

Class outline:

1. Examples of real world applications (as motivation)
2. Simpler example to gain intuition (since not always intuitive the first time)-
3. Back to applications and other examples (next couple classes)

Applications: very useful approach today - some examples - non exhaustive

- Computational Biology ✓
- pattern / speech recognition
- finding shortest path (later)
- cypher to Thomas Jefferson (will mention)

(a) Computational Biology: genome similarity:

- strings from alphabet $\{A, C, G, T\}$

Example: ACGGAT
 CCGCTT

language of DNA - 4 bases

human genome: ~6 Billion bases
gene - different lengths - could be 100,000 bases

- we'll (ask) How "similar" = ?

we'll see why we care and how we do it

- Why? If we understand function of one string (mouse genome), can teach us about function of other string (human genome) which know less about.

- How? (ask)

we'll look for Longest Common Subsequence (must appear same order, not necessarily consecutive)

- number changes from one to other small; allowed to insert blank, change characters
- longest common subsequence

Example: ACGGAT
CCGCTT

(English: CAT
CHAT)

Question: longest common subsequence length? (ask)

Answer: 3 CGT.

Efficiency: How many possible subsequences? (ask)
small example we can see below but

Let's say string of length $n=500$...

A C C C G G T C G A G T G ...
G T C G T T C G G A A T T ...

Brute force: Consider all possible subsequences of first string and check if also subsequence of second string.

$n=500 : 2^{500}$ Not feasible!

To see: pick first character or do not ...
pick 2nd character or do not ...
pick 3rd character or do not ...

$2 \cdot 2 \cdot 2 \dots 2 = 2^n$
n times

actually needs to multiply by length of 2nd string, since check each subsequence in second

Therefore: need efficient algorithm?

Intuition: We'll save solutions to smaller subproblems; we'll come back to this ...

→ will divide and conquer work?
→ we'll come back to Dynamic Programming solution

No, not in a simple way: could be for this example

ACG GAT
CCG CTT

But generally, not straightforward, eg this won't work:
CGT GAC → Finds CG, not CGT
CCG TTT CGT

Simple example (to build intuition):

-5-

Fibonacci!

(Everyone has heard of? I'm going to teach it a little differently)

First write naive recursive algorithm:

$Fib(n)$

1. if $n == 0$ return \bullet (1 or 0 depending on source)
2. if $n == 1$ return 1
3. else return $Fib(n-1) + Fib(n-2)$

Good algorithm? (ask) Does the job but ...

No! Very wasteful

Ways to see:

1. To compute $Fib(25)$ you compute $Fib(24)$ and $Fib(23)$. To compute $Fib(24)$ you compute $Fib(23)$ again!, and $Fib(22)$. That's a lot of re-computing and a lot of work!
we'll look at a recursion tree in a moment.

