Lecture 4: Dynamic Programming Analysis and Design of Computer Algorithms

#### GONG Xiu-Jun

#### School of Computer Science and Technology, Tianjin University

Email: gongxj@tju.edu.cn

Website: http://cs.tju.edu.cn/faculties/gongxj/course/algorithm Discussion : http://groups.google.com/group/algorithm-practice-team

April 29, 2014

- 1 Aims and Objectives
- 2 What is DP
- **3** 0/1 Knapsack Problem
- 4 Matrix Multiplication Chains
- **5** All Pairs Shortest Path
- 6 Maximum Non-crossing Subset of Nets
- Congest Common Subsequences

- Aims and Objectives
  - Know the definition of DP
    - Optimal substructure and
    - Overlapped subproblems
  - prove DP applicable to a problem
  - Solving a real problem using DP
- Practice
  - 0/1 Knapsack
  - Matrix Multiplication Chains
  - All Pairs Shortest Paths
  - Maximum Non-crossing Subset nets
  - Longest Common Subsequence
  - Hidden Markov Model

### Motivations



What is DP

### An exmple of task sheduling

In a weighted staged directed graph G, we want to search the shortest path from start node s to the terminal node e.



## Definition

- The term dynamic programming was originally used in the 1940s by Richard Bellman (Bellman equation)
  - DP refers specifically to nesting smaller decision problems inside larger decisions.
  - Dynamic was chosen because it sounded impressive, not because it described how the method worked.
  - Programming referred to the use of the method to find an optimal programming.
- A method of solving complex problems by breaking them down into simpler steps (subproblems) in recursive manner.
  - the solution can be produced by combining solutions to subproblems;
  - the solution to each subproblem can be produced by combining solutions to sub-subproblems, etc;
  - the total number of subproblems arising recursively is polynomial.

# Where DP appliable?

- Optimal substructure(Principle of Optimization)
  - A problem is said to have optimal substructure if an optimal solution can be constructed efficiently from optimal solutions to its subproblems.
  - In other words, a problem exhibits optimal substructure if an optimal solution to the problem contains within it optimal solutions to subproblems.
  - Subproblems must be only 'slightly' smaller than the larger problem
    - A constant additive factor
    - A multiplicative factor: D&C
- Overlapping subproblems
  - A problem said to have overlapping subproblems if the problem can be broken down into subproblems which are reused several times
- DP  $\approx$  Recursion + Memorization

What is DP

## Procedures



1	Identifying subproblems	1-2 , 1-3 or 1-4 ?
2	Validating principle of optimization	What happen if not
3	Defining an optimal value function	Assumed that c(i) is the short- est distance from i to e
4	Deriving the recursive equation	$c(i) = \min_{j \in N(i)} \{ c(j) + cost(i,j) \}$
5	Solving the recursive equation	c(7) - c(6) - c(5) - c(4)- c(3) - c(2) -c(1)
6	Tracebacking the optimal solution	P(1) - P(3) -P(5) -P(7)

#### Probelm statements

Given a set of items, each with a mass and a value, determine the number of each item to include in a collection so that the total weight is less than or equal to a given limit(capacity) and the total value is as large as possible.

Formal representation: Given n > 0, c > 0 and two n-tuples of positive numbers:  $(w_1, w_2, \dots, w_n)$  and  $(p_1, p_2, \dots, p_n)$ .

we wish to determine the subset Maximize:  $\sum_{i \in T} p_i$  $T \subset \{1...n\}$  such that Subjects to:  $\sum_{i \in T} w_i \le c$ 

If used a n-tuple of indicator variables  $(x_1, x_2, \cdots, x_n)$ , then

Maximize: 
$$\sum_{i=1}^{n} x_i * p_i$$
  
Subjects to:  $\sum_{i=1}^{n} x_i * w_i \le c$ 

### Solutions

- Brute force: Try all 2<sup>n</sup> possible subsets T
- Greedy
  - choose items maximizing value ?
  - choose items maximizing value/size
  - the second looks much great, but what if items dont exactly fit (non-divisible items)?
- any others?
  - traceback
  - branch and bound
  - ...
- Dynamic programming is one of great choice

0/1 Knapsack Problem

## Step 1. Identifying subproblems



- By fewer items: take an item i without consideration , what changed?
  - the number of items n becomes n-1, and
  - the capacities become c or  $c w_i$ , why?
  - So the subproblem is constrained by two parameters.
- By smaller knapsack capacities, reduce the capacity, what changed?
  - the capacities become smaller c'
  - supposed that this is caused by removing an item  $i c' = c - w_c$  and thus the optimal value f(c) = f(c - v)

