

AOA week 8 Notes

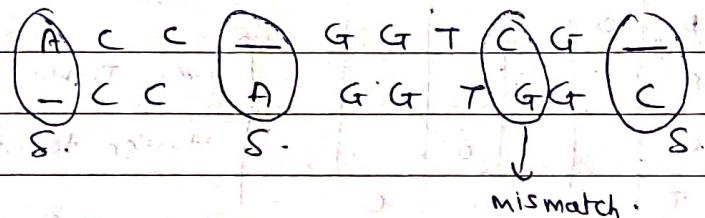
* Sequence Alignment Problem.

DNA strand consists of 4 bases {A, C, G, T}

Comparing DNA strands

$$S_1 = A \text{ C C } G \text{ G G } T \text{ C G }$$

$$S_2 = C \text{ C A } G \text{ G T } G \text{ G C }$$

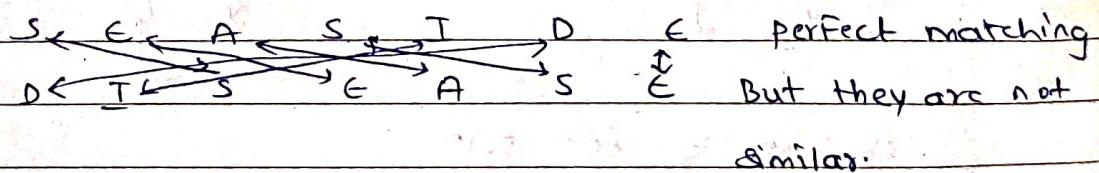


3 gaps 1 mismatch.

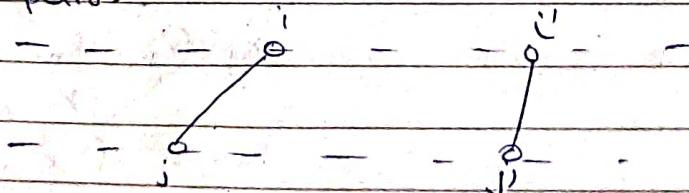
$$X = \{x_1, x_2, \dots, x_m\}$$

$$Y = \{y_1, y_2, \dots, y_n\}$$

A matching is a set of ordered pairs with property that each item occurs at most once.



DEFⁿ: A matching is an alignment if there are no crossing pairs.



$$(i, j) \text{ and } (i', j) \in M$$

$$\& i < i' \Rightarrow j < j'$$

For an alignment M between x and y .

1. we incur a "gap penalty" of s for each gap
 \leftarrow Fixed

2. For each mismatch (of letters $p \neq q$) we incur
a mismatch cost d_{pq}
 \downarrow not fixed.

	A	C	G	T	
A	0				↓ Table Given to us. $d_{xx} = 0$
C		0			
G			0		
T				0	

$d_{xy} \Rightarrow$ +ve int.

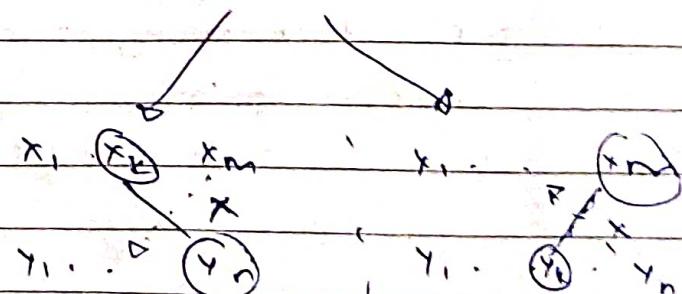
Defn Similarity between strings x and y is the minimum cost of an alignment between x and y .

$$x = \{x_1, \dots, x_m\}$$

$$y = \{y_1, \dots, y_n\}$$

Say M is OPT soln.

$$(x_m, y_n) \in M \quad \text{or} \quad (x_m, y_n) \notin M$$



If (x_m, y_n) exist i.e. if $(x_m, y_n) \in M$
so x_m will not be paired as it will result in cross pair.

