

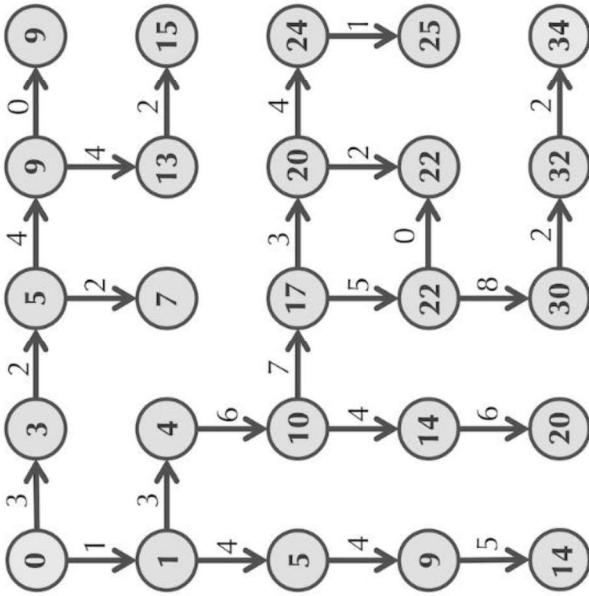
## 22BIO211: Intelligence of Biological Systems - 2

# BACKTRACKING IN ALIGNMENT GRAPHS

Dr. Manjusha Nair M  
Amrita School of Computing, Amritapuri  
Email : [manjushanair@am.amrita.edu](mailto:manjushanair@am.amrita.edu)  
Contact No: 9447745519

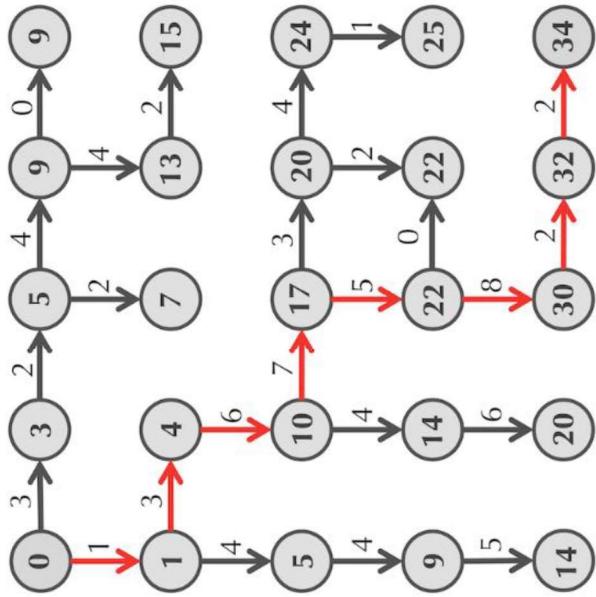
# BackTracking

- Each edge selected by ManhattanTourist.
- To form a longest path, we simply need to find a path from source to sink formed by highlighted edges (more than one such path may exist).
- However, if we were to walk from source to sink along the highlighted edges, we might reach a dead end, such as the node (1, 2).



# Backtracking

- In contrast, every path from the sink will bring us back to the source if we backtrack in the direction opposite to each highlighted edge.
- Following edges from the sink backwards to the source produces the red path in the DAG , which highlights a longest path from the source to sink.
- We can use this backtracking idea to construct an LCS of strings v and w.



# Backtracking - Algorithm

- If we assign a weight of 1 to the edges in  $\text{AlignmentGraph}(v, w)$  corresponding to matches and assign a weight of 0 to all other edges, then  $S_{|v|, |w|}$  gives the length of an LCS.
- The following algorithm maintains a record of which edge was used to compute each value  $s_{i,j}$  by utilizing backtracking pointers, which take one of the three values  $\downarrow$ ,  $\rightarrow$ , or  $\nwarrow$ .
- Backtracking pointers are stored in a matrix *Backtrack*.