 $i,c' = c - w_i$ , and thus the optimal value  $f(c) = f(c - w_i) + v_i$ 

• however, we donn't know which *i*,

• 
$$f(c) = \max_{i=1,w_i < c} f(c - w_i) + v_i$$



0/1 Knapsack Problem

## Step 2. Validating Principle of Optimization



[Optimal substructure] A problem exhibits optimal substructure if an optimal solution to the problem contains within it optimal solutions to subproblems.

- Supposed that (x<sub>1</sub>, x<sub>2</sub>, · · · , x<sub>n</sub>) is the optimal solution of original problem P<sup>0</sup>.
- ► (y<sub>1</sub>, y<sub>2</sub>, · · · , y<sub>n-1</sub>) is the optimal solution of its subproblem P' with items 1, 2, · · · , n 1 and capacities c x<sub>n</sub> \* w<sub>n</sub>.
- We should show that  $(y_1, y_2, \dots, y_{n-1}, x_n)$  is no worse than  $(x_1, x_2, \dots, x_n)$  to  $P^0$
- How? using a "cut-and-paste" technique (homework)



We have known that a subproblem is constrained by two parameters: number of items and capacities.

- ► For a given subproblem with items i, i + 1, ..., n and capacities y, we define its optimal value by f(i, y).
- f(n,0) = 0, and  $f(n, w_n) = p_n$  if  $w_n < c$  else 0
- f(1, c) is our goal, why?
- How to get f(1, c) from f(n, y) ?

0/1 Knapsack Problem

### Step 4. Deriving the recursive equation



To calculate f(i, y) from f(i + 1, y), we only need know if item *i* could be included

- ▶ If  $w_i > y$ , item *i* couldn't be included, f(i, y) = f(i + 1, y).
- If  $w_i \leq y$ , item *i* might be included
  - if so,  $f(i, y) = f(i + 1, y w_i) + p_i$
  - else, f(i, y) = f(i + 1, y), why?
  - finally, we use the larger one.

### Step 5. Solve recursive equation

Now we have the recursive equation

$$f(i,y) = \begin{cases} f(i+1,y) & \text{if } 0 \le y < w_i \\ \max \begin{cases} f(i+1,y-w_i) + p_i \\ f(i+1,y) \end{cases} & \text{if } y \ge w_j \end{cases}$$
(1)

and the initial condition

$$f(n, y) = \begin{cases} p_n & \text{if } y \ge w_n \\ 0 & \text{if } 0 \le y < w_n \end{cases}$$
(2)

How to solve them?

- Recursive version
- Non-recursive version
- Tuple version

#### 0/1 Knapsack Problem

## Recursive version

Algorithm 1: RKnapsack

Algorithm complexity: We use t(n) to denote the algorithm time complexity with n items

- At the best case(step4), t(n) = t(n − 1) + b, so t(n) = Θ(n)
- At the worst case(step6), t(n) = 2t(n − 1) + b, so t(n) = Θ(2<sup>n</sup>)



#### Figure : Recursive Call Relation Tree



If these repetitive calls are saved, the algorithm complexity should be reduced!

#### First attempt: using array

```
template<class T>
 void Knapsack(T p[], int w[], int c, int n, T** f)
[] {// Compute f[i][y] for all i and y.
     // initialize f[n][]
     int vMax = min(w[n]-1,c);
     for (int v = 0; v \le vMax; v \leftrightarrow)
        f[n][v] = 0;
     for (int y = w[n]; y <= c; y++)</pre>
        f[n][y] = p[n];
     // compute remaining f's
    for (int i = n - 1; i > 1; i - 1) {
        yMax = min(w[i]-1,c);
       for (int y = 0; y \le y Max; y \leftrightarrow y
           f[i][y] = f[i+1][y];
        for (int v = wfil: v \leq c: v++)
           f[i][y] = max(f[i+1][y],
                          f[i+1][v-w[i]] + p[i]);
     f[1][c] = f[2][c];
     if (c \ge w[1])
      f[1][c] = max(f[1][c], f[2][c-w[1]] + p[1]);
```