So y_n will not be paired as it will result in gross pair.

$$\text{OPT}(m, n) = \min \left(\begin{array}{l} \text{OPT}(m-1, n-1) + \alpha_{mn}, \\ \text{OPT}(m-1, n) + S, \\ \text{OPT}(m, n-1) + S \end{array} \right)$$

Rec
formula ①

$j \uparrow$	58			
	48			
	38			
	28	← o		
	8	↓		
	0	S	28	38 48

Alignment (x, y)

init $A[i, j] = iS$ for each i

init $A[0, j] = jS$ for each j

For $j=1$ to n

For $i=1$ to m

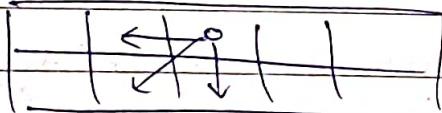
Rec formula ①

Time: $O(m * n)$

Space: $O(m * n)$

Suppose DNA Seq $\approx 3B$ size

Space Req: $3B * 3B$



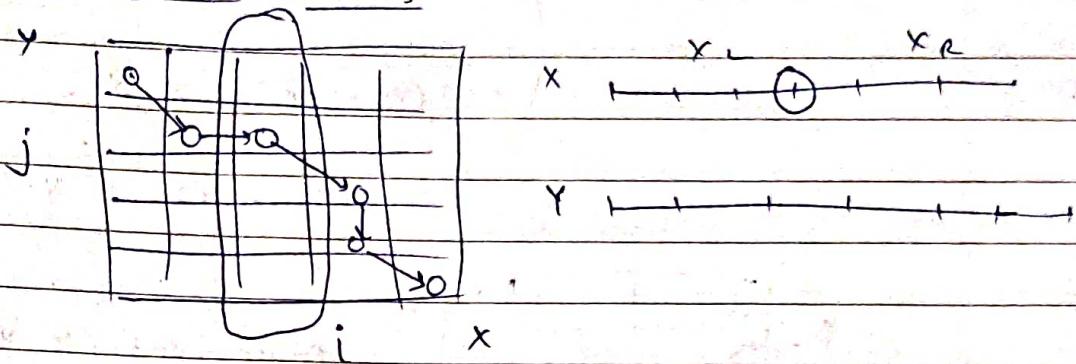
But for calculation of a particular row we only need 1 prev row.

But if we keep 2 rows to solve,

we will get the optimal one but how

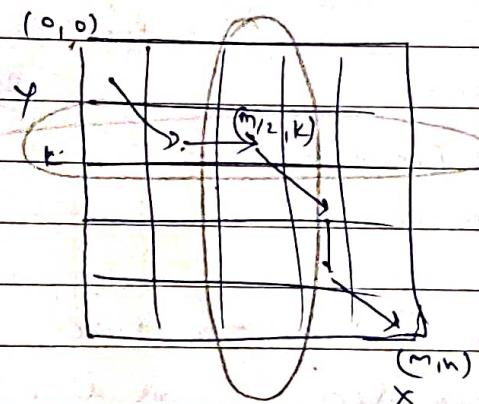
will we find opt soln ~~if~~ as we will not be able to backtrack.

Memory efficient way.



We will use divide + conquer
We will divide x in middle to get x_l and x_r

But how do we know where to divide y



We have x in middle.

We now need to find index 'k' such that
The optimal soln passes through that point

i.e. OPT path passes through index k

where $x = m/2$

$$\text{i.e. } \left(\frac{m}{2}, k \right)$$

Let $f(q, k)$

$g(q, k)$



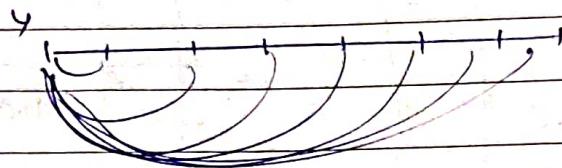
$$\text{OPT } (0,0) \rightarrow (q, k)$$