# Backtracking - Algorithm

```
LCSBackTrack( $v, w$ )
    for  $i \leftarrow 0$  to  $|v|$ 
         $s_{i, 0} \leftarrow \emptyset$ 
        for  $j \leftarrow 0$  to  $|w|$ 
             $s_{0, j} \leftarrow \emptyset$ 
            for  $i \leftarrow 1$  to  $|v|$ 
                for  $j \leftarrow 1$  to  $|w|$ 
                     $match \leftarrow \emptyset$ 
                    if  $v_{i-1} = w_{j-1}$ 
                         $match \leftarrow 1$ 
                    if  $s_{i,j} = s_{i-1,j}$ 
                         $Backtrack_i, j \leftarrow \downarrow$ 
                    else if  $s_{i, j} = s_{i, j-1}$ 
                         $Backtrack_i, j \leftarrow \rightarrow$ 
                    else if  $s_{i, j} = s_{i-1, j-1} + match$ 
                         $Backtrack_i, j \leftarrow \rightarrow$ 
                    return  $Backtrack$ 
     $s_{i, j} \leftarrow \max\{s_{i-1, j}, s_{i, j-1}, s_{i-1, j-1} + match\}$ 
```

# Longest Path in a DAG Problem

- We now need to find a path from the source to the sink formed by the highlighted edges.
- We can use this backtracking idea to construct an LCS of strings  $v$  and  $w$ .
- The algorithm in the next slide solves the Longest Common Subsequence Problem by using the information in *Backtrack*.
- ***OUTPUTLCS(Backtrack, v, i, j)*** outputs an LCS between the  $i$ -prefix of  $v$  and the  $j$ -prefix of  $w$ .
- The initial invocation that outputs an LCS of  $v$  and  $w$  is ***OUTPUTLCS(Backtrack, v, |v|, |w|)***.

# Longest Path in a DAG Problem

```
OutputLCS(backtrack, v, i, j)

if i = 0 or j = 0
    return ""

if backtracki, j = "↓" 
    return OutputLCS(backtrack, v, i - 1, j)

else if backtracki, j = "→" 
    return OutputLCS(backtrack, v, i, j - 1)

else
    return OutputLCS(backtrack, v, i - 1, j - 1) + vi
```

	G	A	A	T	T	C	A	G	T	T	A
0	1	1	1	2	2	3	3	4	4	5	5
G	0	0	0	0	0	0	0	0	0	0	0
G	0	0	0	0	0	0	0	0	0	0	0
A	0	1	1	2	2	3	3	4	4	5	5
T	0	1	2	2	2	3	3	4	4	5	5
C	0	1	2	2	2	3	3	4	4	5	5
G	0	1	2	2	2	3	3	4	4	5	5
A	0	1	2	2	2	3	3	4	4	5	5
A	0	1	2	2	2	3	3	4	4	5	5

# Sequence alignment Example - Backtracking

	G	A	A	T	T	C	A	G	T	T	A
G	0	0	0	0	0	0	0	0	0	0	0
G	0	1	1	1	1	1	1	1	1	1	1
G	0	1	1	1	1	1	1	1	1	1	1
A	0	1	1	2	2	2	2	2	2	2	2
T	0	1	2	2	3	3	3	3	3	3	3
C	0	1	2	3	3	4	4	4	4	4	4
G	0	1	2	2	3	3	4	5	5	5	5
A	0	1	2	3	3	4	5	5	5	5	6

(Seq #1)

(Seq #2)

Alignment:

A  
|  
A

	G	A	A	T	T	C	A	G	T	T	A
G	0	0	0	0	0	0	0	0	0	0	0
G	0	1	1	1	1	1	1	1	1	1	1
G	0	1	1	1	1	1	1	1	1	1	1
A	0	1	1	2	2	2	2	2	2	2	2
T	0	1	2	2	3	3	3	3	3	3	3
C	0	1	2	3	3	4	4	4	4	4	4
G	0	1	2	2	3	3	4	5	5	5	5
A	0	1	2	3	3	4	5	5	5	5	6

