

Greedy Algorithms and Dynamic Programming.

sep 17, 2020

Till now, we have seen problems that have some immediate polynomial time (eg, $O(n^3)$, $O(n^4)$) algorithms and we saw some tools like divide and conquer to make them more efficient.

for example $O(n^2) \rightarrow O(n \log n)$
 $O(n^3) \rightarrow O(n)$

Now, we will see more examples where the obvious algorithm takes exponential time and the challenge is to first design some polynomial time algorithm.

Towards this, the two paradigms Greedy & DP are often helpful.

They will involve more clever and subtle ideas than what we have seen till now.

We will study Greedy & DP simultaneously.

pick up a problem and see whether Greedy works or not, whether DP works or not.

Greedy Algorithm

→ local optimality

→ near-sighted

→ don't care about long-run

Proof of correctness is important.

because correctness is not obvious in most cases.

we will see some basic format of how to argue correctness.

Dynamic Programming

→ Recursion with a better memory management.

or → a systematic approach to search through all possible solutions.

Problem 1

You are allowed to choose five apples from a basket. You want to maximize the total weight of your apples.

Greedy approach works here.

→ Choose heaviest, 2nd heaviest, ..., 5th heaviest

Problem 2 Total weight of the apples ≤ 2 kg.

← Basket has apples of weight

→ 450 gms, 400 gms.

Greedy $4 \times 450 = 1800$ gms.

Alternate $5 \times 400 = 2000$ gms.

Problem 3

← Subsequence Problem.

$S_1 = a c b g a b a g b c a$ sequence

$S_2 = c a b g c$ (subsequence)

Que Given two sequences S_1 & S_2 , whether S_2 is a subsequence of S_1

$|S_1| = n$ $|S_2| = p$

Approach 1 → search through all subsequences of length p and check if one of them is equal to S_2 .

no. of such subsequences = $\binom{n}{p}$

Approach 2: match letter by letter

$S_1 =$ b a c b c b a b c a c b a a
 $S_2 =$ b c b c a

Match the current letter in S_2 with its first occurrence you see in S_1 after the previous matching.

Argument for correctness

S_1 b a c b c b a b c a c b a a
 S_2 b c b c a

We want to argue that the greedy approach works. That is, if S_2 is a subsequence of S_1 then we should be able to see it by matching each subsequent letter of S_2 with its first occurrence in S_1 after the previous matching.

To show that this is indeed true, start with an arbitrary subsequence of S_1 that matches with S_2 . Now, move the matching of first letter of S_2 to its first occurrence in S_1 . You still have the valid subsequence. Repeat the argument for every letter of S_2 one by one.

Dynamic Programming.

Simple Example:

$$F_n = F_{n-1} + F_{n-2}$$

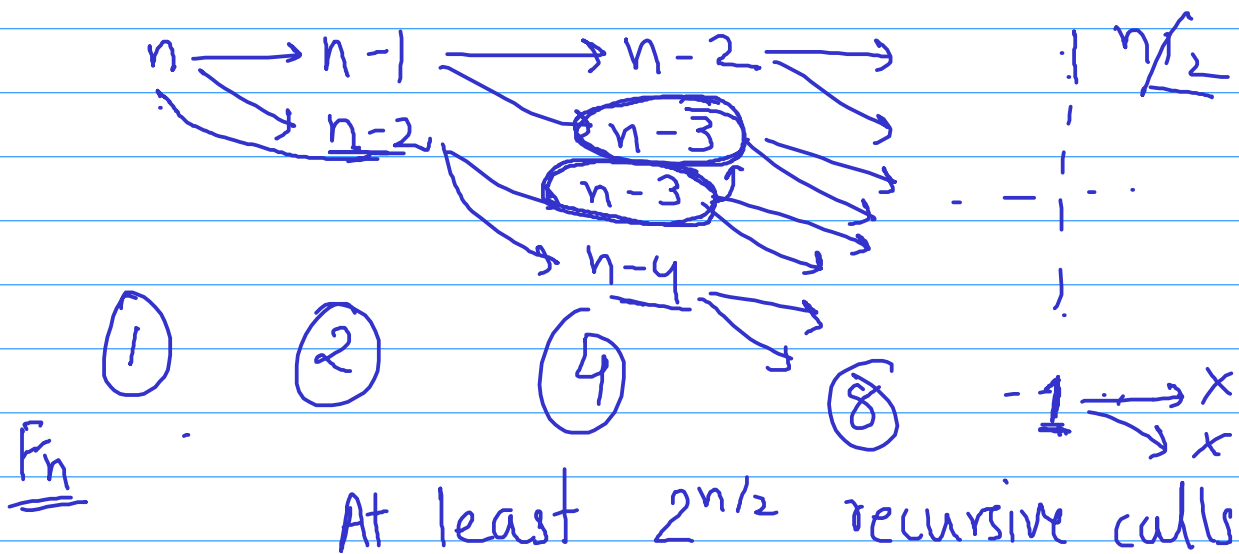
$$F_0 = F_1 = 1$$

Fibonacci (n):

if $n=0$ return 1

if $n=1$ return 1

else return Fibonacci (n-1) + Fibonacci (n-2)



Bad Implementation.

Better Implementation

Array F of length $n+1$.

$$F[0] = 1$$

$$F[1] = 1$$

Bottom-up

for $i=2$ to n

$$F[i] = F[i-1] + F[i-2]$$

end for.

Memoization.

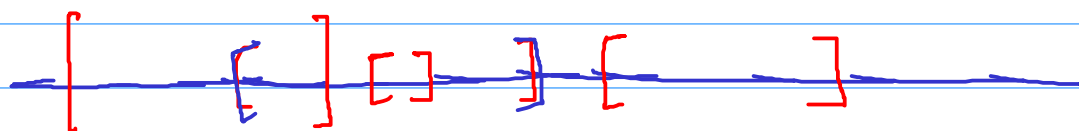
already stored in array.

In dynamic programming, you reduce your problem to one or many subproblems. And if the same subproblem instance is used multiple times then you solve it only once and afterwards, keep using its stored solution. If the total number of **distinct subproblem instances** needed during the course of your algorithm is polynomially bounded then your algorithm is efficient.

For example, in the Fibonacci problem we had n distinct instances of the subproblem $\text{Fibonacci}(n), \text{Fibonacci}(n-1), \dots, \text{Fibonacci}(1)$

Interval Scheduling. [Kleinberg Tardos]
Chapter 4

Resource / server doing computation.



n requests.

$(s_1, t_1), (s_2, t_2), \dots, (s_n, t_n)$

Maximize the number requests you can cater to.
Constraints:

- If you accept a request then you have to allocate the resource for the whole duration of desired interval (No partial allocation)
- The resource can be allocated to only one person at a time.

Greedy Algorithms

Aug 30

→ locally optimal, near-sighted, immediate benefit

→ sometimes intuitively clear, sometimes not.

→ correctness of the algorithm is most often not obvious, and thus requires a concrete argument

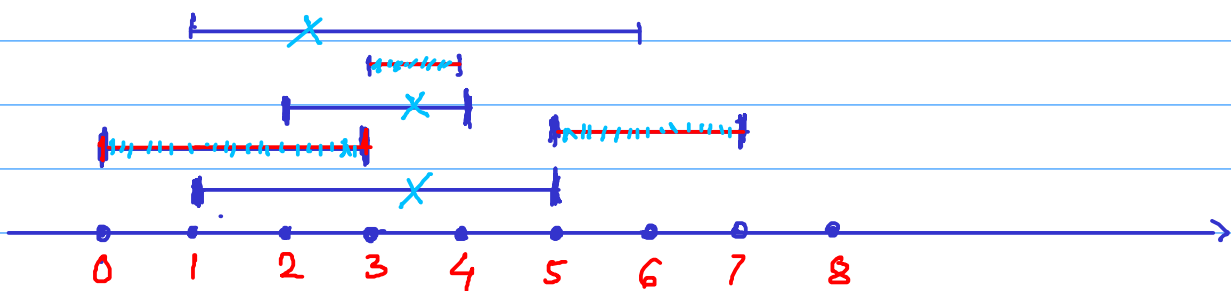
Classic Example: minimum spanning tree



Interval Scheduling: [Kleinberg Tardos Chapter 4]

Given a set of intervals (on real number line)
find the largest subset of disjoint intervals.

Example: → (1, 5), (0, 3), (2, 4), (1, 6), (3, 4), (5, 7)

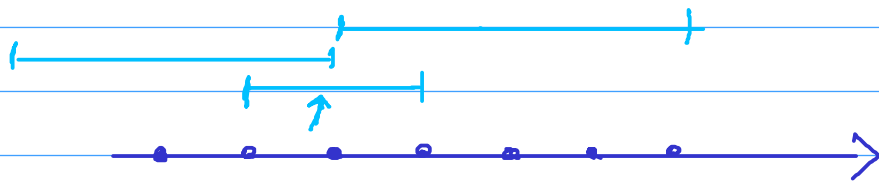


Application: resource allocation requests for fixed intervals of time

- resource can be allocated to only one request at a time
- no partial allocation
- cater to maximum number of requests.

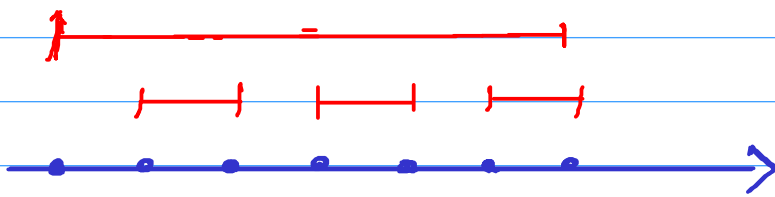
Searching through all possible solutions $\approx 2^n$ solutions

Greedy Strategy 1: keep choosing the smallest interval
 (and removing those which intersect with chosen ones)



Greedy 1
 optimal 2.

