

# STANFORD ALGOS 3

## WEEK 1 LI

### GRAPHS AND THE INTERNET

Claim: the Internet is a graph

[vertices = end hosts and routers, directed edges = direct connections]

BT

Other graphs related to the Internet:

- Web graph
- Social Networks

### INTERNET ROUTING

Suppose Stanford needs to send Cornell some data

Which route to use? We'd want to use the shortest path  
(algorithm)

Issue? Too many vertices  $\nearrow X$ .  
(and changing weighted edges)

Can't use Dijkstra's, but can use Bellman-Ford, which  
is non-controlled and can deal with negative edge lengths  
(but you just said it was useless?)

2)

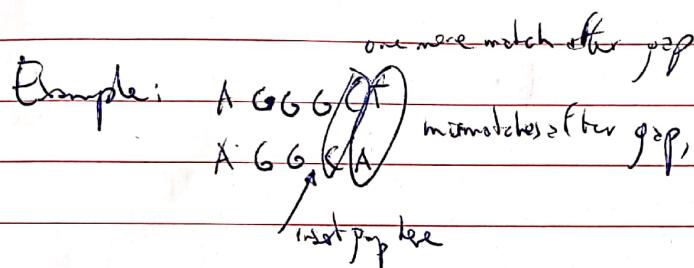
## WEEK 12

### SEQUENCE ALIGNMENT

Input: two strings over the alphabet  $\{A, C, G, T\}$  (or more general)

Problem: How "similar" are they?

- Application examples:
- (1) Extrapolate function of genome substrings
  - (2) Similar genomes to infer proximity in evolutionary tree

Example:  


Penalties for mismatches and gaps!

Key

### PROBLEM

Input: 2 strings over  $\{A, C, G, T\}$ , penalties  $\text{pen}_{\text{gap}} \geq 0$  and  $\text{pen}_A \geq c$

Output: alignment of strings that minimizes the total penalty

Total penalty is called Needleman-Wunsch score

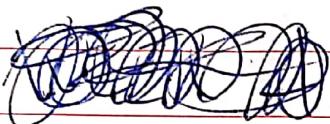
We need an efficient way to find it, and we can't Brute Force it

A  $500 \times 500$  letter string has more alignments than atoms in the universe  
 $(10^{125})$

## Some Algorithm designs.

- Divide and conquer
- Randomized algorithms
- Dyn. Greedy algorithms
- Dynamic Programming.

Greedy algorithms make "myopic" decisions and hope everything works out at the end. Process every state point once

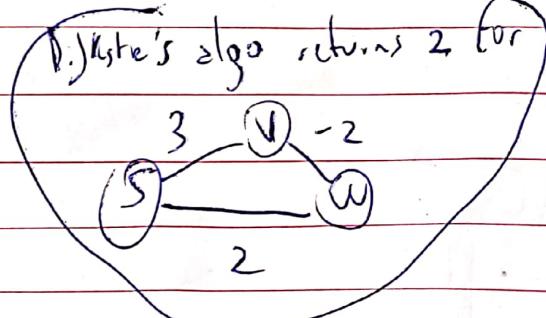


## GREEDY VS DIVIDE & CONQUER

- ① Easy to propose multiple greedy algorithms for problems
- ② And to calculate their running time too
- ③ Hard to establish correctness, even though you might find a correct one you'll have a hard time understanding why.

Dijkstra's algorithm <sup>(5)</sup> is easy to prove-ish with positive lengths

It's easy to prove incorrectness!



Three strategies for PO Correctness:

- ① Induction
- ② Exchange (shouldn't change for an optimal solution)
- ③ Whether it works!

(4)

## 1st WEEK 4

### THE CACHING PROBLEM

- small fast memory (the cache)
- big slow memory
- process sequence of "page requests"
- one fault, what do you evict from cache?

Example: Cache  $\boxed{a \ b \ c \ d}$

Request sequence: c ✓  
d ✓

e ✗ What to do? We have to  
load e, but what do we evict?

Let's say we evict:  $\boxed{e \ b \ c \ d}$  ?

f ✗  
we evict b  $\boxed{e \ f \ c \ d}$

g ✗

b ✗

So in the end we had request sequence cdefab and if we  
we could have gotten only 3 if we ~~hadn't~~ evicted c and d

What eviction algorithm to choose?

### THE OPTIMAL CACHING ALGORITHM

Theorem: (Belady 1960) The "furthest-in-future" algorithm is optimal (but what is it?).

Good benchmark for caching algorithms. LRU should work well enough as long as state exhibits locality of reference.

# WEEK 1 LS

## SCHEDULING

Setup: - one shared resource (example: processor)  
 - many jobs to do

Question: What order do we do the jobs in?

Assume that each job  $j$  has

- weight  $w_j$
- length  $l_j$

Completion Time  $c_j$  of job  $j = \text{sum}$   
 of job lengths up to and including  $j$   
 Example: 3 jobs 1, 2, 3 with  $l_j$  1, 2, 3.  
 $c_j = 1, 3, 6.$

We are trying to optimise the weighted sum of completion times:

$$\min \sum_{j=1}^n w_j c_j$$

Back to our  ~~$(1, 2, 3)$~~   $(l_j (1, 2, 3) w_j (3, 2, 1))$  example:

$$3 + 6 + 15 = 24$$

How to do this in general?

6]

## WEEK 1 L 6

### 1 ALGO FOR SCHEDULING

Intuition  $\rightarrow$  do them in order of smallest length first, <sup>bigger weight</sup> later

These are decent greedy algorithms on their own, better if combined.

Idea: assign jobs to slots, say  $5 \rightarrow 3, 1, 2, 4, 6$  from left to right

1 Guess  $\sigma = w_j - l_j$   $\rightarrow$  Easy to find examples contradicting both  
 Guess  $\sigma' = \frac{w_j}{l_j}$

Always correct!

## WEEK 1 L 7

## AND L 8



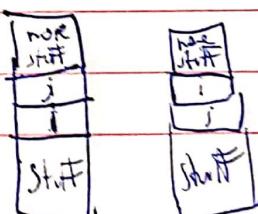
PROOF: by Exchange Argument

Fix  $n$  jobs, proceed by contradiction

Let  $\sigma$  = greedy schedule,  $\sigma^*$  = optimal schedule

Assume:  $\frac{w_i}{l_i}$ 's distinct. Assume you order them  $\frac{w_1}{l_1} > \frac{w_2}{l_2} \dots$

$\sigma$  is  $1, 2, 3, \dots, n$ .  
 If  $\sigma^* \neq \sigma$ , there's consecutive jobs ( $i, j$ ) where  $i >$



What is the effect of swapping these two jobs is that

$c_j$  goes up and  $c_i$  goes down

by ~~l<sub>i</sub> · w<sub>j</sub>~~  
l<sub>i</sub> · w<sub>j</sub>

~~by l<sub>j</sub> · w<sub>i</sub>~~  
l<sub>j</sub> · w<sub>i</sub>

$$\text{As we said, } \frac{w_i}{l_i} < \frac{w_j}{l_j} \Rightarrow l_i w_j > l_j w_i$$

So if  $\sigma^*$  wasn't equal to  $\sigma$ , it would be a contradiction to it being optimal.

