

Dynamic Programming

Sequence Alignment Problem

A DNA strand consists of a string of molecules called bases!

4 Types of bases:

- Adenine A

- Cytosine C

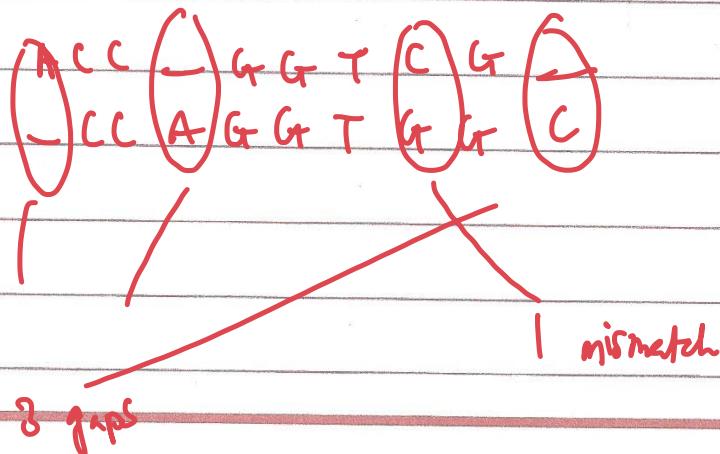
- Guanine G

- Thymine T

Comparing DNA strands

$S_1 = A C C G G T C G$

$S_2 = C C A G G T G G C$



Suppose we have 2 strings X & Y .

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

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

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

perfect
matching

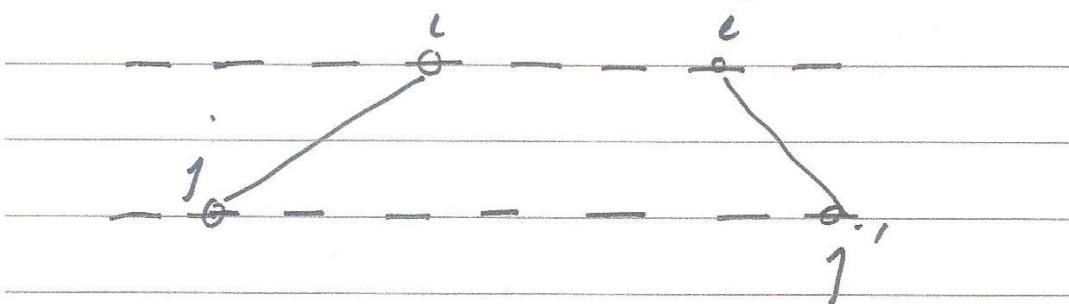
but they are
not similar

SEASIDE
~~SEASIDE~~
DISEASE

Matching is an alignment with no crossing pairs.

Def. A matching is an alignment

if there are no crossing pairs.



$$(i, j), (i', j') \in M$$

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

For an alignment M between X & Y

1- We incur a "gap penalty" of 8 fixed
penalty for each gap.

2- For each mismatch (of letters $p \neq q$) we incur a mismatch cost α_{pq}

Not fixed; this penalty depends on which two characters mismatch

	A	C	G	T	
A	0				Mismatch Matrix
C		0			
G			0		
T				0	

This matrix is given to us.

Def. Similarity between strings X & Y

is the minimum cost of an alignment
between X & Y .

$$X = \{x_1, \dots, x_m\}$$
$$Y = \{y_1, \dots, y_n\}$$

Is this a pair in optimal solution or not

Say M is an opt. solution.

either $(x_m, y_n) \in M$ or $(x_m, y_n) \notin M$

~~Define OPT(~~

$x_1 \dots -$
 $y_1 \dots - \textcircled{-} y_n$

IF (x_m, y_n) is not a pair

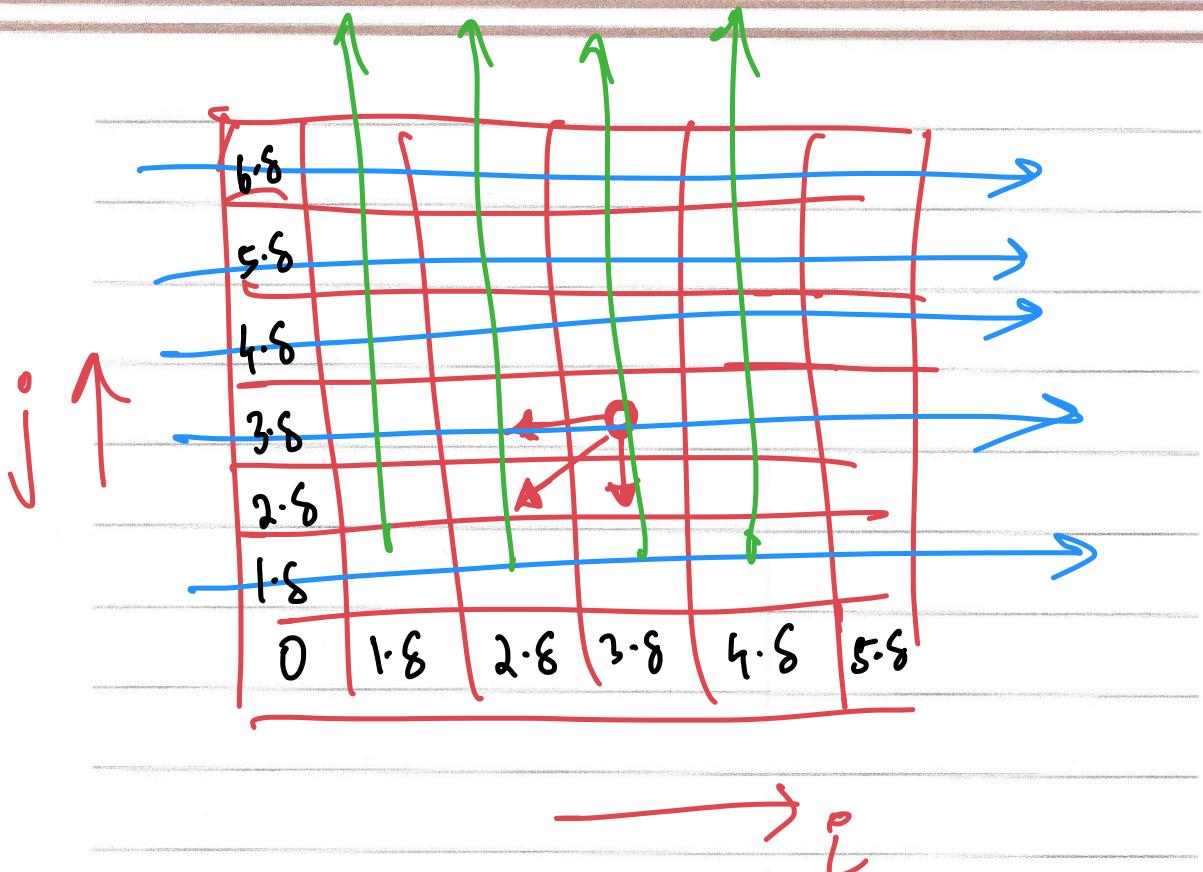
Consider x_m to be paired with some y_i
if this is case y_n cannot be paired
as it would lead to cross pairing

similar is case if y paired then x not considered

In an optimal alignment M , at least one of the following is true:

Formal proof

- 1) - $(x_m, y_n) \in M \rightarrow OPT(m, n) = OPT(m-1, n-1) + \delta_{x_m y_n}$
- 2) - x_m is not matched $\rightarrow OPT(m, n) = OPT(m-1, n) + 8$
- 3) - y_n is not matched $\rightarrow OPT(m, n) = OPT(m, n-1) + 8$



Alignment (x, y)

Initialize $A[i, 0] = i\delta$ for each i

$A[0, j] = j\delta$ for each j

for $j = 1$ to n

 for $i = 1$ to m

 Recurrence Formula ①

Takes $O(n \cdot m)$

Efficient solⁿ

Size of string

Size of string

Both represent size of some input

DNA sequence \approx 3 billion size

Memory waste full way

$$\therefore \text{matrix size } (3 \text{ billion})^2$$

But for calculation we only need 1 previous row

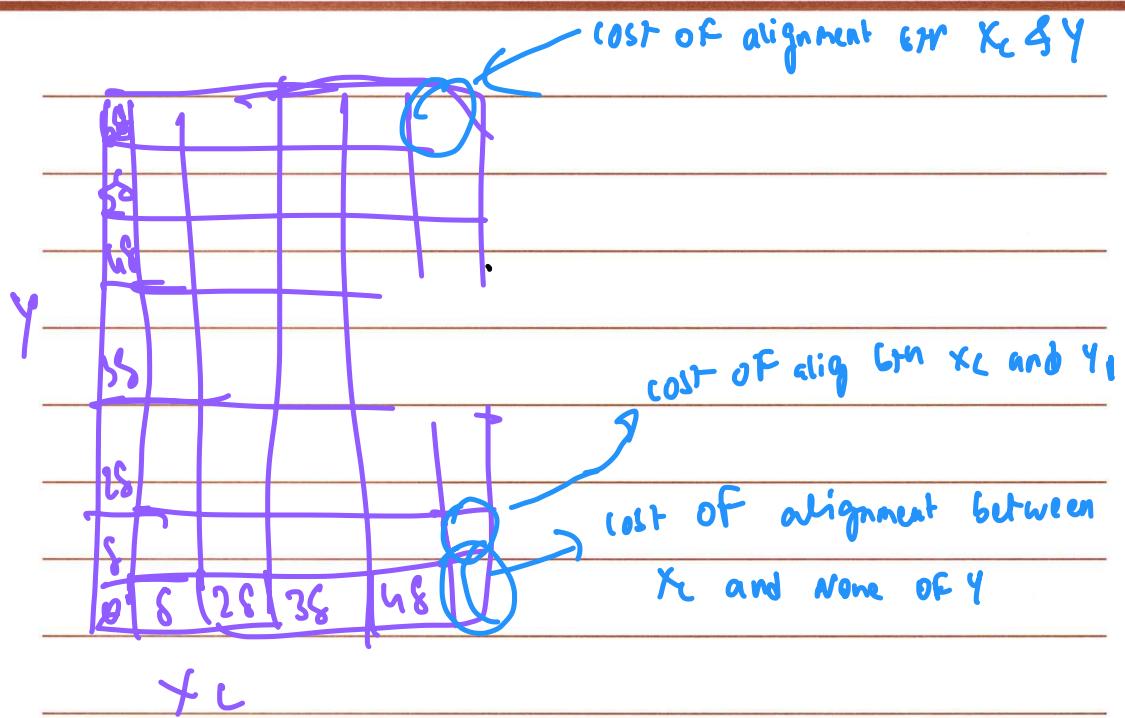
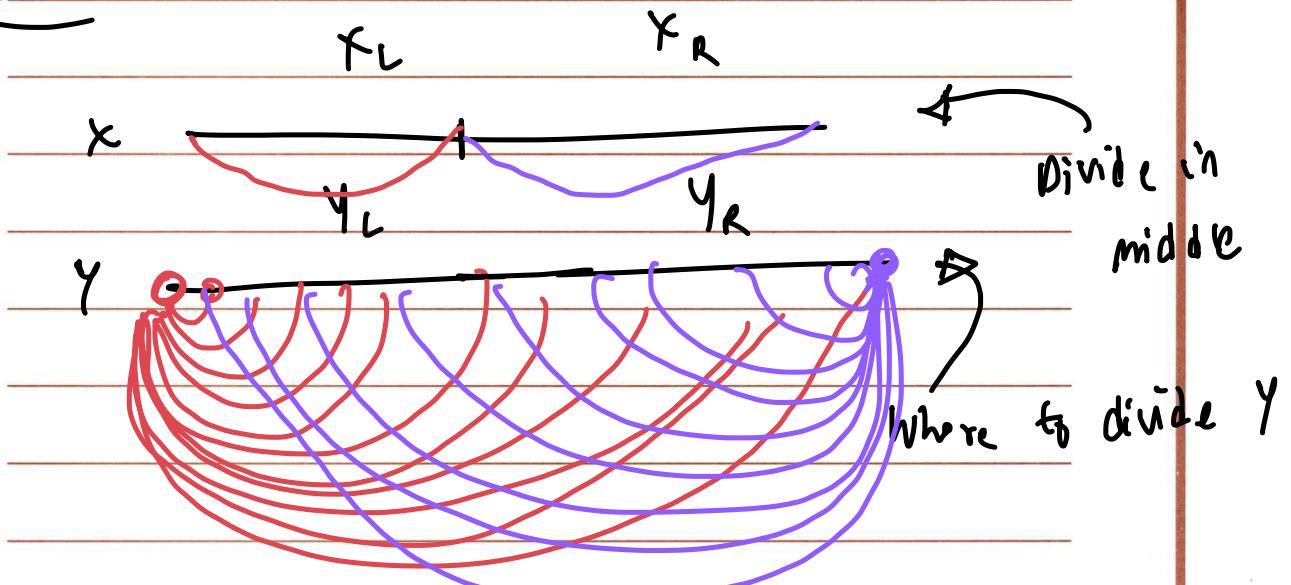
Memory efficient way