$$\text{OPT } (q, k) \rightarrow (m, n)$$

$\therefore k$ is the index
which minimizes

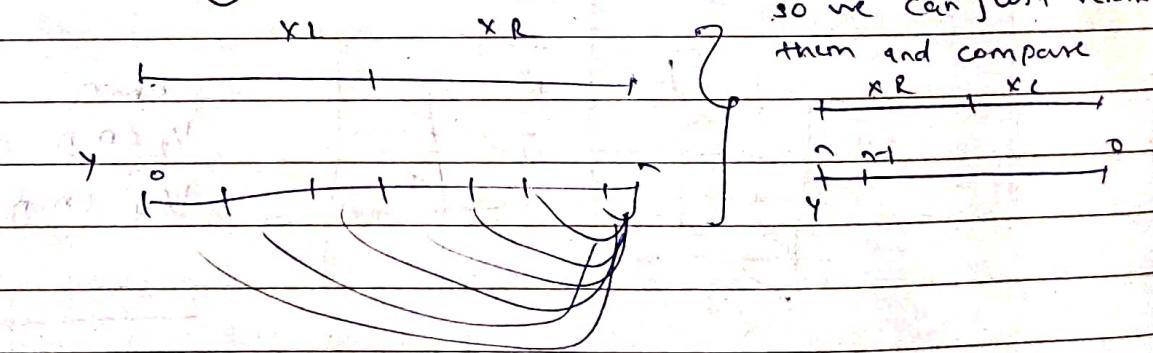
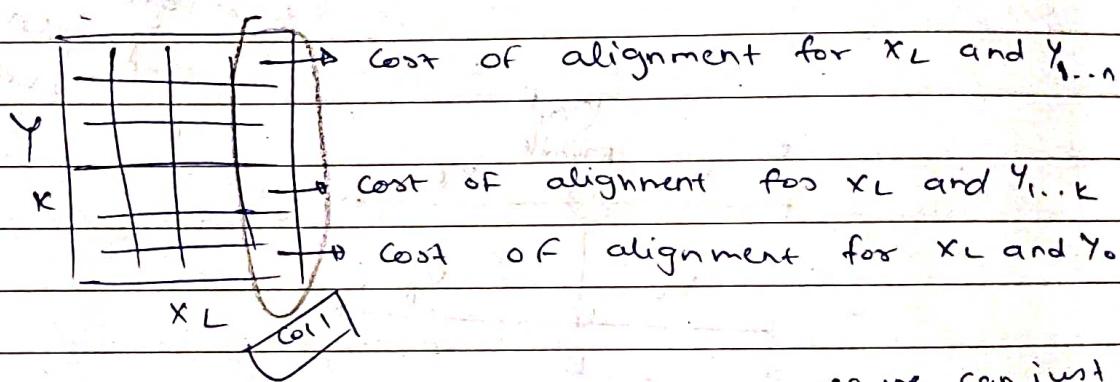
$$\left(f(q, k) + g(q, k) \right)$$

