

Algorithms - Spring '25

Dynamic Pro:

Edit dist

Subset Sum

Trees



Recap

- Back tracking HW - due today
- Next HW - up tomorrow,
due Wed, Feb. 19
(Cover DP)
- Readings set for Friday
and next week

Edit distance:

HUGE in bioinformatics!

One of the basic tools in sequence alignment.

(I have a book with an entire chapter on how to optimize.)

Also: spell checkers, word prediction,

How to begin? (Recursively!)

ALGORITHM
↓
ALTRUISTIC
↓

Start at end, & ask "obvious" question:

Insert, delete, edit

try them all,

~~STOP~~

STOOP

↑
edit Insert

(→) Instead;

insert(0)

STOP

STOOP

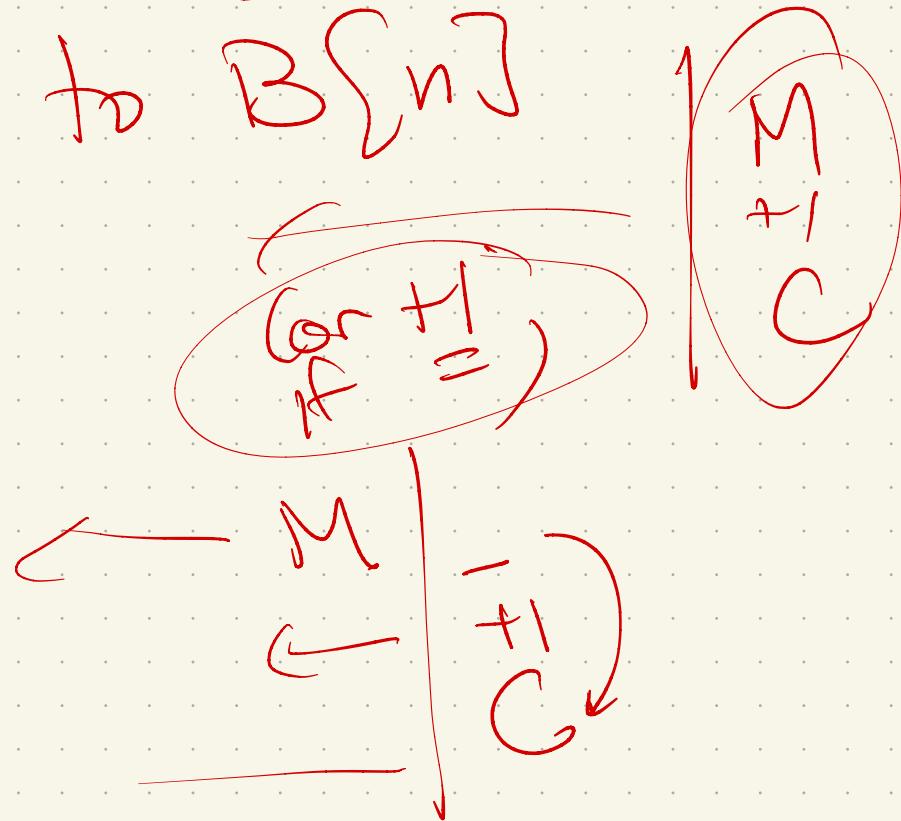
Let's try: A: ALGORITHM

B: ALTRUISTIC

Start at end:

Align: edit $A[m]$
to $B[n]$

Insert



Delete:



Example: TGCATAT
to ATCCGAT

TGCATAT

delete last T

TGCATA

delete last A

TGCAT

insert A at the front

ATGCAT

substitute C for G in the third position

ATCCAT

insert a G before the last A

ATCCGAT

TGCATAT

insert A at the front

ATGCATAT

delete T in the sixth position

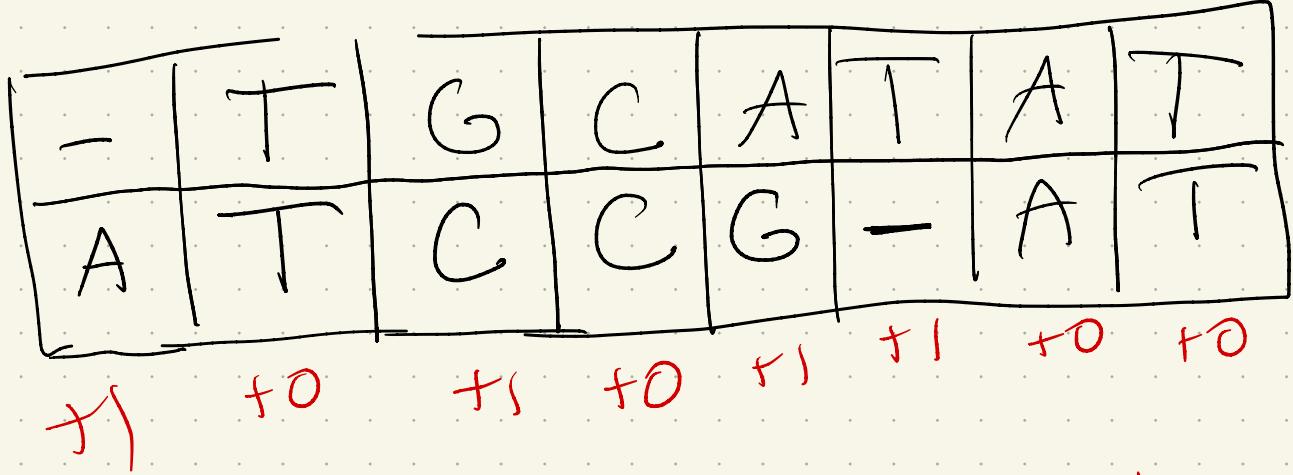
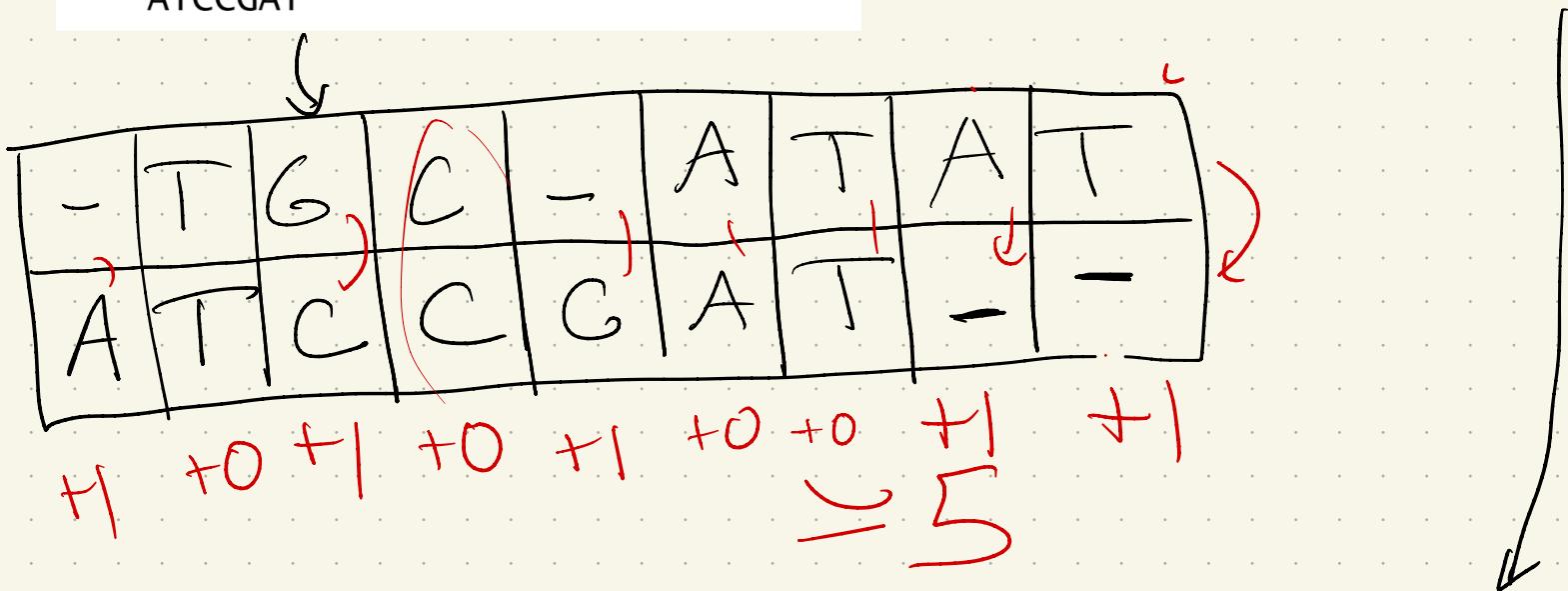
ATGCAAT

substitute G for A in the fifth position

ATGCGAT

substitute C for G in the third position

ATCCGAT



Input: A[1..m]
B[1..n]