6

# Sequence alignment Example - Backtracking

	G	A	A	T	T	C	A	G	T	T	A
G	0	0	0	0	0	0	0	0	0	0	0
G	0	1	1	1	1	1	1	1	1	1	1
G	0	1	1	1	1	1	1	1	1	2	2
A	0	1	2	2	2	2	2	2	2	2	2
T	0	1	2	2	3	3	3	3	3	3	3
C	0	1	2	2	3	3	4	4	4	4	4
G	0	1	2	2	3	3	4	4	4	5	5
											6

(Seq #1)

T A

Alignment:

(Seq #2)

| - A

# Sequence alignment Example - Backtracking

	G	A	A	T	T	C	A	G	T	T	A
G	0	0	0	0	0	0	0	0	0	0	0
G	0	1	1	1	1	1	1	1	1	1	1
G	0	1	1	1	1	1	1	1	2	2	2
A	0	1	2	2	2	2	2	2	2	2	2
T	0	1	2	2	3	3	3	3	3	3	3
C	0	1	2	2	3	3	4	4	4	4	4
G	0	1	2	2	3	3	4	4	4	5	5
											6
											A

	G	A	A	T	T	C	A	G	T	T	A
G	0	1	1	1	1	1	1	1	1	1	1
G	0	1	1	1	1	1	1	1	2	2	2
G	0	1	1	1	1	1	1	1	2	2	2
A	0	1	2	2	2	2	2	2	2	2	2
T	0	1	2	2	3	3	3	3	3	3	3
C	0	1	2	2	3	3	4	4	4	4	4
G	0	1	2	2	3	3	4	4	4	5	5
	0	1	2	2	3	3	4	4	4	5	6
											A

(Seq #1) G A A T T C A G T T A  
(Seq #2) G G A - T C - G - - A

Alignment:

# Scoring Alignment

- To generalize the alignment scoring model
  - we award +1 for matches,
  - also penalize mismatches by some positive constant  $\mu$  (the mismatch penalty)
  - and indels by some positive constant  $\sigma$  (the indel penalty).
- As a result, the score of an alignment is equal to the following expression:
$$S = \# \text{ matches} - \mu \cdot \# \text{ mismatches} - \sigma \cdot \# \text{ indels}$$

# Scoring Alignment

- For example, with the parameters  $\mu = 1$  and  $\sigma = 2$ , the following alignment will be assigned a score of -4.

A	T	-	G	T	T	A	T	A
A	T	C	G	T	-	C	-	C
+1	+1	-2	+1	+1	-2	-1	-2	-1

- $S = 4 - (1 \times 2) - (2 \times 3) = 4 - 2 - 6 = -4$
- Biologists have further refined this cost function to allow for the fact that
  - *some mutations may be more likely than others*
    - which calls for mismatches and indel penalties that differ depending on the specific symbols involved.

# Scoring Alignment

- We will extend the  $k$ -letter alphabet to include the space symbol and then construct a  $(k + 1) \times (k + 1)$  scoring matrix holding the score of aligning every pair of symbols.
- The scoring matrix for comparing DNA sequences ( $k = 4$ ) when all mismatches are penalized by  $\mu$  and all indels are penalized by  $\sigma$  is shown below.

	A	C	G	T	-
A	+1	$-\mu$	$-\mu$	$-\mu$	$-\sigma$
C	$-\mu$	+1	$-\mu$	$-\mu$	$-\sigma$
G	$-\mu$	$-\mu$	+1	$-\mu$	$-\sigma$
T	$-\mu$	$-\mu$	$-\mu$	+1	$-\sigma$
-	$-\sigma$	$-\sigma$	$-\sigma$	$-\sigma$	$-\sigma$

## Scoring matrices for protein sequence comparison

- Scoring matrices for DNA sequence comparison are usually defined only by the parameters  $\mu$  and  $\sigma$
- Scoring matrices for protein sequence comparison weight different mutations differently and become quite involved.
- One example- the PAM<sub>250</sub> scoring matrix (Here, the indel penalty has been set to 8)

## PAM<sub>250</sub> scoring matrix

# Summary

- BackTracking
  - *Algorithm*
- Longest Path in a DAG Problem
- Scoring Alignment
- Scoring matrices for protein sequence comparison
  - $PAM_{250}$  scoring matrix