Greedy Strategy 2: earliest starting time



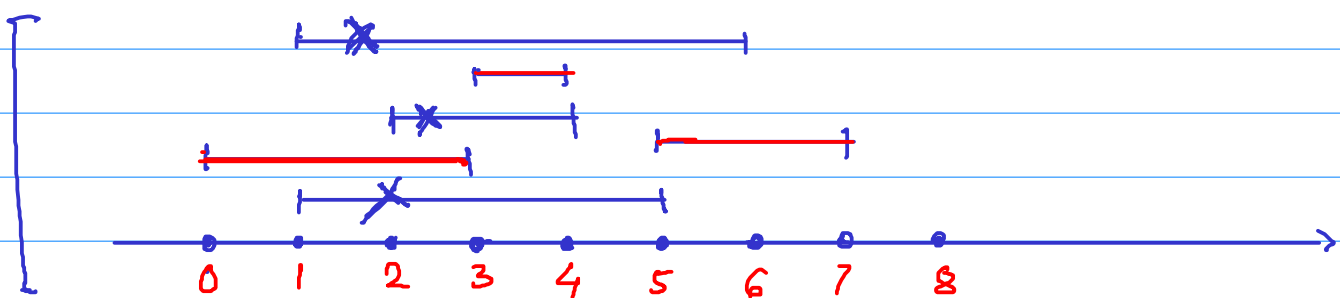
Greedy 1
 optimal 3

Greedy Strategy 3: keep choosing the interval with smallest no. of overlaps.

HW



Greedy Strategy 4: keep choosing the interval with earliest ending time

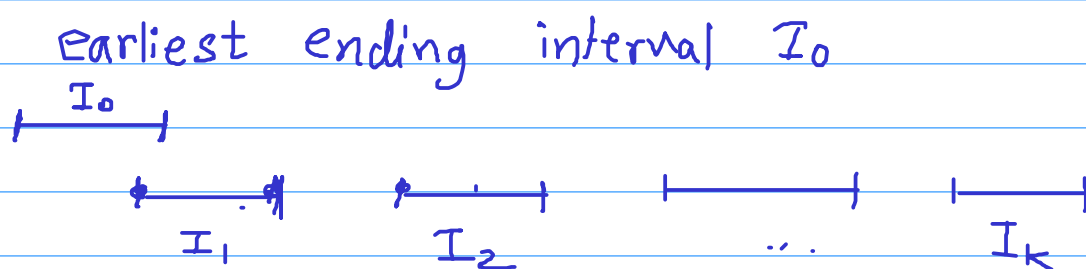


Intuitively, you are maximizing the remaining time, thus maximizing the possible number of requests catered afterwards.

Claim: Greedy strategy 4 always gives an optimal solution

General Framework of argument:

- ① show that there exists an optimal solution that agrees with the first greedy step.
 - ② The rest of the argument will work inductively.
- consider an optimal solution $I_1, I_2, I_3, \dots, I_k$.



swap I_0 with I_1 .

Consider a new solution $I_0, I_2, I_3, \dots, I_k$

Claim It is valid solution.

$$\text{endtime}(I_0) < \text{endtime}(I_1) \leq \text{starttime}(I_2)$$

$\Rightarrow I_0$ does not overlap with I_2, I_3, \dots, I_k

$I_0, I_2, I_3, \dots, I_k$ is an optimal solution agreeing with first greedy step.

Inductive proof based on the number of intervals.

Inductive Hypothesis: Greedy algo works for any instance with up to $n-1$ intervals

Inductive Step: It works for all instances with n intervals

Base case: $n=1$. Obvious.

Input: $\mathcal{I} = \{I_1, I_2, I_3, \dots, I_n\}$

for convenience

assume endtime $(I_1) \leq$ endtime $(I_2) \leq \dots$

Algorithm chooses I_1 .

then removes all intervals overlapping with I_1

$\rightarrow \mathcal{I}' = \mathcal{I} - \{ \text{intervals overlapping with } I_1 \}$

Recursively applying the same algorithm on \mathcal{I}'

By the inductive hypothesis, we know that the algorithm gives optimal solution for \mathcal{I}'

That is,

$$\boxed{\text{Algo}(\mathcal{I}') = \text{OPT}(\mathcal{I}')}_{--}$$

Output: $I_1 + \text{Algo}(\mathcal{I}') = I_1 + \text{OPT}(\mathcal{I}')$

Claim: $(I_1 + \text{any optimal solution for } \mathcal{I}')$
is an optimal solution for \mathcal{I} .

Claim $I_1 + \text{OPT}(\mathcal{I}')$ is an optimal solution for \mathcal{I} .

Proof: Recall that we showed that there exist an optimal solution for \mathcal{I} that contains I_1 .

Let it be $I_1, J_1, J_2, \dots, J_k$

Clearly J_1, J_2, \dots, J_k are disjoint from I_1 .

Hence, $\{J_1, J_2, \dots, J_k\}$ is a valid solution for \mathcal{I}' .

$$\rightarrow |\text{OPT}(\mathcal{I}')| \geq |\{J_1, J_2, \dots, J_k\}|$$

$$\Rightarrow |I_1 + \text{OPT}(\mathcal{I}')| \geq |\{I_1, J_1, J_2, \dots, J_k\}|$$

\downarrow
 optimal for \mathcal{I} . = $|\text{OPT}(\mathcal{I})|$

Pseudocode

Input: $(s_1, f_1), (s_2, f_2), \dots, (s_n, f_n)$

- Sort according to f_j .

$\rightarrow f = -\infty$ (finish time of latest interval selected so far)

```

for (i = 1 to n)
  → if ( $s_i \geq f$ ) then
      select  $(s_i, f_i)$ 
       $f \leftarrow f_i$ 
  
```

HW Assignments

deadlines d_1, d_2, \dots, d_n

time required l_1, l_2, \dots, l_n

Lateness of i -th assignment = $(t_i - d_i)$ if $t_i > d_i$
 \uparrow actual submission

Minimize maximum lateness over all assignments.

- smallest length first?
- earliest deadline first?
- minimum $d_i - l_i$ first?

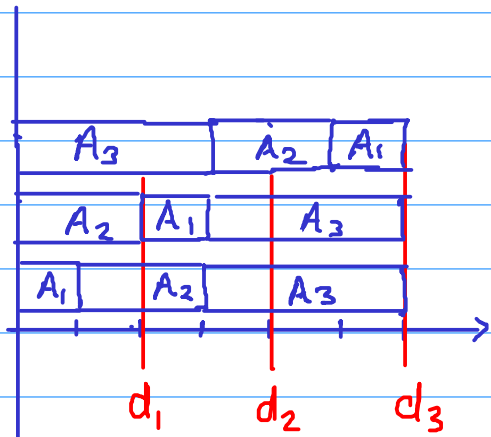
Example :

	A_1	A_2	A_3
length	1	2	3
deadlines	2	4	6

	A_1	A_2	A_3	
Lateness	4	1	0	← $A_3 A_2 A_1$
Latness	1	0	0	← $A_2 A_1 A_3$
Latness	0	0	0	← $A_1 A_2 A_3$

Red shows maximum lateness.

Best order is $A_1 A_2 A_3$.

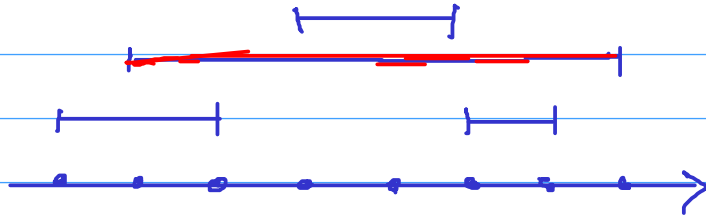


② Interval scheduling: cater to **all requests** with **minimum** number of servers.
 (one server can cater to one request at a time)

Equivalently, finding the minimum number of platforms required for a set of trains stopping at a station.

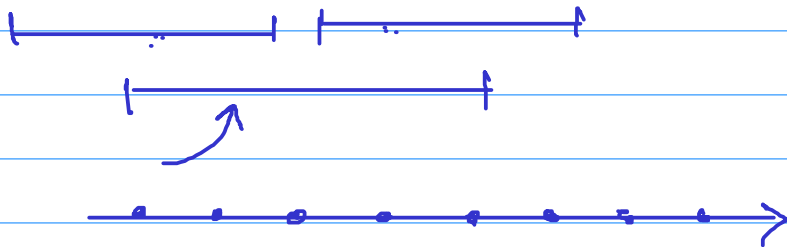
Interval Scheduling : another variant

- select a set of disjoint intervals to maximize the total length of selected intervals.

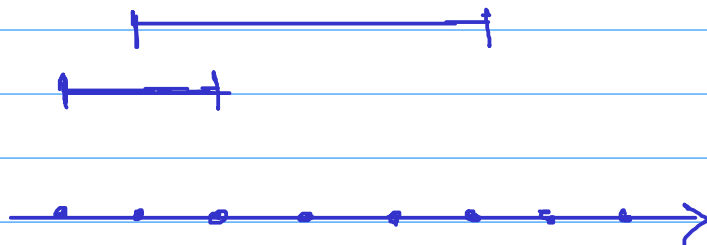


For example, there might be a profit proportional to the duration of use.

Greedy Strategy 1 keep picking the longest length intervals



Greedy Strategy 2 : earliest starting time



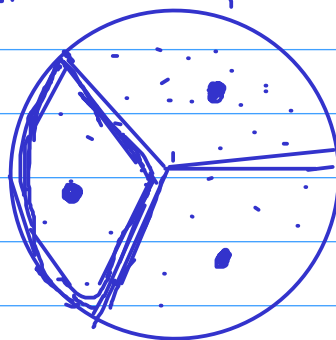
Dynamic Programming:

Recursion with better memory management.

General Idea: try to categorize the possible solutions into different types.

For each type, find the optimal solution via recursion.

