

Design and Analysis of Algorithms

Dynamic Programming (III)

- 1 Return on Investment
- 2 Knapsack Problem
 - Knapsack with Repetition
 - Knapsack without Repetition
- 3 Longest Common Substring
- 4 Edit Distance
- 5 Summary of Dynamic Programming

- 1 Return on Investment
- 2 Knapsack Problem
 - Knapsack with Repetition
 - Knapsack without Repetition
- 3 Longest Common Substring
- 4 Edit Distance
- 5 Summary of Dynamic Programming

Return on Investment

Problem. Given m coins, n projects, and function $f_i(x)$: profit of investing x on the i -th project. Find the optimal investment scheme that maximizes profit..

Solution: a vector (x_1, x_2, \dots, x_n) , x_i : investment on project i

Optimized function: $\max \sum_{i=1}^n f_i(x_i)$

Constraints: $x_1 + x_2 + \dots + x_n = m$, $x_i \in N$

Table: 5 coins on 4 projects

x	$f_1(x)$	$f_2(x)$	$f_3(x)$	$f_4(x)$
0	0	0	0	0
1	11	0	2	20
2	12	5	10	21
3	13	10	30	22
4	14	15	32	23
5	15	20	40	24

Subproblems and Computation Order

Subproblem: defined by k and x

- k : invest on the $1, 2, \dots, k$ projects
- the total investment is less than x

The parameter in matrix multiplication chain is a tuple of index, of the same type

(k, x) are of different types \leadsto 2-dimension dynamic programming

Original problem: $k = n, x = m$

Computation order: $k = 1, 2, \dots, n$; for any $k, x = 1, 2, \dots, m$

- can be implemented by two level loop

Iteration Relation of Optimized Function

Optimized function $F_k(x)$: the maximal profit of investing x coins on the first k projects

Iteration relation: Determine $F_k(x)$ from $F_{k-1}(y \leq x)$

$$F_k(x) = \max_{0 \leq x_k \leq x} \{f_k(x_k) + F_{k-1}(x - x_k)\}, k > 1$$

$$F_1(x) = f_1(x), k = 1 \quad (\text{initial values})$$

Demo of $k = 2$

x	$f_1(x)$	$f_2(x)$	$f_3(x)$	$f_4(x)$
0	0	0	0	0
1	11	0	2	20
2	12	5	10	21
3	13	10	30	22
4	14	15	32	23
5	15	20	40	24

$k = 1$ corresponds to the initial values: $F_1(1) = 11$, $F_1(2) = 12$,
 $F_1(3) = 13$, $F_1(4) = 14$, $F_1(5) = 15$

$$F_2(1) = \max\{f_1(1), f_2(1)\} = 11$$

$$F_2(2) = \max\{f_2(2), F_1(1) + f_2(1), F_1(2)\} = 12$$

$$F_2(3) = \max\{f_2(3), F_1(1) + f_2(2), F_1(2) + f_2(1), F_1(3)\} = 16$$

Similarly, we can compute $F_2(4) = 21$, $F_2(5) = 26$

Memo and Solution

x	$F_1(\cdot) \ s_1(\cdot)$	$F_2(\cdot) \ s_2(\cdot)$	$F_3(\cdot) \ s_3(\cdot)$	$F_4(\cdot) \ s_4(\cdot)$
1	11 1	11 0	11 0	20 1
2	12 2	12 0	13 1	31 1
3	13 3	16 2	30 3	33 1
4	14 4	21 3	41 3	50 1
5	15 5	26 4	43 4	61 1

- $F_k(x)$ records maximized profit of investing x coins on the first k projects
- $s_k(x)$ records the investment on k -th project

$$s_4(5) = 1 \Rightarrow x_4 = 1, s_3(5 - 1) = s_3(4)$$

$$s_3(4) = 3 \Rightarrow x_3 = 3, s_2(4 - 3) = s_2(1)$$

$$s_2(1) = 0 \Rightarrow x_2 = 0, s_1(1 - 0) = s_1(1)$$

$$s_1(1) = 1 \Rightarrow x_1 = 1$$

Solution: $(x_1 = 1, x_2 = 0, x_3 = 3, x_4 = 1), F_4(5) = 61$

Complexity Analysis

Memo table is a matrix of m rows (total number of coins) and n columns (total number of projects), totally mn items:

$$F_k(x) = \max_{0 \leq x_k \leq x} \{f_k(x_k) + F_{k-1}(x - x_k)\}, k > 1$$

$$F_1(x) = f_1(x), k = 1 \quad // \text{initial values}$$

The cost of computing $F_k(x)$: there are possible $x + 1$ different choices of $x_k \Rightarrow$ $x + 1$ times add + x times compare

Total number of add

$$\sum_{k=2}^n \sum_{x=1}^m (x + 1) = \frac{1}{2}(n - 1)m(m + 3)$$

Total number of compare

$$\sum_{k=2}^n \sum_{x=1}^m x = \frac{1}{2}(n - 1)m(m + 1)$$

Time complexity $W(n) = O(nm^2)$, space complexity is $O(mn)$

- 1 Return on Investment
- 2 Knapsack Problem
 - Knapsack with Repetition
 - Knapsack without Repetition
- 3 Longest Common Substring
- 4 Edit Distance
- 5 Summary of Dynamic Programming

Motivation

During a robbery, a burglar finds much more loot than he had expected and has to decide what to take.

- His bag (or “knapsack”) will hold a total weight if at most W pounds.
- He want to figure out the most valuable combination of items he can fit into his bag, **quickly**.

There are two version of this problem:

- with repetition: there are unlimited quantities of each item available
- without repetition: there is one of each item (the bugalar has broken into an art gallery)

Neither version is likely to have a polynomial-time algorithm.

Formal Motivation

If the above motivation seems frivolous

- replace “weight” with “CPU time”
- replace “only W pounds can be taken” with “only W units of CPU times are available”

CPU time can also be replaced by bandwidth

The knapsack problem generalizes a wide variety of resource-constrained selection tasks.

Knapsack with Repetition

Problem. Given n items and a knapsack, item i weighs $w_i > 0$ and has value $v_i > 0$, knapsack has capacity of W

Goal. Fill knapsack so as to maximize total value.