e.g. $x_L \quad x_R$

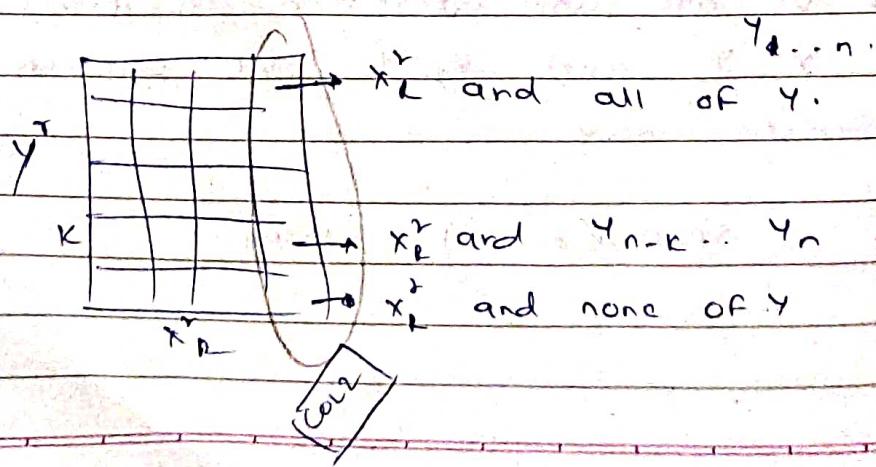


x_L could be solved with y_0

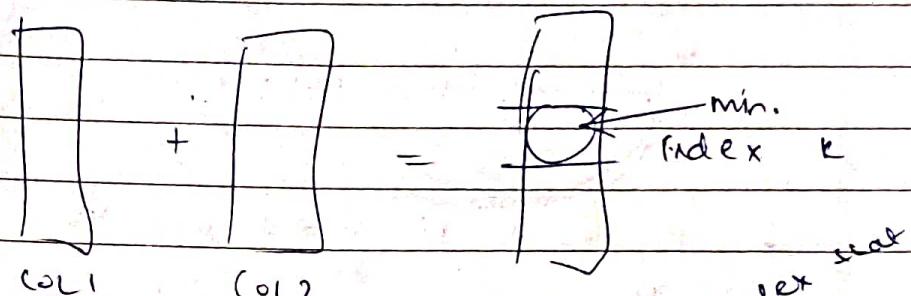
$x_L \quad , \quad " \quad y_1$
 $x_L \quad , \quad " \quad y_n$



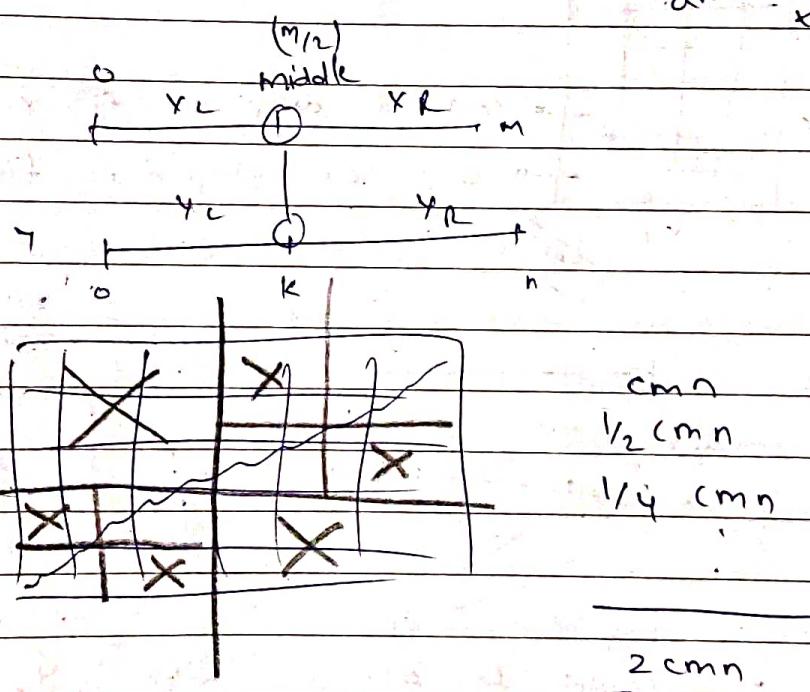
x_R could be aligned with y_n



after getting $\boxed{\text{COL1}}$ and $\boxed{\text{COL2}}$



find the index
give min ans
and split 4 or
that index



Time : $O(2mn)$

As we are just interested in final 2 columns
of $X_L Y$ and $X_R Y$ we use memory
efficient way

Space $O(m+n)$ linear.

memory efficient way

Divide \leftarrow use dynamic programming to
find optimal split point for

Solve

Conquer \leftarrow just concatenate the ans of
subproblems.