Then compare the various types of optimal solutions with each other.



Interval scheduling with maximum total length.

$$\mathcal{I} = \{I_1, I_2, I_3, \dots, I_n\}$$

Two kinds of solutions

→ contain I_1

→ don't contain I_1

Can we find optimal from both the kinds recursively?

$$\mathcal{I}' = \mathcal{I} - \underline{\text{overlap}(I_1)}$$

$$\underline{\text{OPT}(\mathcal{I}')} + I_1$$

$$|\mathcal{I}'' = \mathcal{I} - I_1 \quad \underline{\text{OPT}(\mathcal{I}'')}$$

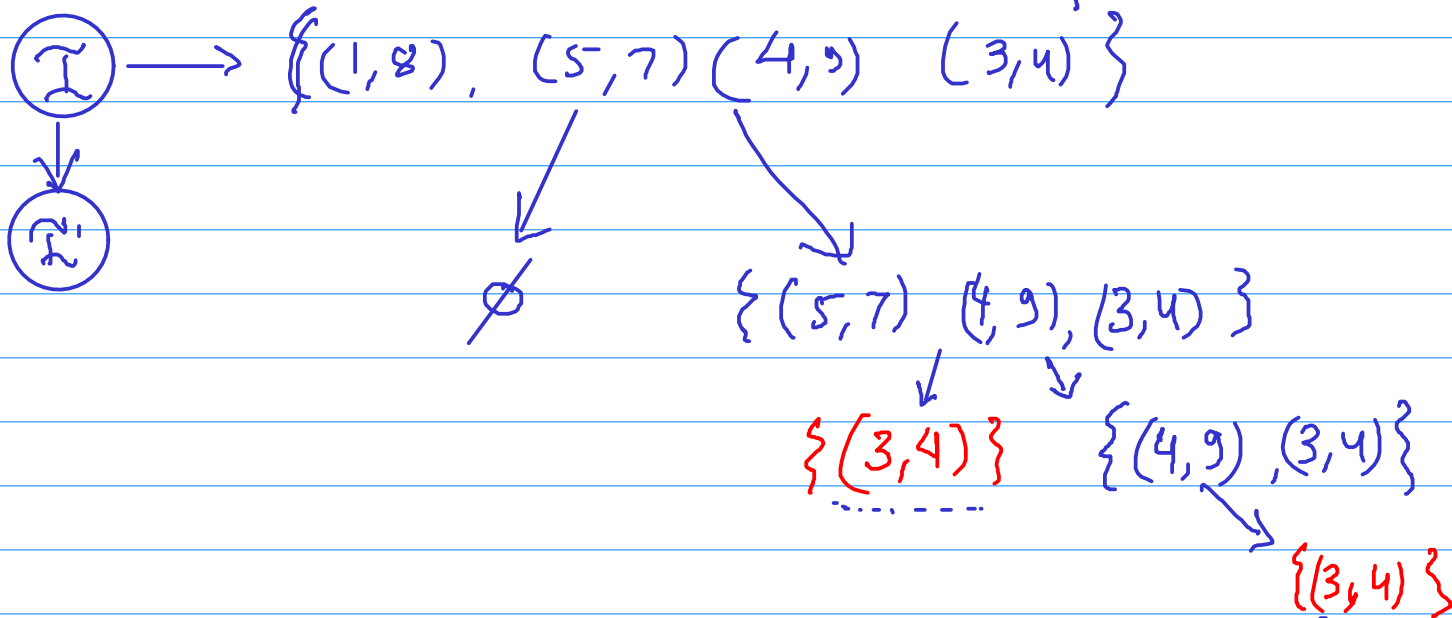
$$\text{OPT}(\mathcal{I}) = \max \left\{ I_1 + \text{OPT}(\mathcal{I}'), \text{OPT}(\mathcal{I}'') \right\}$$

Example $\mathcal{I} = \{ (0,5), (1,8), (5,7), (4,9), (3,4) \}$

$$I_1 = (0,5)$$

$$\mathcal{I}' = \{ (5,7) \}$$

$$\mathcal{I}'' = \{ (1,8), (5,7), (4,9), (3,4) \}$$



Nothing clever here. We are simply trying to go over all possible solutions recursively.

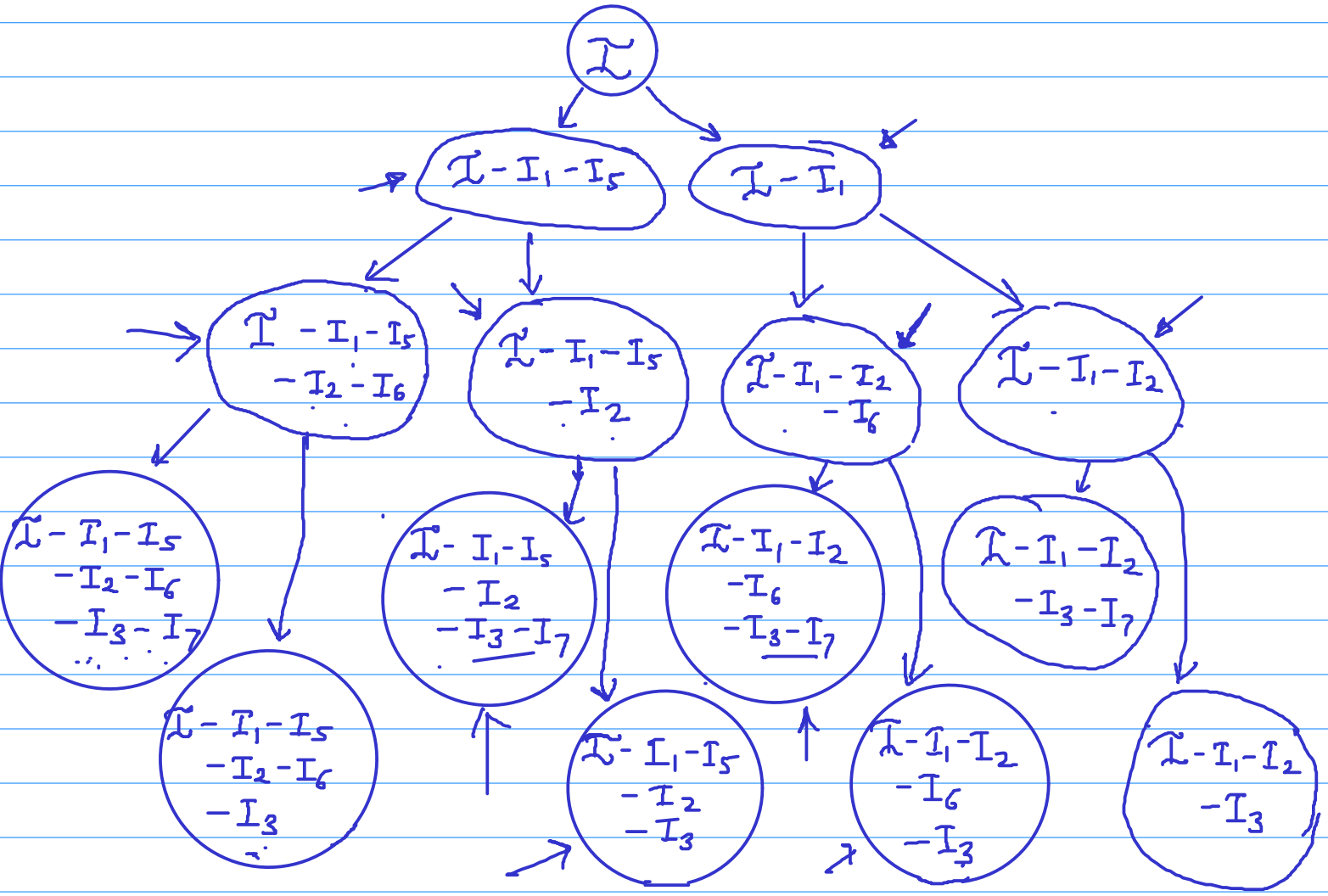
No. of recursive calls seems to be growing exponentially because each call makes two new recursive calls.

Efficient implementation possible if no. of total distinct recursive calls is small.

$$\mathcal{I} = \{ \underbrace{(1,3)}_{I_1}, \underbrace{(11,13)}_{I_2}, \underbrace{(21,23)}_{I_3}, \underbrace{(31,33)}_{I_4}, \underbrace{(2,4)}_{I_5}, \underbrace{(12,14)}_{I_6}, \underbrace{(22,24)}_{I_7}, \underbrace{(32,34)}_{I_8} \}$$