8]

## WEEK 1 L10

### MST PROBLEM.

Goal: Connect a bunch of vertices together as cheap as possible

Applications: clustering, networking

Blazingly fast Greedy Algorithms:

- Prim's Algorithm
- Kruskal's Algorithm

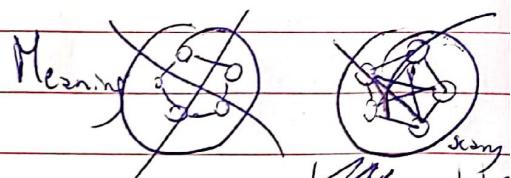
$(O(m \log(n)))$ , really fast, almost linear!

### DEFINITION OF PROBLEM

Input: undirected graph, (we won't talk about directed graph branching)  
 $G(V, E)$  cost  $c_e$  for each edge

- Assume adjacency list representation, no need to exclude negative edges

Output: minimum ~~cost~~ tree  $T \subseteq E$  that spans all vertices



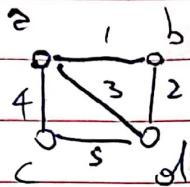
### STANDING ASSUMPTIONS

#1 Input graph  $G$  is connected (which is easy to check in pre processing)  
with DFS

#2 Edge costs are distinct. (Prim + Kruskal remains correct though)

# WEEK 1 L11

## EXAMPLE OF PRIM-KRUSCAL



Like with Dijkstra's, we'll grow our tree vertex by vertex

~~Starting from b → (b, e) → (b, d) → (b, d, c)~~

~~Starting from b →  $b_1$  →  $b_1 + bd$  →  $b_1 + bd + ec$~~

## DEFINITION OF PRIM'S MST ALGO

- Initialize  $X = \{s\}$ , ( $s$  chosen arbitrarily,  $s \in V$ )

- $T = \emptyset$  (initially formed by tree w.r.t.  $\emptyset$ )

- while  $X \neq V$ :

- let  $e = (u, v)$  be the cheapest edge between  $(X)$  and  $(V - X)$
- add  $v$  to  $X$
- add  $e$  to  $T$

Theorem: PRIM'S MST ALGO is CORRECT

Proof part I: Computes a spanning tree  $T^*$

Proof part II:  $T^*$  is an MST

10

# WEEK 1 L12

## CUTS

Claim: Prim's algo outputs a spanning tree  $T^*$ .

Definitions: a cut of a graph  $G = (V, E)$  is a partition of  $V$  into 2 non-empty sets.

Empty cut Lemma: a graph is not connected  $\Leftrightarrow$  ∃ cut  $A, B$  with

Double crossing Lemma: Suppose a cycle  $C \subseteq E$  crosses a cut; it has to cross it twice

Lonely cut corollary: if  $e$  only  $\in E$ , not in any cycle



## PROOF:

- ① It maintains invariant that  $T$  spans  $X$
- ② Can't get stuck with  $X \neq V$ , otherwise the cut  $(X, V-X)$  would be empty. By Empty cut lemma, the input graph  $G$  is disconnected.

*there no Cs to proceed on*

- ③ No cycles ever get created in  $T$  because you only ever choose one edge leading to a new vertex that gets added to  $X$

# WEEK 1 L13

L11

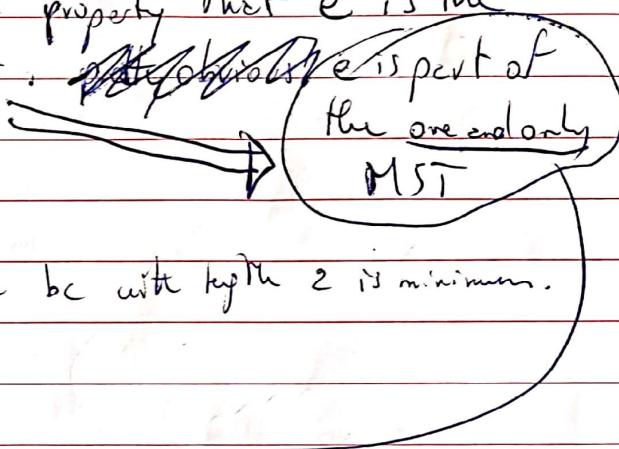
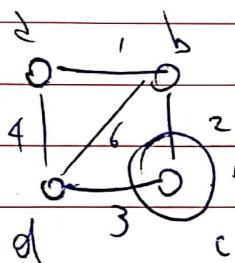
## Correctness CORRECTNESS OF PRIM PROOF PART 2

Theorem 2: Prim's Algo computes minimum-cost  $T^*$

Key Question: When is it "safe" to include an edge in the tree-sofar?

CUT PROPERTY: consider an edge  $e$  of  $G$ :

Suppose there is a cut  $(A, B)$  with the property that  $e$  is the cheapest edge of  $G$  that crosses it.



← two edges crossing the cut, edge  $bc$  with weight 2 is minimum.

→ PROOF IN L14

12)

## WEEK 1 L105

### RUNNING TIME OF PRIM'S ALGO

- Init  $X = \{s\}$
- $T = \emptyset$
- while  $X \neq V$ 
  - let  $e$  be cheapest  $x \in X$   ~~$\rightarrow$~~   $(x)(N-x)$  edge
  - add  $e$  to  $T$ , add  $v$  to  $X$

Running Time of naive implementation :

- $O(n)$  iterations
- $O(m)$  time per iteration

$O(nm)$  this is ~~nothing~~ nothing to sneeze at, but ~~looks~~ worse than  $O(mn)$

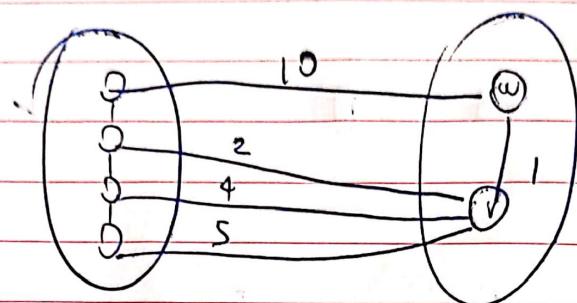
(so we do better?)

Yes, with heaps, just like with Djikstra's

Heapify runs in  $O(m\log n)$ , allows us to run Prim's in  $O(n\log n)$

## WEEK 1 L16

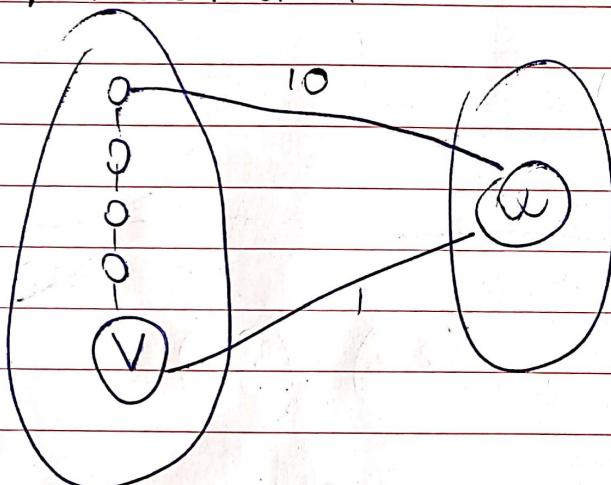
## MAINTAINING INVARIANT #2

~~Keep min heap~~

$$\text{Key}[u] = 10$$

$$\text{Key}[v] = 2$$

After one iteration:



$$\text{Key}[u] = 1$$

We need to recompute some keys every step of the min-while loop

Pseudocode: when  $v$  colored to  $X$ :

- foreach edge  $(v, w) \in E$ 
  - - if  $w \in V - X$  (the keys might have changed for these):
    - Delete  $w$  from heap
    - recompute  $\text{Key}[w] := \min\{\text{Key}[w], c_{vw}\}$
    - re-insert  $w$  into heap.