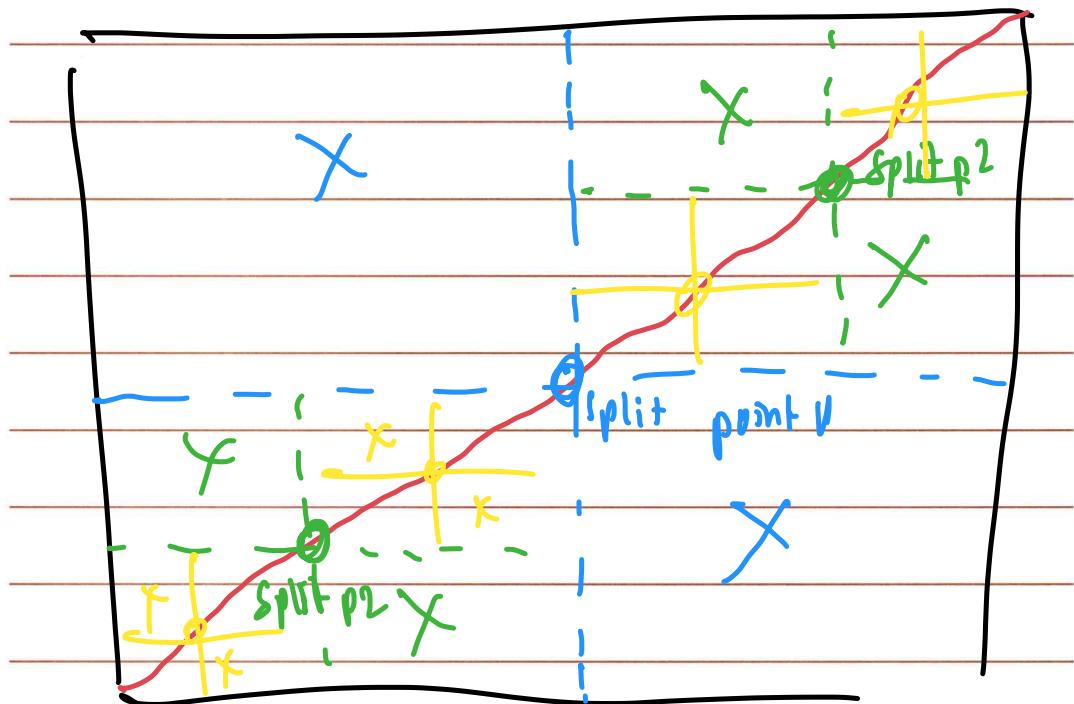
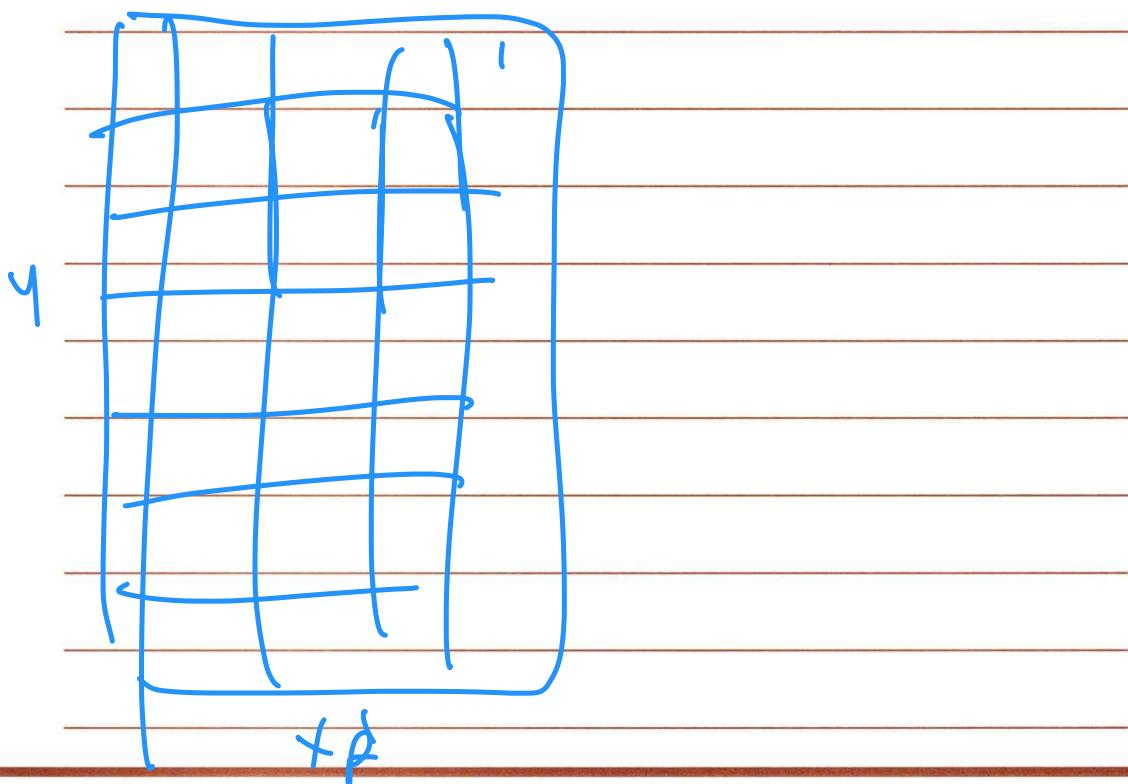
But if we only keep two rows,
we won't be able to backtrack the
answer (alignment)

TRY TO FIND ALIGNMENT WITH MEMORY EFFICIENT WAY

Solution



DP is at divide strip



C_{mn}

$C_{mn}/2$

$C_{mn}/4$

!