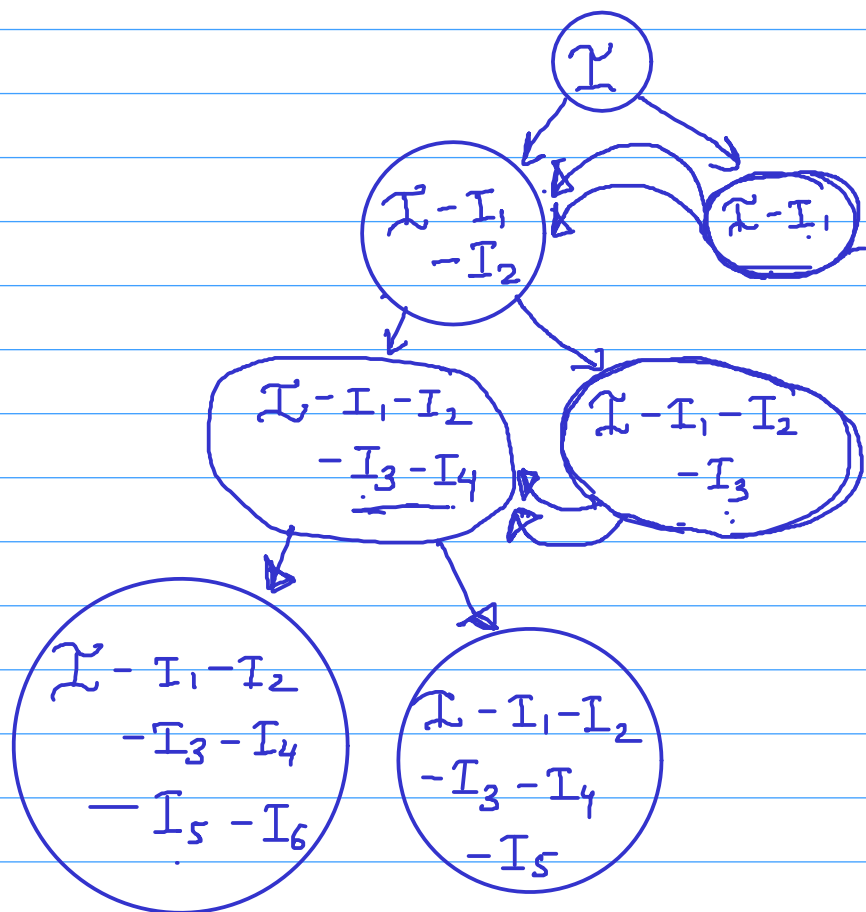
HW work out the recursion tree and figure out whether the number of recursive calls is growing exponentially or polynomially

$$\mathcal{I} = \left\{ \begin{array}{cccc} (1, 3) & (11, 13) & (21, 23) & (31, 33) \\ I_1 & I_2 & I_3 & I_4 \\ (2, 4) & (12, 14) & (22, 24) & (32, 34) \\ I_5 & I_6 & I_7 & I_8 \end{array} \right\}$$



If we have $2n$ intervals, we will have at least 2^n distinct recursive calls or subproblems

$$\mathcal{I} = \left\{ \begin{array}{cccccc} (1,3), & (2,4), & (11,13), & (12,14), & (21,23), & (22,24), \\ I_1 & I_2 & I_3 & I_4 & I_5 & I_6 \\ \\ (31,33), & (32,34) & & & & \\ I_7 & I_8 & & & & \end{array} \right\}$$



Only $2n$ distinct recursive calls or subproblems.

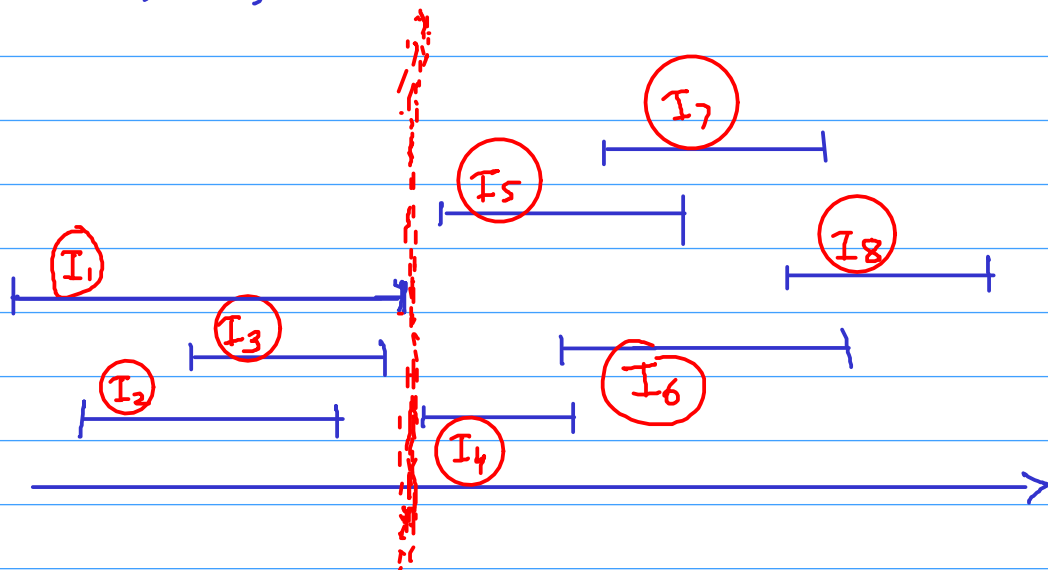
It seems if the intervals are arranged in a particular order, the number of distinct subproblems will be small.

Possible orders:

- ① starting time
- ② ending time

sort the intervals in increasing order of starting time

$I_1, I_2, I_3, \dots, I_n$



Claim: The input set of intervals for any recursive call will look like

$$\rightarrow \{I_j, I_{j+1}, I_{j+2}, \dots, I_n\}$$

(as opposed to an arbitrary subset of intervals)

This is because when we remove intervals overlapping with I_k , the remaining set is simply

all intervals with starting time $>$ end-time(I_k).

No. of distinct subproblems $\leq n$.

Recursive Implementation

$Opt \leftarrow$ array with all zeros.

$Opt[n] \leftarrow \text{length}(I_n)$

Output $ALG(i)$;

$ALG(j)$: // $ALG(j)$ computes optimal solution for $\rightarrow \{I_j, I_{j+1}, \dots, I_n\}$

if $Opt[j] > 0$ return $Opt[j]$

// means already solved and stored.

else

$Opt[j] \leftarrow \max \left\{ \begin{array}{l} ALG(j+1) \leftarrow \\ \text{length}(I_j) + ALG(P(j)) \end{array} \right.$

return $Opt[j]$

least index k
such that

$\text{start}(I_k) > \text{end}(I_j)$

Iterative Implementation:

$Opt \leftarrow$ array with all zeros.

$Opt[n] \leftarrow \text{length}(I_n)$

for ($j = n-1$ to 1)

$Opt[j] \leftarrow \max \left\{ \begin{array}{l} Opt[j+1] \\ \text{length}(I_j) + Opt[P(j)] \end{array} \right.$

HW

Add code to compute the optimal set of intervals.

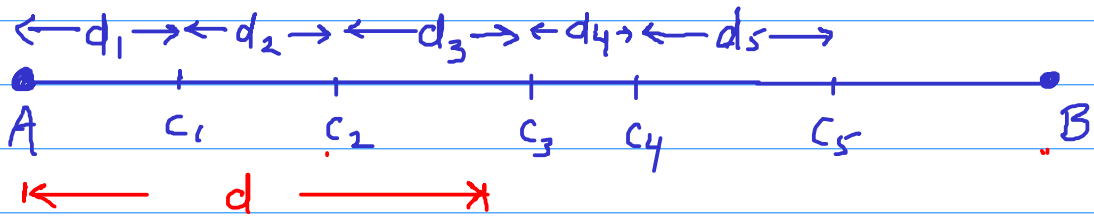
Conclusion: Order of processing the input is important.

Intervals: Order of starting time / ending time

Sequences: left to right

Try to ensure that no. of distinct subproblems is small.

HW 1



Maximum travel in a day - d.

Night stay prices - P_1, P_2, \dots

Minimize total stay cost during the journey.

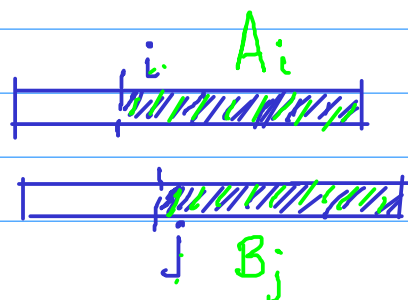
$O(n)$

no. of distinct subproblems = n .

HW 2

$\rightarrow A = \overset{2}{b} \overset{5}{a} \overset{3}{c} \overset{1}{b} \overset{2}{c} \overset{5}{b} \overset{3}{a} \overset{1}{b} \overset{3}{a} \overset{1}{c} \overset{2}{b} \overset{4}{a} \overset{1}{b}$
 $\rightarrow B = \overset{2}{b} \overset{5}{c} \overset{3}{b} \overset{1}{c} \overset{2}{a}$

Match B inside A with minimum cost



no. of distinct subproblems $O(mn)$

Subset Sum problem

Given set of integers $A = \{a_1, a_2, a_3, \dots, a_n\}$

is there a subset with sum zero?

Example:

↓ ↓ ↓ ↓ ↓
2, -5, 1, 7, -4, -6

{ 2, 7, -5, -4 }

Solutions $\begin{cases} \rightarrow \text{containing } a_n \\ \rightarrow \text{not containing } a_n \end{cases}$

is there a subset of { 2, -5, 1, 7, -4 } with sum zero?

is there a subset of { 2, -5, 1, 7, -4 } with sum six?

Subproblem $A_j = \{a_1, \dots, a_j\}$, number N

is there a subset of A_j with sum N .

No. of distinct subproblems = $n \times \sum_i |a_i|$

ALG(j, N):

\rightarrow ALG(j-1, N) OR ALG(j-1, N-a_j)

Pseudo polynomial time.

A polynomial time algorithm is supposed to take time $\text{poly}(n, \text{no. of bits in } a_1, a_2, \dots)$.

Subset Sum Pseudocode

$S \leftarrow$ Boolean two-dim array of size $n \times (\sum |a_i| + 1)$

$S[j, k]$ will denotes whether there is a subset of first j numbers with sum = k .

say, range of j $1 \leq j \leq n$.

and range of k $neg \leq k \leq pos$

↓
sum of all
negative numbers

↘
sum of all
positive numbers

Initialization:

$$S[1, k] \leftarrow \begin{cases} \text{True} & \text{for } k = a_1 \\ \text{False} & \text{for all other valuse of } k. \end{cases}$$

for ($j = 2$ to n)

for ($neg \leq k \leq pos$)

$$S[j, k] \leftarrow S[j-1, k] \text{ OR } S[j-1, k-a_j]$$

To construct the desired set

Assume $S[n, 0]$ is true.

sum $\leftarrow 0$

for ($j = n$ to 1)

if $S[j-1, \text{sum}] = \text{true}$, don't take a_j

elseif $S[j-1, \text{sum}-a_j] = \text{true}$, take a_j and
sum $\leftarrow \text{sum} - a_j$

l_1, l_2, \dots, l_k

$$\text{Variance} = \underbrace{\text{avg}(l_1^2, l_2^2, \dots, l_k^2)} - \underbrace{(\text{avg}(l_1, l_2, \dots, l_k))^2}_{\substack{\uparrow \\ \text{fixed.}}}$$

Arrange the words in lines to minimize the **sum of squares of the slacks** of all lines.