Edit(,)

= min {

+ Base cases!

this way:

$$Edit(i, j) = \begin{cases} i & \text{if } j = 0 \\ j & \text{if } i = 0 \\ \min \left\{ \begin{array}{l} Edit(i, j - 1) + 1 \\ Edit(i - 1, j) + 1 \\ Edit(i - 1, j - 1) + [A[i] \neq B[j]] \end{array} \right\} & \text{otherwise} \end{cases}$$

So: what's our "memory" data structure?

Then, our algorithm:

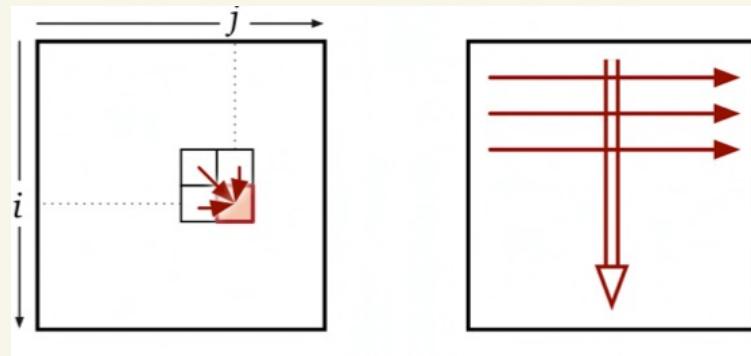
- start w/ base case
(row + column)
- Fill in :

$$Edit(i, j) = \begin{cases} i & \text{if } j = 0 \\ j & \text{if } i = 0 \\ \min \left\{ \begin{array}{l} Edit(i, j - 1) + 1 \\ Edit(i - 1, j) + 1 \\ Edit(i - 1, j - 1) + [A[i] \neq B[j]] \end{array} \right\} & \text{otherwise} \end{cases}$$

Result:

```
EDITDISTANCE( $A[1..m], B[1..n]$ ):  
    for  $j \leftarrow 0$  to  $n$   
         $Edit[0,j] \leftarrow j$   
    for  $i \leftarrow 1$  to  $m$   
         $Edit[i,0] \leftarrow i$   
        for  $j \leftarrow 1$  to  $n$   
             $ins \leftarrow Edit[i,j-1] + 1$   
             $del \leftarrow Edit[i-1,j] + 1$   
            if  $A[i] = B[j]$   
                 $rep \leftarrow Edit[i-1,j-1]$   
            else  
                 $rep \leftarrow Edit[i-1,j-1] + 1$   
             $Edit[i,j] \leftarrow \min \{ins, del, rep\}$   
    return  $Edit[m,n]$ 
```

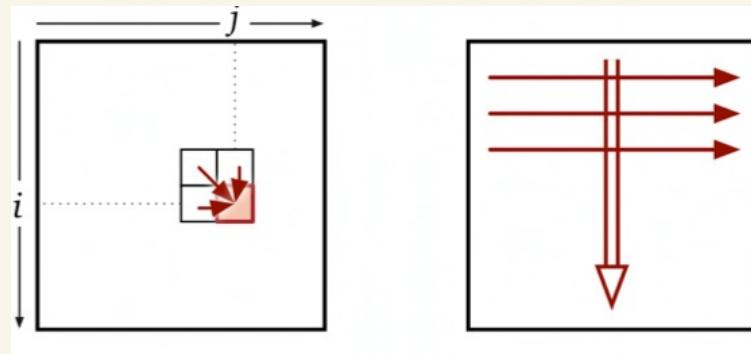
Picture :



Result:

```
EDITDISTANCE( $A[1..m], B[1..n]$ ):  
    for  $j \leftarrow 0$  to  $n$   
         $Edit[0,j] \leftarrow j$   
    for  $i \leftarrow 1$  to  $m$   
         $Edit[i,0] \leftarrow i$   
        for  $j \leftarrow 1$  to  $n$   
             $ins \leftarrow Edit[i,j-1] + 1$   
             $del \leftarrow Edit[i-1,j] + 1$   
            if  $A[i] = B[j]$   
                 $rep \leftarrow Edit[i-1,j-1]$   
            else  
                 $rep \leftarrow Edit[i-1,j-1] + 1$   
             $Edit[i,j] \leftarrow \min \{ins, del, rep\}$   
    return  $Edit[m,n]$ 
```