memory wasteful
way

memory efficient way

Time

$O(mn)$

$O(1)$

Space

$O(mn)$

$O(m+n)$

In D&C we solve for OPT SOLN so

we directly get OPT SOLN after

concat

① Matrix chain multiplication

$$A = A_1 \cdot A_2 \cdot A_3 \cdots \cdots \cdots \cdot A_n$$

e.g

$$B \cdot C \cdot D$$

$$B = 2 \times 10$$

$$C = 10 \times 50$$

$$D = 50 \times 20$$

$$B \cdot C \cdot D = (B \cdot C) \cdot D \Rightarrow 3000 \text{ operations}$$

OR

$$B \cdot (C \cdot D) \Rightarrow 10400 \text{ operations}$$

$$((A_i \cdots A_k) (A_{k+1} \cdots A_j))$$

Finally

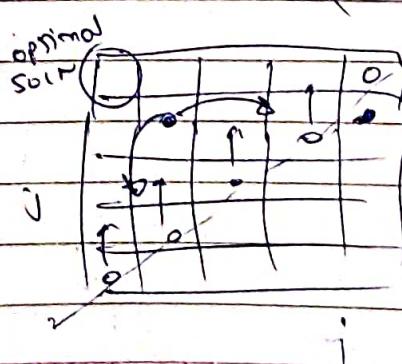
$$\left(\begin{matrix} 1 \text{ matrix} \\ \vdots \\ 1 \text{ matrix} \end{matrix} \right)$$

We need to find where to split

$\text{OPT}(i, j)$ = optimal cost of multiplying matrices
i through j

$$\text{OPT}(i, j) = \min_{i \leq k \leq j} \left(\text{OPT}(i, k) + \text{OPT}(k+1, j) \right) + (r_i c_k c_j)$$

formula



final product cost

For $i=1$ to n $\text{OPT}(i, i) = 0$

~~For $j=1$ to i~~

For $j=2$ to n

For $i=j-1$ to 1 ^{dec by -1}

$O(n)$

$O(n)$

Rec formula ② $\rightarrow O(n)$

end for

end for

$O(n^3)$ Time

Efficient soln.

* Shortest path problem Dynamic Programming:

Assumption: NO -ve cost cycles

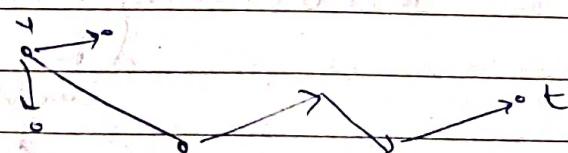
There can be -ve cost edges.

If G has no -ve cycles, then there is a shortest path from s to t that is simple and hence atmost ' $n-1$ ' edges.

$\text{OPT}(i, v) \Rightarrow \min$ cost from v to t

With atmost ' i ' edges

Problem statement via $\text{OPT}(n-1, s)$

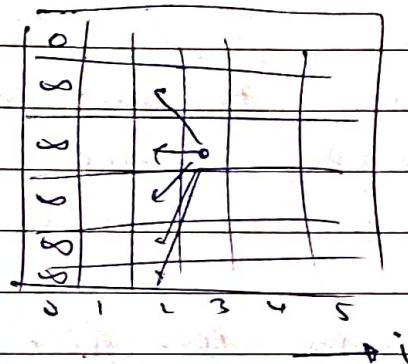


$$\text{OPT}(i, v) = \min_{\text{For } w \in \text{adj}(v)} \left\{ (vw + \text{OPT}(i-1, w)) \right\}$$

But by the recurrence relation we are only considering paths of exactly len ' i '
 Path can also be atmost ' i '
 so they can also be less than ' i '

$$OPT(i, v) = \min \left(OPT(i-1, u) \right)$$

Min for
 $\{ w \in \text{Adj}(v) \}$
 $w \rightarrow OPT(i-1, w)$



Bellman Ford Algo.

shortest-path (G, s, t)