Puzzle find integers a, b, c such that

$$a + b + c = 14$$

and $a^2 + b^2 + c^2$ is minimized.

Answer:

What do you observe?

Balanced Margins

Idea 1: fit as many words as you can

→ can be very unbalanced

Greedy Idea:

compute the average slack per line. Go line by line and try to keep the slack for each line as close as possible to the average slack.

doesn't give an optimal solution.

is

$\overbrace{5}$	5		2	<u>is</u>	$\frac{11}{3}$
$\underbrace{\quad}$	1		4		
\rightarrow	5		5	51, <u>45</u>	<u><u>$= 3.66$</u></u>

Dynamic Programming.

- Try to categorize the set of all possible solutions

Each solution is a partition of words into k lines.

$$n = n_1 + n_2 + n_3 + \dots + n_k$$

Categories:

$n_k = 1$	last line has 1 word.
$n_k = 2$	last line has 2 words.
\vdots	
$n_k = n-1$	last line has $n-1$ words.

Assuming last line has words w_{p+1}, \dots, w_n , (i.e. $n-p$ words) can you compute the optimal solution via a subproblem?

$$\text{OPT}(w_1, w_2, \dots, w_p) + \text{slack}(w_{p+1}, w_{p+2}, \dots, w_n)^2$$

$$\text{OPT}(n) = \min \left\{ \begin{array}{l} \text{OPT}(n-1) + [W - w_n]^2 \\ \text{OPT}(n-2) + [W - w_n - w_{n-1} - 1]^2 \\ \text{OPT}(n-3) + [W - w_n - w_{n-1} - w_{n-2} - 2]^2 \\ \vdots \\ \cdot \\ \cdot \end{array} \right.$$

Running time $O(n^2)$

Pseudocode for computing the optimal value and the optimal solution.

$S \leftarrow$ array of length $n+1$

// $S[j]$ denotes the minimum sum of squares of slacks for first j words, for $1 \leq j \leq n$

$N \leftarrow$ array of length $n+1$

// $N[j]$ denotes the index of the first word in the last line in the optimal arrangement of first j words.

$s[0] \leftarrow 0$; $s[j] \leftarrow \infty$ for $j > 0$; $N[0] \leftarrow 0$

for ($j = 1$ to n)

for ($r = j$ to 1)

slack $\leftarrow W - (w_r + w_{r+1} + \dots + w_j + j - r)$

if (slack ≥ 0 and $s[j] > s[r-1] + (\text{slack})^2$)

$s[j] \leftarrow s[r-1] + (\text{slack})^2$

$N[j] \leftarrow r$

Optimal Arrangement:
Run Arrange(n).

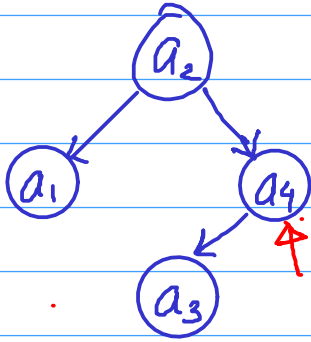
```
Arrange(j) {  
  Arrange(N(j)-1)  
  print words from N(j) to j  
  in a new line.  
}
```

Optimal Binary Search Tree

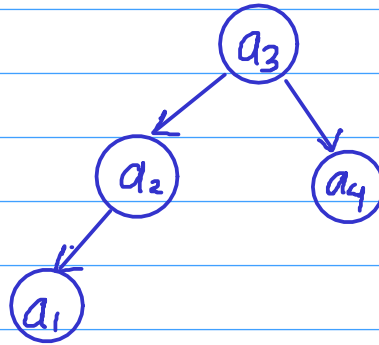
Sep 7

$$a_1 < a_2 < a_3 < a_4$$

$$\begin{array}{r} 7 \times 2 \\ 5 \times 1 \\ 4 \times 3 \\ 6 \times 2 \\ \hline 43 \end{array}$$



$$\begin{array}{r} 7 \times 3 \\ 5 \times 2 \\ 4 \times 1 \\ 6 \times 2 \end{array}$$

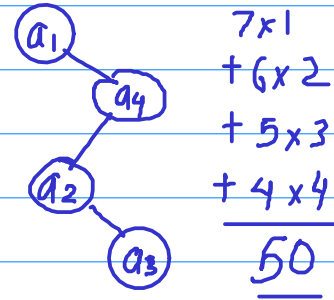


suppose we know how frequent are the search queries for each of the elements.

Say

frequencies

a_1	\rightarrow	7
a_2	\rightarrow	5
a_3	\rightarrow	4
a_4	\rightarrow	6



Total search cost
43

First binary tree

second binary tree

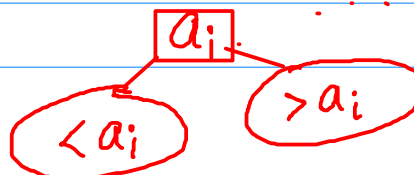
47

Compute the binary search tree with **minimum** total search cost $\sum_i f_i h_i$

where $f_i \rightarrow$ frequency (a_i) $h_i \rightarrow$ depth (a_i)

Greedy Idea 1:

maximum frequency element \rightarrow root

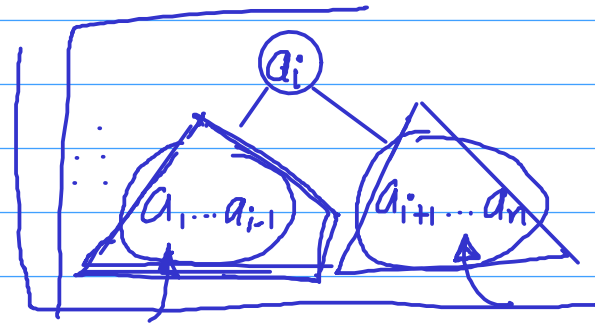


Dynamic Programming

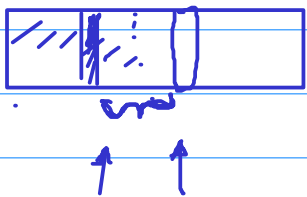
$$a_1 < a_2 < a_3 \dots < a_n$$

Possible solutions : every valid binary search tree.

Categories :
 Root $\leftarrow a_1$
 Root $\leftarrow a_2$
 ...
 Root $\leftarrow a_n$

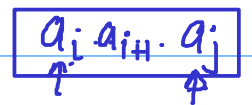
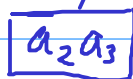
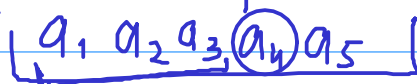


$$\text{Tree}(a_1, \dots, a_n) = \min_{1 \leq i \leq n}$$



$$\text{Tree}(a_1, \dots, a_{i-1}) + \text{Tree}(a_{i+1}, \dots, a_n) + \sum_{j=1}^n f_j \cdot 1$$

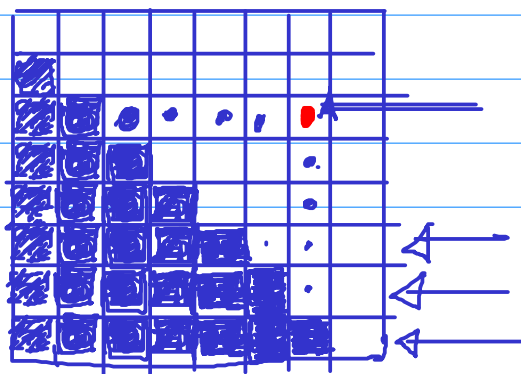
How many distinct subproblems = $\binom{n}{2}$



$$\text{Tree}(i, j) = \min_{\substack{k \\ i \leq k \leq j}} \{ \text{Tree}(i, k-1) + \text{Tree}(k+1, j) \} + \sum_{k=i}^j f_k$$

$j > i$

$i \downarrow$



OPT \leftarrow 2 dim array.

$O(n^3)$

OPT $[i, i] \leftarrow f_i$ for each i

for ($i = n-1$ to 1)

for ($j = i+1$ to n)

$$\text{OPT}[i, j] = \sum_{k=i}^j f_k + \min_{\substack{k: \\ i \leq k \leq j}} \{ \text{OPT}[i, k-1] + \text{OPT}[k+1, j] \}$$

↑
Optimal cost
for $(a_i a_{i+1} \dots a_j)$

Add/modify code to compute the optimal solution.

Sequence Alignment

Computational biology, Spell checking

Similarity between two strings

Needleman and Wunsch defined a notion of similarity

Example

$UGCTGACU$
 $\rightarrow GAATGCA$

Given Gap Penalty δ
 Mismatch cost (for each pair) α_{xy}

$UGC-TG-ACU$
 $-GAATGCA--$
 $5\delta + \alpha_{CA}$

Total cost = sum of the gap and mismatch costs.

Find an alignment of the two strings with **minimum total cost**.

Input: x_1, x_2, \dots, x_m and y_1, y_2, \dots, y_n
 $\delta, \{\alpha_{x_i y_j}\}$

Categories of solutions.