Picture :



Back to an example:

	A	L	G	O	R	I	T	H	M
A	0 → 1 → 2 → 3 → 4 → 5 → 6 → 7 → 8 → 9								
L	1	0 → 1 → 2 → 3 → 4 → 5 → 6 → 7 → 8							
T	2	1	0 → 1 → 2 → 3 → 4 → 5 → 6 → 7						
R	3	2	1	1 → 2 → 3 → 4 → 4 → 5 → 6					
U	4	3	2	2	2 → 3 → 4 → 5 → 6				
I	5	4	3	3	3	3 → 4 → 5 → 6			
S	6	5	4	4	4	3 → 4 → 5 → 6			
T	7	6	5	5	5	4	4	5	6
I	8	7	6	6	6	5	4	5	6
C	9	8	7	7	7	7	6	5	6
	10	9	8	8	8	8	7	6	6

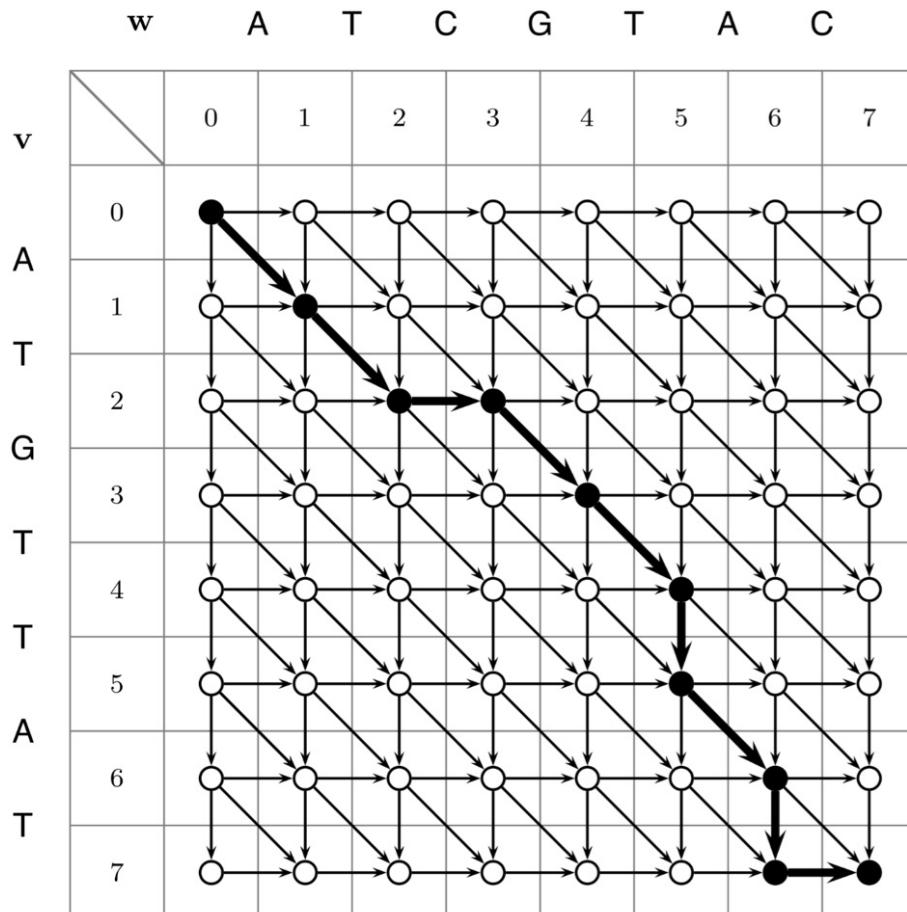
The memoization table for $\text{Edit}(\text{ALGORITHM}, \text{ALTRUISTIC})$

A	L	G	O	R	I	T	H	M	
A	L	T	R	U	I	S	T	I	C

$$\text{Edit}(i, j) = \begin{cases} i & \text{if } j = 0 \\ j & \text{if } i = 0 \\ \min \left\{ \begin{array}{l} \text{Edit}(i, j - 1) + 1 \\ \text{Edit}(i - 1, j) + 1 \\ \text{Edit}(i - 1, j - 1) + [A[i] \neq B[j]] \end{array} \right\} & \text{otherwise} \end{cases}$$

Another (from Bioinformatics book):

$v =$	0	1	2	2	3	4	5	6	7	7
	A	T	-	G	T	T	A	T	-	
$w =$	A	T	C	G	T	-	A	-	C	
	0	1	2	3	4	5	5	6	6	7



\searrow	\swarrow	\rightarrow	\searrow	\swarrow	\downarrow	\searrow	\downarrow	\rightarrow
A	T	-	G	T	T	A	T	-
A	T	C	G	T	-	A	-	C

Question:

Can we do better?

