

Data Mining and Machine Learning

Assignment Project Exam Help

Sequence Analysis & Dynamic Programming

<https://powcoder.com>

Add WeChat powcoder

Peter Jančovič



Objectives

- To consider data mining for sequential data
- To understand Dynamic Programming (DP)
- Using DP to compute distance between sequences
- To understand what is meant by:
 - An alignment path
 - The DP recurrence equation
 - The distance matrix
 - The accumulated distance matrix
 - The optimal path



Sequences

- Sequences are common in real applications:
 - DNA analysis in bioinformatics and forensic science
 - Sequences of the letters A, G, C and T
 - Signature recognition biometrics
 - Words and text
 - Spelling and grammar checkers, author verification,...
 - Speech, music and audio
 - Speech/speaker recognition, speech coding and synthesis
 - Electronic music
 - Radar signature recognition...



Mining sequential data

- Sequences may not be amenable to human interpretation (complexity, dimension, quantity)
- Need for automated sequential data retrieval/mining
- For clustering and other tools, the fundamental requirement is for a measure of the distance between two sequences

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



Basic definitions

- In a typical sequence analysis application we have a basic alphabet consisting of N symbols

Assignment Project Exam Help

$$A = \{\alpha_1, \dots, \alpha_n, \dots, \alpha_N\}$$

<https://powcoder.com>

- Examples:
 - In text A is the set of letters plus punctuation plus ‘white space’
 - Bioinformatics $A = \{A, G, C, T\}$ (elements of DNA sequences)

Add WeChat powcoder



Sequences of continuous variables

- In some applications, elements of a discrete sequence are taken from a continuous vector space, rather than a finite set
- Sequences of continuous variables can be dealt with in two ways:
 - Directly
 - Vector quantization (VQ):
 - Represent space as a set of K centroids:
 - Replace each data point by its closest centroid



Distance between sequences (1)

- Sequences from the alphabet $\mathbf{A} = \{A, B, C, D\}$
- How similar are the sequences:
 - $S_1 = ABCD$
 - $S_2 = ABD$
- Intuitively S_2 is obtained from S_1 by deleting C
- Alternatively S_2 is obtained from S_1 by substituting D for C and then deleting D



Distance between sequences (2)

- Or S_2 was obtained from S_1 by deleting ABCD and inserting ABC
- ... **Assignment Project Exam Help**
- First explanation is intuitively ‘correct’, final explanation is intuitively ‘wrong’. Why?
- We favour the simplest explanation, involving the minimum number of insertions, deletions and substitutions
- ...but maybe not always



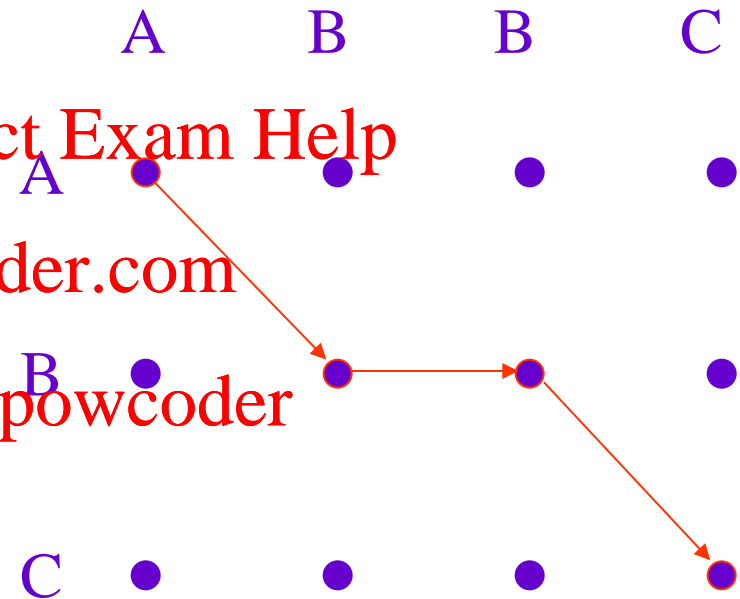
Distance between sequences (3)

- Consider:
 - $S_1 = AABC$
 - $S_2 = SABC$
 - $S_3 = PABC$
 - $S_4 = ASCB$
- If these sequences were typed then maybe S_2 is closer to S_1 than S_3 is, because A and S are adjacent on a keyboard
- Similarly S_4 is close to S_2 because letter-swapping ($SA \rightarrow AS$ etc) is a common typographical error



