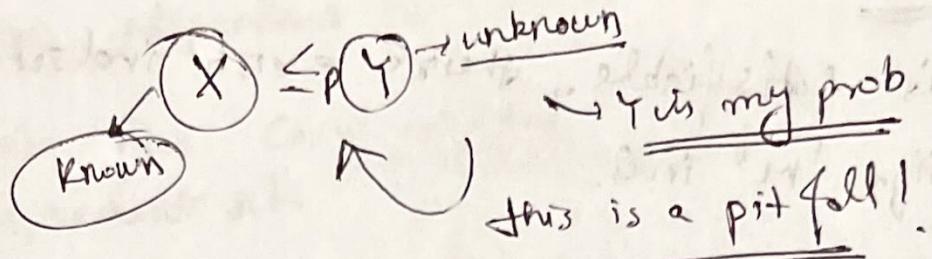
Hence, proved.

∴ 3-SAT problem is NP-complete.

$$\begin{array}{c} \text{3-SAT} \leq_p \text{IS}(D) \\ \text{IS}(D) \leq_p \text{VC}(D) \end{array} \quad \left\{ \begin{array}{c} \text{by transitive property} \\ \text{by transitive property} \end{array} \right\} \Rightarrow \text{3-SAT} \leq_p \text{VC}(D)$$

Is(D) is at least as hard as 3-SAT

Pitfalls \rightarrow



\Rightarrow In absence of ideas, reduce from 3-SAT.

(e.g. we should see which prob is most suitable to reduce from).

Traveling Salesmen Problem (TSP):-

Tour a simple cycle that visits every vertex exactly once.

Defn \rightarrow

Given n cities $\{1, \dots, n\}$, a set of non-negative distances d_{ij} b/w every pair of cities & a budget B , is there a tour of deryth $\leq B$?

Equivalently, is there a permutation π s.t.

① $\pi(1) = \pi(n+1) = 1$; that is, we start & end at city 1.

② total distance travelled satisfies

$$\sum_{i=1}^n d_{\pi(i)\pi(i+1)} \leq B$$

App's & Google street view car.

Simple instances of TSP :-

- Asymmetric: $d_{ij} \neq d_{ji}$ → directed graphs
- Symmetric: $d_{ij} = d_{ji}$ → undirected graphs
- Metric: satisfy the triangle inequality $d_{ij} \leq d_{ik} + d_{kj}$
- Euclidean: e.g. cities are in \mathbb{R}^2 hence city i corresponds to pt. (x_i, y_i) , then,

$$d_{ij} = \sqrt{(x_i - x_j)^2 + (y_i - y_j)^2}$$

using
approximate
algo, we can
solve few
type of instances
but not for
asymmetric

⇒ A related problem of TSP(D) :-

Hamiltonian cycle: Given a graph $G = (V, E)$, is there a simple cycle that visits every vertex exactly once?

Claim →

① Hamiltonian cycle is NP-complete.

Proof: Reduction from 3-SAT.

Claim →

② Symmetric TSP(D) is NP-complete.

Proof: Redⁿ from undirected Hamiltonian cycle.

\Rightarrow graph for Travel cycle is unweighted.

Ham. cycle \Rightarrow tour of length n in unweighted graph!

Q how to prove TSP \in NP (define a certifier)

$\Rightarrow \theta: (u, \pi)$?

(\hookrightarrow permutation of vertices)

This is a polynomial time certifier.

\hookrightarrow TSP is in NP!

Now, to prove it is NP complete, Ham cycle \leq_p TSP(θ)

To prove

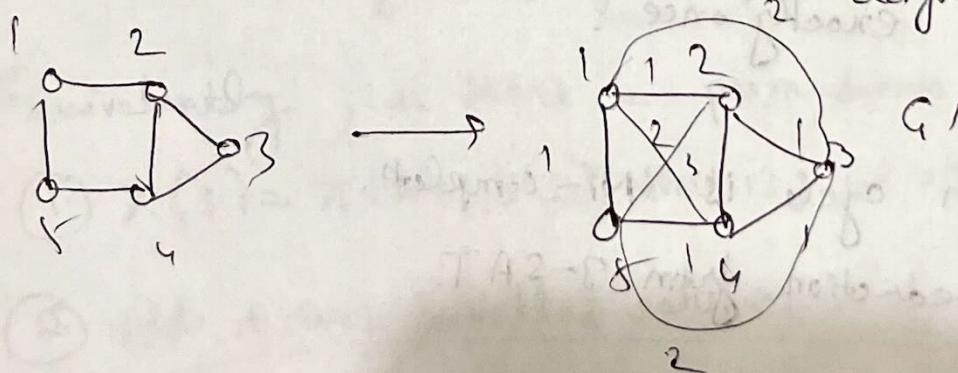
Ham cycle is already a decision problem, so, we don't need to convert it.