$\leq (n-1)$  inserts  
 $\leq (n-1)$  Extract mins  
 $\leq m$  Delete-Insert combos

N of operations at most  $O(m \log(n))$  in

[4]

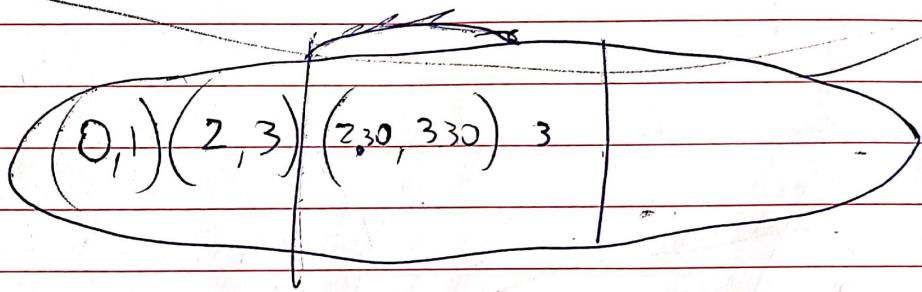
## Prob scribbling

$s_i, t_i, \dots$

$$(0, 1) (1, 2) (2, 3) (4, 5) (4, 6)$$

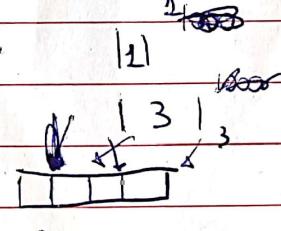
$$(0, 1) (1, 3) (2, 5) (4, 4, 3, 0) (4, 3, 0, 5)$$

Delta not earliest start time



\*

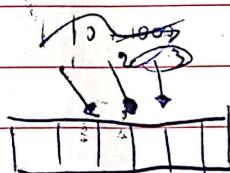
earliest finish



$$C_j - d_j$$

$$(2, 2) (3, 2)$$

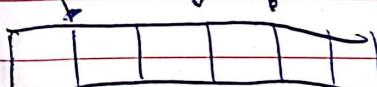
Delta not necessary earlier



$$(p_j + p_{j-1} + p_{j-2}) - d_j$$

$$1 \quad 5 \quad 1$$

$$d_1 \quad d_2 \quad d_3$$

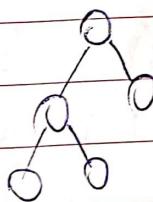
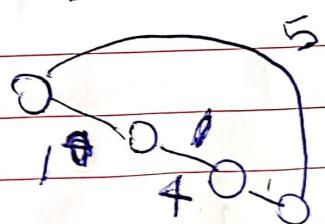
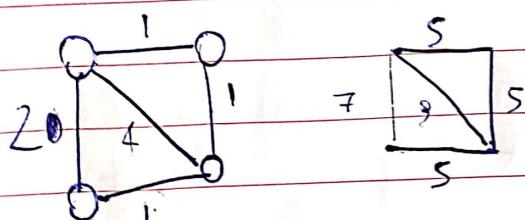
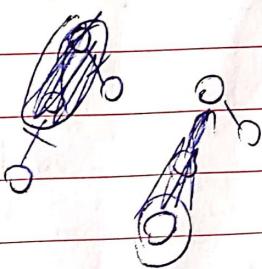
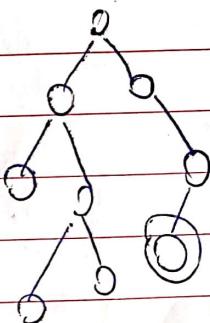
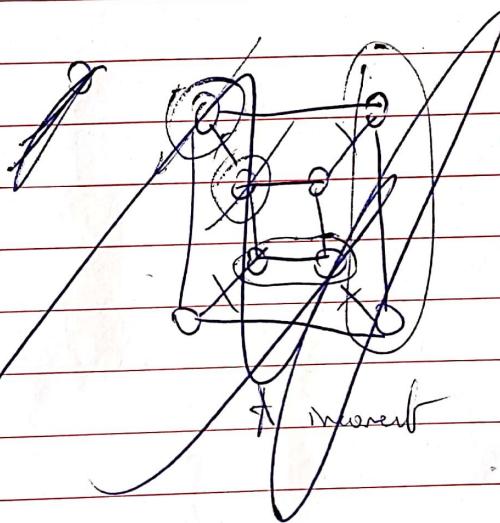
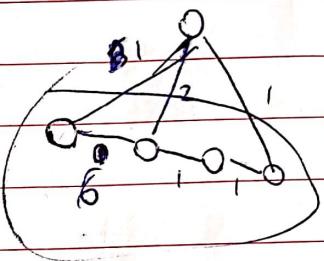
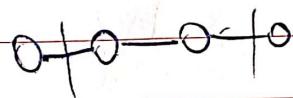
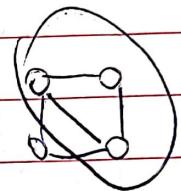
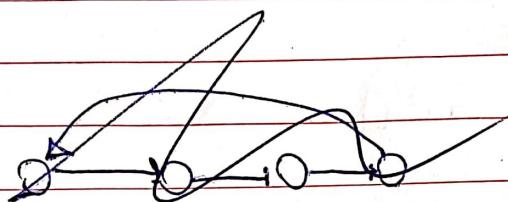
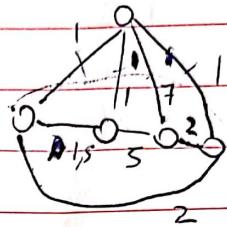
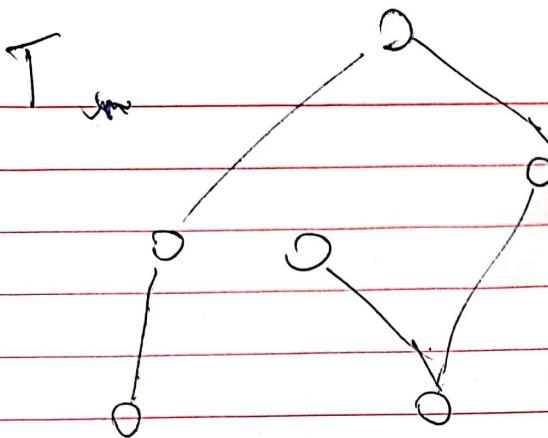


$$p_j - d_j = x_j, \quad p_j - d_j = \cancel{p_j - d_j}$$

$$= (p_j + p_{j-1} + p_{j-2}) - d_j - d_{j-1} + d_{j-2}$$

$$p_j + x_j + (d_j - d_j) \xrightarrow{\text{max}}$$

$$= p_j + x_j - d_j \text{ rel.} = \cancel{p_j} x_j + x_j + \cancel{(d_j)}$$



# WEEK 2 L1

## KRUSCAL'S MST ALGORITHM

- Uses union-find data structure

Cut Property: if  $e$  is the cheapest edge

crossing any cut  $(A, B)$  of  $G(V, E)$  then  ~~$e$  belongs to  $T$~~   $e \in T$

the  
mst

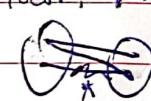
- Sort edges in order of increasing cost (reverse edges 1, 2, 3, ..., m so that)  $\text{cost}[1] < \text{cost}[2] < \dots < \text{cost}[m]$
- $T = \emptyset$
- For  $i = 1$  to  $m$ :
  - if  $T \cup \{e_i\}$  has no cycles
  - add  $e_i$  to  $T$
- return  $T$

## WEEK 2 L2: Proving Correctness of Kruskal

Theorem: Kruskal's algo is correct

- ① No cycles, we check in loop
- ②  $T^*$  is connected (when you're adding edges with one vertex outside of  $X$ )
- ③ more rigorous proof of ②:

fix cut  $(A, B)$ , since  $G$  connected,  $\exists$  one edge, at least, Kruskal will include first edge crossing  $(A, B)$  that it sees.