① $\boxed{\dots\dots\dots} x_m$
 $\boxed{\dots\dots\dots} y_n$

② $\boxed{\dots\dots\dots x_{m-1}} x_m$
 $\boxed{\dots\dots\dots y_n}$

③ $\boxed{\dots\dots\dots x_m}$
 $\boxed{\dots\dots\dots y_{n-1}} y_n$

$$OPT(m, n) = \min \begin{cases} \alpha_{x_m y_n} + OPT(m-1, n-1) \\ \delta + OPT(m-1, n) \\ \delta + OPT(m, n-1) \end{cases}$$

$x_1 \dots x_{m-1}$
 \uparrow
 $y_1 \dots y_{n-1}$

$OPT(i, j)$ no. of distinct subproblems $m \times n$

Implementation

$A \leftarrow$ 2D array $m \times n$

// $A[i, j]$ denotes the minimum cost of alignment for $x_1 x_2 \dots x_i$ and $y_1 y_2 \dots y_j$

$A[i, 0] \leftarrow \delta_i$ for each i

$A[0, j] \leftarrow \delta_j$ for each j

for ($i = 1$ to n)

for ($j = 1$ to n)

$$A[i, j] = \min \left\{ \begin{aligned} &\alpha_{x_i y_j} + A[i-1, j-1], \\ &\delta + A[i-1, j], \\ &\delta + A[i, j-1] \end{aligned} \right.$$

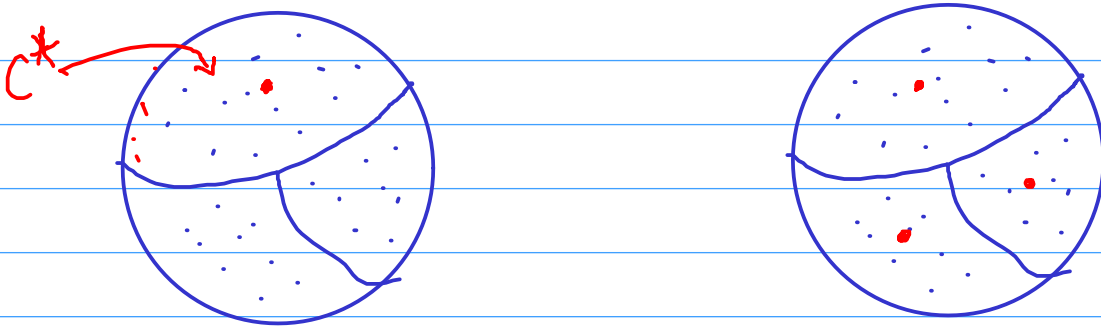
HW space $O(mn)$

Can you get space $O(m+n)$
 and time $O(mn)$

Divide and Conquer

Summarizing Greedy and Dynamic Programming

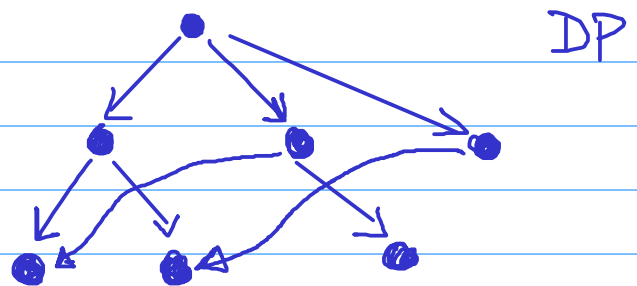
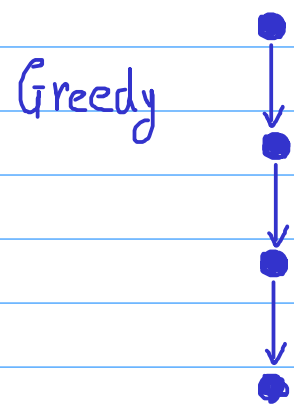
- Dividing the set of possible solutions into multiple categories.



Greedy: there must be an optimal solution in a certain category C^* (greedy choice)

DP: Will take best of the optimal solutions from each category.

To compute the optimal solution from a chosen category \rightarrow smaller subproblem (Recursion)



No. of distinct subproblems should be small.

Data Compression: Coding

assign a fixed length 0-1 string to each character.

a → 00001
b → 00010
⋮

5 bit encoding can work for up to 32 characters.

Is a smaller length encoding possible?

With fixed length - not possible

Variable length encoding

- can be more efficient when no. of characters is not 2^k .
- can use smaller length codes for more frequent characters.

Example: Morse Code (dots and dashes and spaces)

e • z
t - q
a •-
⋮
• -
⋮

Problem

• - aa
 eta
 aet
 etet

Solution: Gap after every character

Prefix Code

Def: For any two different characters x and y
 $c(x)$ should not be a prefix of $c(y)$.

Ex $x \rightarrow 00$
 $y \rightarrow 001$ } Not a prefix code

Ex $a \rightarrow 0$
 $b \rightarrow 10$
 $c \rightarrow 11$ } Prefix code.

abaca \rightarrow $\underbrace{0100110}$

Claim: For a prefix code, any 0-1 string is unambiguously decodable.

Just scan left to right, as soon as the current substring matches one of the codewords, output the corresponding character.

Optimal Prefix Codes

	Freq		Code 1	Code 2
	0.4	A	00	0
	0.4	T	01	10
\rightarrow	0.1	P	10	110 \leftarrow
\rightarrow	0.1	G	11	111 \leftarrow

	<u>2</u>	<u>0.4×1</u>	Avg Bit length $\frac{180}{180}$ bits.
		$+ 0.4 \times 2$	
		$+ 0.1 \times 3 + 0.1 \times 3$	
		$= 1.8$	
100 char	200 bits		

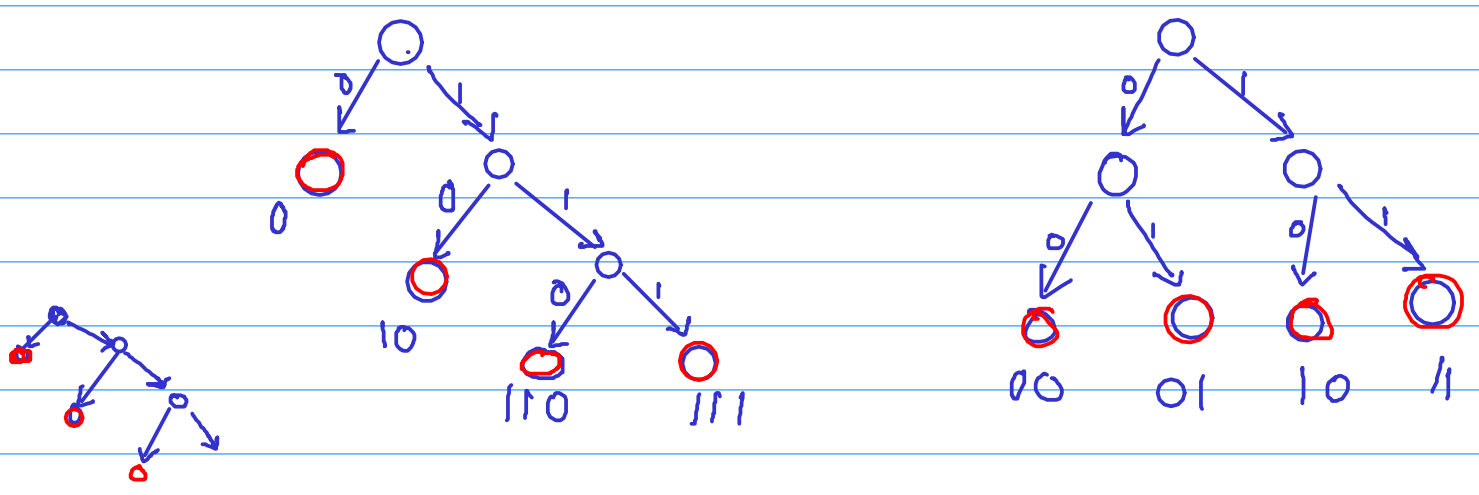
Prob Given frequencies f_1, f_2, \dots, f_n for n characters find a prefix code that minimizes $\sum f_i l_i$.

$(\sum f_i = 1)$

avg encoding length for i -th character.
 length of the encoding

Observation: Each prefix code corresponds to a binary tree.

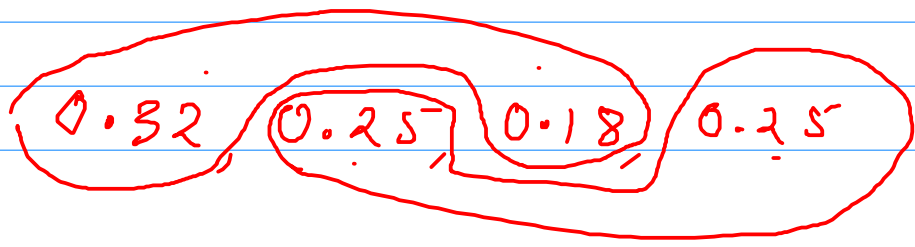
codewords correspond to leaves of the tree.



Approach 1: assign 0 to highest frequency.

Approach 2: assign 0 to highest freq if it is above some threshold.

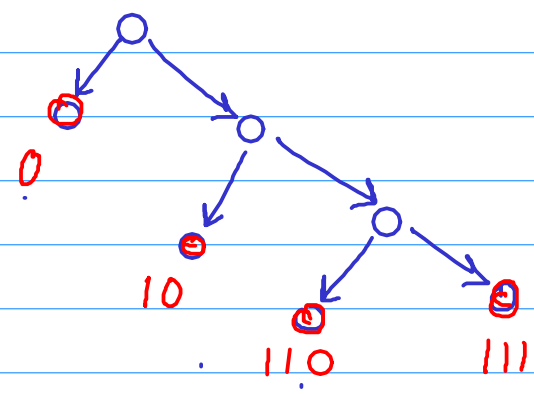
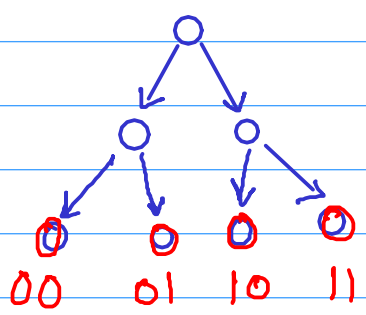
Approach 3: Do a balanced division of frequencies into two parts.