A really good question!

Lots of attention in
bioinformatics

Clever divide and conquer
can reduce space.

↳ but will give #, not
sequence, w/out some
nice tricks

Subset sum (revisited)

Key takeaway (I think):

Sometimes, our backtracking recurrences can be memoized

(Note: Sometimes, they can't!

Think n queens.)

Recall:

Given a set $X[1..n]$ of numbers + a target T ,
find a subset of X whose sum is $= T$.

Ch2 solution

«Does any subset of X sum to T ?»

SUBSETSUM(X, T):

```
if  $T = 0$ 
    return TRUE
else if  $T < 0$  or  $X = \emptyset$ 
    return FALSE
else
     $x \leftarrow$  any element of  $X$ 
    with  $\leftarrow$  SUBSETSUM( $X \setminus \{x\}, T - x$ )    «Recurse!»
    wout  $\leftarrow$  SUBSETSUM( $X \setminus \{x\}, T$ )    «Recurse!»
    return (with  $\vee$  wout)
```

«Does any subset of $X[1..i]$ sum to T ?»

SUBSETSUM(X, i, T):

```
if  $T = 0$ 
    return TRUE
else if  $T < 0$  or  $i = 0$ 
    return FALSE
else
    with  $\leftarrow$  SUBSETSUM( $X, i - 1, T - X[i]$ )    «Recurse!»
    wout  $\leftarrow$  SUBSETSUM( $X, i - 1, T$ )    «Recurse!»
    return (with  $\vee$  wout)
```

The recursion
(Note: same thing as code!!)

$$SS(i, t) = \begin{cases} \text{TRUE} & \text{if } t = 0 \\ \text{FALSE} & \text{if } t < 0 \text{ or } i > n \\ SS(i + 1, t) \vee SS(i + 1, t - X[i]) & \text{otherwise} \end{cases}$$

$$SS(i, t) = T \text{ or } F$$

$0 \leq i \leq n$ $0 \leq t \leq T$

So: another 2-d table!

To decide:

$$SS(i, t) = \begin{cases} \text{TRUE} & \text{if } t = 0 \\ \text{FALSE} & \text{if } t < 0 \text{ or } i > n \\ SS(i+1, t) \vee SS(i+1, t-X[i]) & \text{otherwise} \end{cases}$$

↑ ↑
 look at these 2
 cells.

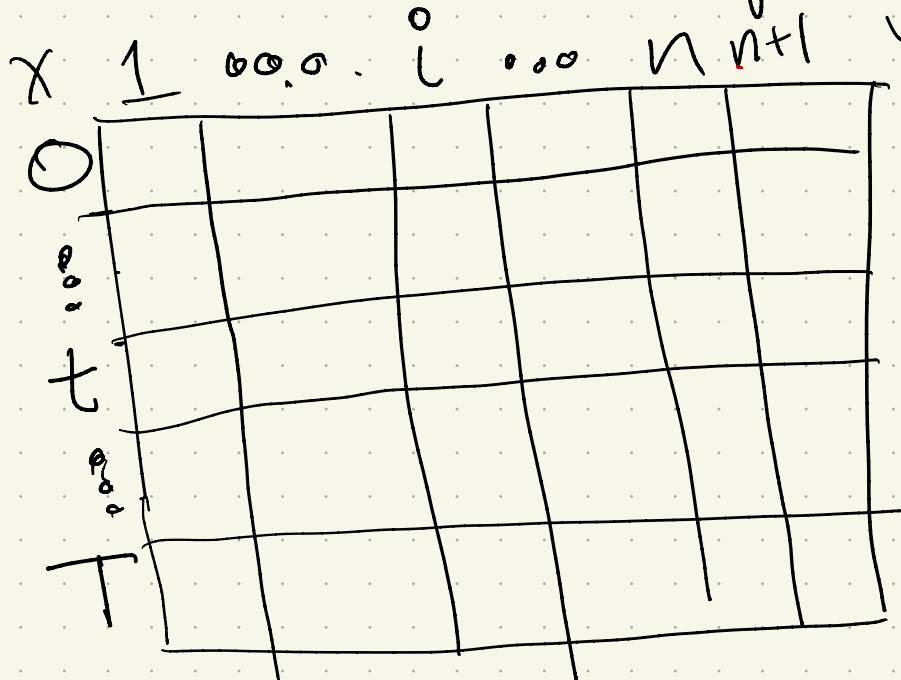
One note: if $t - X[i] < 0$, wasting time!
 Equivalent to:

$$SS(i, t) = \begin{cases} \text{TRUE} & \text{if } t = 0 \\ \text{FALSE} & \text{if } i > n \\ SS(i+1, t) & \text{if } t < X[i] \\ SS(i+1, t) \vee SS(i+1, t-X[i]) & \text{otherwise} \end{cases}$$

Now - need to code this:

$$SS(i, t) = \begin{cases} \text{TRUE} & \text{if } t = 0 \\ \text{FALSE} & \text{if } i > n \\ SS(i+1, t) & \text{if } t < X[i] \\ SS(i+1, t) \vee SS(i+1, t-X[i]) & \text{otherwise} \end{cases}$$

How should our loops go?



Fill:

Hs code :

```
FASTSUBSETSUM( $X[1..n]$ ,  $T$ ):  
     $S[n+1, 0] \leftarrow \text{TRUE}$   
    for  $t \leftarrow 1$  to  $T$   
         $S[n+1, t] \leftarrow \text{FALSE}$   
    for  $i \leftarrow n$  downto 1  
         $S[i, 0] = \text{TRUE}$   
        for  $t \leftarrow 1$  to  $X[i]-1$   
             $S[i, t] \leftarrow S[i+1, t]$     «Avoid the case  $t < 0$ »  
        for  $t \leftarrow X[i]$  to  $T$   
             $S[i, t] \leftarrow S[i+1, t] \vee S[i+1, t-X[i]]$   
    return  $S[1, T]$ 
```

Correctness :

Time/Space Analysis :

Note:

How big is this, & is it even a good idea??

Input: number T and array $X[1..n]$

table has a column for every number $1..T$.

How bad?

Well, X could be a list of 5000 #s, but T could be in the millions!

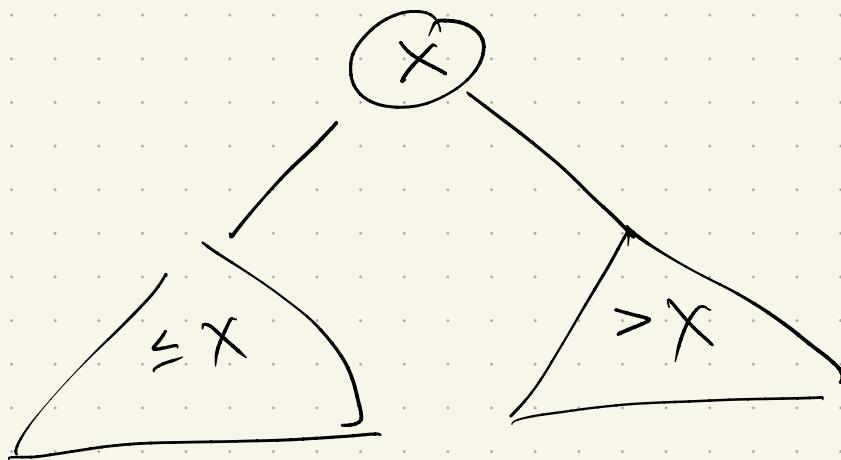
(lots of empty columns, many of which are impossible to hit!)

Balanced search trees (again)

Recall:

What is the "best" one?

Recap:



Time to search for k in T

=

Goal: Given frequencies, built
best BST for those frequencies.

Example:

f: 100, 1, 1, 2, 8

A: 1, 2, 3, 4, 5

assume
sorted

Many BSTs: Which is best?

Construction methods we've studied
in data structures:

Here: Given $X[1..n]$
 $F[1..n]$

element $X[i]$ will have
 $F[i]$ searches.

Intuitively - want higher $F[i]$
 to be closer to the root!