Table: knapsack instance, $W = 11$

i	1	2	3	4	5
v_i	1	6	18	22	28
w_i	1	2	5	6	7

- Greedy by value (maximum v_i first): $\{5, 2 \times 2\}$ has value 40
- Greedy by weight (minimum w_i first): $\{1 \times 11\}$ has value 11
- Greedy by ratio (maximum ratio v_i/w_i first): $\{5, 2 \times 2\}$ has value 40

Observation. None of greedy algorithms is optimal.

Modeling

Solution vector: $x = (x_1, x_2, \dots, x_n) \in (\mathbb{Z}^+)^n$, x_i is the number of item i

Optimized goal: $\max \sum_{i=1}^n v_i x_i$

Constraint: $\sum_{i=1}^n w_i x_i \leq W, x_i \in \mathbb{N}$

- linear programming: find min or max of optimized function with linear constraints
- integer programming: linear programming when x_i are non-negative integers

As always, the main question in dynamic programming is:

what are subproblems

- It usually takes a little experimentation to figure out exactly what works.

Dynamic Programming: False Start

We can shrink the original problem in two ways:

- ① smaller knapsack capacities $w \leq W$
- ② fewer items (for instance, items $1, 2, \dots, j$ for $j \leq n$)

Def. $K(j)$ = maximum value achievable with items $1, \dots, j$ with weight limit W .

Case 1. K does not select item j .

- K selects best of $\{1, 2, \dots, j-1\} \rightsquigarrow$ satisfy optimal substructure property (proof via exchange argument)

Case 2. K selects item j

- We don't know the consequence of selecting item j , cause it will change weight limit of subproblems \rightsquigarrow cannot make a decision

Need more subproblems!

Dynamic programming: Adding a New Variable

Def. $K_j(w) = \max$ value of choosing from items $\{1, \dots, j\}$ with weight limit w .

Case 1. K does not select item j

- K selects best of $\{1, 2, \dots, j-1\}$ using weight limit w .

Case 2. K selects item j (at least 1)

- New weight limit $= w - w_j$.
- K selects best of $\{1, 2, \dots, j\}$ using this new weight limit (cause we allow repetition)

Both cases satisfy optimal substructure property (proof via exchange argument)

Wrap it Up

Subproblem: defined by two variables j and w

- j : select from subset of $\{1, 2, \dots, j\}$
- w : limit on weight

$K_j(w)$: maximum value achievable of selecting from the first j items with weight limit w

Computation order: $j = 1 \rightarrow n$; for any k , $w = 1 \rightarrow W$

$$\begin{cases} \underline{K_j(w) = \max\{K_{j-1}(w), K_j(w - w_j) + v_j\}} \\ K_0(w) = 0, 0 \leq w \leq W, K_j(0) = 0, 0 \leq j \leq n \\ K_1(y) = \left\lfloor \frac{W}{w_1} \right\rfloor v_1, K_j(w) = -\infty, w < 0 \end{cases}$$

- $K_j(w - w_j) + v_j$: maximum value when selecting at least one j -th item

Pseudocode of Knapsack

Algorithm 1: Knapsack($n, W, w_1, \dots, w_n; v_1, \dots, v_n$)

```
1: for  $w = 0$  to  $W$  do  $K_0(w) \leftarrow 0$ ;  
2: for  $j = 1$  to  $n$  do  $K_j(0) \leftarrow 0$ ;  
3: for  $w = 0$  to  $W$  do  $K_1(w) \leftarrow \lceil W/w_1 \rceil v_1$ ;  
4:  $K_j(w) = -\infty, w < 0$ ;  
5: for  $j = 1$  to  $n$  do  
6:   for  $w = 0$  to  $W$  do  
7:      $K_j(w) = \max\{K_{j-1}(w), K_j(w - w_j) + v_j\}$   
8:   end  
9: end
```

- Bottom-up approach

Table: knapsack instance, $n = 4, W = 10$

i	1	2	3	4
v_i	1	3	5	9
w_i	2	3	4	7

Computation process of $K_j(w)$ (hint: how to fill the matrix)

- left to right, top to down
- top to down, left to right

$j \backslash w$	1	2	3	4	5	6	7	8	9	10
1	0	1	1	2	2	3	3	4	4	5
2	0	1	3	4	4	6	6	7	9	9
3	0	1	3	5	5	6	8	10	10	11
4	0	1	3	5	5	6	9	10	10	12

A Remark

Alternative optimization function: like ROI problem

$$K_j(w) = \max_{0 \leq x_j \leq \lfloor w/w_j \rfloor} \{K_j(w - x_j \cdot w_j) + x_j \cdot v_j\}$$

- Pros: more intuitive and easy to understand
- Cons: complexity of computing $K_j(w)$ depends on w , in contrast to the original representation which only requires one comparison.

Lesson

The design of optimized function is vital

Trace Function

$s_j(w)$: the biggest item number in solution $K_j(w)$

$$s_j(w) = \begin{cases} s_{j-1}(w) & K_{j-1}(w) > K_j(w - w_k) + v_k \\ j & K_{j-1}(w) \leq K_j(w - w_k) + v_k \end{cases}$$

$$s_1(w) = \begin{cases} 0 & w < w_1 \\ 1 & w \geq w_1 \end{cases}$$

Trace function is used to trace solution and output the detailed information

Algorithm 2: TraceSolution($s[n, W]$)

Input: table $s_j(w)$, $j \in [n]$, $w \in [W]$

Output: solution vector x_1, x_2, \dots, x_n

```
1: for  $i \leftarrow 1$  to  $n$  do  $x_i \leftarrow 0$ ;  
2:  $w \leftarrow W$ ,  $k \leftarrow n$ ;  
3:  $x_k \leftarrow 0$  ;  
4: while  $s_k(w) = k$  do                                //continue select  $k$ -th item  
5:    $w \leftarrow w - w_k$ ;  
6:    $x_k \leftarrow x_k + 1$ ;  
7: end  
8: if  $s_k(w) \neq 0$  then  $k \leftarrow k - 1$ , goto 4;        //trace next item  
9: else finishes tracing;
```

Trace Solution

Table: $s_j(w)$

$j \backslash w$	1	2	3	4	5	6	7	8	9	10
1	0	1	1	1	1	1	1	1	1	1
2	0	1	2	2	2	2	2	2	2	2
3	0	1	2	3	3	3	3	3	3	3
4	0	1	2	3	3	3	4	3	4	4

- $s_4(10) = 4 \Rightarrow x_4 \geq 1$
- $s_4(10 - w_4) = s_4(3) = 2 \Rightarrow x_4 = 1, x_3 = 0, x_2 \geq 1$
- $s_2(3 - w_2) = s_2(0) = 0 \Rightarrow x_2 = 1, x_1 = 0$

Solution: $x_1 = 0, x_2 = 1, x_3 = 0, x_4 = 1$, max profit is 12.