- ④ every edge in  $T^*$  is justified by the Cut Property (it really is an MST)

On the contrary, since  $T \cup e(u, v)$  has no cycle, there is no  $u, v$  path in  $T$ .

which shows that there exists a cut  $A, B$  with  $u \in A, v \in B$ .

⑤ by  $T^*$  being connected, no edges crossing  $(A, B)$  were previously considered by Kruskal, so  $(u, v)$  is the cheap

## WEEK 2 L2

### KRUSKAL IMPLEMENTATIONS

#### NAIVE UNION-FIND

- sort edges in order of increasing cost
- $T = \emptyset$
- for  $i = 1 \rightarrow m \rightarrow O(m)$  iterations
  - if  $T \cup \{i\}$  has no cycles
    - add  $i$  to  $T$
- return  $T$

why not  $m$ ?

at most  $n^2$

$$O(m \log(n)) \quad \log m \leq \log(n)$$

→ How much work do we do every iteration?

- $O(\log n)$ ? Checking for a cycle boils down to finding if there is a  $(v_i, v_j)$  path in  $T$ . DFS or BFS,  $O(n)$ , if you find a path you're done.

~~Running time~~

$$\rightarrow O(mn) + O(m \log(n)) = O(mn)$$

(CAN WE DO BETTER?)

Plan: data structure for  $O(1)$  - time

Cycle checks  $\Rightarrow O(m \log(n))$  time +  $O(m) = O(m \log(n))$  time

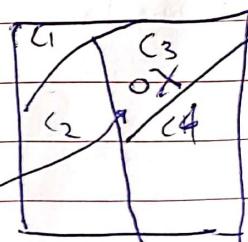
### THE UNION-FIND DATA STRUCTURE

Raison d'être of a union-find data structure: maintain partition of a set of objects

If we have it a value, we want to know what partition that value is in

$find(x)$ : return partition  $x$  belongs to

returns (3)



$Union(C_i, C_j)$ : fuse groups containing  $C_i$  and  $C_j$  into a single one

If it is useful for Kruskal's algorithm:

objects = vertices

adding new edge  $(u, v)$  to  $T \oplus \leftrightarrow$  fusing partitions (components) of  $u, v$

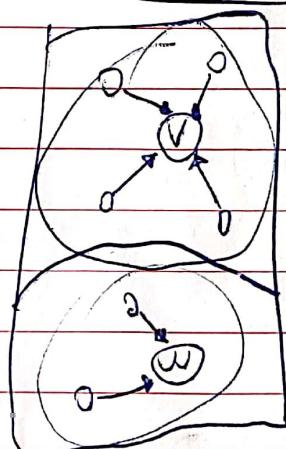
## WEEK 2 (3)

### ■ BASICS OF UNION-FIND

Goal #1 - maintain one linked structure per connected component of  $(V, T)$

- Each component has an arbitrary leader vertex

Invariant: each vertex points to the leader of its ~~component~~  
"name" of a component is inherited from leader vertex



Key point: given edge  $(u, v)$ , can check if  $u, v$  is already in the same component in  $O(1)$  time:  
 $\{ \text{Find}(u) = \text{Find}(v) \}$ .

### ■ HOW TO MAINTAIN THE INVARIANT (when trying)

Note: when new edge  $(u, v)$  added to  $T$ , connected components of  $u$  and  $v$  merge

Question: At most how  $O(n)$  leader pointer updates for  $n$  rats  
wif. we could change  $w \xrightarrow{\text{points to}} v$ ?

We're happy already with  $O(m \log n)$  work. Every time two sets merge we  
at most double and it can happen  $\log_2 n$  times

# WEEK 2 @ LS

## CLUSTERING

MSTs have multiple greedy algorithms which are correct.

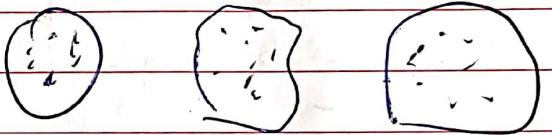
Internal Goal: cluster connected objects into coherent groups

Assumption: ① as input, given a (dis)similarity measure, a distance  $d(p, q)$  between each pair of objects

$$\textcircled{2} \quad d(p, q) = d(q, p) \quad (\text{symmetric})$$

Examples: Euclidean distance, genome similarity, etc.

Goal: Some cluster  $\Leftrightarrow$  nearby



## MAX-SPACING K-CLUSTERING

Assume: you want  $K$  clusters in total

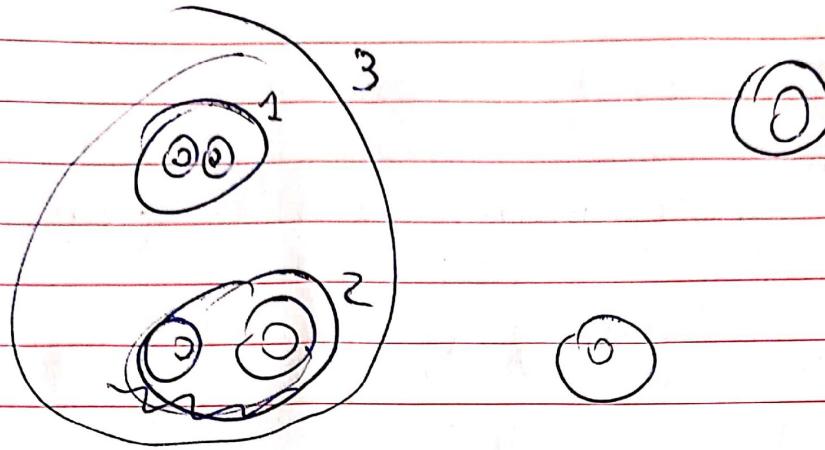
Two points  $p \in q$  separated if they are assigned to different clusters.

Definition: The spacing of a  $K$ -clustering is the minimum distance after which two objects end up in different clusters. The bigger, the less clusters there are.

Problem statement

Problem statement: Given a distance measure  $d$  and  $K$  ~~connected them~~ compute the  $K$  clustering with maximum spacing.

## A GREEDY ALGORITHM



How to decrease # clusters? 1, 2, 3, keep increasing distances

### Pseudocode:

- Initially, each point in separate clusters
- Repeat until only K clusters:
  - let  $p, q =$  closest pair of separated points
  - merge clusters containing  $p, q$  into a single one

Reminds you of Kruskal's, but aborted early

Objects are the vertices, distances are like edge costs.

This type of fusing clusters is called single link clustering.