I(P2)

$$G = (V, E) \xrightarrow[\text{Convert to}]{} \left(G' = (V', E', w), B \right)$$

budget

Set G' for a Ham cycle (\Leftrightarrow G' for a tour of depth B).



Complete Graph formation

The vertices in G' remain the same in G . The ~~edges~~ in G is replicated in G' with weight 2.

Now, in the TSP, we need to know the weight (paths) of all possible vertices, \therefore we make all possible edges (even those not in G). Now G' becomes a complete graph. Further, we set the wts. of those new edges to an arbitrary large #. To keep it clean, we set it to 2.

We set the budget to be ' n ', since a ham cycle is a tour of length n , $\therefore \boxed{B=n} \rightarrow$ Red in polytime!

Equivalence of statements

\Rightarrow Forward direction

If G has a ham cycle, G' has a tour of length ' n '! The edges of the tour in G' will be exactly the same as in ham cycle G . The weight of the tour would be ' n ' since each of these edges also appear in G & its edge wt in $G' = 1$.

\Leftarrow Reverse 2

Given a tour of length ' n ' in G' , exhibit a ham cycle in G .
 G' is a complete graph, \therefore if it has at least one hamiltonian cycle. Given a tour of length ' n ', it means that every edge wt in the ham cycle must be one, \therefore a tour of length ' n ' in G' form a ham cycle in G . (This is because all the ' n ' edges in the tour in G' also appear in G .)

(*) Approximation algos

An α -approx. algo for an optimization problem
 is a poly-time algo that, for all instances of
 the problem, produces a solⁿ whose value is
 within a factor of α of the value of the optimal
 solⁿ.

Remark: α is the approx ratio / factor.

→ for min. prob, $\alpha \geq 1$

→ for max prob, $\alpha \leq 1$.

To prove \nexists 2 approx algo to prove TSP.

→ we'll prove by contradiction

\exists 2 approx algo for TSP

Assume that \exists 2 approx algo for TSP

if $A(G) \leq 2^n$, then G has a tour cycle
 if $A(G) > 2^n$, then G does not have a tour
 cycle

An approx. algo is an efficient algo! we observe that
 above is a tour, then we'll get one once in
 poly time, which is not possible unless P = NP.

$\therefore n-1 + (n!) \geq 2^{n+1}$ we conclude that
 there does not exist a 2-approx. algo.

(*) It is possible to provide very good approximate
 algo to euclidean TSP, which is not for symmetric TSP

Linear Programming

(a)

LP describes a broad class of optimization problems where both the constraints and the objective function are linear functions.

Applications

- 1) Resource allocation
- 2) Production planning
- 3) Military strategy forming
- 4) Error correction
- 5) Graph theoretic problems.

It can establish existence of polynomial time (efficient) algs.
Guide the design of approximation algs for computationally hard problems.

Duality allows to unify abstract views of seemingly unrelated results and is useful in algo design.

* GOALS

- ① Structure of an LP.
- ② Geometry of the solution in small dimensions.
- ③ LP solvers as a black box.
- ④ Integer LP and NP-completeness.

Q A boutique chocolatier has 2 products: chocolate & truffles
 their profit is: \$1/c, \$6/T
 They can produce ≤ 400 boxes/day
 daily demand is limited to
 1) ≤ 200 boxes of C 2) ≤ 300 boxes of T

What are the optimal levels of production?

\Rightarrow Variables: $x_1 = \#$ of boxes of chocolate
 $x_2 = \#$ of boxes of truffles
 $x_1 \geq 0, x_2 \geq 0$ (special constraints)

CONSTRAINTS: $x_1 + x_2 \leq 400$
 $x_1 \leq 200$
 $x_2 \leq 300$

Objective function: Maximize $x_1 + 6x_2$

Maximize
 $x_1 \geq 0, x_2 \geq 0$
 subject to
 $x_1 + 6x_2$
 $x_1 \leq 200$
 $x_2 \leq 300$
 $x_1 + x_2 \leq 400$

general problem

Ques.