#### Comments

- The algorithm saves all possible values using an array f, so that every f(i, y) is calculated only once.
- It needs Θ(nc) extra space.
- Its time complexity is Θ(nc).
  - It is not polynomial: to describe c need log<sub>2</sub> c bits
  - but pseudo-polynomial: exponential dependence on numerical inputs
- Its disadvantage
  - The capacity c must an integer
  - The complexity might still be very hight when c is large enough, for instance c = 2<sup>n</sup>

### Second attempt: using a tuple



- Calculate f values using array
  - $f(5,0) = 0, \dots, f(5,4) = 6, \dots, f(5,10) = 6$
  - $f(4,0) = 0, \dots, f(4,4) = 6, \dots, f(4,9) = 10$ , f(4,10) = 10•  $\dots$

Save step points only for each i

- Define a tuple (a, b), where a = y and b = f(i, y)
- (*a*, *b*) corresponds an optimal loading with capacity *a* and value *b*
- Put all tuples into a set P<sub>i</sub>
- $P_n$  can be got easily.  $P_n = \{(0,0), (w_n, p_n)\}$
- *P*<sub>1</sub> contains our goal! Why?

## Tuple method: principle

• Let 
$$Q = \{(s, t) | w_i \le s < c, (s - w_i, t - p_i) \in P_{i+1}\}$$

- Q corresponds  $P_i$  in which item *i* has been selected
- $P_{i+1}$  corresponds  $P_i$  in which item *i* has not been selected
- thus,  $P_i = Q \bigcup P_{i+1}$
- Merge Q and  $P_{i+1}$  to get  $P_i$ 
  - remove dominated tuples ( a tuple (a, b) is dominated by (u, v) if a > u, but b > v)
  - remove repeated tuples
  - remove over capacity tuples in which a > c
- ► The complexity is still O(2<sup>n</sup>). The number of elements in P<sub>i</sub> increase exponentially at worst case

Dynamic Progra	amming		0/1 Knap	osack Prob	lem					
Tuple method:an example										
	n=5, c=10 p= w= X	$ \begin{array}{c} 6 & 3 \\ 2 & 2 \\ x_1 & x_2 \end{array} $	5 6 $2$ $x_3$	4 5 ×4	6 4 x <sub>5</sub>					
1. P(5)=[(0,0)	(4,6)]		Q=[(5	5,4),(9	,10)]					
2. P(4)=[(0,0)	( <b>4,6)</b> ,(9,10	)]	Q=[(6	6,5),(1	0,11)]					
$\begin{array}{c} P(3) = [(0, 0)] \\ 3. \\ (10, 11)] \end{array}$	), <b>(4, 6)</b> , (9	9, 10),	Q=[(2	2,3),(6	,9)]					
4. P(2)=[(0,0) (6,9), (9,10)	(2,3), ), (10,11)]	(4,6),	Q=[(2 (8,15)	2,6), ]	(4,9),	(6,12),				
$5. \begin{array}{c} P(1) = [(0,0)] \\ (6,12), \ \textbf{(8,1)} \end{array}$	, (2,6), <b>5)</b> ]	(4,9),								

The optimal value is 15 and the optimal solution is [1,1,0,0,1] by tracebacking

GONG Xiu-Jun

Lecture 4: Dynamic Programming

## Problem Description

- Two matrices A×B with dimension (m,n) and (n,q) takes mnq multiplications.
- Let's consider three matrices:  $A \times B \times C$  with dimension (100,1), (1,100) and (100,1). The product can be done:
  - $(A \times B) \times C$  takes 20000 multiplications.
  - $A \times (B \times C)$  takes 200 multiplications.
  - Any other way? No
    - Associative law :  $(A \times B) \times C = A \times (B \times C)$
    - Commutative law:  $A \times B \neq B \times A$
  - the order of multiplications makes a big difference in the final running time!
- Matrix Multiplication Chains. Suppose that we multiply q matrices
  - $M(1,q) = M_1 \times M_2 \times \cdots \times M_q$
  - There are q + 1 ( $r_1, r_2, \dots, r_q, r_{q+1}$ )numbers(parameters) to constrain their dimensions. Why?