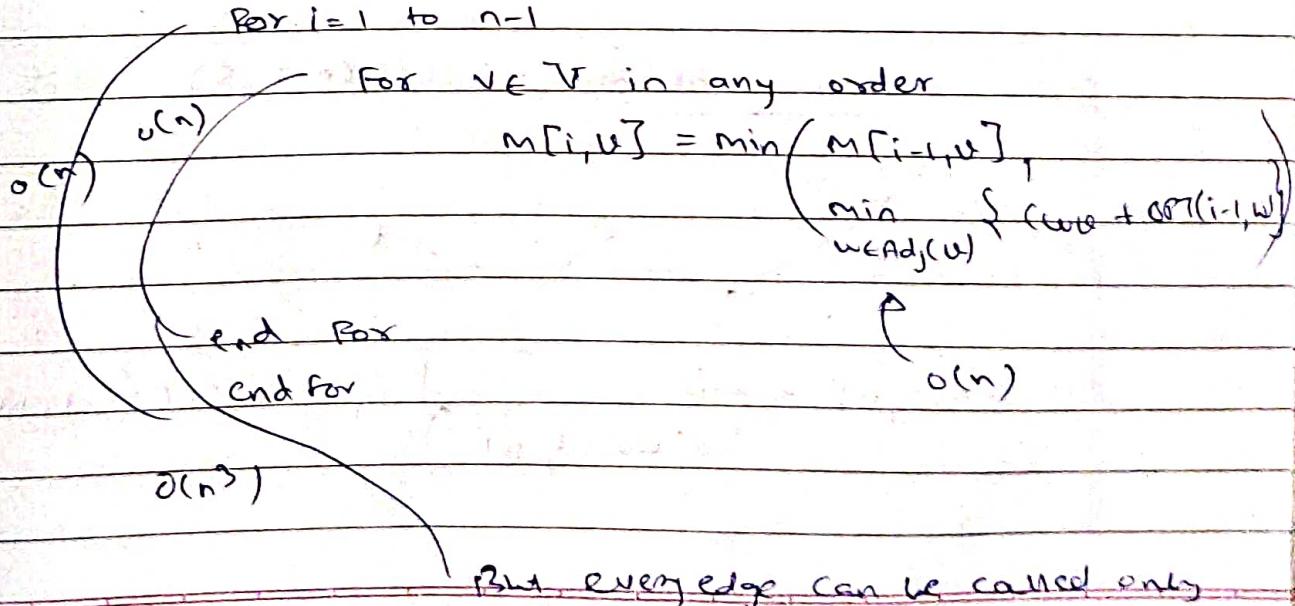
$n = \text{no. of nodes in } G$

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

For $i = 1$ to $n-1$

 For $v \in V$ in any order

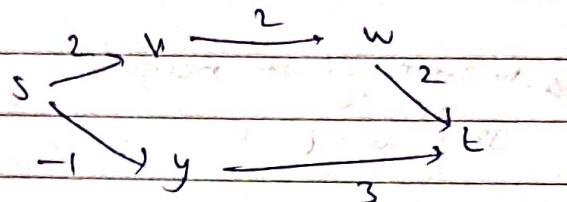
$$M[i, v] = \min \left(M[i-1, v], \min_{w \in \text{Adj}(v)} \{ w \rightarrow OPT(i-1, w) \} \right)$$



$O(mn)$

at max $m = n^2$

e.g.:



2 col's same.

t	0	0	0	0		STOP HERE
y	00	3	3	3		
w	∞	2	2	2		
v	∞	∞	4	4		
s	∞	∞	2	2		
	0	1	2	3		

1 COLUMN

(Search + with almost 1 edge)

t-t 0

y-t 3

w-t 2

v-t ∞

s-t ∞

2 COLUMN

(Search + with almost 2 edges)