Complexity Analysis

The above algorithm solves the knapsack problem with n items and maximum weight W in $\Theta(nW)$ time and $\Theta(nW)$ space.

According to the optimization function

$$K_j(w) = \max\{K_{j-1}(w), K_j(w - w_j) + v_j\}$$

- Memo computation: takes $O(1)$ time per table entry, there are $\Theta(nW)$ table entries
- Trace back: at most $\Theta(n + W)$ steps (think why?)

The total time complexity and space complexity are $O(nW)$

Remarks

- Not polynomial in input size, cause for integer W , binary representation requires $\log W$ bit, thus input size is n and $\log W \leftarrow$ **super-polynomial**

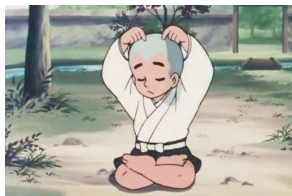
A Second Thought

Do we really have to use 2-dimension dynamic programming?



A Second Thought

Do we really have to use 2-dimension dynamic programming?

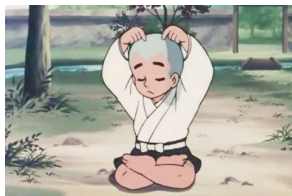


Consider only put restriction on weight, define:

$K(w)$ = maximum value achievable with a knapsack with limit w

A Second Thought

Do we really have to use 2-dimension dynamic programming?



Consider only put restriction on weight, define:

$K(w)$ = maximum value achievable with a knapsack with limit w

How to express this in terms of smaller subproblems?

- If the optimal solution to $K(w)$ includes item i , then removing this item leaves an optimal solution to $K(w - w_i)$. In other words, $K(w) = K(w - w_i) + v_i$, for some i .
- We don't know which i , so we need to try all possibilities.

Iteration Relation

The algorithm now writes itself \leadsto incredibly simple and elegant

Algorithm 3: Knapsack(P, n, W)

```
1:  $K(0) \leftarrow 0$ ;  
2: for  $w = 1$  to  $W$  do  
3:    $K(w) = \max_{i:w_i \leq w} \{K(w - w_i) + v_i\}$   
4: end  
5: return  $K(W)$ ;
```

The algorithm fills in a one-dimension table of length $W + 1$, in left-to-right order

- each entry can take up to $O(n)$ time to compute \Rightarrow overall running time is $O(nW)$.

Think over It

As always, there is an underlying dag. Try constructing it, and you will be rewarded with a startling insight

- this particular variant of knapsack boils down to finding the longest path in a dag.



你品 你细品

Knapsack without Repetition

What if repetitions are not allowed?

Our earlier subproblems now become complete useless.

- For instance, knowing that the value $K(w - w_j)$ does not help to make further decision, cause we don't know whether or not item j has already got used up in this partial solution.

We must refine the subproblem to carry additional information about the items being used as before \leadsto add another parameter $0 \leq j \leq n$:

$K_j(w)$ = maximum value using items $\{1, \dots, j\}$ and weight limit w

The answer we seek is $K_n(W)$.

Iteration Relation

How to express $K_j(w)$ in terms of smaller subproblems?

Quite simple: either item j is needed to achieve the optimal value or it isn't needed.

$$K_j(w) = \max\{K_{j-1}(w - w_j) + v_j, K_{j-1}(w)\}$$

This algorithm fills out a 2-dimension table, with $W + 1$ rows and $n + 1$ columns. Each table entry takes just constant time. The running time remains the same: $O(nW)$.

Algorithm 4: Knapsack(P, n, W)

```
1:  $K_0(w) \leftarrow 0$  for  $w \in [0, W]$ ,  $K_j(0) = 0$  for  $j \in [0, n]$ ;  
2: for  $j = 1$  to  $n$  do  
3:   for  $w = 1$  to  $W$  do  
4:      $K_j(w) = \max\{K_{j-1}(w - w_j) + v_j, K_{j-1}(w)\}$   
5:   end  
6: end  
7: return  $K_n(W)$ ;
```

Memoization

In dynamic programming, we write out a recursive formula that express large problems in terms of smaller ones and then use it to fill a table of solution values in a bottom-up manner, from smaller subproblem to largest.

The formula also suggests a recursive algorithm.

As we saw earlier that naive recursion can be terribly inefficient, because it solves the same subproblems over and over again.

What about a more intelligent recursive implementation? One that remembers its previous invocations and thereby avoids repeating them?

Memoization

On the knapsack problem (with repetitions), algorithm would use hash table to store $K(\cdot)$ that had already been computed.

- At each recursive call requesting some $K(w)$, the algorithm would first check if the answer was already in the table and then would proceed to its calculation only if it wasn't.
- This trick is called *memoization*.

Complexity: recursive algorithm never repeats a subproblem \leadsto running time is $O(nW)$, just like dynamic program.

- However, the constant factor in the big- O notation is substantially larger because of the overhead of recursion.

In some cases, memoization pays off.

- Dynamic programming automatically solves every subproblem that could *conceivably be needed*, while memoization only ends up solving the ones that are actually needed.

Extension of Knapsack Problem

Decision version of knapsack problem is \mathcal{NP} -COMPLETE.

There exists a poly-time algorithm that produces a feasible solution that has value within 1% of optimum.