## Identifying subproblems

 The number of orders of multiplications equal to the number of parenthesization

$$P(q) = \begin{cases} 1 & \text{if } q = 2 \ T(q) \approx O(4^{q} q^{\frac{3}{2}}) \\ \sum_{k=1}^{q-1} p(k) p(q-k) & \text{if } q \ge 2 \end{cases} \approx O(2^{q})$$

- In fact, a particular parenthesization can be represented by a binary tree
  - The leaves corresponds to matrices
  - The root corresponds to the final product
  - Interior nodes are intermediate products
  - If a tree to be optimal, its subtrees must also be optimal.
- The binary tree representation suggests that
  - A subtree corresponds to a subproblem in MPC , so
  - MPC satisfies the principle of optimization

## Define the optimal function

Define c(i,j) is the cost of  $M_i \times M_{i+1} \times \cdots \times M_j$ 

- c(0,0) = 0
- c(i,i) = 0

• 
$$c(i, i+1) = r_i * r_j * r_{j+1}$$

- ► c(1, q) is our goal
- The recursive equation is

$$c(i,j) = \min_{i \le k < j} \{ c(i,k) + c(k,j) + r_k * r_j * r_{j+1} \}$$
(3)

- Let kay(i,j) is the number minimizing above equation.
  - kay(i,i)=0
    kay(i,i+1)=i

### An exmple

Suppose that q =5 and r =(10 , 5 , 1 , 10 , 2 , 10), find the optimal orders of multiplications

	$M(1,5) = M(1,2) \times M(3,5)$	$M(3,5) = M(3,4) \times M(5,5)$
c(1,5)=min	${c(1,1)+c(2,5)+500, c(1,2)+c(3,5)+100, }$	c(1,5)=190, kay(1,5)=2
	c(1,3)+c(4,5)+1000, $c(1,4)+c(5,5)+200\}$	c(1,3)=150, kay(1,3)=2 c(1,4)=90, kay(1,4)=2
c(2,5)=min	$ \begin{array}{l} \{c(2,2)+c(3,5)+50,\\ c(2,3)+c(4,5)+500,\\ c(2,4)+c(5,5)+100\} \end{array} $	c(2,5)=90, kay(2,5)=2
c(3,5)=min	$ \begin{array}{l} \{c(3,3)+c(4,5)+100,\\ c(3,4)+c(5,5)+20\}\\ min\{300,40\}=\!40 \end{array} $	c(3,5)=40,kay(3,5)=4
c(2,4)=min	$ \begin{array}{l} \{c(2,2)+c(3,4)+10,\\ c(2,3)+c(4,4)+100\}\\ \min\{30,150\}=30 \end{array} $	c(2,4)=30,kay(2,4)=2

Matrix Multiplication Chains

#### Recursive Algorithm for MPC

```
1
    int RC(int i, int j)
2
    \{// Return c(i, i) and compute kay(i, i)=kay[i]
          ][i].
 3
    // Avoid recomputations, check if already
          computed
        if (c[i][j] > 0) return c[i][j];
4
 5
    // c[i][j] not computed before, compute now
 6
       if (i == i) return 0: // one matrix
7
       if (i == i - 1) {// two matrices
8
                          kav[i][i+1] = i;
                                                        1
                                                           void Traceback(int i, int j, int
 9
                          c[i][j] = r[i]*r[i+1]*r[
                                                                  **kay)
                                i+21:
                                                           ſ
                                                        2
10
                          return c[i][i];}
                                                        3
                                                               if (i == j) return;
11
    // more than two matrices
                                                        4
                                                               Traceback(i, kay[i][j], kay);
12
    // set u to mini term for k = i
                                                        5
                                                               Traceback(kay[i][j]+1, j, kay
13
        int u = RC(i,i) + RC(i+1,j) + r[i]*r[i
                                                                    ):
             +1] *r[i+1];
                                                        6
                                                               cout << "Multiply...M..." << i <<</pre>
14
        kav[i][i] = i:
                                                                      ","" << kay[i][j];
15
    // compute remaining min terms and update u
                                                        7
                                                               cout << "uanduMu" << (kay[i][
16
        for (int k = i+1; k < j; k++) {
                                                                    i]+1) << ",..." << i
17
           int t = RC(i,k) + RC(k+1,i) + r[i]*r[k]
                                                        8
                                                                    << endl:
                 +1] *r[j+1];
                                                        9
                                                           }
18
           if (t < u) {// smaller min term
19
                        \mathbf{u} = \mathbf{t};
20
                        kay[i][j] = k;
21
           3
22
        c[i][j] = u;
23
       return u:
24
    }
```