$$Fib(25) = Fib(24) + Fib(23)$$

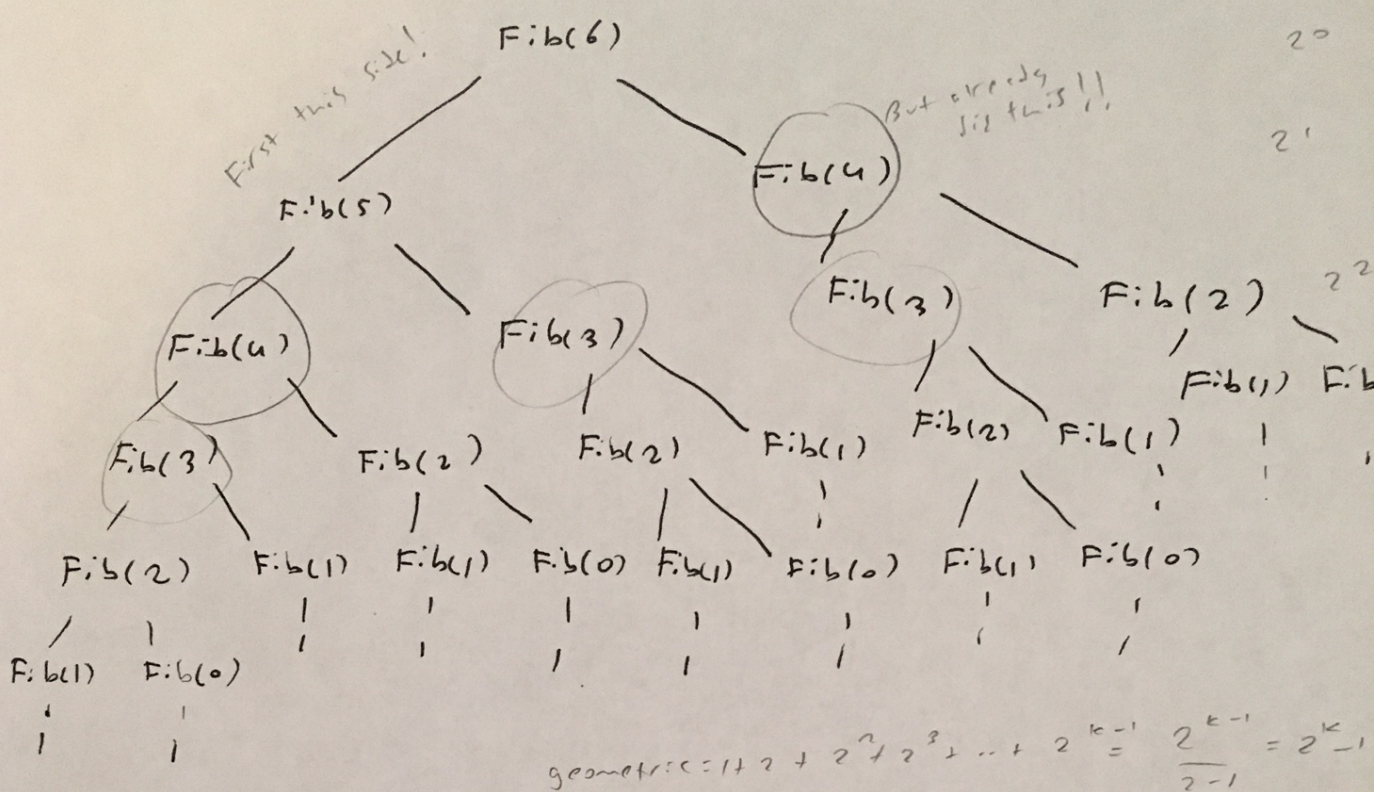
$$Fib(24) = Fib(23) + Fib(22)$$

2. Remember from divide & conquer how to write out recursion:

$$- T(n) = T(n-1) + T(n-2) + \text{const} \quad \text{to combine, adding}$$
$$\geq 2T(n-2)$$

- each time subtract 2 and multiply by 2 ($\frac{n}{2}$ times)
- $$T(n) \geq 2^{\frac{n}{2}} \quad (\text{run time at least } 2^{\frac{n}{2}})$$

3. To get more intuition, recursion tree for $n=6$:
(ask to do)



- Fib(4) computes twice!
- Fib(3) three times!
- (imagine if we started from Fib(50)!!)

Bad: keep repeating computations. (if I knew answer to Fib(4), wouldn't have to keep repeating whole subtree!)

4. See online: order of computing:

www.cs.usfca.edu/visualization/DPFib.html

www.usfca.edu/ngalles/visualization/DPFib.html
using ngalles (example for 6)

Summary so far:

Fib problem has two important properties:

1. overlapping subproblems (lots!)
 2. solution to big problem can be constructed from solutions to subproblems
- type of problem well solved by dynamic programming

Dynamic Programming Fibonacci: (What to do?)

- 7 -

main idea: save in dictionary subproblems already solved so don't need to recompute (array or hash table).

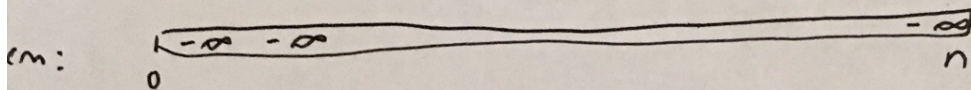
called memoization: from memo pad or memory

clear
hash, memo

memoized Dynamic Programming algorithm:

a. initialization:

Let mem be a new array with values initialized to $-\infty$



mem is dictionary
each time compute
Fib number, first
check if already
in dictionary

b. $\text{Fib}(n)$ // memoized dynamic programming version

1. if $\text{mem}[n] \geq 0$
2. return $\text{mem}[n]$ } if already previously computed and in memo pad

3. if $n=0$ return 1
4. if $n=1$ return 1
5. else $f = \text{Fib}(n-1) + \text{Fib}(n-2)$ } otherwise compute and save value