Variants of Knapsack

- Knapsack with constraint on item number: maximum number of i -th item is n_i
 - 0 – 1 Knapsack: $x_i = 0, 1; i \in [n]$
- Multi-Knapsack: m knapsack, the weight limit of knapsack i is $W_i, i \in [m]$.
- 2-dimension Knapsack: each item with weight w_i and volume $t_i, i \in [n]$, the weight limit is W , the volume limit is V

- 1 Return on Investment
- 2 Knapsack Problem
 - Knapsack with Repetition
 - Knapsack without Repetition
- 3 Longest Common Substring
- 4 Edit Distance
- 5 Summary of Dynamic Programming

Longest Common Substring

Let $X = (x_1, x_2, \dots, x_m)$ and $Z = (z_1, z_2, \dots, z_n)$ be two strings. Z is a **substring** of X if there exists an index sequence of strict increasing order (i_1, \dots, i_k) such that $z_k = x_{i_k}$ for all $k \in [n]$.

Common substring of X and Y : the substring of both X and Y .

Problem. Find the longest string of $X = (x_1, x_2, \dots, x_m)$ and $Y = (y_1, y_2, \dots, y_n)$.

Example

- $X : A \text{ } \color{red}{B} \text{ } \color{red}{C} \text{ } \color{red}{B} \text{ } D \text{ } \color{red}{A} \text{ } B$
- $Y : \color{red}{B} \text{ } D \text{ } \color{red}{C} \text{ } A \text{ } \color{red}{B} \text{ } \color{red}{A}$

LCS: $\color{red}{B} \text{ } \color{red}{C} \text{ } \color{red}{B} \text{ } \color{red}{A}$, length is 4

Brute Force Algorithm

Assume $m \leq n$, $|X| = m$, $|Y| = n$

Brute force algorithm: for each substring of X , check if the substring appears in Y

Complexity analysis

- check each substring takes $O(n)$
 - think how? hint: sequentially scan two strings using two pointers (after each comparison, at least one point moves forward, thus the maximum number of comparison is $2n$)
- there are totally 2^m substrings in X

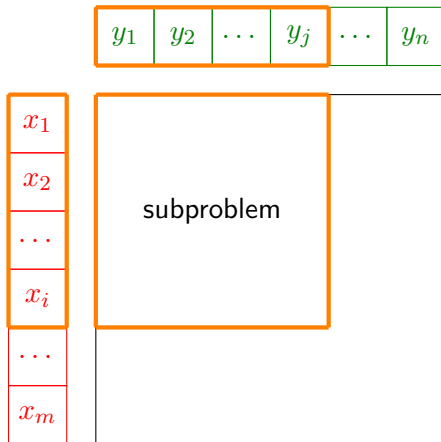
Complexity: $O(n2^m)$

Dynamic Programming: Subproblem

Introduce i and j to define subproblem

X right boundary is i , Y right boundary is j

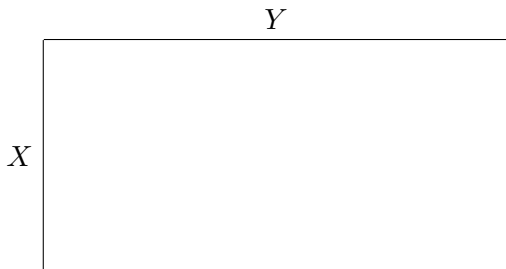
$X_i = (x_1, x_2, \dots, x_i)$, $Y_j = (y_1, y_2, \dots, y_j)$



Relations Between Problems and Subproblems

$$X_m = (x_1, x_2, \dots, x_m), Y_n = (y_1, y_2, \dots, y_n) \\ Z_k = (z_1, z_2, \dots, z_k) = \text{LCS}(X_m, Y_n)$$

Consider the following cases:



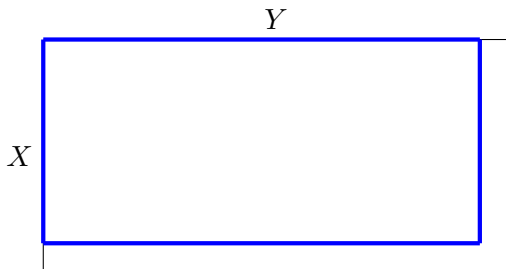
Relations Between Problems and Subproblems

$$X_m = (x_1, x_2, \dots, x_m), Y_n = (y_1, y_2, \dots, y_n)$$

$$Z_k = (z_1, z_2, \dots, z_k) = \text{LCS}(X_m, Y_n)$$

Consider the following cases:

- $x_m = y_n \Rightarrow z_k = x_m = y_n, Z_{k-1} = \text{LCS}(X_{m-1}, Y_{n-1})$



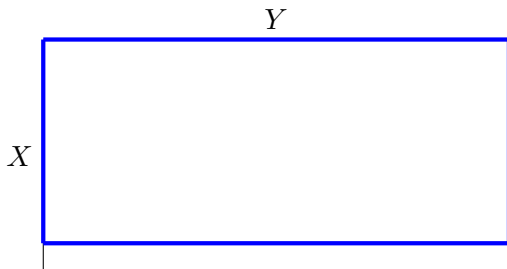
Relations Between Problems and Subproblems

$$X_m = (x_1, x_2, \dots, x_m), Y_n = (y_1, y_2, \dots, y_n)$$

$$Z_k = (z_1, z_2, \dots, z_k) = \text{LCS}(X_m, Y_n)$$

Consider the following cases:

- $x_m = y_n \Rightarrow z_k = x_m = y_n, Z_{k-1} = \text{LCS}(X_{m-1}, Y_{n-1})$
- $x_m \neq y_n$ (the following cases occur)
 - $z_k \neq x_m \Rightarrow Z_k = \text{LCS}(X_{m-1}, Y_n)$

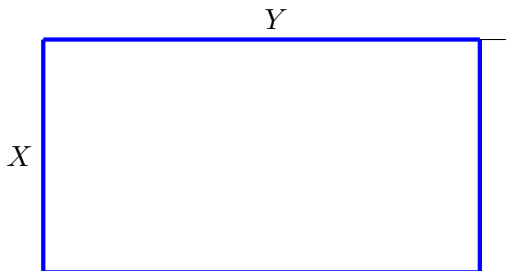


Relations Between Problems and Subproblems

$$X_m = (x_1, x_2, \dots, x_m), Y_n = (y_1, y_2, \dots, y_n) \\ Z_k = (z_1, z_2, \dots, z_k) = \text{LCS}(X_m, Y_n)$$

Consider the following cases:

- $x_m = y_n \Rightarrow z_k = x_m = y_n, Z_{k-1} = \text{LCS}(X_{m-1}, Y_{n-1})$
- $x_m \neq y_n$ (the following cases occur)
 - $z_k \neq x_m \Rightarrow Z_k = \text{LCS}(X_{m-1}, Y_n)$
 - $z_k \neq y_n \Rightarrow Z_k = \text{LCS}(X_m, Y_{n-1})$

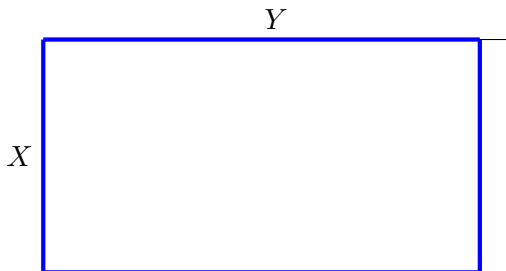


Relations Between Problems and Subproblems

$$X_m = (x_1, x_2, \dots, x_m), Y_n = (y_1, y_2, \dots, y_n) \\ Z_k = (z_1, z_2, \dots, z_k) = \text{LCS}(X_m, Y_n)$$

Consider the following cases:

- $x_m = y_n \Rightarrow z_k = x_m = y_n, Z_{k-1} = \text{LCS}(X_{m-1}, Y_{n-1})$
- $x_m \neq y_n$ (the following cases occur)
 - $z_k \neq x_m \Rightarrow Z_k = \text{LCS}(X_{m-1}, Y_n)$
 - $z_k \neq y_n \Rightarrow Z_k = \text{LCS}(X_m, Y_{n-1})$



satisfy optimal sub-structure

Optimized Function and Iteration Relation

Optimized function: $L(i, j)$

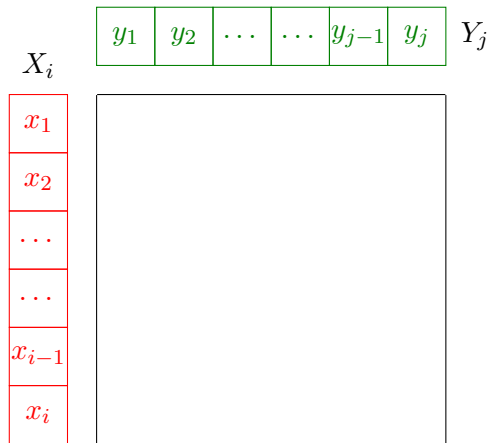
- LCS length of $X_i = (x_1, x_2, \dots, x_i)$ and $Y_j = (y_1, y_2, \dots, y_j)$

Iteration relation

$$L(i, j) = \begin{cases} 0 & i = 0 \vee j = 0 \\ L(i-1, j-1) + 1 & i, j > 0 \wedge x_i = y_j \\ \max\{L(i, j-1), L(i-1, j)\} & i, j > 0 \wedge x_i \neq y_j \end{cases}$$

Indicator Function

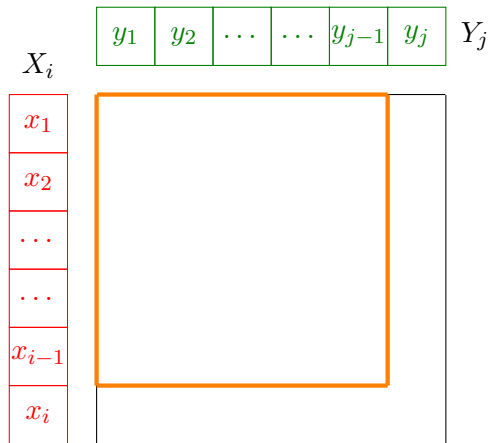
Indicator function $s(i, j)$ with values: $\nwarrow, \leftarrow, \uparrow$



Indicator Function

Indicator function $s(i, j)$ with values: $\nwarrow, \leftarrow, \uparrow$

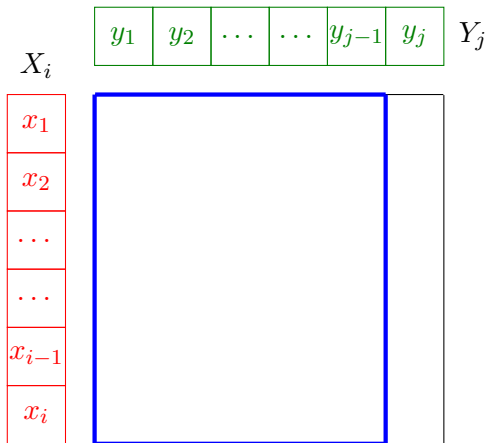
- $L(i, j) = L(i - 1, j - 1) + 1$: \nwarrow



Indicator Function

Indicator function $s(i, j)$ with values: $\nwarrow, \leftarrow, \uparrow$

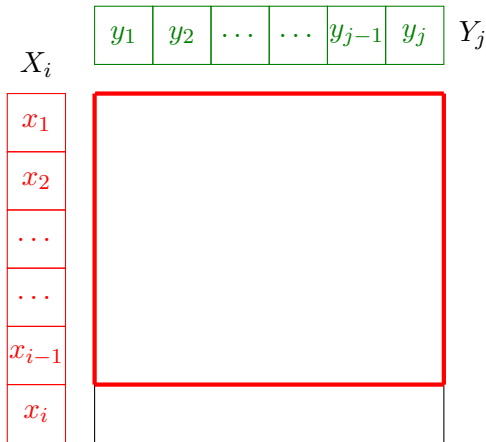
- $L(i, j) = L(i - 1, j - 1) + 1$: \nwarrow
- $L(i, j) = L(i, j - 1)$: \leftarrow



Indicator Function

Indicator function $s(i, j)$ with values: $\nwarrow, \leftarrow, \uparrow$

- $L(i, j) = L(i - 1, j - 1) + 1$: \nwarrow
- $L(i, j) = L(i, j - 1)$: \leftarrow
- $L(i, j) = L(i - 1, j)$: \uparrow