,

$2 C_{mn}$

memory wasteful

Space: $3B + 3B$

Time: C_{mn}

memory efficient

Space: $3B + 2$

Time: $2 C_{mn}$

Matrix Chain Multiplication

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

e.g.:

$$B \cdot C \cdot D$$

$$B \in 2 \times 10$$

$$C \in 10 \times 50$$

$$D \in 50 \times 20$$

$$m \begin{bmatrix} & \\ & \end{bmatrix}_n + \begin{bmatrix} & \\ & \end{bmatrix}_n$$

operations = $m * n * k$

$$\therefore B \cdot C \cdot D = B \cdot (C \cdot D) \Rightarrow 10,400 \text{ operations}$$

OR

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

$$\left((A_1 \dots A_K) \times (A_{K+1} \dots A_j) \right)$$

↓ ↓

Final step (| matrix | matrix)

We need to find where we need to split

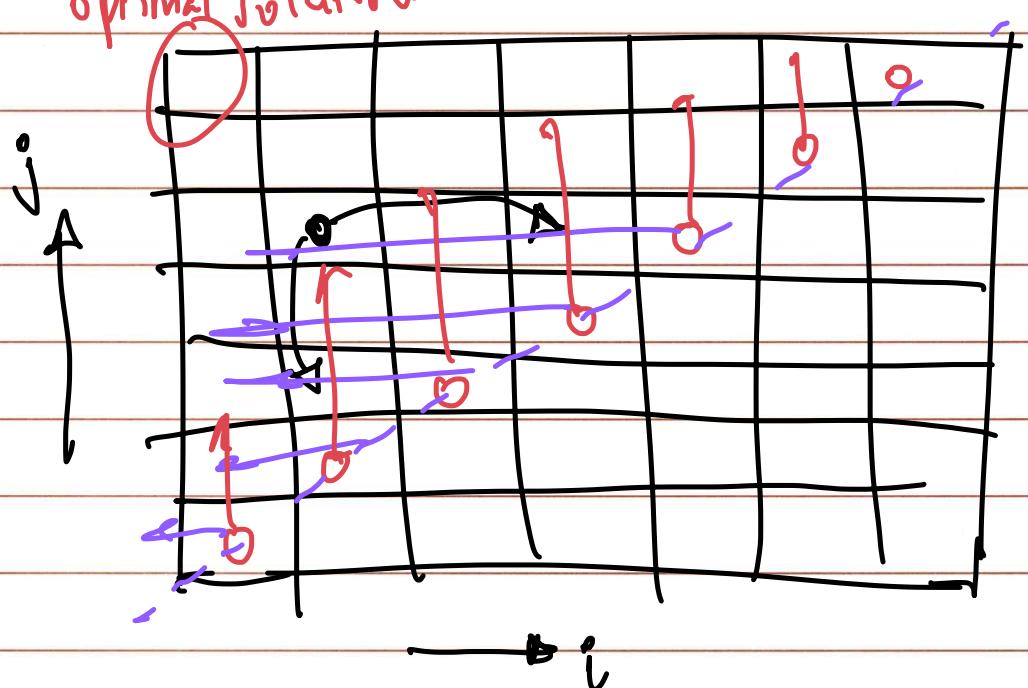
$OPT(i, j)$ = Optimal cost of multiplying matrices
i through j

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

+ Final product cost
($R_i C_k L_j$)

Rec Formula ②

Optimal solution



for $i = 1$ to n

$$OPT(i, i) = 0$$

For $j = 2$ to n

for $i = j-1$ to 1 by -1

use rec formula ②

$O(n)$

$O(n^2)$

$O(n)$

end for

end for

$O(n^2)$ Time

Efficient solution.

Shortest Path Problem

Dynamic Programming

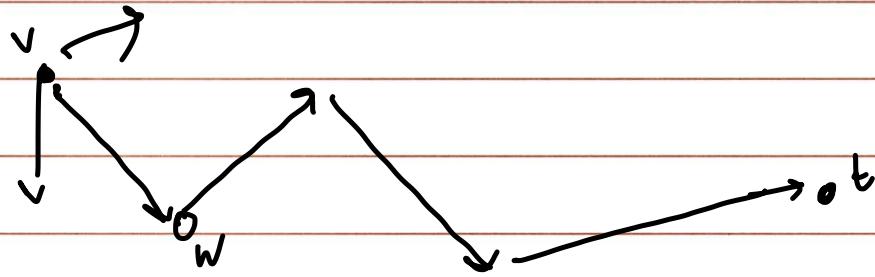
- Assumption : No -ve cost cycles

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 at most $n-1$ edges

$\text{OPT}(i, v)$ denotes the min cost of a v to t
path using at most i edges

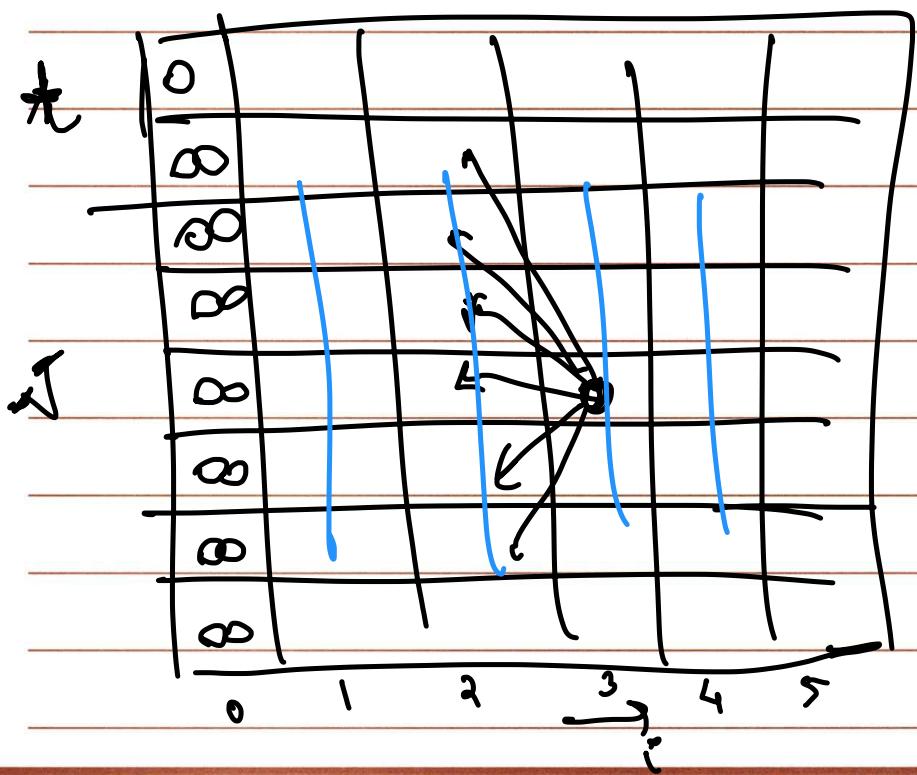
$$\text{OPT}(n-1, S)$$



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

This path should also have fewer edges than i

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



Apply deCursion Formula:

Bellman-Ford Alg.