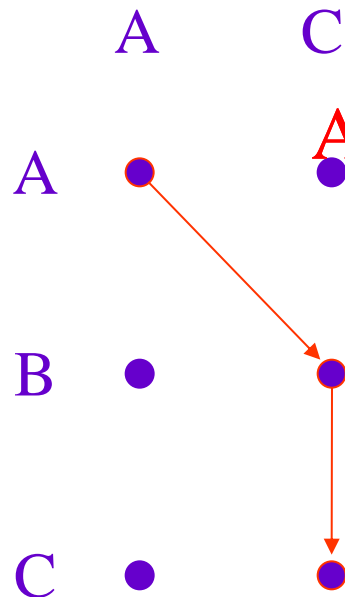
Alignments

- Relationship between two sequences can be expressed as an alignment between their elements
- Insertion (w.r.t. ABC) is a horizontal step

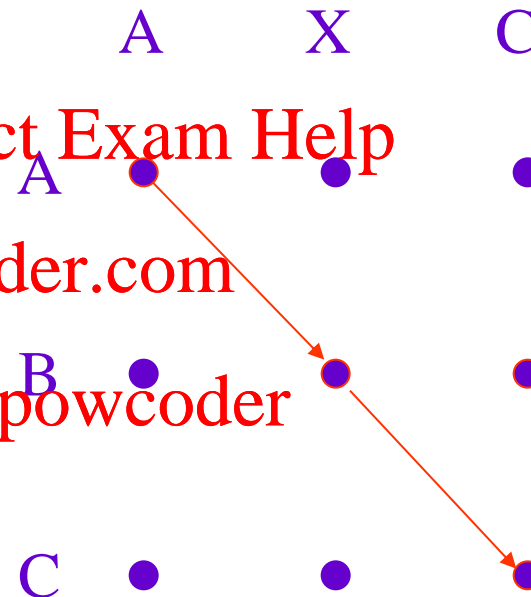


Alignment: deletion and substitution

Deletion is a
vertical step



Substitution or perfect
alignment are diagonal steps



Assignment Project Exam Help

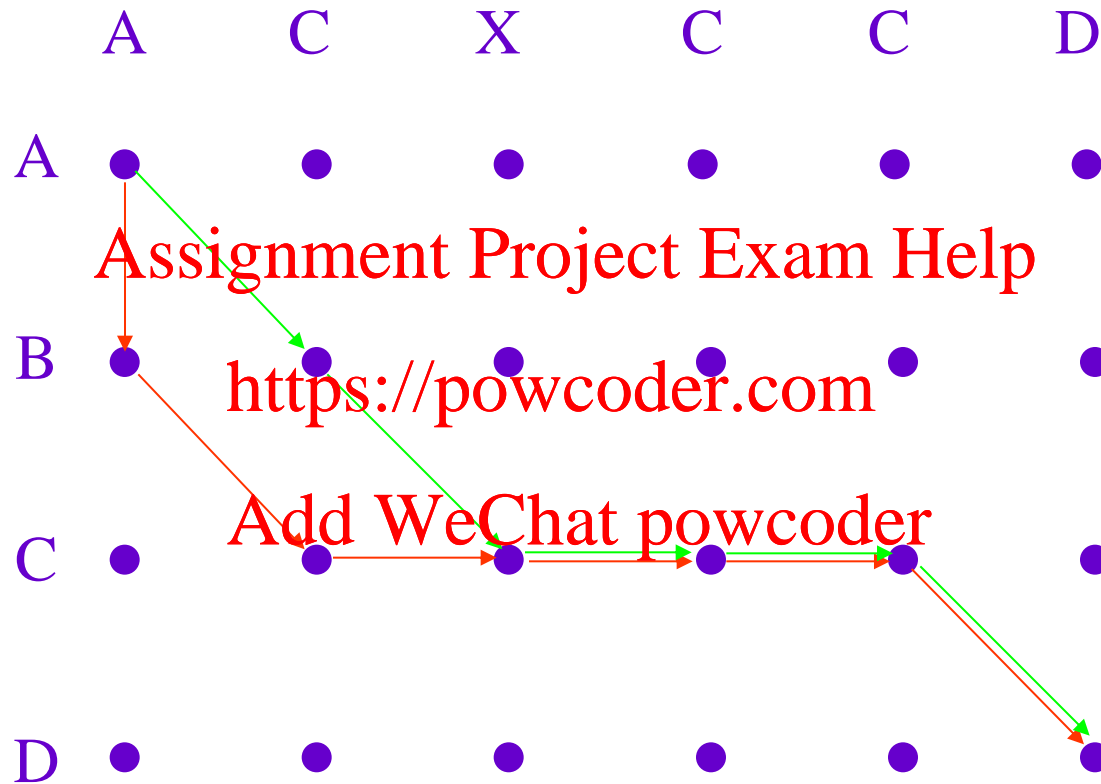
<https://powcoder.com>

Add WeChat powcoder

N.B: Edits described relative to vertical string



Alternative alignment paths



Which alignment path is best?