Pseudocode of LCS

Algorithm 5: $\text{LCS}(X[m], Y[n])$

```
1:  $L(i, 0) \leftarrow 0, i \in [m], L(0, j) \leftarrow 0, j \in [n];$ 
2: for  $i \leftarrow 1$  to  $m$  do
3:   for  $j \leftarrow 1$  to  $n$  do
4:     if  $X[i] = Y[j]$  then
5:        $L(i, j) = L(i - 1, j - 1) + 1, s(i, j) \leftarrow (\nwarrow)$ 
6:     else if  $L(i - 1, j) \geq L(i, j - 1)$  then
7:        $L(i, j) \leftarrow L(i - 1, j), s(i, j) \leftarrow (\uparrow)$  else
8:        $L(i, j) \leftarrow L(i, j - 1), s(i, j) \leftarrow (\leftarrow)$ 
9:     end
10:  end
11: end
```

Algorithm 6: TrackLCS(s, m, n)

Output: LCS of X and Y

```
1: while  $m \neq 0 \wedge n \neq 0$  do
2:   if  $s(m, n) = (\nwarrow)$  then
3:     output  $X[m]$ ;  $m = m - 1, n = n - 1$ , continue;
4:   end
5:   if  $s(m, n) = (\uparrow)$  then
6:      $m = m - 1$ , continue;
7:   end
8:   if  $s(m, n) = (\leftarrow)$  then
9:      $n = n - 1$ , continue;
10:  end
11: end
```

Demo of Indicator Function

$$X = (A, B, C, B, D, A, B), Y = (B, D, C, A, B, A)$$

	1	2	3	4	5	6
1	$s[1, 1] = \uparrow$	$s[1, 2] = \uparrow$	$s[1, 3] = \uparrow$	$s[1, 4] = \nwarrow$	$s[1, 5] = \leftarrow$	$s[1, 6] = \nwarrow$
2	$s[2, 1] = \nwarrow$	$s[2, 2] = \leftarrow$	$s[2, 3] = \leftarrow$	$s[2, 4] = \uparrow$	$s[2, 5] = \nwarrow$	$s[2, 6] = \leftarrow$
3	$s[3, 1] = \uparrow$	$s[3, 2] = \uparrow$	$s[3, 3] = \nwarrow$	$s[3, 4] = \leftarrow$	$s[3, 5] = \uparrow$	$s[3, 6] = \uparrow$
4	$s[4, 1] = \uparrow$	$s[4, 2] = \uparrow$	$s[4, 3] = \uparrow$	$s[4, 4] = \uparrow$	$s[4, 5] = \nwarrow$	$s[4, 6] = \leftarrow$
5	$s[5, 1] = \uparrow$	$s[5, 2] = \uparrow$	$s[5, 3] = \uparrow$	$s[5, 4] = \uparrow$	$s[5, 5] = \uparrow$	$s[5, 6] = \leftarrow$
6	$s[6, 1] = \uparrow$	$s[6, 2] = \uparrow$	$s[6, 3] = \uparrow$	$s[6, 4] = \nwarrow$	$s[6, 5] = \uparrow$	$s[6, 6] = \nwarrow$
7	$s[7, 1] = \uparrow$	$s[7, 2] = \uparrow$	$s[7, 3] = \uparrow$	$s[7, 4] = \uparrow$	$s[7, 5] = \uparrow$	$s[7, 6] = \uparrow$

$$\text{Solution: } \text{LCS} = (X[2], X[3], X[4], X[6]) = (B, C, B, A)$$

Complexity Analysis

Computation of optimized function

- Initialization: $O(m + n)$
- Computation: in each loop, require ≤ 2 times comparison, complexity is $\Theta(mn)$

Computation of indicator function

- Computation: $\Theta(mn)$
- Trace solution: $\Theta(m + n)$ (reduce the size of X or/and Y by 1 in each step)

Overall time complexity: $\Theta(mn)$

Space complexity: $\Theta(mn)$

Further Discussion

Standard LCS problem

- Dynamic programming: $\Theta(nm)$
- Generalized suffix tree: $\Theta(n + m)$

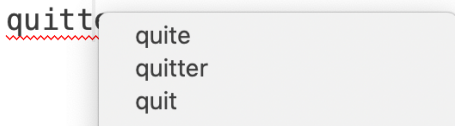
Generalized LCS problem: find LCS for k strings with length n_1, \dots, n_k

- k -dimension Dynamic programming: $\Theta(n_1 \cdots n_k)$
- Generalized suffix tree: $\Theta(n_1 + \cdots + n_k)$

- 1 Return on Investment
- 2 Knapsack Problem
 - Knapsack with Repetition
 - Knapsack without Repetition
- 3 Longest Common Substring
- 4 Edit Distance
- 5 Summary of Dynamic Programming

Motivation of String Similarity

When a spell checker encounters a possible misspelling, it looks in its dictionary for other words that are close by.



Q. What is the appropriate notion of closeness or similarity for two strings?

Edit Distance

Edit distance. [Levenshtein 1966, Needleman-Wunsch 1970]

Given two strings x and y , after a sequence of operations (replace, insert, delete), change y to x . The minimal number of operations is called the edit distance of between x and y , write as $\Delta(x, y)$.

capture similarity between two strings

Justify the definition: satisfy three rules of distance

- Non-negative: $\Delta(x, y) \geq 0$. $\Delta(x, y) = 0$ iff $x = y$
- Symmetric: $\Delta(x, y) = \Delta(y, x)$ (just reverse the operation)
- Triangle inequality: $\forall x, y, z, \Delta(x, z) + \Delta(z, y) \geq \Delta(x, y)$

How to Compute Edit Distance

Sequence alignment. A natural measure of edit distance is the extent to which they can be aligned, or matched up.

- an alignment is simply a way of writing the strings one above the other, allow adding \perp

S	\perp	N	O	W	Y
---	---------	---	---	---	---

S	U	N	N	\perp	Y
---	---	---	---	---------	---

1 mismatches, 2 gap

\perp	S	N	O	W	\perp	Y
---------	---	---	---	---	---------	---

S	U	N	\perp	\perp	N	Y
---	---	---	---------	---------	---	---

1 mismatches, 4 gap

- \perp indicates a “gap”: can be placed in either string — interpreting as delete or insert
- Cost of an alignment is the number of columns in which the letters differ.

$$\text{cost} = \underbrace{\sum_{x_i \neq y_i} \text{diff}(i, j)}_{\text{mismatch}} + \underbrace{\sum_{x_i \text{ unmatched}} \alpha + \sum_{y_j \text{ unmatched}} \beta}_{\text{gap}}$$

Insight of Edit Distance

Edit distance between two strings is the cost of their best alignment.