6. $\text{mem}[n] = f$ } save value in memo pad

7. return f

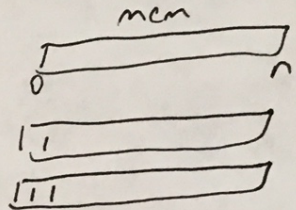
Another Fibonacci Dynamic Programming

algorithm: Bottom up:

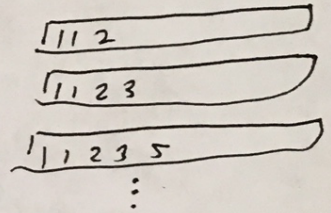
(point to tree; recursive starts at $Fib(n)$ and works down; here start from bottom - more like you first learned $Fib!$)

$Fib(n)$ // bottom up dynamic programming

1. Let $mem[0 \dots n]$ be a new array
2. $mem[0] = 1$
3. $mem[1] = 1$
4. For $i = 2$ to n



5. $mem[i] = mem[i-1] + mem[i-2]$
6. Return $mem[n]$



(could also just keep around last two in table)

same as before

Question: Same or different from previous memoized algorithm?

Answer: Easy to see (one for loop) time is $\Theta(n)$
Equivalent solution to memoization (same things happen in same order)

(see online: www.cs.usfca.edu/visualization/DPFib.html)

Going back to Longest Common Subsequence

problem: (DNA similarity)

(start to get intuition; will continue this and another problem next time)

Example: CCGCTT X_6
ACGGAT Y_5

Longest common subsequence: CGT

$$LCS(6, 6) = 3$$

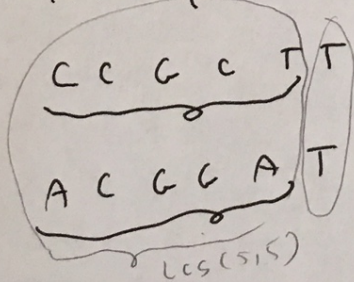
// length of longest common subsequence

↑ ↑
i j
CCGCTT ACGGAT

How to formulate recursion:

→ as in Fib, we'll try and start from largest sequence, and formulate recursion from smaller subproblems:

→ Look at full sequences $i = j = 6$:



ideas for formulating recursion?

In this example: $X_6 = Y_6 = T$:

solution is that of smaller subproblem length 5, plus 1 (for the TT)

$x_i = x_j :$

CCGCT
ACGGA

(T)
(T)

$i = j = 6$

$LCS(6,6) = LCS(5,5) + 1 = 3$

CCGCT T CGT
ACGGA

= 2 CG

Longest common subsequence has to include (T,T); otherwise we could always append (T,T), finding longer sequence - contradiction.

$x_i \neq y_j$: a little more tricky:

Let's say we have:

CCGCT
ACGGA

(C)
(T)

$i = j = 6$

they can't both be in LCS; what to do?

$LCS[6,6] = \max(LCS[5,6], LCS[6,5]) = 3$

CCGCT CGT
ACGGA

= 3 CGT

CCGCTC
ACGGA

= 2 CG

Summary:

$LCS[i,j] = \begin{cases} 0 & i=0 \text{ or } j=0 \\ LCS[i-1, j-1] + 1 & i, j > 0 \text{ } \underline{x_i = y_j} \\ \max(LCS[i-1, j], LCS[i, j-1]) & i, j > 0 \text{ } \underline{x_i \neq y_j} \end{cases}$

Dynamic Programming 2

-0-

Reminder previous class:

Longest common subsequence:

C C G C T C
A C G G A T

answer: CGT

Recursive solution:

$$LCSC[i, j] = \begin{cases} 0 & i=0 \text{ or } j=0 \\ LCSC[i-1, j-1] + 1 & i, j > 0, x_i = y_j \\ \max(LCSC[i-1, j], LCSC[i, j-1]) & i, j > 0, x_i \neq y_j \end{cases}$$

$i=0$ or $j=0$

$i, j > 0, x_i = y_j$

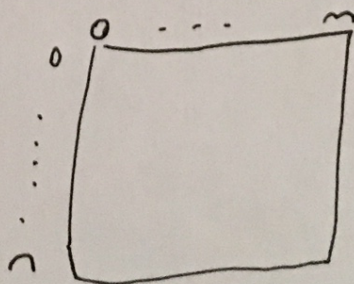
$i, j > 0, x_i \neq y_j$

Dynamic programming solution:

→ Define table $C[0..m, 0..n]$

$n = x.length$

length of first
subsequence



$m = y.length$

length of second
subsequence

→ Either memoize solutions to subproblems not yet computed, or compute solutions to sub-problems bottom up.