The Distance Matrix

- Let d be a metric, so $d(A,B)$ is the distance between the alphabet symbols A and B
- Examples:
 - $d(A,B) = 0$ if $A = B$, otherwise $d(A,B) = 1$
 - In typing, $d(A,B)$ might indicate how unlikely it is that A would be mistyped as B
 - For continuous valued sequences d could be Euclidean distance, or City Block distance, or L_∞ distance



Notation

- Suppose we have an alphabet:

$$\mathbf{A} = \{\alpha_1, \dots, \alpha_n, \dots, \alpha_N\}$$

- The distance matrix for \mathbf{A} is an $N \times N$ matrix

$$D = [D_{m,n}], \quad 1 \leq m, n \leq N$$

Add WeChat powcoder

where $D_{m,n} = d(\alpha_m, \alpha_n)$ is the distance between the m^{th} and n^{th} alphabet symbols

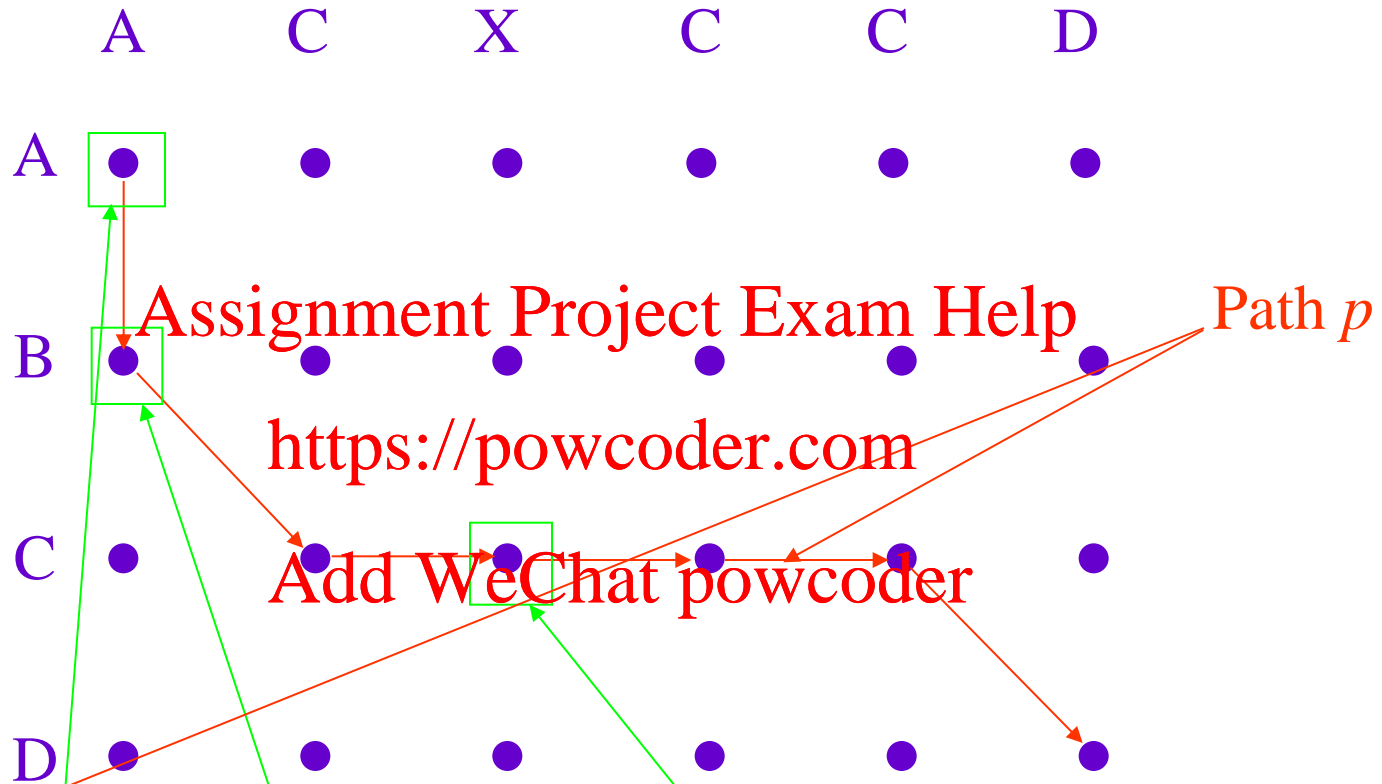


The Accumulated Distance

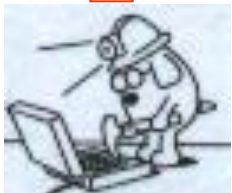
- Consider two sequences:
 - $S_1 = ABCD$
 - $S_2 = ACXCCD$
- For an alignment path p between S_1 and S_2 the accumulated distance between S_1 and S_2 , denoted by $AD_p(S_1, S_2)$, is the sum over all the nodes of p of the corresponding distances between elements of S_1 and S_2



