TODAY: Dynamic Programming III (of 4)
- subproblems for strings
- parenthesization
- edit distance (& longest common subseq.)
- knapsack
- pseudopolynomial time

\* 5 easy steps to dynamic programming:
   ① define subproblems     count # subprobs.
   ② guess (part of solution)    count # choices
   ③ relate subprob. solutions   compute time/subprob.
   ④ recurse + memoize      time = time/subprob.
  OR build DP table bottom-up    · # subprobs.
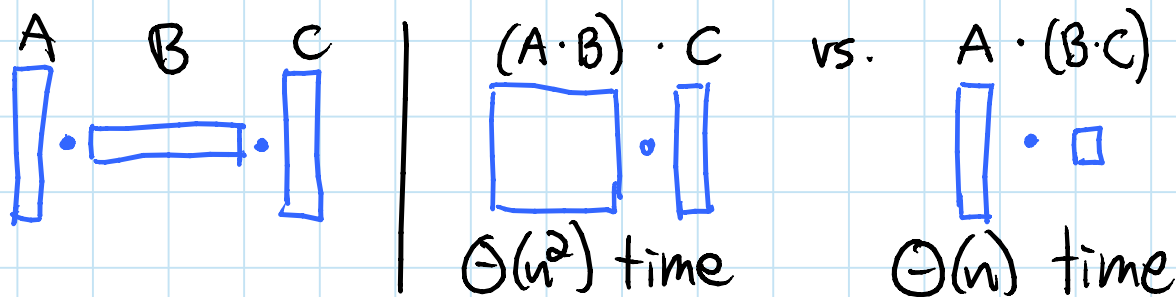      ~ check subprobs. acyclic/topological order
   ⑤ solve original problem: = a subproblem
  OR by combining subprob. solutions (⇒ extra time)

- problems from L20 (text justification, Blackjack)
are on sequences (words,      cards)

\* useful subproblems for strings/sequences $x$:
      ┌→ - suffixes $x[i:]$    } $\Theta(|x|)$ ← cheaper ⇒
  L20     - prefixes $x[:i]$         use if possible
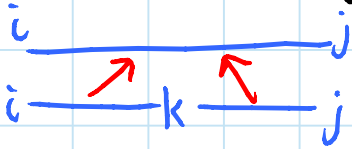         - substrings $x[i:j]$   } $\Theta(|x|^2)$

Parenthesization: optimal evaluation of an associative expression $A[0] \cdot A[1] \cdot \cdots \cdot A[n-1]$
— e.g. multiplying rectangular matrices

A   B   C     $(A \cdot B) \cdot C$   vs.   $A \cdot (B \cdot C)$

$\Theta(n^2)$ time     $\Theta(n)$ time

② <u>guessing</u> = outermost multiplication: $(\cdots)(\cdots)$
⟹ #choices = $O(n)$     $\uparrow$   $\uparrow$    $k-1$   $k$

① <u>subproblems</u> = ~~prefixes & suffixes?~~ <span style="color:red">**NO**</span>
      = cost of substring $A[i:j]$
⟹ # subproblems = $\Theta(n^2)$

③ <u>recurrence</u>:
— $DP[i,j] = \min \left( \begin{array}{l} DP[i,k] + DP[k,j] + \text{cost of} \\ (A[i]\cdots A[k-1]) \cdot (A[k]\cdots A[j-1]) \\ \text{for } k \text{ in range}(i+1, j) \end{array} \right)$

<span style="color:red">DAG:</span>   $i \xrightarrow{} j$    $i \xrightarrow{} k \xleftarrow{} j$

— $DP[i, i+1] = \emptyset$
⟹ cost per subproblem = $\Theta(j-i) = O(n)$

④ <u>topological order</u>: increasing substring size
— total time = $\Theta(n^3)$

⑤ <u>original problem</u> = $DP[0, n]$
<span style="color:green">(& use parent pointers to recover parens.)</span>

Edit distance: (used for DNA comparison, diff, CVS/SVN/···, spellchecking (typos), plagiarism detection, etc.)

given two strings $x$ & $y$, what's the cheapest possible sequence of character edits to transform $x$ into $y$?

insert $c$    delete $c$    replace $c \to c'$

- cost of edit depends only on characters $c, c'$
- e.g. in DNA, $C \to G$ common mutation ⇒ low cost
- cost of sequence = sum of costs of edits