Observations:

low frequency \rightarrow higher length
 high frequency \rightarrow lower length

Approach 1: for highest frequency character, assign '0' i.e. length one codeword.



$(0.25, 0.25, 0.25, 0.25)$

$(0.28, 0.24, 0.24, 0.24)$

2

$$0.28 \times 1 + 0.24 \times 2 + 0.24 \times 3 + 0.24 \times 3 = 2.2$$

$(0.37, 0.21, 0.21, 0.21)$

2

2

= 2.05

$(0.35, 0.35, 0.15, 0.15)$

1

2

= 1.95

Assign '0' if frequency higher than certain threshold.

Doesn't work.

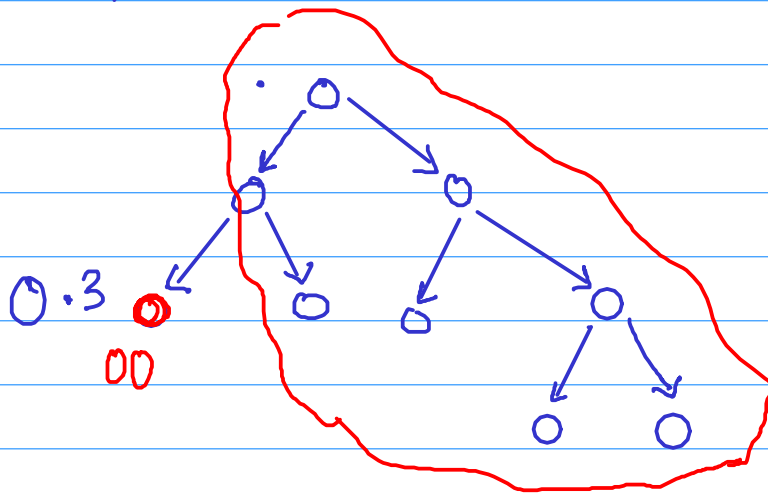
Cannot decide l , just by looking at f_1 .

Suppose there is some way to decide the encoding length for highest frequency character.

say length 2. '00'

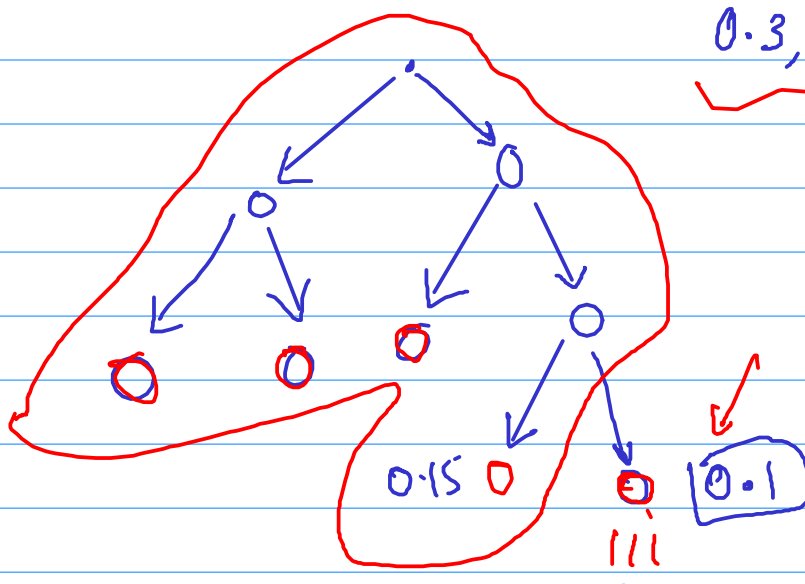
0.3, 0.25, 0.2, 0.15, 0.1

can we reduce the rest of the encodings to a subproblem?



Not able to frame it as a smaller instance of the same problem.

→ Suppose we can fix the length for the least frequent character.



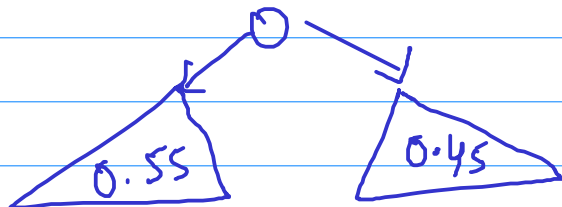
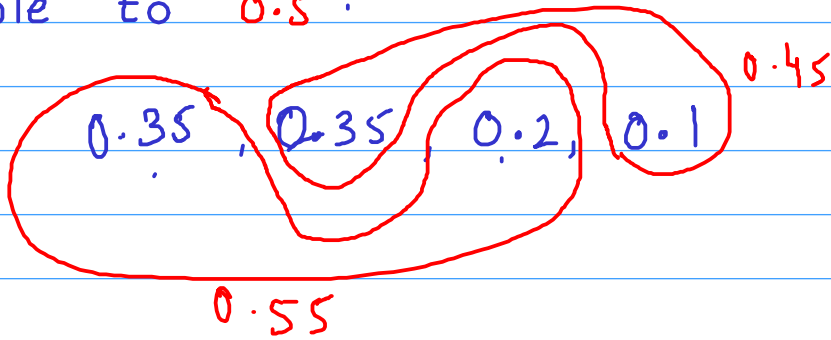
0.3, 0.25, 0.2, 0.15, 0.1

same issue.

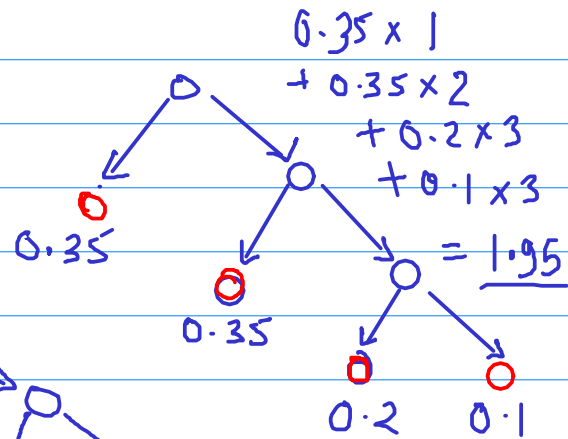
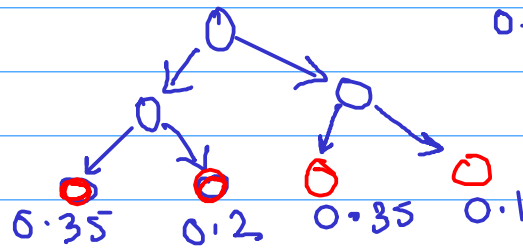
Shannon and Fano (1940's)

Balanced Partition

Divide the list into two parts such that total frequency of each part is as close as possible to 0.5.

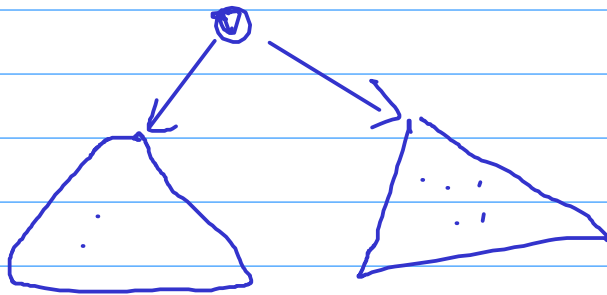


2



By balanced partition approach we are getting average encoding length = 2. But there is a better solution with 1.95.

→ Try various possibilities for the partition?

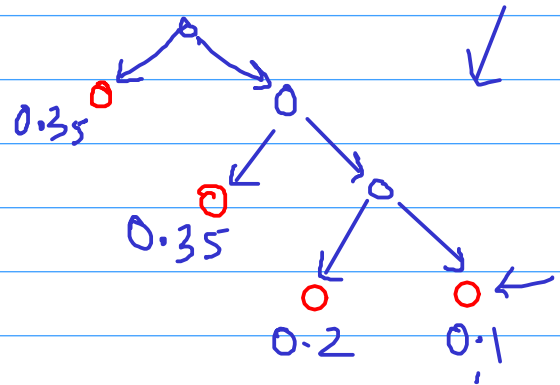
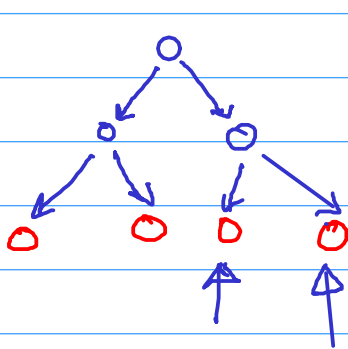


Exponentially many possibilities.

Huffman [1952]

Obs 1: If you fix a binary tree, then there is a natural way to map characters to leaves

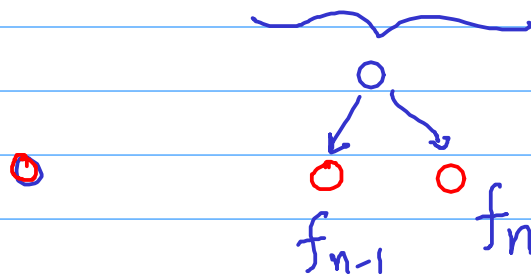
0.35, 0.35, 0.2, 0.1



Lowest frequency character has the largest depth.