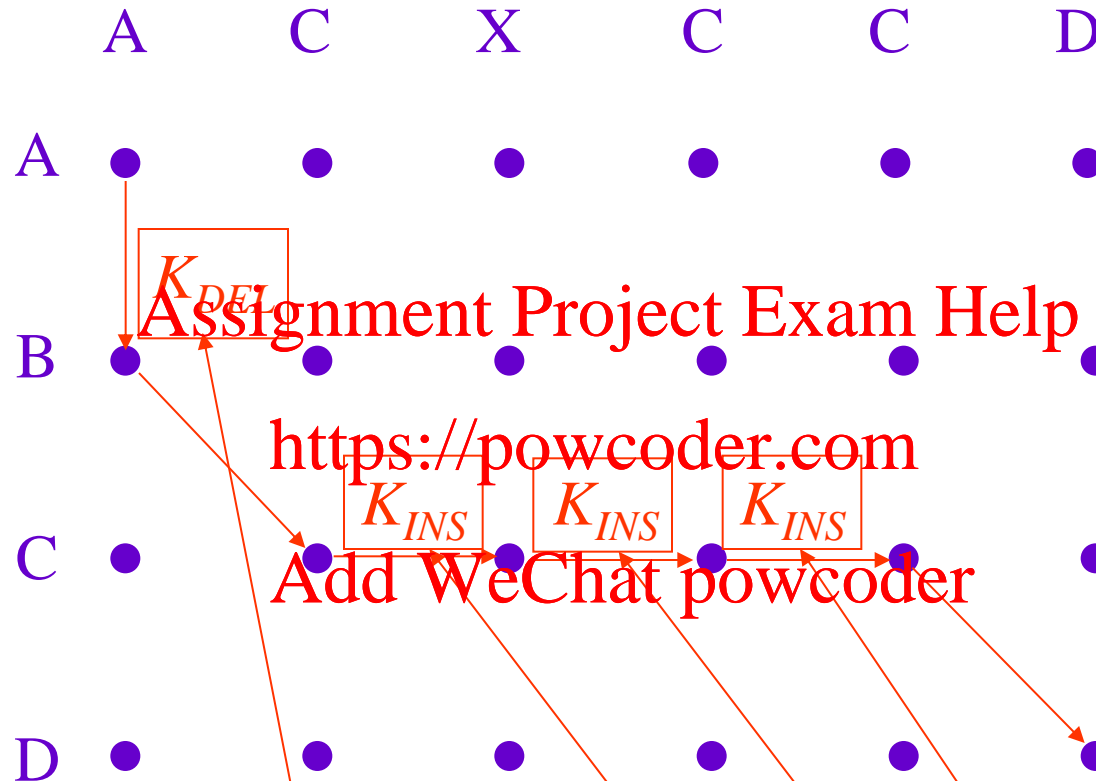
Accumulated distance along p



$$AD_p(S_1, S_2) = d(A, A) + d(A, B) + d(C, C) + d(C, X) + d(C, C) + d(C, C) + d(D, D)$$



Accumulated distance (continued)



$$AD_p(S_1, S_2) = d(A, B) + d(C, X)$$

$$AD_p(S_1, S_2) = K_{DEL} + d(A, B) + K_{INS} + K_{INS} + K_{INS} + d(C, X)$$



Optimal path and DP distance

- Optimal path is path with minimum accumulated distance

- Formally the optimal path is \hat{p}

where: <https://powcoder.com>

$$\hat{p} = \arg \min_p AD_p(S_1, S_2), \text{ or } AD_{\hat{p}}(S_1, S_2) = \min_p AD_p(S_1, S_2)$$

- The DP distance, or accumulated distance $AD(S_1, S_2)$ between S_1 and S_2 is given by:

$$AD(S_1, S_2) = AD_{\hat{p}}(S_1, S_2)$$



Calculating the optimal path

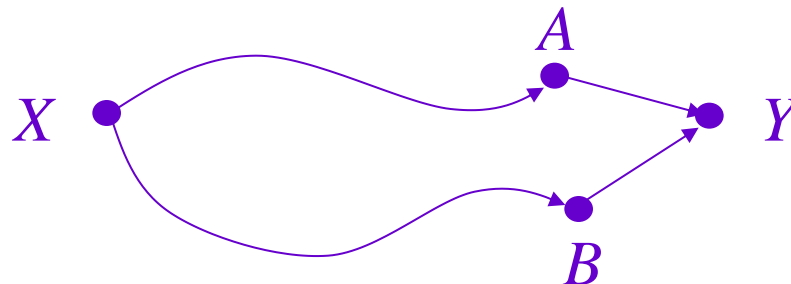
- Given
 - the distance matrix D ,
 - the insertion penalty K_{INS} , and
 - the deletion penalty K_{DEL}^*
- How can we compute the optimal path between two (potentially very long) sequences S_1 and S_2 ?