### Revision recursive algorithm



The above figure suggest that c(i,j) can be calculated in an iterative miner

$$c(i, i+s) = \min_{i \le k < s} \{c(i, k) + c(k, j) + r_i * r_k * r_{i+s+1}\}$$
  
s = 1, 2, \dots, q

#### Iterative algorithm for MPC

```
void MatrixChain(int r[], int q, int **c, int **kay)
 1
    {//Compute costs and kay for all Mij's.
 2
    //initialize c[i][i],c[i][i+1],and kay[i][i+1]
 3
 4
       for (int i = 1; i < q; i++) {</pre>
 5
           c[i][i] = 0;
 6
          c[i][i+1] = r[i]*r[i+1]*r[i+2];
 7
           kav[i][i+1] = i:
 8
           }
 9
       c[q][q] = 0;
    //compute remaining c's and kay's
10
11
       for (int s = 2; s < q; s++)
12
           for (int i = 1; i \le q - s; i++) {
13
              // min term for k = i
              c[i][i+s] = c[i][i] + c[i+1][i+s]
14
15
                           + r[i]*r[i+1]*r[i+s+1];
16
              kav[i][i+s] = i;
17
              // remaining mini terms
18
              for (int k = i+1; k < i + s; k++) {
19
                 int t = c[i][k] + c[k+1][i+s]
                         + r[i]*r[k+1]*r[i+s+1]:
20
21
                 if (t < c[i][i+s]) \{// smaller mini term
22
                    c[i][i+s] = t;
23
                    kav[i][i+s] = k;
24
                 }
25
              3
26
    3
```

#### Problem statement

▶ Input: Given a directed graph *G* = (*V*, *E*) and a matrix (*a*<sub>*ij*</sub>) where

$$V = \{1, 2, \cdots, n\} \text{with edge weight function } W : E \to R$$
$$a_{ij} = \begin{cases} w(i, j) & \text{if } (i, j) \in E \\ 0 & \text{if } i = j \\ \infty & \text{otherwise} \end{cases}$$

- **Output**: A  $n \times n$  matrix of shortest-path lengths c(i, j)
- Assumption: No negative-weight cycles

#### Subprobem identification by edges

- Define d<sup>m</sup><sub>ij</sub> = weight of a shortest path from i to j that only uses <u>at most</u> m edges
- ▶ d<sup>n−1</sup><sub>ij</sub> is our goal
- We have known that

• 
$$d^0_{ij}=0$$
 if  $i=j$  , and  $\infty$  if  $i
eq j$ 

• 
$$d_{ij}^{1} = 0$$
 if  $i = j$ , and  $a_{ij}$  if  $i \neq j$ 

#### Theorem

For 
$$m = 1, 2, \cdots, n-1$$
, we have

$$d_{ij}^{m} = \min_{1 \le k < m} \{ d_{ik}^{m-1} + a_{kj} \}$$

## Proof



Running time  $O(n^4)$  - similar to n runs of Bellman-Ford algorithm

## Subprolems identification by intermediate vertices

- Define d<sup>m</sup><sub>ij</sub> = weight of a shortest path from i to j that only uses intermediate vertices from set {1, , m} or
- Define d<sup>m</sup><sub>ij</sub> = weight of a shortest path from i to j that the orders of intermediate vertices is no larger than m
- d<sup>n</sup><sub>ij</sub> is our goal
- We have known that