- Finding the edit distance is equivalent to finding the optimal alignment.

Edit distance is so named because it can also be thought of as the minimum number of *edits* — insertion, deletions, and substitutions — needed to transform the first string to the second.

- The above example: insert 'U', substitute 'O' \rightarrow 'N', and delete 'W'

In general, there are so many possible alignments between two strings \leadsto it would be terribly inefficient to search through all of them for the best one.

A Dynamic Programming Solution

When solving a problem by dynamic programming, the most crucial question is

What are the subproblems?

As long as they are chosen so as to have the **optimal substructure**, it is easy to write the algorithm: iteratively solve one subproblem after the other, in order of increasing order.

Goal. Finding the edit distance $E(m, n)$ between two strings $x[1 \dots m]$ and $y[1 \dots n]$.

Subproblem. Looking at the edit distance between some *prefix* of $x[1 \dots i]$ and some prefix of $y[1 \dots j]$, call the subproblem $E(i, j)$.

E	X	P	O	N	E	N	T	I	A	L
P	O	L	Y	N	O	M	I	A	L	

subproblem $E(7, 5)$

Structure of Problem

We need somehow express $E(i, j)$ in terms of smaller subproblems. Analyze the best alignment between $x[1 \dots i]$ and $y[1 \dots j]$: their rightmost column can only be one of three things:

$$\begin{array}{ccc} x[i] & \perp & x[i] \\ \perp & y[j] & y[j] \end{array}$$

Case 1a. leave x_i unmatched

- pay gap for x_i + min cost of aligning $x[i - 1]$ and $y[j]$.

Case 1b. leave y_j unmatched

- pay gap for y_j + min cost of aligning $x[i]$ and $y[j - 1]$.

Case 2. M matches $x_i - y_j$.

- pay (mis)match for $x_i - y_j$ + min cost of aligning $x[i - 1]$ and $y[j - 1]$.

optimal substructure property (proof via exchange argument)

Iteration Relation for Optimized Function

Optimized function: $E(i, j)$ — edit distance between $x[1, \dots, i]$ and $y[1, \dots, j]$

Initial values: $E(i, 0) = i, E(0, j) = j$

Iteration relation. We have expressed $E(i, j)$ in terms of three *smaller* subproblems $E(i-1, j), E(i, j-1), E(i-1, j-1)$.

- We have no idea which of them is the right one, so we need to try them all and pick the best

$$E(i, j) = \min\{1+E(i-1, j), 1+E(i, j-1), \text{diff}(i, j)+E(i-1, j-1)\}$$

$$\text{diff}(i, j) = \begin{cases} 0 & x[i] = y[j] \\ 1 & x[i] \neq y[j] \end{cases}$$

Computation Order

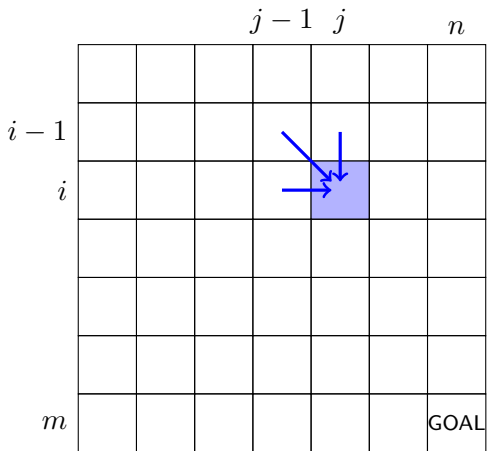
The answers to all the subproblems $E(i, j)$ form a 2-dimensional table.

What order should these subproblems be solved?

Any order is fine, as long as $E(i - 1, j)$, $E(i, j - 1)$ and $E(i - 1, j - 1)$ are handled before $E(i, j)$.

- 1 fill in the table one row at a time, from top row to bottom row, and moving left to right across each row
- 2 or fill in the table column by column

Both methods would ensure that by the time we get around to compute a particular table entry, all the other entries we need are already filled in.