*If K_{DEL} and K_{INS} are not defined you should assume that they are zero



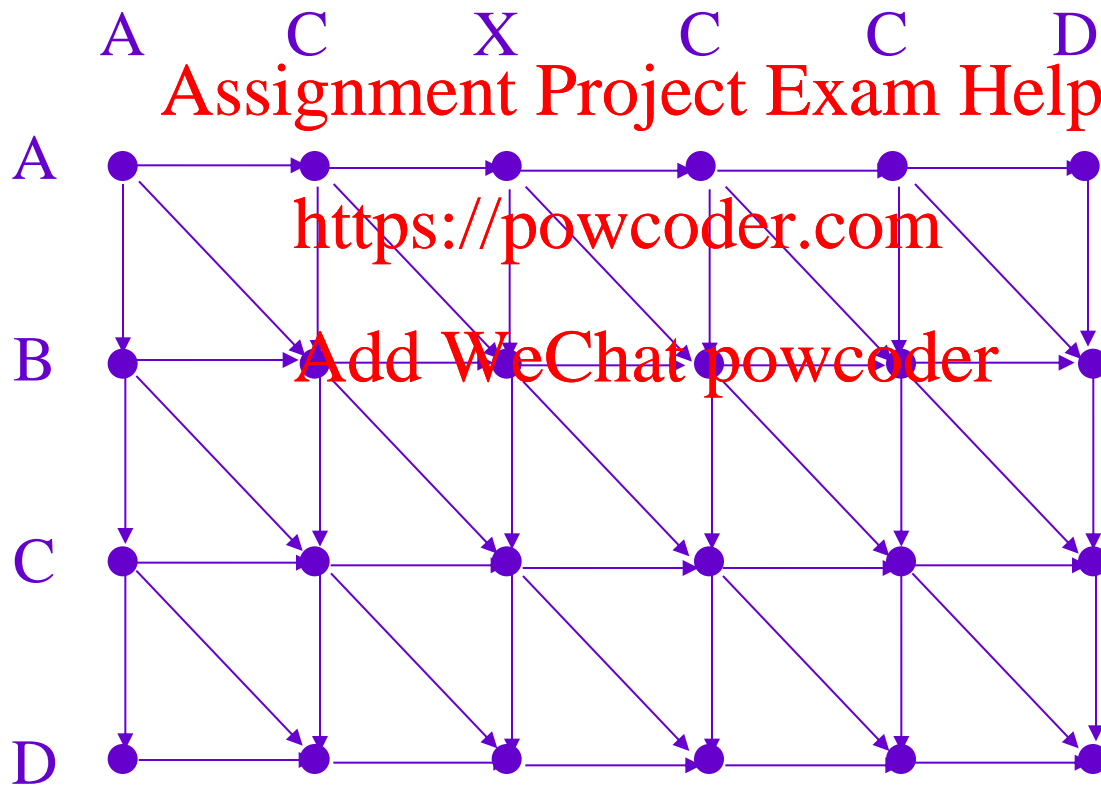
Dynamic Programming (DP)

- Optimal path calculated using Dynamic Programming (DP), based on principle of optimality
- If paths from X to Y go through A or B immediately before Y , optimal path from X to Y is best of:
 - Best path from X to A plus cost to go from A to Y
 - Best path from X to B plus cost to go from B to Y