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 $n-1$

$O(n)$

$O(n)$

for $v \in V$ in any order

$$M[i, v] = \min(M[i-1, v], \min_{w \in \text{Adj}(v)} (M[i-1, w] + C_{vw}))$$

end for

end for

$O(n)$

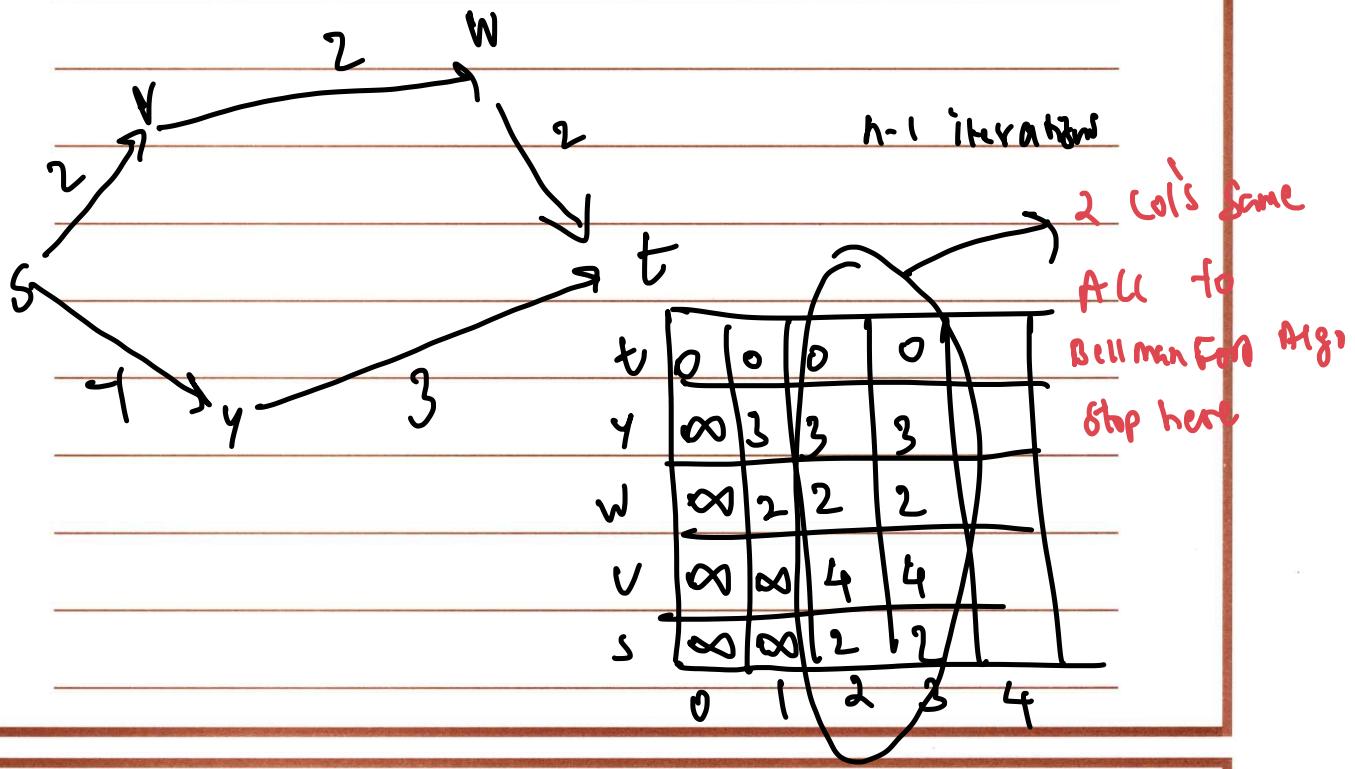
$O(n^3)$

But each edge can
only be called once

$O(n \cdot m)$

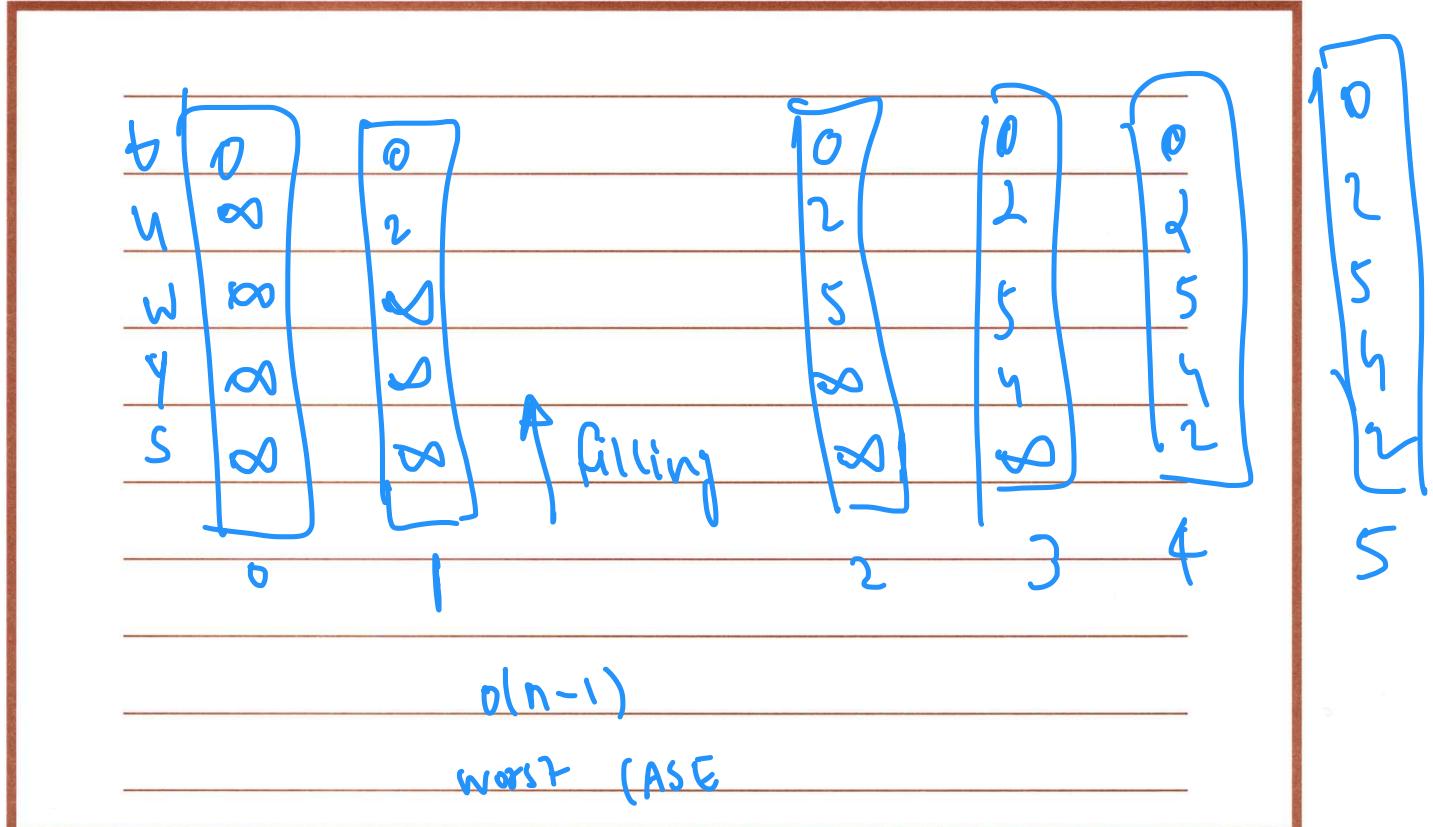
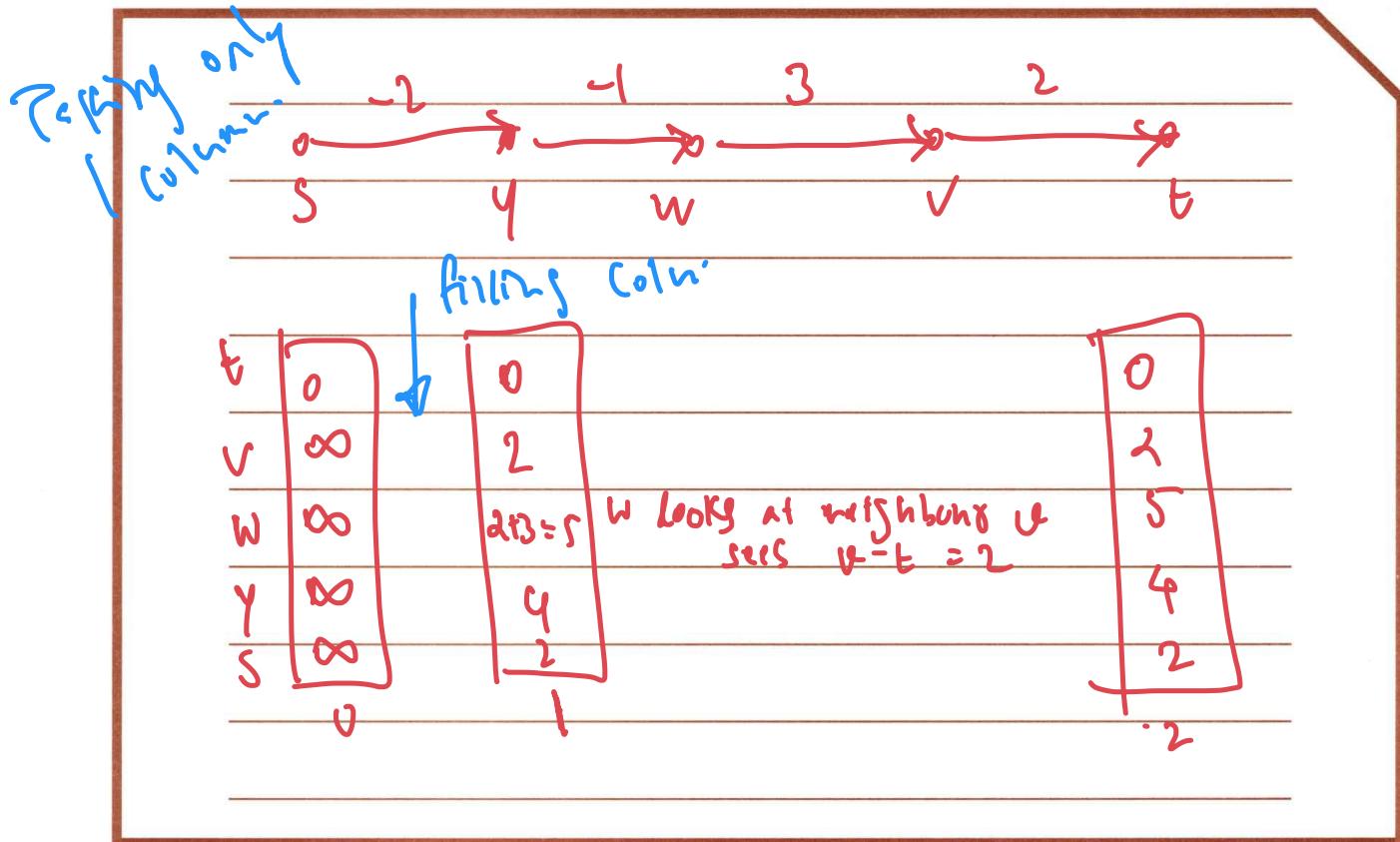
(Same as that of
Dijkstra's)

Efficient solution



i column (reach t with i edge)	2 column (reach t with 2 edge)
$t-t$ 0	$t-t$ 0
$y-t$ 3	$y-t$ 3
$w-t$ 2	$w-t$ 2
$v-t$ ∞	$v-t$ 4
$s-t$ ∞	$s-t$ 2

At any time atmost only last 2 col's are required.



What is the performance of filling

BFS order from T

As its neighbours will convey info
then their neighbours

↓

so on

We found shortest distance

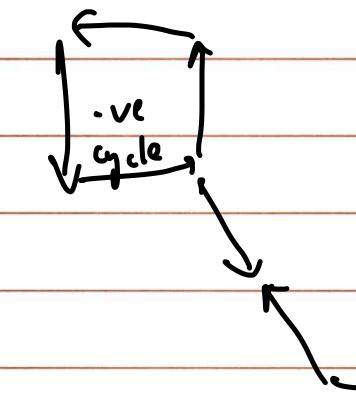
for shortest path.

We just keep track of parent i.e. neighbour
that give shortest distance to t.