Track Solution

	x	S	N	O	W	Y
y	0	1	2	3	4	5
S	1	0	1	2	3	4
U	2	1	1	2	3	4
N	3	2	1	2	3	4
N	4	3	2	2	3	4
Y	5	4	3	3	3	3

Track Solution

	x	S	N	O	W	Y
y	0	1	2	3	4	5
S	1	0	1	2	3	4
U	2	1	1	2	3	4
N	3	2	1	2	3	4
N	4	3	2	2	3	4
Y	5	4	3	3	3	3

$$E(5, 5) \leftarrow E(4, 4) + 0$$

Y

Y

Track Solution

	x	S	N	O	W	Y
y	0	1	2	3	4	5
S	1	0	1	2	3	4
U	2	1	1	2	3	4
N	3	2	1	2	3	4
N	4	3	2	2	3	4
Y	5	4	3	3	3	3

$$E(5, 5) \leftarrow E(4, 4) + 0$$

$$E(4, 4) \leftarrow E(4, 3) + 1$$

W Y

\perp Y

Track Solution

	x	S	N	O	W	Y
y	0	1	2	3	4	5
S	1	0	1	2	3	4
U	2	1	1	2	3	4
N	3	2	1	2	3	4
N	4	3	2	2	3	4
Y	5	4	3	3	3	3

$$E(5, 5) \leftarrow E(4, 4) + 0$$

$$E(4, 4) \leftarrow E(4, 3) + 1$$

$$E(4, 3) \leftarrow E(3, 2) + 1$$

O W Y

N ⊥ Y

Track Solution

	x	S	N	O	W	Y
y	0	1	2	3	4	5
S	1	0	1	2	3	4
U	2	1	1	2	3	4
N	3	2	1	2	3	4
N	4	3	2	2	3	4
Y	5	4	3	3	3	3

$$E(5, 5) \leftarrow E(4, 4) + 0$$

$$E(4, 4) \leftarrow E(4, 3) + 1$$

$$E(4, 3) \leftarrow E(3, 2) + 1$$

$$E(3, 2) \leftarrow E(2, 1) + 0$$

N O W Y

N N \perp Y

Track Solution

	x	S	N	O	W	Y
y	0	1	2	3	4	5
S	1	0	1	2	3	4
U	2	1	1	2	3	4
N	3	2	1	2	3	4
N	4	3	2	2	3	4
Y	5	4	3	3	3	3

$$E(5, 5) \leftarrow E(4, 4) + 0$$

$$E(4, 4) \leftarrow E(4, 3) + 1$$

$$E(4, 3) \leftarrow E(3, 2) + 1$$

$$E(3, 2) \leftarrow E(2, 1) + 0$$

$$E(2, 1) \leftarrow E(1, 1) + 1$$

⊥ N O W Y

U N N ⊥ Y

Track Solution

	x	S	N	O	W	Y
y	0	1	2	3	4	5
S	1	0	1	2	3	4
U	2	1	1	2	3	4
N	3	2	1	2	3	4
N	4	3	2	2	3	4
Y	5	4	3	3	3	3

$$E(5, 5) \leftarrow E(4, 4) + 0$$

$$E(4, 4) \leftarrow E(4, 3) + 1$$

$$E(4, 3) \leftarrow E(3, 2) + 1$$

$$E(3, 2) \leftarrow E(2, 1) + 0$$

$$E(2, 1) \leftarrow E(1, 1) + 1$$

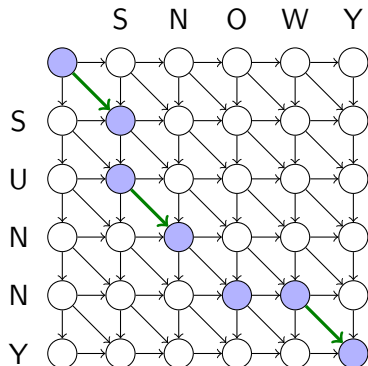
$$E(1, 1) \leftarrow E(0, 0) + 0$$

S ⊥ N O W Y

S U N N ⊥ Y

The Underlying DAG

Every dynamic programming has an underlying dag structure.



Set all edge lengths to 1 except the green ones

Final answer is the **shortest path** from $E(0,0)$ and $E(m,n)$

- move down: delete
- move right: insert
- move diagonal: match or substitution

By altering the weights on the DAG, we can allow generalized forms of edit distance: insertion, deletion, and substitution have different associated costs.

Algorithm 7: SequenceAlignment($x[m], y[n]$)

```
1: for  $i = 0$  to  $m$  do  $E(i, 0) = i$ ;  
2: for  $j = 0$  to  $n$  do  $E(0, j) = j$ ;  
3: for  $i = 1$  to  $m$  do  
4:   for  $j = 1$  to  $n$  do  
5:      $E(i, j) \leftarrow \min\{1 + E(i - 1, j), 1 + E(i, j -$   
         $1), \text{diff}(i, j) + E(i - 1, j - 1)\}$   
6:   end  
7: end  
8: return  $E(m, n)$ ;
```

There are totally mn subproblems, each subproblem requires constant time $\Rightarrow \Theta(mn)$ time and $\Theta(mn)$ space.

- 1 Return on Investment
- 2 Knapsack Problem
 - Knapsack with Repetition
 - Knapsack without Repetition
- 3 Longest Common Substring
- 4 Edit Distance
- 5 Summary of Dynamic Programming

How to Find Subproblems

Finding the right subproblems takes creativity and experimentation.

But there are a few standard choices that seem to arise repeatedly in dynamic programming.

One-Dimension Dynamic Programming

The input is x_1, x_2, \dots, x_n . A subproblem is x_1, x_2, \dots, x_i

x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8	x_9	x_{10}
-------	-------	-------	-------	-------	-------	-------	-------	-------	----------

The number of subproblems is therefore $O(n)$.

Examples

- shortest path in dag
- longest increasing subsequence
- max interval sum
- image compression

Two-Dimension Dynamic Programming: Type 1

The input is x_1, x_2, \dots, x_n . A subproblem is x_i, \dots, x_j

x_1 x_2 x_3 x_4 x_5 x_6 x_7 x_8 x_9 x_{10}

The number of subproblems is therefore $O(n^2)$.

Examples

- matrix multiplication chain
- optimal binary search tree

Two-Dimension Dynamic Programming: Type 2

The input is x_1, x_2, \dots, x_n and y_1, \dots, y_n . A subproblem is x_1, \dots, x_i and y_1, \dots, y_j

x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8	x_9	x_{10}
y_1	y_2	y_3	y_4	y_5	y_6	y_7	y_8		

The number of subproblems is therefore $O(mn)$.

Examples

- return on investment
- knapsack problem
- longest common substring
- edit distance

Another Important Characterization of DP

The computational complexity of DP algorithm not only depends on the number of subproblems, but also depends on the complexity of the recursive relation, i.e., the extent that a problem relates to its subproblems.

We capture such dependence as *locality*.

Case 1: depend on linear number of subproblems

- shortest path in DAG, longest increasing sequence, maximum interval sum, matrix multiplication chain, optimal binary search tree,

Case 2: depend on constant number of subproblems

- knapsack problem, longest common substring, edit distance

Greedy vs. Dynamic Programming

DP is mainly an optimization over plain recursion.

- Wherever we see a recursive solution that has repeated calls for the same inputs, we can optimize it using DP.
 - simply store the results of subproblems so that we do not have to re-compute them when needed later
- This simple optimization reduces time complexity from exponential to polynomial.
- **Example.** A simple recursive solution for Fibonacci numbers leads to exponential time complexity. But, if we optimize it by storing solutions of subproblems, time complexity reduces to linear.

We can think of Dynamic Programming as finding a DAG in a huge recursion tree (iterative approach) or travel the recursion tree with memo (recursive approach).

Greedy Algorithm vs. Dynamic Programming

Optimality: make choice seems best at the moment in the hope to obtain global optimal solution, rigorous proof is needed

Memorization: efficient in terms of memory as it never look back or revise previous choices

Fashion: computes its solution by making its choices in a serial forward fashion, never looking back or revising previous choices.

Optimality: make decision at each step considering current problem and solution to previously solved subproblem to calculate optimal solution; optimality is automatically guaranteed since DP actually considers **all possible cases** and then choose the best.

Memorization: requires DP table for memorization

Fashion: computes its solution by synthesizing from smaller optimal sub solutions: bottom up or top down