DP – step 1

- Step 1: draw the trellis of all possible paths

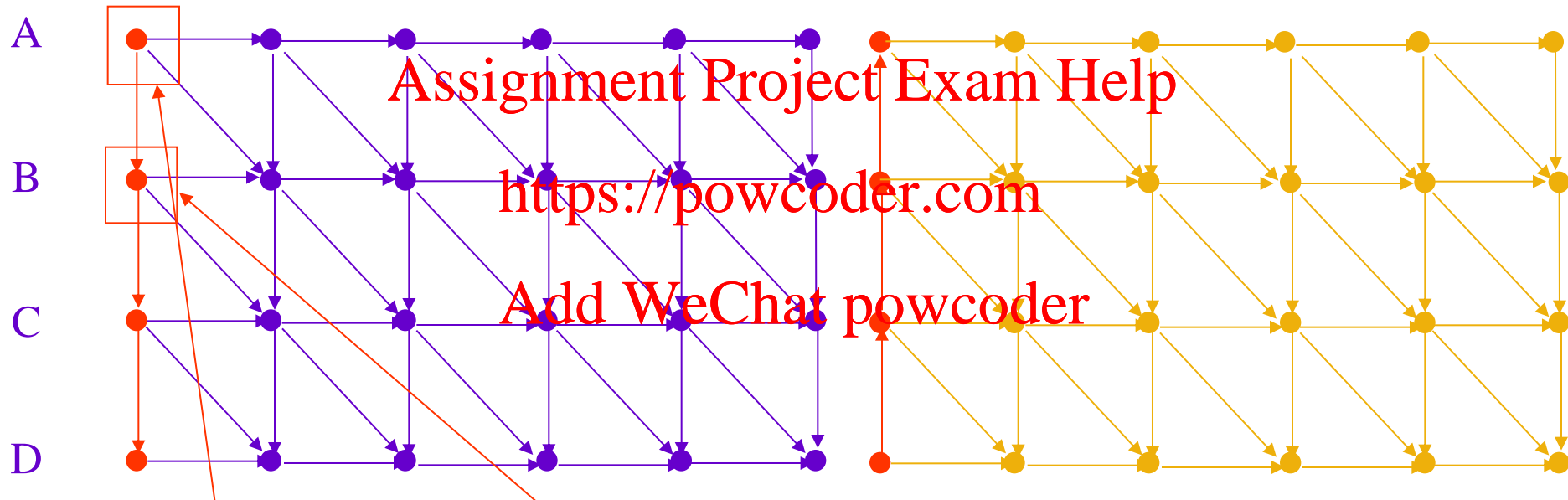


DP – forward pass – initialisation

Accumulated distance matrix

A C X C C D

Path matrix



Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

$$ad(1,1)=d(A,A)$$

$$ad(2,1)=ad(1,1)+d(2,1)+K_{\text{DEL}}$$

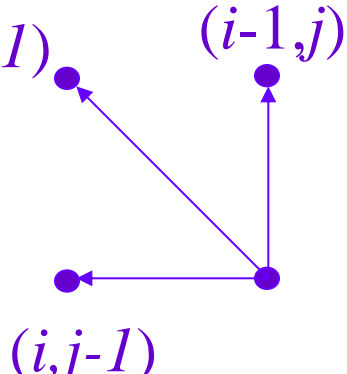


$ad(i,j)$

- $ad(i,j)$ is the sum of distances along the best (partial) path from (1,1) to (i,j)
- Calculated using the principle of optimality

$$ad(i,j) = \min \begin{cases} ad(i-1,j) + K_{DEL} + d(i,j) \\ ad(i,j-1) + K_{INS} + d(i,j) \end{cases}$$