∴ Shortest path always found in ~~$O(n^2)$~~ $n-1$
iterations

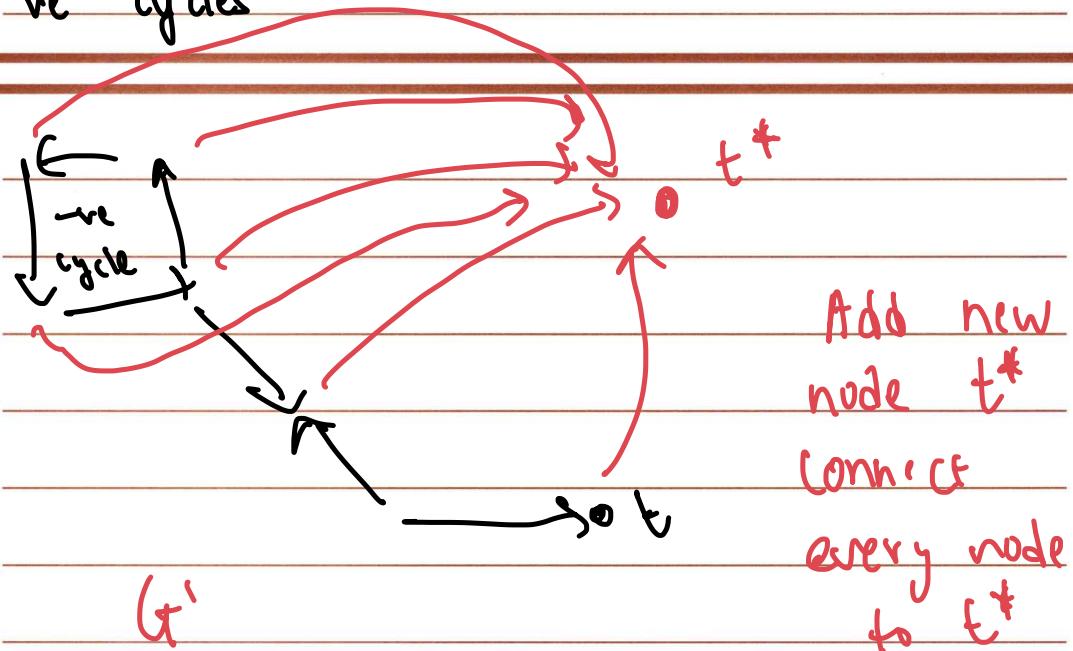


If after $n-1$ iterations value goes down then
indicates presence of -ve cycle



In this case running it after $n-1$ times wont change $(v)^n$ value as -ve cycle is not on path to t

Then how to find node from where to run BellmanFord and check if there are -ve cycles



Now run Bellman Ford for t^*

Bellman Ford

$\Theta(m \cdot n)$

If no -ve cost edges
can use BF | Dijkstr'a's

Dijkstra's
 $\Theta(n \log n)$

Discussion 7

If -ve cost edges use
Bellman
ford

- When their respective sport is not in season, USC's student-athletes are very involved in their community, helping people and spreading goodwill for the school. Unfortunately, NCAA regulations limit each student-athlete to at most one community service project per semester, so the athletic department is not always able to help every deserving charity. For the upcoming semester, we have S student-athletes who want to volunteer their time, and B buses to help get them between campus and the location of their volunteering. There are F projects under consideration; project i requires s_i student-athletes and b_i buses to accomplish, and will generate $g_i > 0$ units of goodwill for the university. Our goal is to maximize the goodwill generated for the university subject to these constraints. Note that each project must be undertaken entirely or not done at all -- we cannot choose, for example, to do half of project i to get half of g_i goodwill.
- Suppose you are organizing a company party. The corporation has a hierarchical ranking structure; that is, the CEO is the root node of the hierarchy tree, and the CEO's immediate subordinates are the children of the root node, and so on in this fashion. To keep the party fun for all involved, you will not invite any employee whose immediate superior is invited. Each employee j has a value v_j (a positive integer), representing how enjoyable their presence would be at the party. Our goal is to determine which employees to invite, subject to these constraints, to maximize the total value of invitees.
- You are given a set of n types of rectangular 3-D boxes, where the i^{th} box has height $h(i)$, width $w(i)$ and depth $d(i)$ (all real numbers). You want to create a stack of boxes which is as tall as possible, but you can only stack a box on top of another box if the dimensions of the 2-D base of the lower box are each strictly larger than those of the 2-D base of the higher box. Of course, you can rotate a box so that any side functions as its base. It is also allowable to use multiple instances of the same type of box.

In some case Bellman Ford works better than Dijkstra's

as Bellman Ford can be easily parallelized

As Bellman Ford only looks at its neighbors but Dijkstra's need the whole G \therefore it needs info centrally whereas Bellman Ford does not need info centrally.

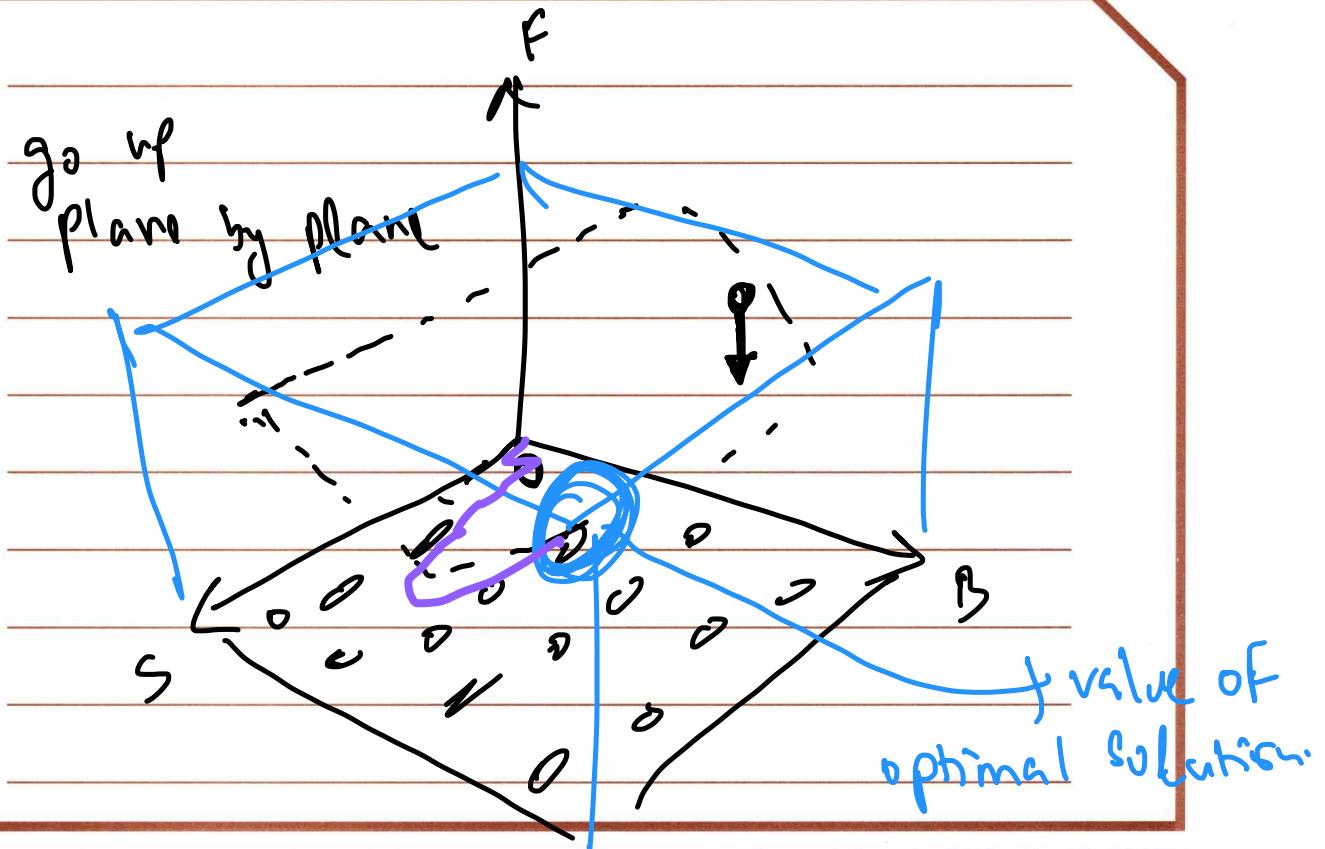
(pn cost fastest
I/O cost
communication cost slower)