## WEEK 2 L6

### CORRECTNESS CLAIM

Theorem: single-link clustering finds the max-spacing K-clustering

Proof: Let  $C_1, \dots, C_k$  = greedy clustering with spray  $S$   
 Let  $\hat{C}_1, \dots, \hat{C}_n$  = arbitrary other clustering

Need to Show: spray of  $\hat{C}_1, \dots, \hat{C}_n$  is  $\leq S$

### CORRECTNESS PROOF

Case 1: If  $\hat{C}_i = C_i$ , spray  $\hat{S} \leq S$

Case 2: Otherwise, can find 2 point pair  $p, q$  such that:

(A)  $p, q$  in the same greedy cluster  $C_i$

(B)  $p, q$  in different clusters  $\hat{C}_i, \hat{C}_j$

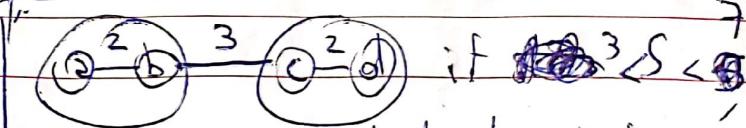
Property of greedy algorithm: if two points  $x, y$  were merged,

$d(x, y) \leq S$  (since  $S$  is ever increasing)

Easy Case: if  $p, q$  were merged,  $S \geq d(p, q) \geq$  the spray in  $\hat{C}_1, \dots, \hat{C}_n$

Tricky Case:  $p, q$  indirectly merged

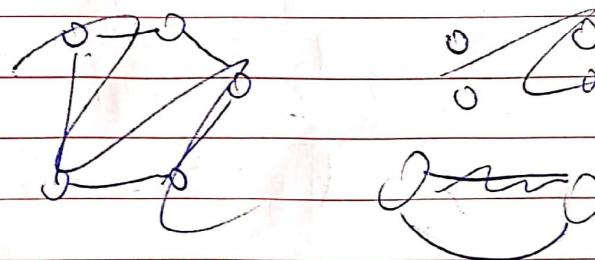
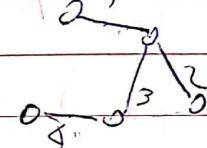
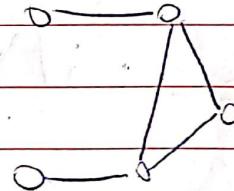
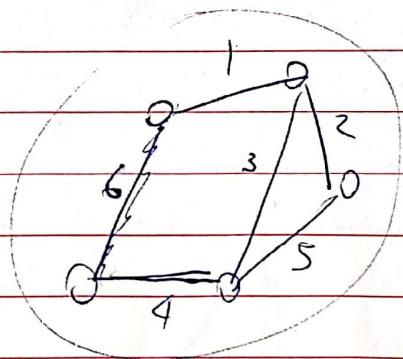
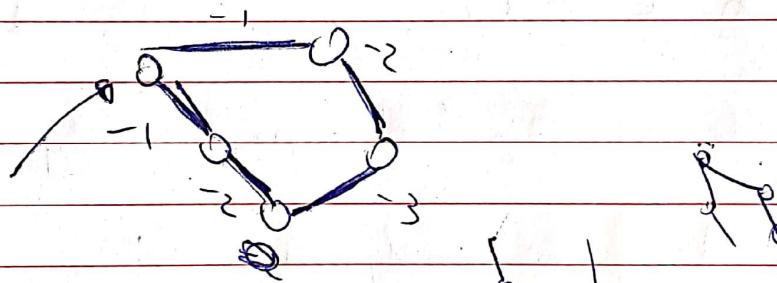
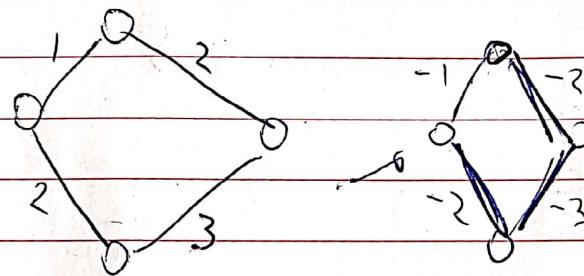
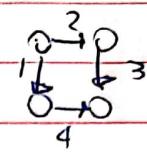
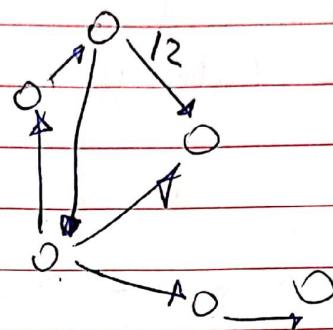
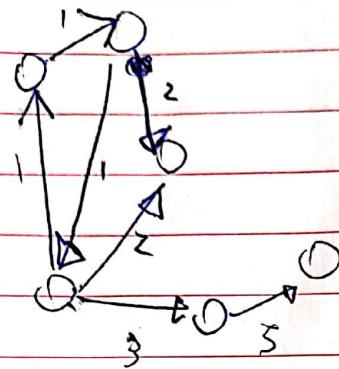
Since  $p \in \hat{C}_i$  and  $q \notin \hat{C}_i$   $\exists$  consecutive pair  $z_j, z_{j+1}$  with  $z_j \in \hat{C}_i, z_{j+1} \notin \hat{C}_i$



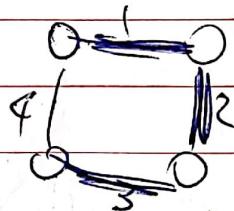
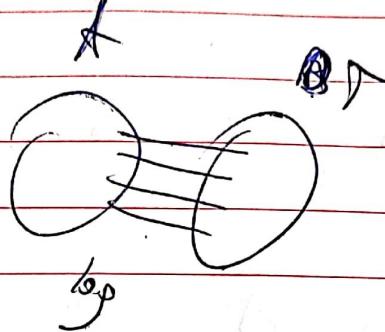
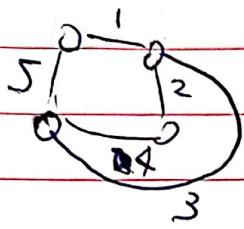
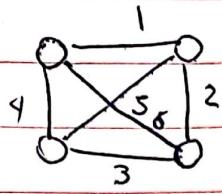
$\Rightarrow$  end  $b$  and  $d$  end up in same group without ever having been merged directly

21

## Problem scribbling



~~edge~~



\* charge.

then O(1) to Red lookup

# WEEK 3 LI

## BINARY CODES

Binary code: maps each character of an alphabet  $\Sigma$  to a binary string.

Example:  $\Sigma = \{a - z\}$ ,  $\overset{a}{00000}, \overset{b}{00001}, \dots$

### AMBIGUITY

Hypothetical:  $\Sigma = \{A, B, C, D\}$ . Fixed length encoding is  $(0, 010)$

(cant be  $\{0, 01, 10, 1\}$ )! 001 could be AB or AAD

## PREFIX-FREE CODES

Problem: with variable length codes it's never clear where one character ends and another begins

Solution: prefix-free codes - make sure that for every pair  $(i, j) \in \Sigma^*$ , neither of the encodings is a prefix of the other