Obs 2: The leaf with the largest depth must have a sibling which is a leaf.

$$f_1 \geq f_2 \geq \dots \geq f_{n-1} \geq f_n$$

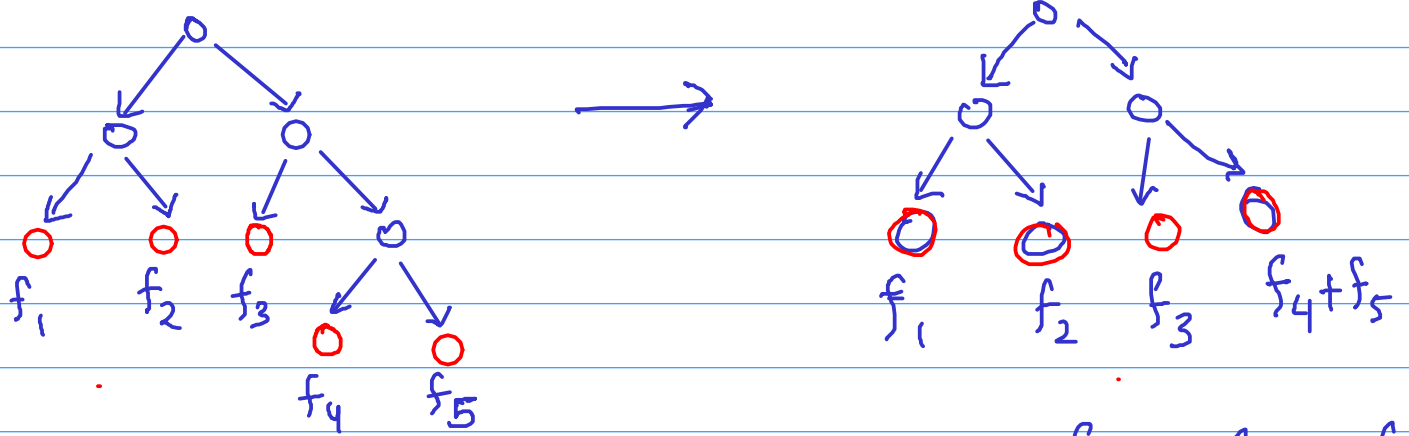


Claim: \downarrow

Lowest and second lowest frequency characters can be mapped to two largest depth siblings.

Now, the rest of the tree can be found recursively.

f_1, f_2, f_3, f_4, f_5 (in decreasing order)

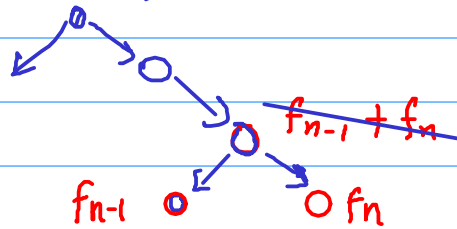


$$\text{cost}_0: 2f_1 + 2f_2 + 2f_3 + 3f_4 + 3f_5 \quad \text{cost}_1: 2f_1 + 2f_2 + 2f_3 + 2f_4 + 2f_5$$

$$\boxed{\text{cost}_0 = \text{cost}_1 + f_4 + f_5}$$

→ for any binary tree with n leaves where f_{n-1} and f_n are siblings, there is a corresponding binary tree with $n-1$ leaves with labelings

$f_1, f_2, f_3, \dots, f_{n-2}, f_{n-1} + f_n$



Algorithm:

Input: $f_1 \geq f_2 \geq \dots \geq f_{n-1} \geq f_n$ (n char)

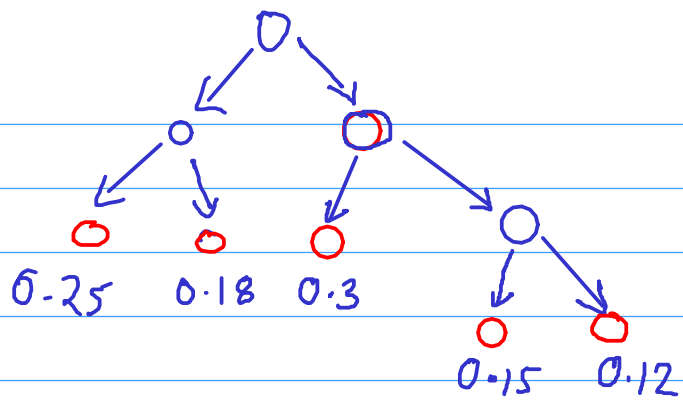
Recursively compute optimal binary tree for

→ $f_1, f_2, \dots, f_{n-2}, f_{n-1} + f_n$ ($n-1$ char)

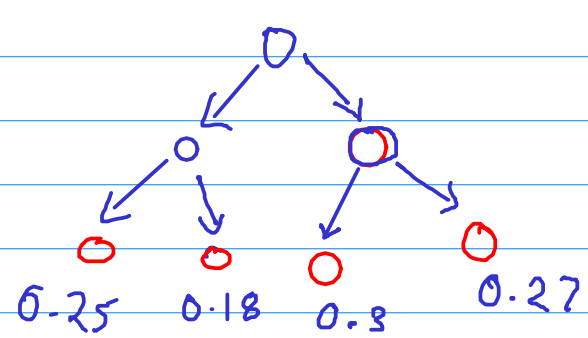
For the leaf labeled $f_{n-1} + f_n$, add two children with label f_{n-1}, f_n .

Example

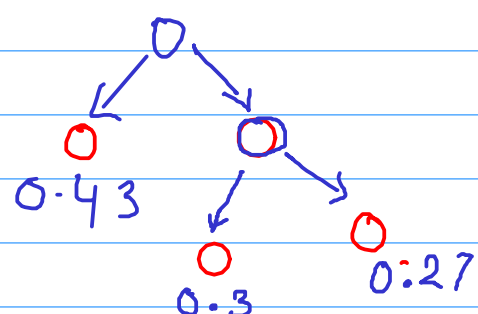
0.3, 0.25, 0.18, 0.15, 0.12



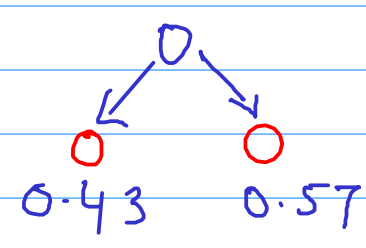
0.3, 0.25, 0.18, 0.27



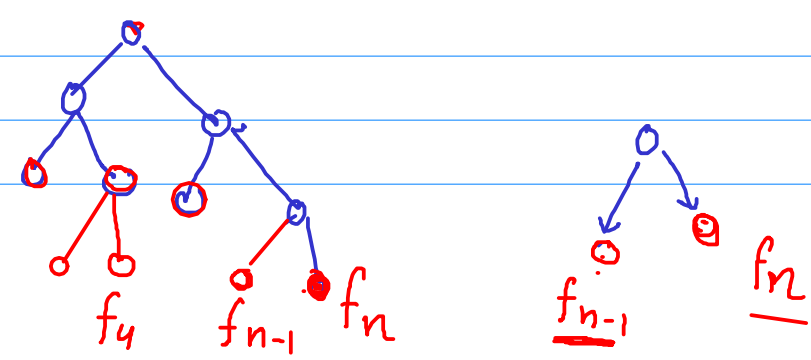
0.3, 0.27, 0.43



0.43, 0.57



$\sum f_i l_i$



Proof of Correctness

① There is an optimal solution where f_{n-1} and f_n are siblings.

② There is a one-to-one correspondence between

$A = \{ \text{full binary trees with } n \text{ leaves labeled } f_1, f_2, \dots, f_n \}$
 where f_n and f_{n-1} are siblings

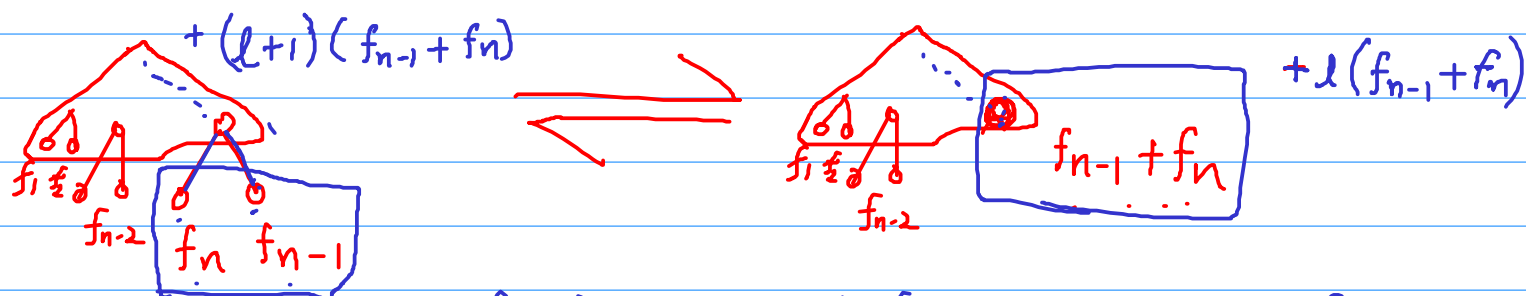
$B = \{ \text{full binary trees with } n-1 \text{ leaves labeled } f_1, f_2, f_3, \dots, f_{n-2}, \underline{f_{n-1} + f_n} \}$

Say, $A = \{ T_1, T_2, T_3, \dots, T_N \}$
 $B = \{ R_1, R_2, R_3, \dots, R_N \}$

Remove leaves labeled f_{n-1}, f_n .
 label their parent $f_{n-1} + f_n$

$T_j \longleftrightarrow R_j$

For the leaf labeled $f_{n-1} + f_n$,
 add two children labeled f_{n-1}, f_n



Moreover $\text{cost}(T_j) = \text{cost}(R_j) + f_{n-1} + f_n$

Thus, R_j is optimal in $B \iff T_j$ is optimal in A .

Huffman Codes:

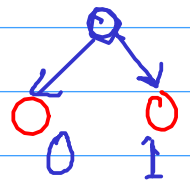
Text over some large alphabet size \longrightarrow space efficient bit representation

- Que: Can we apply this technique when the data is already in 0/1 bit representation?

Suppose there is data where 0 is much more frequent than 1.

b a g ... a a

001000110101000011100110000010



$\rightarrow a$	}	000 ... 2/8	\longrightarrow	00
b		001 3/16	\longrightarrow	100
c	
h		111 1/16	\longrightarrow	1101

- There are many other data compression techniques

- adaptive
- Algebraic