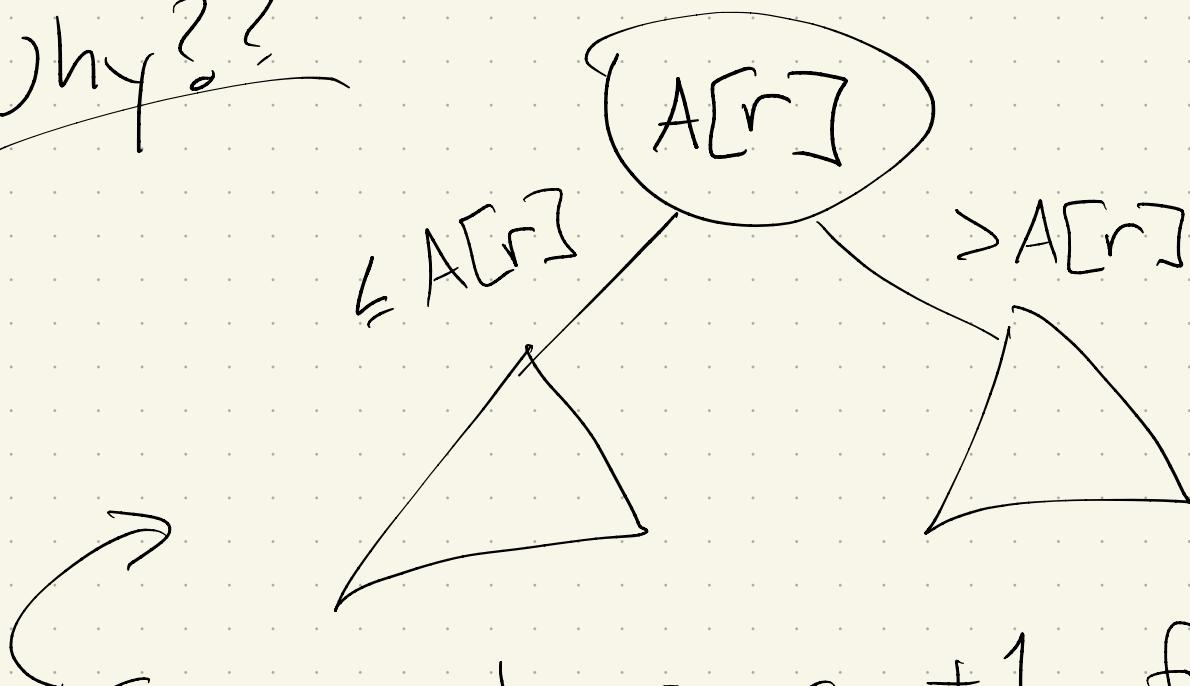
Last chapter:

$$\begin{aligned} Cost(T, f[1..n]) = & \sum_{i=1}^n f[i] + \sum_{i=1}^{r-1} f[i] \cdot \# \text{ancestors of } v_i \text{ in } left(T) \\ & + \sum_{i=r+1}^n f[i] \cdot \# \text{ancestors of } v_i \text{ in } right(T) \end{aligned}$$



$$OptCost(i, k) = \begin{cases} 0 & \text{if } i > k \\ \sum_{j=i}^k f[j] + \min_{i \leq r \leq k} \left\{ OptCost(i, r-1) + OptCost(r+1, k) \right\} & \text{otherwise} \end{cases}$$

Why??



Every node pays +1 for
the root, because search
path must compare to it.

So: we're regrouping!

$$\sum_{i=0}^{n-1} F[i] \cdot (\text{depth in tree})$$

$$= \sum_{\substack{\text{levels } i \\ \text{in tree}}} (\text{sum of frequencies of nodes in level } i \text{ or deeper})$$

$$OptCost(i, k) = \begin{cases} 0 & \text{if } i > k \\ \sum_{j=i}^k f[j] + \min_{i \leq r \leq k} \left\{ OptCost(i, r-1) + OptCost(r+1, k) \right\} & \text{otherwise} \end{cases}$$

Use this to build the "best" tree.

Choose root.