Example:  $\Sigma = \{0, 10, 110, 1101\}$

## BUT WHY WOULD WE WANT TO DO THIS?

	Frequency		
A	60%	00	0
B	25%	01	10
C	10%	10	110
D	5%	11	11

→ Less space to store,  
on average 1.55 bits

# WEEK 3 L 2

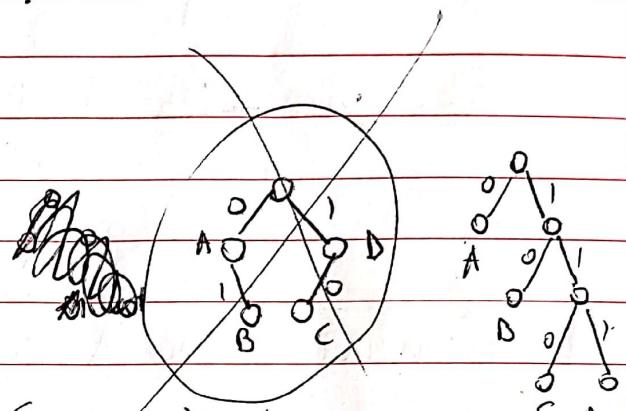
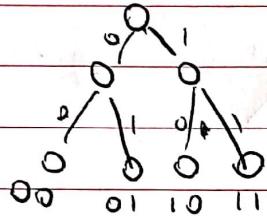
## HUFFMAN CODES, PROBLEM DEFINITION

### CODES AS TREES

Goal: best binary prefix-free encoding

(useful to think of it as a binary tree)

Example:



$$\{00, 01, 10, 11\} = \{A, B, C, D\} \neq \{0, 01, 10, 11\}, \text{ ambiguous}$$

You need  $i$  endpoints / leaves to represent  $i$  values

To decode  $i$ : repeatedly follow path from the root until you hit a leaf (so  $\log(i)$  time?)

### Problem definition

Input: probability  $p_i$  for each character  $i \in \Sigma$

Notation: if  $T$  is the binary tree with leaves  $\leftrightarrow$  symbols of  $\Sigma$ ,  
 $L(T) = \sum_{i \in \Sigma} p_i \cdot (\text{depth of } i \text{ in } T)$

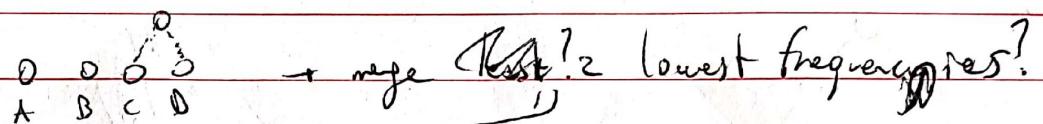
Output:  $T$  with minimum  $L(T)$

# A GREEDY ALGORITHM

Question: best approach to building  $T$ ?

Natural idea: ① Partition  $\Sigma$  into  $\Sigma_1, \Sigma_2$  each with  $\approx$  of total frequency,  
② put one side on the left the other on the right

Huffman's (optimal) idea: build tree bottom up using successive merges



→ WHAT GREEDY APPROACH TO USE?

Question: which pair of symbol is "safe" to merge?

Observation: final encoding length of  $T \in \Sigma = \sum p_i l_i$  of mergers its subtree endev.

Recurse? Now?

Suppose: 1st iteration of algorithm merges symbols  $a, b$  with  $p_{ab} = p_a + p_b$

Idea: Replace symbols  $a, b$  by meta symbol  $ab$  with  $p_{ab} = p_a + p_b$

- I proved my ht!

Huffman's Algorithm ( $\Sigma$ ):

- If  $|\Sigma| = 2$ , return

- Let  $a, b \in \Sigma$  have the smallest frequencies

- Let  $\Sigma' = \Sigma$  with  $a, b$  replaced by  $(ab)$ , with  $p(ab) = p(a) + p(b)$

- Recursively compute  $T'$  (for  $\Sigma'$ )

- Extend  $T'$  (with leaves  $\leftrightarrow \Sigma'$ ) to tree  $T$  with leaves  $\leftrightarrow \Sigma$  by splitting leaf  $ab$  back into  $a$  and  $b$ .

- Return  $T$

Input

A B C D E F  
3 2 6 8 2 6

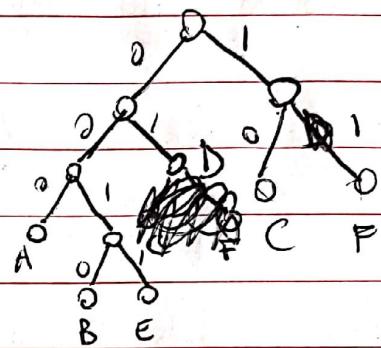
(if you don't like to use fractions)

Step 1 : Merge B and E : A ~~D E~~ C D ~~F~~  
3 4 6 8 6

Step 2 : Merge A and BE → A(BE) C D ~~F~~  
7 12 6 8 6

Step 3 : Merge C and F → (A(BE)) (C F) ~~D~~  
7 12 8

Step 4 : Merge A(BE) and D → (A(BE)) (C F) ~~D~~ =



A : 000

D : 01

B : 0010

E : 0011

C : 0010

F : 11

# WEEK 3 CS16

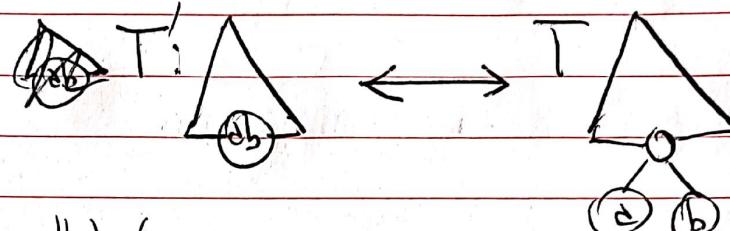
## PROOF OF CORRECTNESS OF HUFFMAN'S

We prove it by induction on  $n = |\Sigma|$

Base case, when there are 2 elements, algorithm outputs the optimal tree.

Inductive hypothesis: Algorithm does solve smaller subproblems optimally.

Consider



$a, b$  = symbols with smallest frequencies

$$p(a+b) = p(a) + p(b)$$

~~$$L(T) - L(T') = \text{Offset} + p_a \times (\text{depth in } T - p_a) + p_b \times (\text{depth in } T - p_b)$$

$$= \text{Offset} + (p_a + p_b) = p_a + p_b$$~~

$$\begin{aligned}
 L(T) - L(T') &= p_a [\text{a's depth in } T] + p_b [\text{b's depth in } T] - p_a - p_b \\
 &= p_a(d+1) + p_b(d+1) - (p_a + p_b)d = p_a + p_b
 \end{aligned}$$

## NOTES ON RUNNING TIME

Naive implementation is  $O(n^2)$

A heap is perfect for this, (extract min twice + reinsert new next symbol)  $\cdot n$ .

-  $O(n \log(n))$

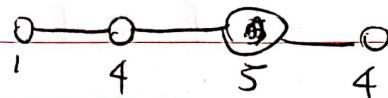
BUT WAIT 2 QUESTIONS DOES IT IN WRS +  $O(n)$  time

## WEEK 3 L7

# DYNAMIC PROGRAMMING

**PROBLEM:** FINDING WIS, WEIGHT INDEPENDENT SETS

Input: path graph  $G(V, E)$ , non-negative ~~edges~~ vertex weights



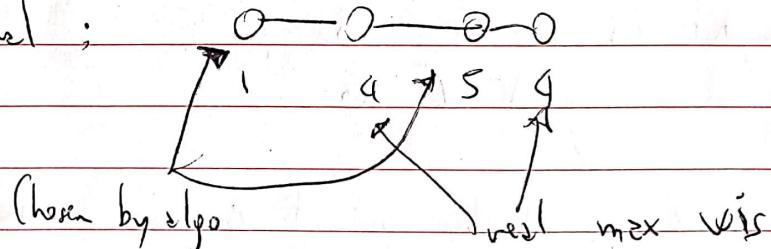
Desired output is a subset of non adjacent vertices of maximum total weight.