• 
$$d_{ij}^0 = aij$$
  
•  $d_{ik}^{k-1} = d_{ik}^k$   
•  $d_{jk}^{k-1} = d_{jk}^k$ 

#### Theorem

For 
$$m = 1, 2, \cdots$$
, n-1, we have

$$d_{ij}^m = \min\{d_{ij}^{m-1}, d_{im}^{m-1} + d_{mj}^{m-1}\}$$

### Proof



Running time  $O(n^3)$  Known as Floyd-Warshall algorithm

#### Iterative algorithm for ASAP

```
template < class T>
 1
 2
     void AdjacencyWDigraph<T>::AllPairs(T **c,
          int **kay)
                                                       1
     {// All pairs shortest paths.
 3
 4
      // Compute c[i][j] and kay[i][j] for all i
                                                       2
            and i.
        // initialize c[i][j] = c(i, j, 0)
                                                       3
 5
        for (int i = 1: i <= n: i++)
 6
                                                       4
 7
           for (int j = 1; j <= n; j++) {</pre>
 8
               c[i][i] = a[i][i];
                                                       5
 9
              kav[i][i] = 0;
10
                                                       6
11
        for (int i = 1; i \le n; i++)
12
           c[i][i] = 0;
                                                       7
                                                           }
13
                                                       8
14
        // compute c[i][j] = c(i, j, k)
                                                       9
15
        for (int k = 1; k \le n; k++)
                                                      10
16
           for (int i = 1: i <= n: i++)
17
               for (int j = 1; j <= n; j++) {</pre>
                                                      11
18
                 T t1 = c[i][k];
                 T t2 = c[k][j];
19
                                                      12
20
                 T t3 = c[i][j];
21
                  if (t1 != NoEdge && t2 != NoEdge 13
                         & &
                                                      14
22
                     (t3 == NoEdge || t1 + t2 < t3
                           )) {
                                                      15
23
                       c[i][i] = t1 + t2;
                                                      16
                       kay[i][j] = k;
24
                                                      17
25
                  }
                                                      18
26
                                                      19
```

 ${\rm Kay}[i][j]$  is used to store the largest vertex in the shortest path from i to j, so that traceback the shortest path

```
void outputPath(int **kay, int i,
     int i)
{// Actual code to output i to j
     path.
   if (i == i) return:
   if (kay[i][j] == 0) cout << j <<</pre>
         ·..·;
   else {outputPath(kay, i, kay[i][j
        1):
         outputPath(kay, kay[i][j],
               i);}
template < class T>
void OutputPath(T **c, int **kay, T
     NoEdge.
                           int i, int
{// Output shortest path from i to j
   if (c[i][j] == NoEdge) {
      cout << "There_is_no_path_from</pre>
           ..." << i << "...to..."
            << j << endl;
      return: }
   cout << "The_path_is" << endl;</pre>
   cout << i << '''';
   outputPath(kay,i,j);
```

#### An example

	Adjacent matrix				Distance matrix of shortest path					matrix kay <sub>ij</sub>					
0	1	2	3	4	0	1	4	4	8	0	0	2	3	4	
3	0	1	2	3	3	0	1	5	9	0	0	0	3	4	
2	2	0	1	2	2	2	0	1	8	0	0	0	0	4	
5	5	3	0	1	8	8	9	0	1	5	5	5	0	0	
4	4	2	3	0	8	8	2	9	0	3	3	0	3	0	

- we can traceback the shortest paths from matrix (kay<sub>ij</sub>), for example, from 1 to 5
  - kay(1,5)=4 , kay(4,5) =0 , 4  $\rightarrow$  5 is a sub-path
  - kay(1,4)=3, kay(3,5)=0,  $3 \rightarrow 4$  is a sub-path
  - kay(1,3)=2, kay(2,3)=0,  $2 \rightarrow 3$  is a sub-path
  - kay(1,2)=0 ,  $1 \rightarrow 2$  is a sub-path
  - The final shortest path is  $1 \rightarrow 2 \rightarrow 3 \rightarrow 4 \rightarrow 5$

## Problem statement



• A 1-1 map  $(i, c_i)$  is called a subnet

- ► Two subnets (i, c<sub>i</sub>) and (j, c<sub>j</sub>) are non-crossed if i < j then c<sub>i</sub> < c<sub>j</sub>
- ▶ The set  $MNS(i,j) = \{(u, c_u) | u \le i, c_u \ne j\}$  is called non-crossing set if for  $\forall (p, c_p)$  and  $\forall (q, c_q) \in MNS(i, j)$  then  $(p, c_p)$  and  $(q, c_q)$  are non-crossed
- Our goal is to find a MNS(n,n) with the maximum number of elements

#### Derive the recursive equation

- Our goal is to maximize size(n,n)
- We have known that

$$size(1,j) = \begin{cases} 0 & \text{if } j < c_1 \\ 1 & \text{if } j \ge c_1 \end{cases}$$
(4)

- We want to know the relationship between size(i,j) and size(i-1,j)
  - if  $j < c_i$  then  $(i, c_i) \notin MNS(i 1, j)$  , thus size(i,j)=size(i-1,j)
  - if  $j \ge c_i$  , there are two cases
    - put (i, c<sub>i</sub>) into MNS(i-1,j), but (i, c<sub>i</sub>) might cross with items in MNS(i-1,j) or result in smaller size, we have size(i,j)=size(i-1,j)
    - put (i, c<sub>i</sub>) into MNS(i-1,j), no crossing, then c<sub>i-1</sub> must be less than c<sub>i</sub> - 1, else crossed with (i, c<sub>i</sub>), thus we have size(i,j)=size(i-1,c<sub>i</sub> -1)+1
    - finally, we choose the maximum of them!