t-t 0

y-t 3

w-t 2

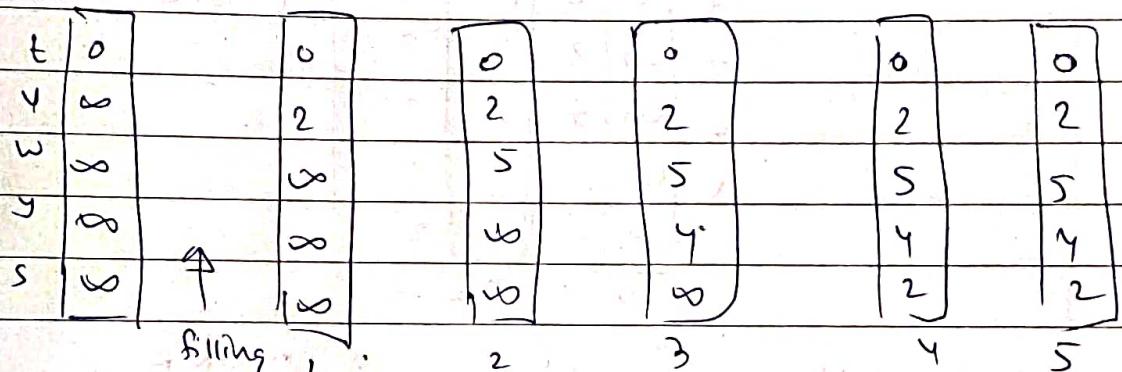
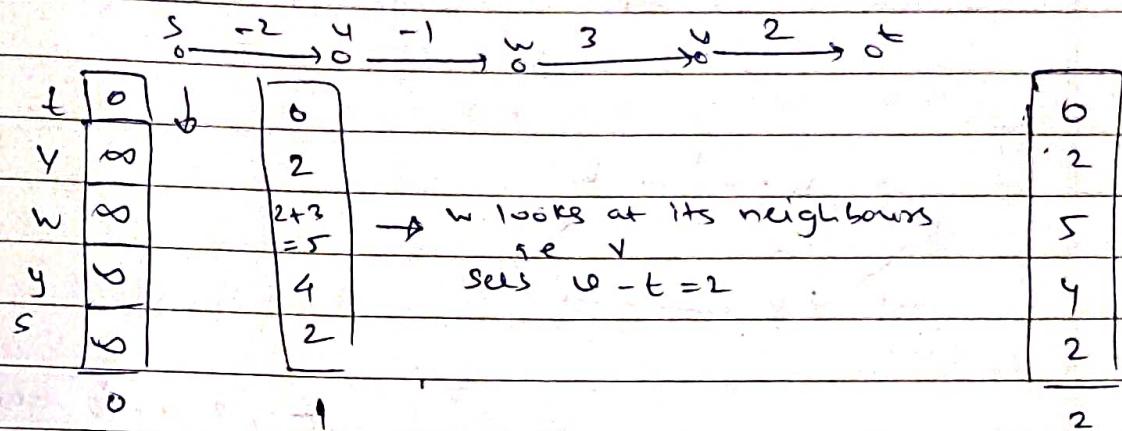
v-t 4

s-t 2

At any time only last 2 col's are in consideration

i.e. cols can be judged from Col 2 only.

Trying to solve with only 1 column.



$O(n-1)$ wont come

what is preference of filling?

BFS order from T

as its neighbours will convey info
then their neighbours

so on.

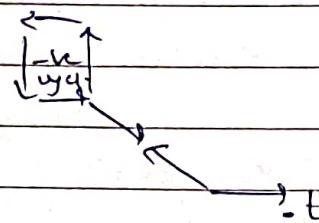
we found shortest dist

for shortest path we just keep a track of parent
i.e. neighbour that give shorter distance to t.

\therefore shortest path always found in ' $n-1$ ' iterations

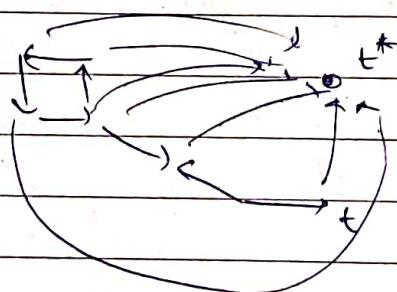
If after $n-1$ iterations value go down this indicates presence of -ve cycle.

Q: How to FIND -ve CYCLE



In this running it after $n-1$ times won't change column value as -ve cycle is not on path to t

To FIND CYCLE



Add new node t^*
connect every node to t^*
run bellmanford for t^*

Bellman Ford
 $O(mn)$

Dijkstras
 $O(m \log n)$

IF -ve cost edges can only use BellmanFord

IF no -ve cost edges

can use Bellmanford + Dijkstras

Some cases Bellmanford faster than Dijkstras
as it can be easily parallelized.