Brute Force is unfeasible as it requires exponential time

## **GREEDY APPROACH**

Ideas: Iteratively choose max weight vertex not adjacent to any previously chosen vertex

Non optimal:



## **D&C APPROACH**

Ideas: split in two and then combine

but recursive subproblems might conflict



## WEEK 3 L8

### THE OPTIMAL SUBSTRUCTURE OF LIS

Last lecture we tried to guess a GT quickly, but often

to find a good one it is necessary to reason about the structure of what an optimal solution would have.

We can narrow down the set we have to brute force through

Notation: let  $S \subseteq U$  be a max weight IS  
let  $v_n$  = last vertex of path

Case 1:  $v_n \notin S$ .  $G' = G - \{v_n\}$

After deleting  $\{v_n\}$ ,  $S$  is still the max in  $G'$

Case 2: if  $v_n \in S$ ,  $S - \{v_n\}$  is the max in  $G''$   
 $(v_n \text{ deleted})$

Upshot: A max weight IS must be either:

- a max weight IS of  $G'$  or
- ←  $v_n + \text{ a max weight of } G'$

Corollary: if we knew whether or not  $v_n$  was in the max, we could recursively compute the max weight IS of  $G''$  and be done.

Could we try both possibilities and return the better solution?

YES

# WEEK 3 L9

## WHAT TIME WOULD THE RECURSIVE ALGO RUN IN?

$\Theta(n)$ !

worst case scenario each iteration removes one vertex, best case tree.

But it runs  $n$  times, sometimes on the same subproblems more than once.

To eliminate redundancy, cache every subproblem in a hashtable to get  $O(1)$ -time lookup later on.

Or better: bottom-up iterative algorithm: instead of plucking a dot!

$G_i = \text{1st } (i) \text{ vertices of } G$

Populate left to right with  $A[i] = \text{value of the maxis of } G_i$

$$A[i] = \max(A[-1], A[i-2] + w_i)$$

Runtime is  $O(n)$ !

OPTIMAC

L10

We can trace back through the array to reconstruct the optimal solution.

Let  $A$  = filled in array

10	4	7	7	...	18	F
0	1	2	3		$n$	

Runs in  
 $O(n)$

- Let  $S = \emptyset$

- While  $i \geq 1$  (from right to left)

- If  $A[i-1] \geq A[i-2] + w_i$

- Decrease  $i$  by 1

- Else

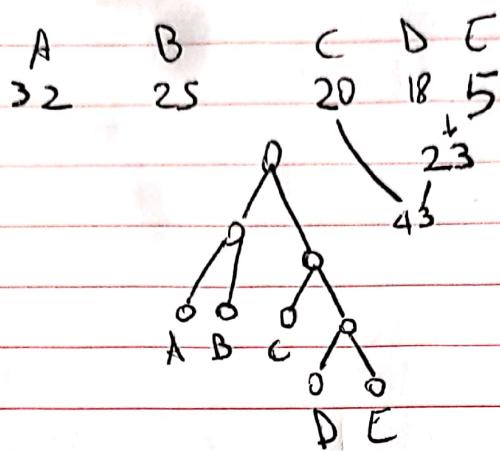
- Else  
- Else  
- Add  $v_i$  to  $S$ ,  $i = 2$   
Return  $S$

## WEEK 3 L11

### PRINCIPLES OF DYNAMIC PROGRAMMING

- (1) Identify a small number of subproblems (example: MDPs  
for  $G_i = 0, 1, \dots, n$ ,  
not  $2^n$ )
- (2) Can quickly solve larger subproblems given the solutions to "smaller subproblems" (usually through recurrence)
- (3) After solving all subproblems, quickly compute the final solution.

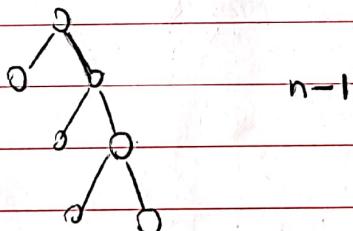
# Problem Subbling



$$32 \cdot 2 + 25 \cdot 2 + 20 \cdot 2 + 18 \cdot 3 + 5 \cdot 3 =$$

$$\begin{array}{r}
 64 + \\
 50 + \\
 40 + \\
 54 + \\
 15 = \\
 \hline
 \end{array}$$

$$223 \cdot \frac{1000}{100} = 2230$$



$n=1$



0.33, 0.33, 0.33

0.4, 0.45, 0.



# WEEK 4 L1

## THE KNAPSACK PROBLEM

- ① Input: n items, each with value:
  - = value  $v_i$  (non-negative)
  - = size  $w_i$

As much value as possible, size limited set  $S \subseteq \{1, 2, 3, \dots, n\}$

## DEVELOPING A GA FROM OPT. SUBSTRUCTURE

Let  $S$  be a max-value solution to an instance of knapsack.

Delete item  $n$ .

(case 1: item  $n \notin S$ ,  $S$  still optimal for  $n-1$  items)

(case 2: item  $n \in S$ ,  $S$  still optimal for first  $n-1$  items with capacity decreased by  $w_n$ )

# WEEK 4 L2

## RESUMING FROM LAST TIME

Notation:  $V_{i,x}$  = value of solution that uses first  $i$  items, size  $\leq x$

$$V_{i,x} = \max \{ V_{(i-1),x}, V_i + V_{(i-1), x-w_i} \}$$

+ Recurrence

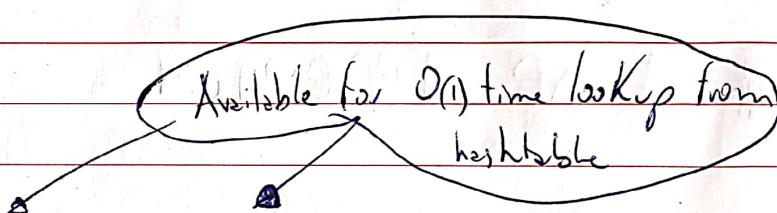
- Let  $A = 2D$  Array

- Initialize  $A[0,x] = 0$  for  $x = 0, 1, \dots, W$

- For  $i = 1, 2, \dots, n$

- For  $x = 0, 1, \dots, W$

$$A[i,x] = \max \{ A[i-1,x], A[i-1,x-w_i] + V_i \}$$



Return  $A[n,W]$

RUNNING TIME IS  $O(nW)$

Example:

$$W = 6$$

$V_1 = 3, w_1 = 4$	$V_2 = 2, w_2 = 3$	$V_3 = 4, w_3 = 2$	$V_4 = 4, w_4 = 3$
6	5	4	3
0 3 3 7 8	0 3 3 6 8	0 0 2 4 4	0 0 0 4 0
0 3 3 4 4	0 0 2 4 4	0 0 0 4 0	0 0 0 0 0
x 0 0 0 0	x 0 0 0 0	x 0 0 0 0	x 0 0 0 0
i: 0 1 2 3 4			

Then: Start from  $(4,6) = 8$ .

Remove item 4.  $6 - w_4 = 3$ .

What now?