$$size(i,j) = egin{cases} size(i,j-1) & ext{if } j < c_i \ \max\{size(i-1,j), size(i-1,c_i-1)+1\} & ext{if } j \geq c_i \end{cases}$$

### Problem statements

#### Definition 1: Subsequence

Given a sequence  $X = x_1 x_2 \cdots x_m$ , another sequence  $Z = z_1 z_2 \cdots z_k$  is a subsequence of X if there exists a strictly increasing sequence  $x_1 x_2 \cdots x_k$  of indices of X such that for all j=1,2,...k, we have  $x_{i_i} = z_j$ .

Example 1: If X=abcdefg, Z=abdg is a subsequence of X.

#### Definition 2: Common subsequence

Given two sequences X and Y, a sequence Z is a common subsequence of X and Y if Z is a subsequence of both X and Y.

Example 2: X=abcdefg and Y=aaadgfd. Z=adf is a common subsequence of X and Y

## Definitions

#### Definition 3: Longest common subsequence:LCS

A longest common subsequence of X and Y is a common subsequence of X and Y with the longest length.

- Longest common subsequence may not be unique, for example, strings both acd and abd are LCS of abcd and acbd
- LCS has been successfully applied to many fields
  - Bioinformatics, e.g. long preserved regions in genomes
  - file comparison, e.g. diff
- Solutions
  - Brute force approach:  $O(n2^m)$
  - DP approach: O(nm)

## DP approach for LCS

- Let's consider the prefixes of x and y
  - x[1..i] ith prefix of x[1..m]
  - y[1..j] jth prefix of y[1..n]
- ► Subproblem: define c[i,j] = |LCS(x[1..i], y[1..j])|
- c[m,n] is our goal
- We have known that c[1,1]=1 if  $x_1 = y_1$  otherwise 0
- ► To derive recursive relationship, we use the following theorem Theorem: Let X=x1x2···xm and Y=y1y2···yn be two sequences, and Z=z1z2···zk be any LCS of X and Y, then
- If x<sub>m</sub> = y<sub>n</sub>, then z<sub>k</sub> = x<sub>m</sub> = y<sub>n</sub> and Z[1..k-1] is an LCS of X[1..m-1] and Y[1..n-1].
- ② If  $x_m \neq y_n$ , then  $z_k \neq x_m$  implies that Z is an LCS of X[1..m-1] and Y.
- If  $x_m \neq y_n$ , then  $z_k \neq y_n$  implies that Z is an LCS of X and Y[1..m-1].

#### **Recursive equation**

By the theorem, we can easily get the recursive equation

$$c[i,j] = \begin{cases} 0 & \text{if } i=0 \text{ or } j=0\\ c[i-1,j-1]+1 & \text{if } x[i]=y[j]\\ \max\{c[i-1,j],c[i,j-1]\} & \text{otherwise} \end{cases}$$
(6)

```
LCS(X,Y,m,n,b)
 1
 2
    for i=1 to m do
                                                         PrintLCS(b,X,i,j)
                                                     1
 3
             c[i, 0] = 0;
                                                     2
                                                         i = m
 4
    for j=0 to n do
                                                     3
                                                         j=n;
 5
             c[0, i]=0;
                                                     4
                                                         if i==0 or i==0 then exit:
 6
    for i=1 to m do
                                                     5
                                                         if b[i,j]==1 then
 7
             for j=1 to n do
                                                     6
                                                         ſ
 8
             {//b[i,j] stores the directions.
                                                     7
                                                                   i = i - 1;
                                                     8
 9
             if x[i] ==v[j] then
                                                                   j = j - 1;
10
                      c[i,j]=c[i-1,j-1]+1;
                                                     9
                                                                   print x[i];
                     b[i,j]=1; //1-diagonal,
11
                                                    10
                                                       }
             else if c[i-1,j]>=c[i,j-1] then
12
                                                   11 if b[i,j]==2 i=i-1
                              c[i,j]=c[i-1,j]
13
                                                   12 if b[i,j]==3 j=j-1
14
                              b[i, j]=2;//2-up,
                                                   13
                                                       Goto Step 3.
15
                      else c[i,j]=c[i,j-1]
16
                               b[i,j]=3; //3-
                                                                       Print LCS algorithm
                                     forward.
17
             }
```