<https://powcoder.com>
 Add WeChat powcoder



- Forward path matrix records local optimal paths

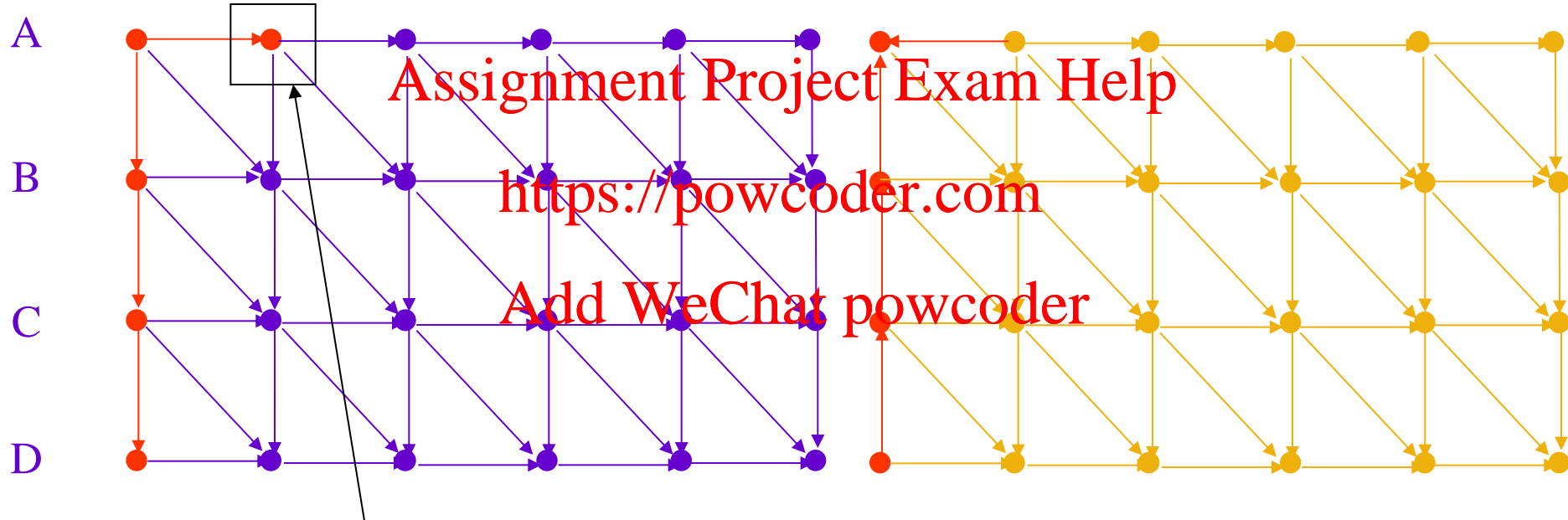


DP – forward pass – continued

Accumulated distance matrix

A C X C C D

Path matrix



$$ad(1,2) = ad(1,1) + d(A,C) + K_{INS}$$

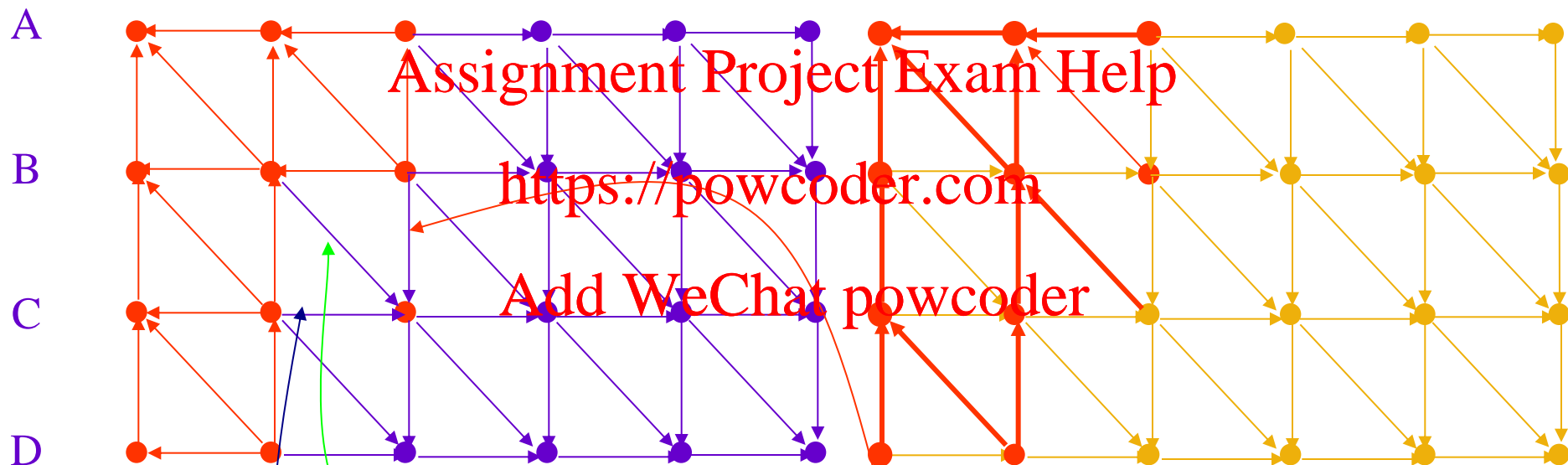


DP – forward pass – continued

Accumulated distance matrix

A C X C C D

Path matrix



$$ad(3,3) = \min \begin{cases} ad(2,3) + K_{DEL} + d(C, X) \\ ad(2,2) + d(C, X) \\ ad(3,2) + K_{INS} + d(C, X) \end{cases}$$