Remove item 3.  $3 - w_3 = 1$

0, 0, 0

Optimal: {8, {item3, item4}}

# WEEK 4 L4

## SEQUENCE ALIGNMENT

Example: AGGGCT  
A G G - C A

Needleman-Wunsch score, total penalty =  $\alpha_{gap} + \alpha_{mismatch}$

Input: two strings over some alphabet  $\Sigma$ , penalties  $\alpha_{gap}$  and  $\alpha_{mismatch}$

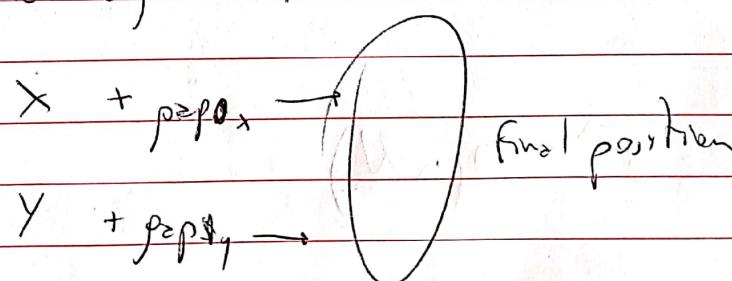
~~try different gaps~~, Insert gaps to equalize lengths of the string

Output: Alignment with minimum penalty

## A DP APPROACH

Identify subproblems from optimal substructure

An optimal alignment:



Either  $p_{gapX} = p_{gapY} = 0$  or  $p_{gapX} = 1^{\text{st}}$  or  $p_{gapY} = 5^{\text{th}}$

- (1)  $x_m \& y_n$
- (2)  $x_m \& p_{gap}$
- (3)  $y_n \& p_{gap}$

## WEEK 4 L5

### RECURSION (APPLYING)

$P_{ij}$  = penalty of optimal alignment of  $x_i$  and  $y_j$

Recurrence: for all  $i = 1, \dots, m$  and  $j = 1, \dots, n$

$$P_{ij} = \min \begin{cases} (1) \alpha_{x,y} + P_{i-1, j-1}, & P_{i-1, j-1} \\ (2) \alpha_{gap} + P_{i-1, j}, & P_{i-1, j} \\ (3) \alpha_{gap} + P_{i, j-1} & \end{cases}$$



$P_{i,0}$  and  $P_{0,j}$  have penalty  $i \cdot \alpha_{gap}$ .

### ALGORITHM PSEUDO CODE (Similar to knapsack)

$A = 2D$  Array

$$A[i,0] = A[0,j] = i \cdot \alpha_{gap}, \forall i \geq 0$$

For  $i = 1$  to  $m$ :

For  $j = 1$  to  $n$ :

$$A[i,j] = \min \left\{ \begin{array}{l} A[i-1, j-1] + \alpha_{x,y} \\ A[i-1, j] + \alpha_{gap} \\ A[i, j-1] + \alpha_{gap} \end{array} \right.$$

All available for  $O(1)$ -time lookup.

Running time is  $O(mn)$

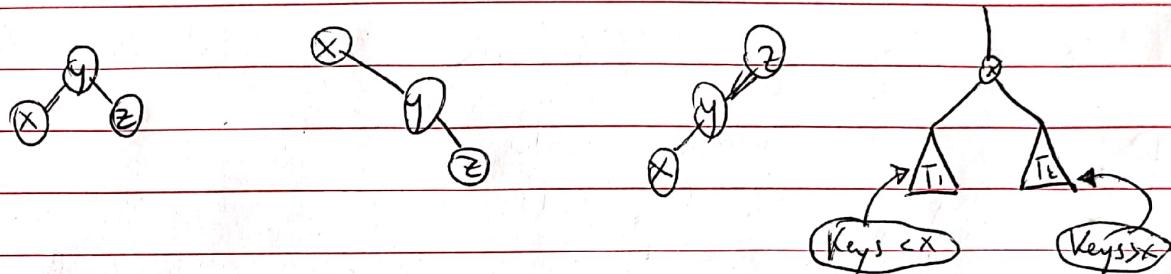
YOU CAN RECONSTRUCT THE SOLUTION in  $O(m+n)$ , From  $A[i,j]$

- IF  $A[i,j]$  filled by (1), match  $x_i$  &  $y_j$  and go to  $A[i-1, j-1]$
- IF  $A[i,j]$  filled by (2), match  $x_i$  with a gap and go to  $A[i-1, j]$
- IF  $A[i,j]$  filled by (3), match  $y_j$  with a gap and go to  $A[i, j-1]$

## WEEK 4 L6

### OPTIMAL BINARY SEARCH TREES

For a given set of keys there are lots of valid search trees



Question: What is the "best" search tree for a given set of keys?

A good answer: balanced search tree, like a red-black tree.

An RBT has height  $\geq \log(n)$  so  $\Theta(\text{height}) = \Theta(\log(n))$

### EXPLOITING NON-UNIFORMITY

Question: Suppose we have keys  $x < y < z$  and we know that

80% of searches are for x

10% of searches are for y

10% of searches are for z

Therefore we want to optimize the search time;



Same as Huffman codes.

The cost changes because we're looking for min depth

instead of  $\Sigma$  (the whole weighted sum)

Also contains symbols only it loses  $\rightarrow$  search tree property

# WEEK 4 C7

## OPTIMAL SUBSTRUCTURE OF BSTs

Inputs: frequencies  $p_1, p_2, \dots, p_n$  for items  $1, 2, \dots, n$ .

Goal: Compute a valid search tree that minimizes the weighted average search time

Ideally, we'd want the <sup>least</sup> frequently searched items at the bottom.

Both bottom up and top down are faulty: choosing a <sup>chosen</sup> root may not always be the most efficient everything up.

What if we calculated it w/ ~~multiple~~ possibilities?

for <sup>an optimal</sup> generic BST:



$T_1$  and  $T_2$  are optimal for respectively items  $\{1, \dots, r\}$  and  $\{r+1, \dots, n\}$

## WEEK 4 C8

### OPTIMAL SUBSTRUCTURE LEMMA

#### RELEVANT SUBPROBLEMS

Let  $\{1, \dots, n\}$  = original keys. For which subsets  $S \subseteq \{1, \dots, n\}$  might we need to compute the optimal BST for  $S$ ?

$$S = \{i, i+1, \dots, i-1, j\} \text{ for every } i \leq j$$

So we'll recurse on that:

### THE ALGORITHM

Important: solve smallest subproblems first.

DPBST (~~A<sub>(20, 2)</sub>, frequencies[ ]~~) :

~~A<sub>(20, 2)</sub> = opt BST  
value of tree (1, ..., j)~~

- A = 2D Array, with  $A(i, j)$  representing opt BST value for items  $\{1, \dots, j\}$
- For  $s = 0$  to  $n-1$  :
  - For  $i = 1$  to  $n$ 
    - $A[i, i+s] = \min_{r=i, \dots, i+s} \left( \sum_{k=i}^{i+s} p_k + A[i, r-1] + A[r+1, i+s-1] \right)$

$O(n^2)$  subproblems,  $O(s-i)$  times,  $O(n^3)$  overall, which is slow

but we can optimize to  $O(n^2)$  if we really want to  
(By remembering previous subproblems and avoiding repetition)

## PROBLEM SCRIBBLING

41

$$\{ \textcircled{1,2}, 1, \textcircled{1,5} \}$$

3 3 6? 5,  
1 2 3

J

Table will still be filled okay

	1	2	3	4	5	6	7
④	5	40	8	4	10	<del>10</del>	23
11	26	11					

Reuse in O(nm) time

(mn?) like subshy?

MWIS only shows 2 values even  $O(1)$

SA uses ~~multiple~~ 2 columns  $O(n)$

BST needs the whole book