⑥

1. a set of variables
2. a set of linear constraints on the variables (equities / inequalities)
3. a linear objective function to maximize (or minimize)

O/P An assignment of real values to the variables such that

1. the constraints are satisfied.
2. the objective fn is maximized (or minimized)

* Terminology

1. Solution : An assignment of real values to the variables.
2. Possible solⁿ : solⁿ that satisfies all the constraints (including the special ones)
3. Feasible region : the set of all feasible solutions
4. Optimal solution : a feasible solⁿ that maximizes (or minimizes) the objective function if the LP is a maximization (minimization) LP
5. Cost/value of solⁿ : the value of the objective fn for the solⁿ.
6. Optimal value of LP : the value of an optimal soln.

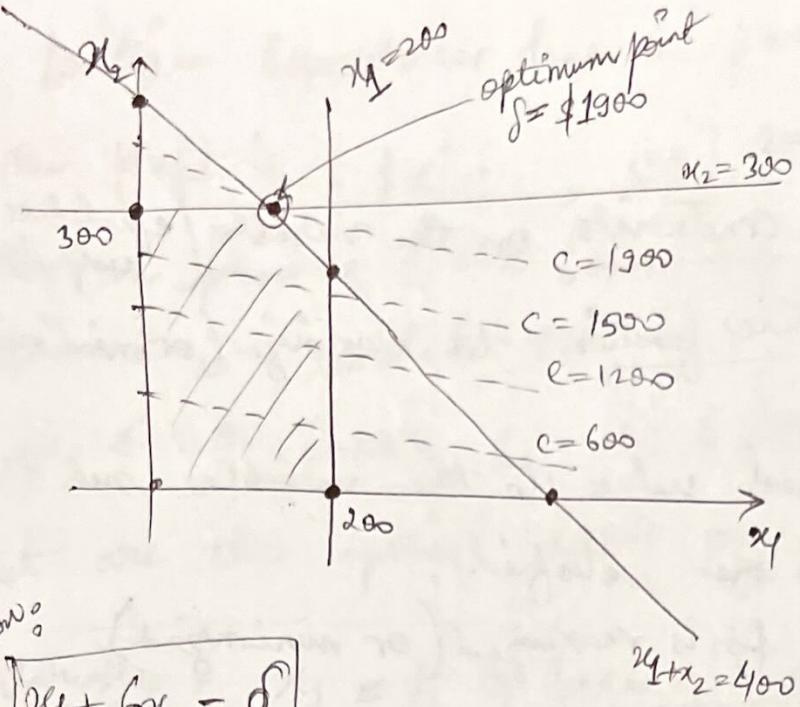
Let's find the feasible region now!

• convex

↓
it will always be convex! (by construction)

• not convex

Means that if you draw a line segment, the line must fully lie in the set.



objective function:

$$2x_1 + 6x_2 = 0$$

defines a line on the plane that contains all points that achieves a profit of S .

Now, we need to find the optimal solⁿ in this feasible region.

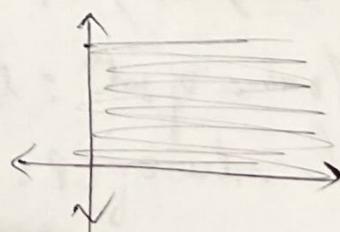
* Theorem

The optimum is achieved at a vertex of the feasible region

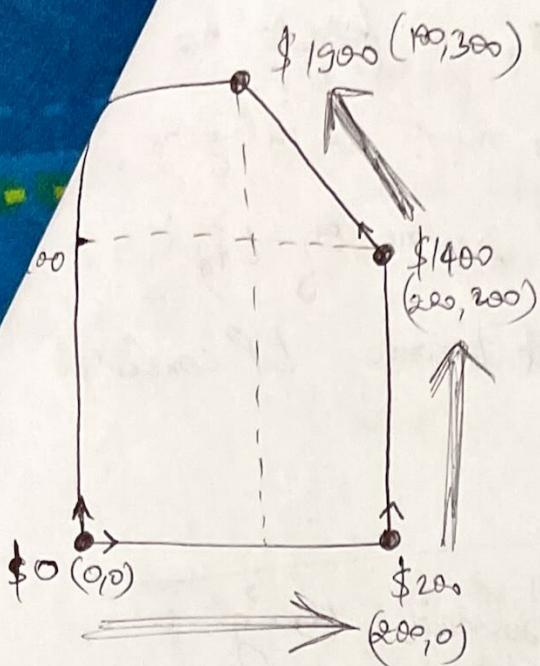
EXCEPTIONS 1. The linear program is infeasible (no solution)
eg: $x \leq 1, x_1, x_2$

2. The optimum value is unbounded

$$\begin{array}{ll} \max & 2x_1 + x_2 \\ x_1 \geq 0, x_2 \geq 0 & \end{array}$$



* Solving LPs : the simplex method [Dantzig]



It is a hill climbing problem.