① S - Students
 B - buses
 F - projects

In 0/1 knapsack
 we had 1 resource (W)
 Here we have
 2 resources (S, B)

$\text{OPT}(i, S, B)$ = optimal value of the soln
 for projs 1..i with
 S students and B buses

$$\text{OPT}(i, S, B) = \max \left(G_i + \text{OPT}(i-1, S-S_i, B-B_i), \text{OPT}(i-1, S, B) \right)$$



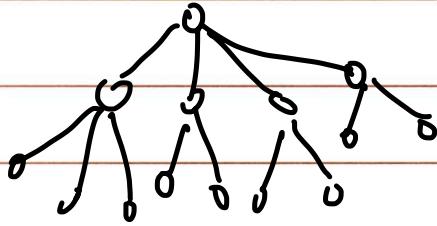
$\Theta(B * F * S)$
 $O(F) \rightarrow$ complete to Down Parse

not efficient

because of F and S

BFS \rightarrow $F \ 2^{\log_2 B} \ 2^{\log_2 S}$

②



$\text{OPT}(i) = \max \text{fun-factor for the subtree rooted at node } i$

$$\text{OPT}(i) = \max \left(v_i + \sum_{\substack{g \in \text{grandchildren}(i) \\ \sum_{c \in \text{children}(i)}}} \text{OPT}(g) \right)$$

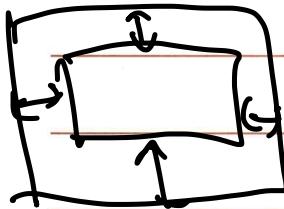
$\Theta(n)$ Time

③ Box Type

height $h_i = 2$	$2 \times (3 \times 5)$
width $w_i = 3$	$3 \times (2 \times 5)$
depth $d_i = 5$	$5 \times (2 \times 3)$
	↑ ↑ ↑ height depth width

Instead of n boxes

we can consider $3 \times n$ boxes that
cannot rotate



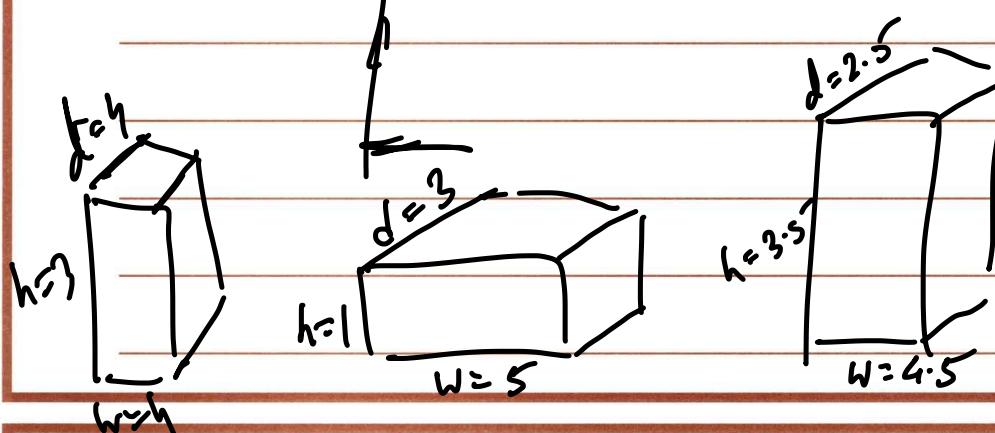
Can sort boxes based on
Base Area, depth, width.

Suppose we sort acc. to Base Area

$h(j) =$ Height of tallest stack of boxes i...j

$$H(j) = \max [H(j-1), \max_{i \leq k \leq j} (h_j + H(k))]$$

AND
 $w_k > w_j$
 AND
 $d_k > d_j$



$$H(1) = 3 \leftarrow$$

$$H(2) = 3$$

$$H(3) = 3.5 + 3 = \cancel{6.5}$$

\uparrow
 h_j $H(3)$ came from $H(2)$

$H(2)$ came from $H(1)$

~~Box 3~~ compatible with ~~Box 2~~

Box 3 and Box 1 not compatible

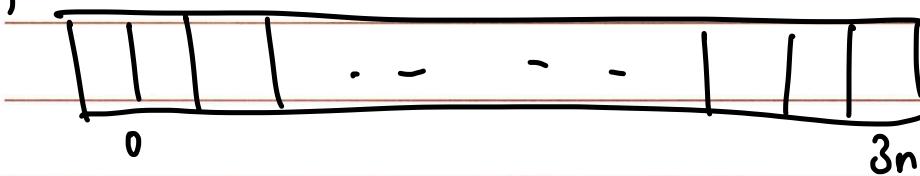
\therefore above formula is wrong

$H(j)$ = height of tallest stack of boxes $1..j$ with
f at the top of stack.

$$H(j) = h_j + \max \left[H(i) \right]$$

$i < j$ AND
 $w_j < w_i$ AND
 $d_j < d_i$

$H()$



Need to search for highest stack when
we are done filling out H array.

Takes $O(n^2)$