→ Run time: $\Theta(mn)$

mn subproblems (at most)

constant computation each

→ We'll write out full bottom up solution (similar to animation we saw)

www.cs.usfca.edu/wgalles/visualization/DPFib.html

Note: main properties that allow DP solution: (ask)

- overlapping subproblems.
- solution to big problem constructed from solutions to smaller subproblems.

Bottom-up from book:

H2X4
- 0 11 -

LCS-LENGTH(x, y)

1. $m = x.length$

2. $n = y.length$

3. let $b[1..m, 1..n]$ and $c[0..m, 0..n]$ be new tables.

4. for $i = 1$ to m

5. $c[i, 0] = 0$

6. for $j = 0$ to n

7. $c[0, j] = 0$

8. For $i = 1$ to m

9. For $j = 1$ to n

10. if $x_i == y_j$:

11. $c[i, j] = c[i-1, j-1] + 1$

12. $b[i, j] = \uparrow$

13. elseif $c[i-1, j] \geq c[i, j-1]$

14. $c[i, j] = c[i-1, j]$

15. $b[i, j] = \uparrow$

16. else $c[i, j] = c[i, j-1]$

17. $b[i, j] = \leftarrow$

18. return c and b

key to reconstructing
right answer
of characters

saves table of length

initialize

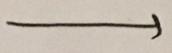
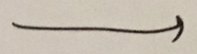
run time:
 $\Theta(mn)$

Example:

CHAT xi

CAT 95

Filled: left to right, from top



:

		J	C	H	A	T
		0	1	2	3	4
i	0	0	0	0	0	0
C	1	0 ↖ 1	← 1	← 1	← 1	
A	2	0 ↑ 1	↑ 1	↖ 2	← 2	
T	3	0 ↑ 1	↑ 1	↑ 2	↖ 3	

Printing: ↖ print(i-1, J-1) // recursive follow arrow
 xi // print out xi

↑ print(i-1, J) // recursive

→ print(i, J-1) // recursive

C A T

(Todo yourself: memoized version + see also animations.)

→ We've structured as large subproblem composed of smaller subproblems.

→ Is our recursive solution efficient?? (ask)

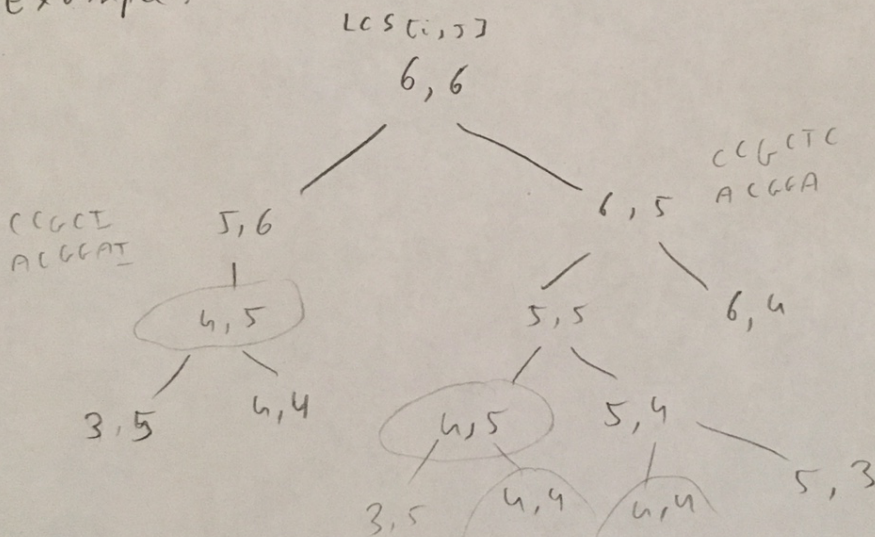
Answer: No. Only if memoized previous solutions (or build bottom up, saving answers to smaller subproblems). → just like in fib!

→ Are there overlapping subproblems? (ask)

Answer: yes.

CCGCTC
ACCGAT

Example:



Dynamic Programming
(we'll write out algorithm next class)

CAF
CHAT

indices start from
-1 + 1 = 0

→ show animation:

[www.cs.usfca.edu / ngeller /
visualization / DPLCS.html](http://www.cs.usfca.edu/~ngeller/visualization/DPLCS.html)

gets stuck in
recur ON