- if insert & delete cost 1, replace costs ∅, min. edit distance equivalent to finding longest common subsequence

  sequential but not necessarily contiguous

- e.g.:  H I E R O G L Y P H O L O G Y    } ⟶ HELLO
  vs.  M I C H A E L A N G E L O

Subproblems for multiple strings/sequences: combine suffix/prefix/substring subproblems
- multiply state spaces
- still polynomial for $O(1)$ strings

# Edit distance DP:

① <u>subproblems</u>: $c(i,j) =$ edit-distance$(x[i:], y[j:])$
for $0 \leq i < |x|,\ 0 \leq j < |y|$

$\Rightarrow \Theta(|x| \cdot |y|)$ subproblems

② <u>guess</u> whether to turn $x$ into $y$
- $x[i]$ deleted
- $y[j]$ inserted
- $x[i]$ replaced by $y[j]$    } 3 choices

③ <u>recurrence</u>: $c(i,j) = \max \{$
$\quad \text{cost}(\text{delete } x[i]) + c(i+1, j)$  if $i < |x|,$
$\quad \text{cost}(\text{insert } y[j]) + c(i, j+1)$  if $j < |y|,$
$\quad \text{cost}(\text{replace } x[i] \to y[j]) + c(i+1, j+1)$
$\qquad\qquad\qquad$ if $i < |x|$ & $j < |y| \}$

- base case: $c(|x|, |y|) = \emptyset$
$\Rightarrow \Theta(1)$ time per subproblem

④ <u>topological order</u>: DAG in 2D table:  $\longleftarrow j \longrightarrow$
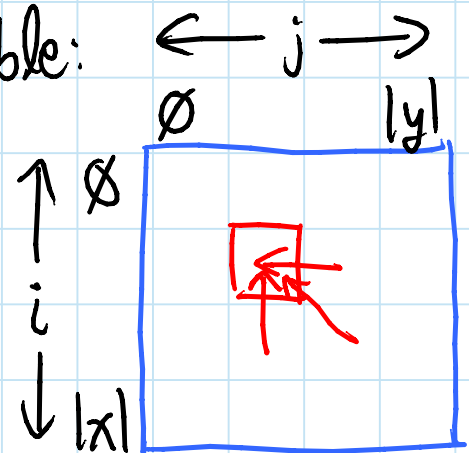- bottom up <u>OR</u> right to left
- only need to keep last
  2 rows/columns
$\Rightarrow$ linear space
- total time $= \Theta(|x| \cdot |y|)$

⑤ <u>original problem</u>: $c(\emptyset, \emptyset)$

# Knapsack of size $S$ you want to pack
  - item $i$ has integer <u>size</u> $s_i$ & real <u>value</u> $v_i$
  - <u>goal</u>: choose subset of items of max. total value
         subject to total size $\leq S$

<u>First attempt:</u>
  ① <s><u>subproblem</u> = value for suffix $i$:</s>   WRONG
  ② <u>guessing</u> = whether to include item $i$
  $\Rightarrow$ #choices = 2
  ③ <u>recurrence</u>:
     $-DP[i] = \max(DP[i+1], v_i + DP[i+1]$ if $\boxed{s_i \leq S}$?!)
     $-$ not enough information to know whether
        item $i$ fits $-$ how much space is left?
                                                   GUESS!
<u>Correct:</u>
  ① <u>subproblem</u> = value for suffix $i$:
                 given knapsack of size $X$
     $\Rightarrow$ # subproblems = $O(n S)$      (!)
  ③ <u>recurrence</u>:
     $-DP[i, X] = \max(DP[i+1, X],$
                       $v_i + DP[i+1, X - s_i]$ if $s_i \leq X)$
     $- DP[n, X] = \emptyset$
     $\Rightarrow$ time per subproblem = $O(1)$
  ④ <u>topological order</u>: for $i$ in $n, \ldots, 0$: for $X$ in $0, \ldots, S$
     $-$ total time = $O(n S)$
  ⑤ <u>original problem</u> = $DP[0, S]$
     (& use parent pointers to recover subset)

<u>Polynomial time</u> = polynomial in <u>input size</u>
— here $\Theta(n)$ if number $S$ fits in a word
  $O(n \lg S)$ in general
— $S$ is <u>exponential</u> in $\lg S$ (not polynomial)

<u>Pseudopolynomial time</u> = polynomial in
the problem size AND the <u>numbers</u> in input
here: $S, s_i, s_n, v_i's$
— $\Theta(n S)$ is pseudopolynomial

$\Longrightarrow$

<u>Remember</u>:  polynomial — GOOD
exponential — BAD
pseudopoly. — SO SO