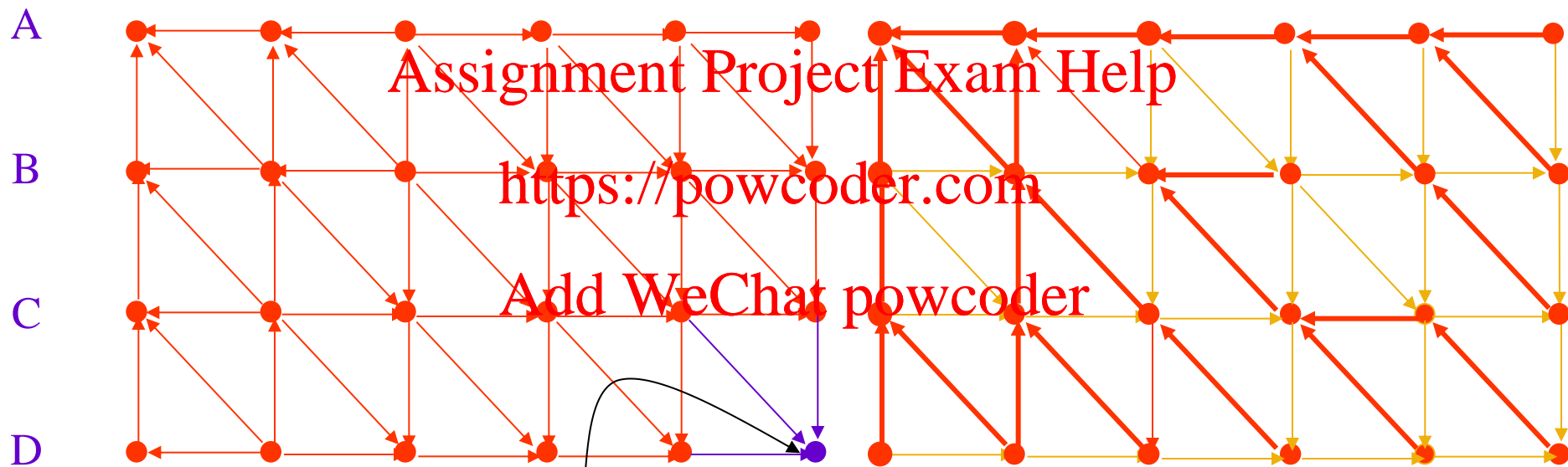


DP – forward pass – continued

Accumulated distance matrix

A C X C C D

Path matrix



$$AD(S_1, S_2) = ad(4, 6) = \min \begin{cases} ad(3, 6) + K_{DEL} + d(D, D) \\ ad(3, 5) + d(D, D) \\ ad(4, 5) + K_{INS} + d(D, D) \end{cases}$$

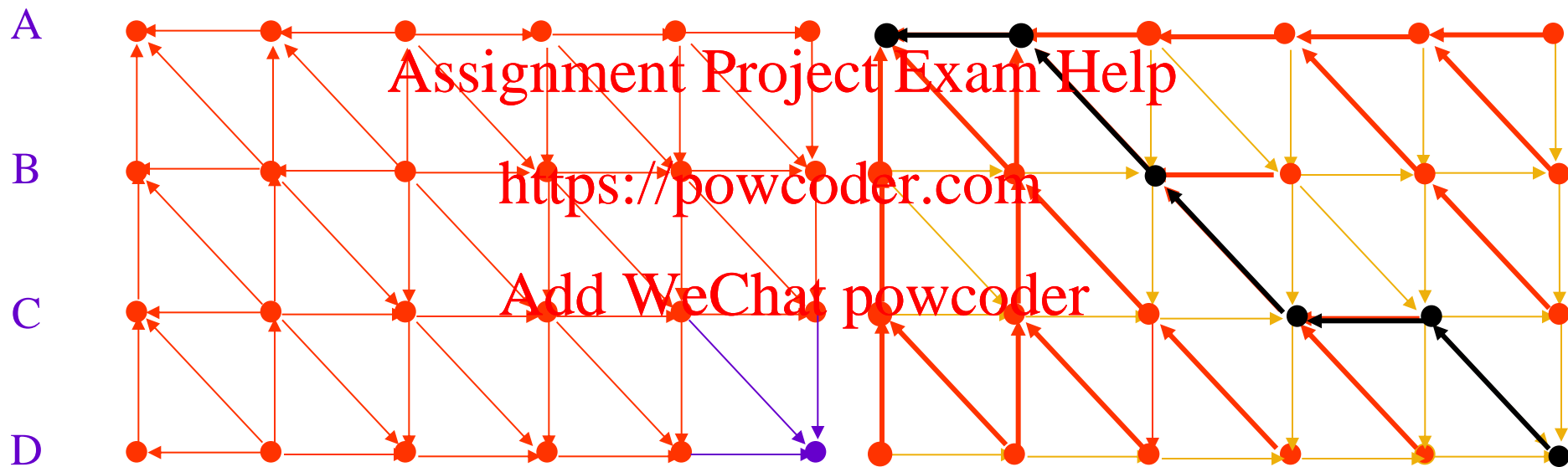


DP – forward pass – continued

Accumulated distance matrix

A C X C C D

Path matrix



Optimal path obtained by tracing back through path matrix, starting at the bottom right-hand corner



Summary

- Introduction to sequence analysis
- Dynamic Programming (DP) and the principle of optimality
- Computing the accumulated distance using DP
 - Distance matrix, Accumulated distance matrix, Path matrix, and Optimal path
- Recovering the optimal path