(c)

- ① Find the vertex to start with, "first vertex". let's say $(0,0)$ in the feasible region
- ② Move towards the direction of improvement, i.e., an adjacent vertex that \uparrow the value of the current vertex.
- ③ keep repeating this until you're stuck. Nowhere else to go where the value of the objective fn can be maximized further.

The vertex at which the algorithm terminates, becomes the global optima because of the convex shape of the feasible region.

→ The simplex algorithm always moves in the direction of improvement. It is the most imp alg of the previous century!

* History on LP solvers

1. Simplex Method

- fast in practice
- exponential worst case performance

2. Ellipsoid Method

- provably poly-time algorithm
- slow in practice

3. Interior point Method

- polynomial time algorithm
- fast in practice.

4. Interior point methods is a field of active research.

In the problem of chocolates & truffles, we forgot to put a constraint that the # of boxes must be integers. If we do so, we are no longer in the realm of LP, but we enter integer LP (ILP). This is because LP consists of solutions containing real values.

Also, this problem could have been solved using a greedy approach / without LP. Prof used this just to motivate the problem & drew a graph (represent visually) on the plane

→ ILP problems are not solvable in polynomial time!

(d)

alternative proof 28th April 2022
that \$1900 is optimal

Max
 $x_1, x_2 \geq 0$
subject to

$$\begin{aligned} & x_1 + 6x_2 \\ & (x_1 \leq 200) .0 \\ & (x_2 \leq 300) .5 \\ & (x_1 + x_2 \leq 400) .1 \end{aligned}$$

PRIMAL LP

$$\begin{aligned} 0.2x_1 + 5x_2 + x_1 + x_2 &\leq 200.0 + 300.5 + 400.1 \\ x_1 + 6x_2 &\leq 1900 \end{aligned}$$

How did we come up with these numbers?

- Multiply these constraints with ~~some~~ different non-negative real #s.
so that they don't change the direction of the inequality
- Add the LHS and RHS of these constraints separately
- Make the LHS similar to the objective function.

~~Maximize~~

$$\begin{aligned} & (x_1 \leq 200) . y_1 \\ & (x_2 \leq 300) . y_2 \\ & (x_1 + x_2 \leq 400) . y^3 \end{aligned}$$

s.t. $y_1, y_2, y^3 \geq 0$

$$y_1 x_1 + y_2 x_2 + y^3 x_1 + y^3 x_2 \leq 200 y_1 + 300 y_2 + 400 y^3$$

$$(y_1 + y^3)x_1 + (y_2 + y^3)x_2 \leq 200 y_1 + 300 y_2 + 400 y^3 \quad \textcircled{1}$$

now, we must upper bound our objective function.

If we can upper bound it by the LHS of ①, then we can also upper bound it by the RHS of ①.

$$\therefore x_1 + 6x_2 \leq (y_1 + y^3)x_1 + (y_2 + y^3)x_2$$

To make sure that this inequality holds? -

$$y_1 + y_3 \geq 1 \quad \text{and} \quad y_2 + y_3 \geq 6$$

This becomes:-

DUAL LP

min

$$y_1 \geq 0$$

$$y_2 \geq 0$$

$$y_3 \geq 0$$

$$200y_1 + 300y_2 + 400y_3$$

s.t.

$$y_1 + y_3 \geq 1 \quad \text{and} \quad y_2 + y_3 \geq 6$$

minimization

OBSERVATIONS

- 1) If the primal LP has m constraints, there are m variables in the dual.
- 2) ~~number of variables~~ The # of variables in the primal = # of constraints in the dual.
- 3) For every primal variable, we get a dual constraint.
- 4) If the primal is a maximization, then the dual is a minimization and vice versa.

Let U_p be the cost of any solution to the primal LP

Let U_D be the cost of any solution to the dual LP.

WEAK DUALITY

$$U_p \leq U_D \text{ by construction}$$

Cost of optimal soln

If the dual is maximization, then $U_p \geq U_D$