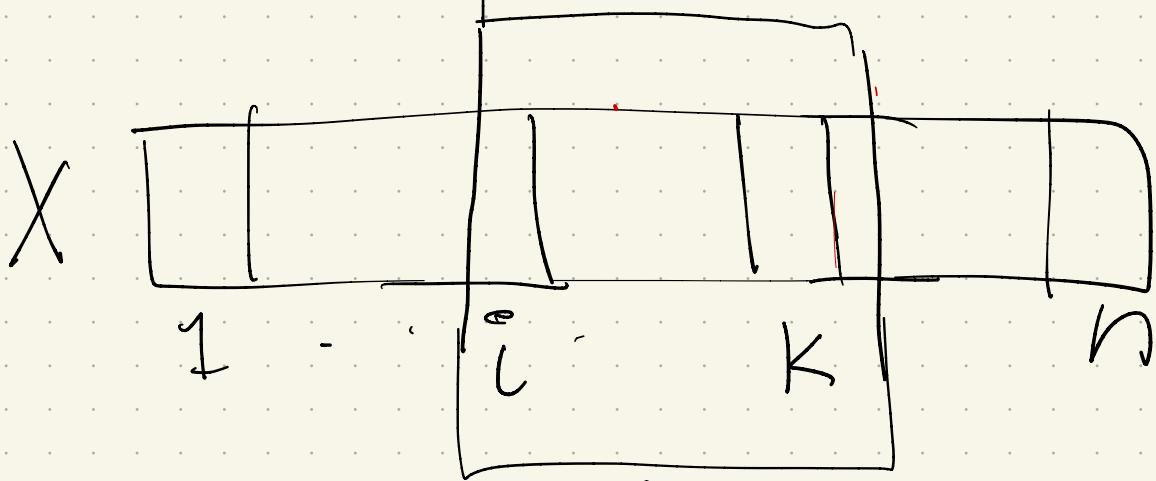
Recursively find best left subtree, + best right subtree.

(Note: try all roots in backtracking!)

How to memoize?

$$OptCost(i, k) = \begin{cases} 0 & \text{if } i > k \\ \sum_{j=i}^k f[i] + \min_{i \leq r \leq k} \left\{ OptCost(i, r-1) + OptCost(r+1, k) \right\} & \text{otherwise} \end{cases}$$

Remember Input:



build best tree here
Everyone here pays $\sum_{j=i}^k f[i]$,
so first precompute &
store these sums.

Time / space:

Let $F[i][k] = \sum_{j=i}^k f[j]$

Now:

$$OptCost(i, k) = \begin{cases} 0 & \text{if } i > k \\ \sum_{j=i}^k f[j] + \min_{i \leq r \leq k} \left\{ OptCost(i, r-1) + OptCost(r+1, k) \right\} & \text{otherwise} \end{cases}$$

\downarrow

$OptCost(i, k) = \left\{ \begin{array}{l} 0 \\ F[i][k] + \end{array} \right.$

Memoize: $0 \leq i \leq k \leq n$

So: 2D table!

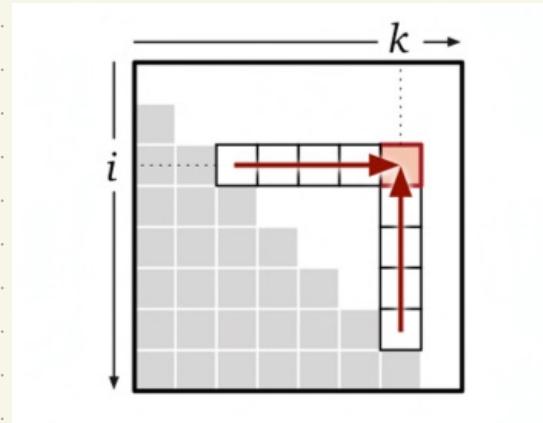
Each $Opt[i][k]$ needs:

- $F[i][k]$
- and

i					
:					
n					

Hs picture (prether):

$$OptCost(i, k) = \begin{cases} 0 & \text{if } i > k \\ F[i, k] + \min_{i \leq r \leq k} \left\{ OptCost(i, r - 1) + OptCost(r + 1, k) \right\} & \text{otherwise} \end{cases}$$



So:

```
OPTIMALBST( $f[1..n]$ ):  
    INITF( $f[1..n]$ )  
    for  $i \leftarrow 1$  to  $n + 1$   
         $OptCost[i, i - 1] \leftarrow 0$   
    for  $d \leftarrow 0$  to  $n - 1$   
        for  $i \leftarrow 1$  to  $n - d$     «...or whatever»  
            COMPUTE $OptCost(i, i + d)$   
    return  $OptCost[1, n]$ 
```

Time:

Space: