

# Sequence Assembly-II

Hammad Naveed

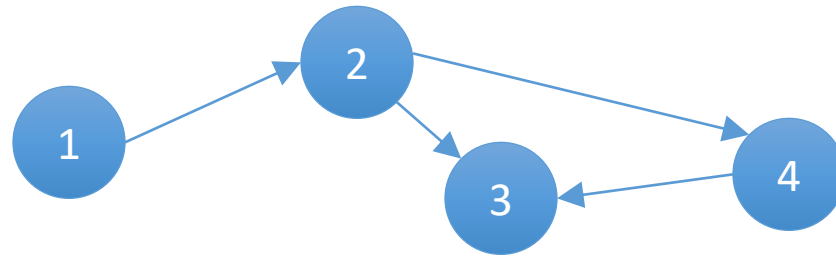
hammad.naveed@nu.edu.pk

# Graph Basics

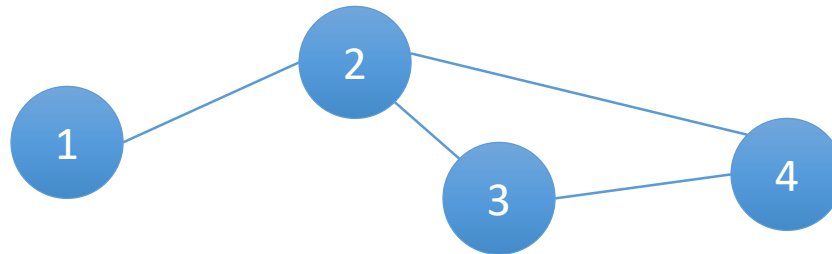
- A graph ( $G$ ) consists of vertices ( $V$ ) and edges ( $E$ )

$$G = (V, E)$$

- Edges can either be *directed* (*directed graphs*)

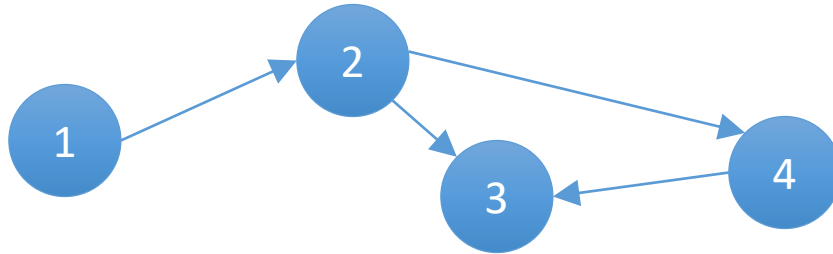


- or *undirected* (*undirected graphs*)



# Vertex degrees

- The *degree* of a vertex: the # of edges incident to that vertex
- For directed graphs, we also have the notion of
  - *indegree*: The number incoming edges
  - *outdegree*: The number of outgoing edges



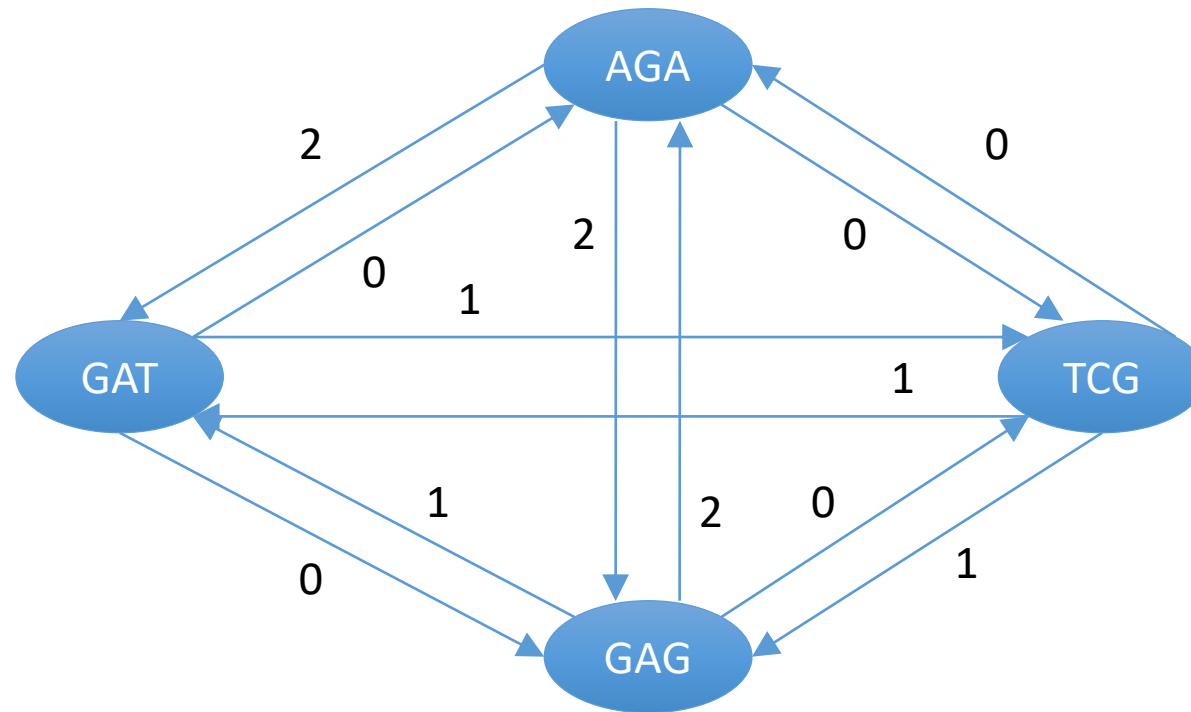
$degree(v_2) = 3$   
 $indegree(v_2) = 1$   
 $outdegree(v_2) = 2$

# Overlap graph

- For a set of sequence reads  $S$ , construct a directed weighted graph  $G = (V, E, w)$ 
  - with one vertex per read ( $v_i$  corresponds to  $s_i$ )
  - edges between all vertices (a *complete* graph)
  - $w(v_i, v_j) = \text{overlap}(s_i, s_j) = \text{length of longest suffix of } s_i \text{ that is a prefix of } s_j$

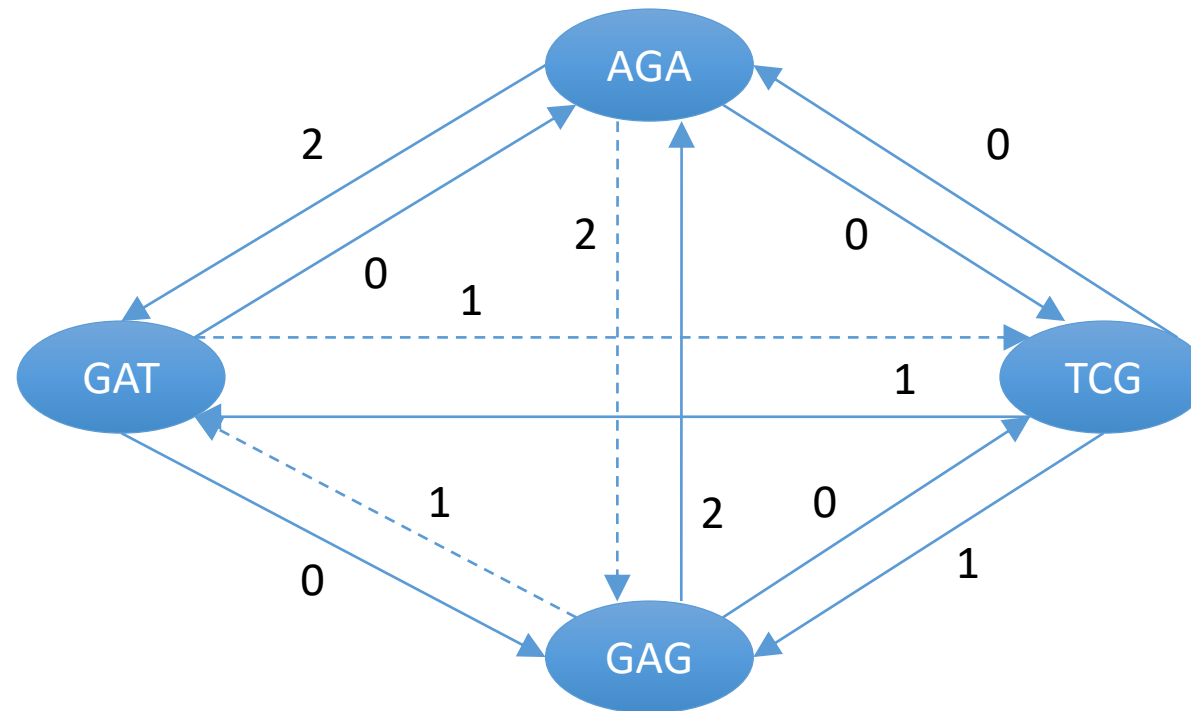
# Overlap graph example

- Let  $S = \{\text{AGA}, \text{GAT}, \text{TCG}, \text{GAG}\}$



# Assembly as Hamiltonian Path

- *Hamiltonian Path*: path through graph that visits each vertex exactly once

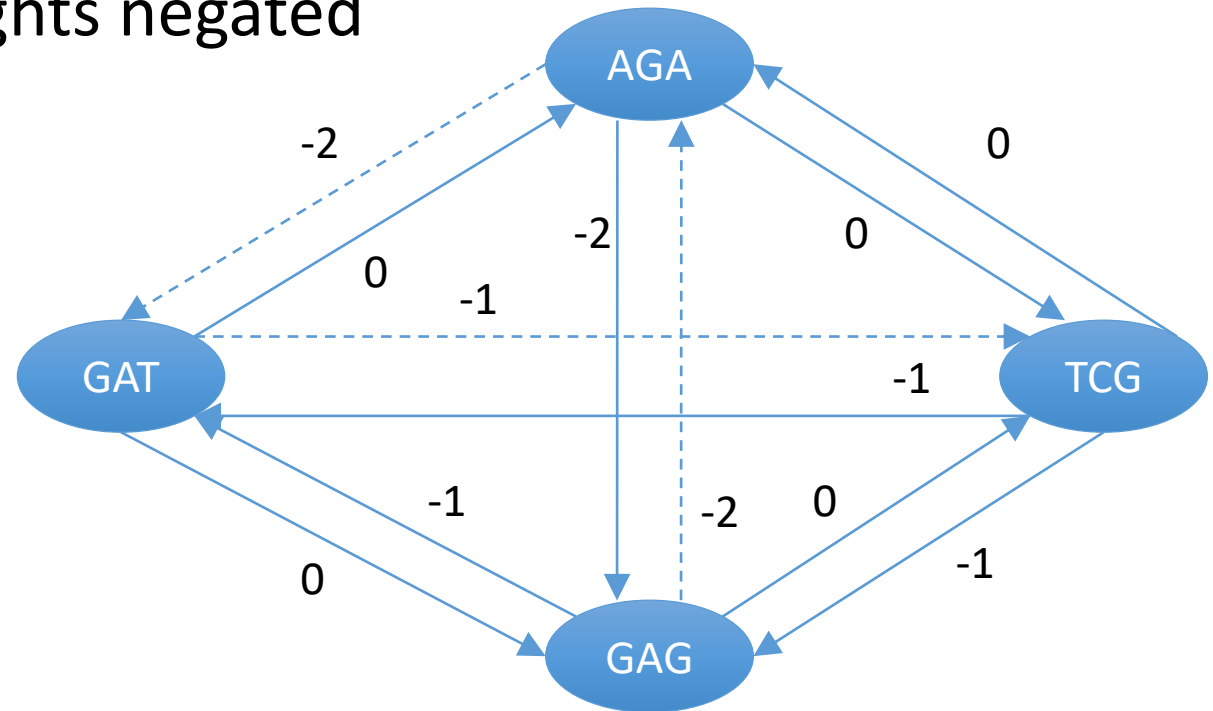


Path: AGAGATCG

# Shortest superstring as TSP

- minimize superstring length  $\rightarrow$  minimize hamiltonian path length in overlap graph with edge weights negated

Path: GAGATCG  
Path length: -5  
String length: 7



- This is essentially the Traveling Salesman Problem (also *NP*-complete)

# Greedy Algorithms

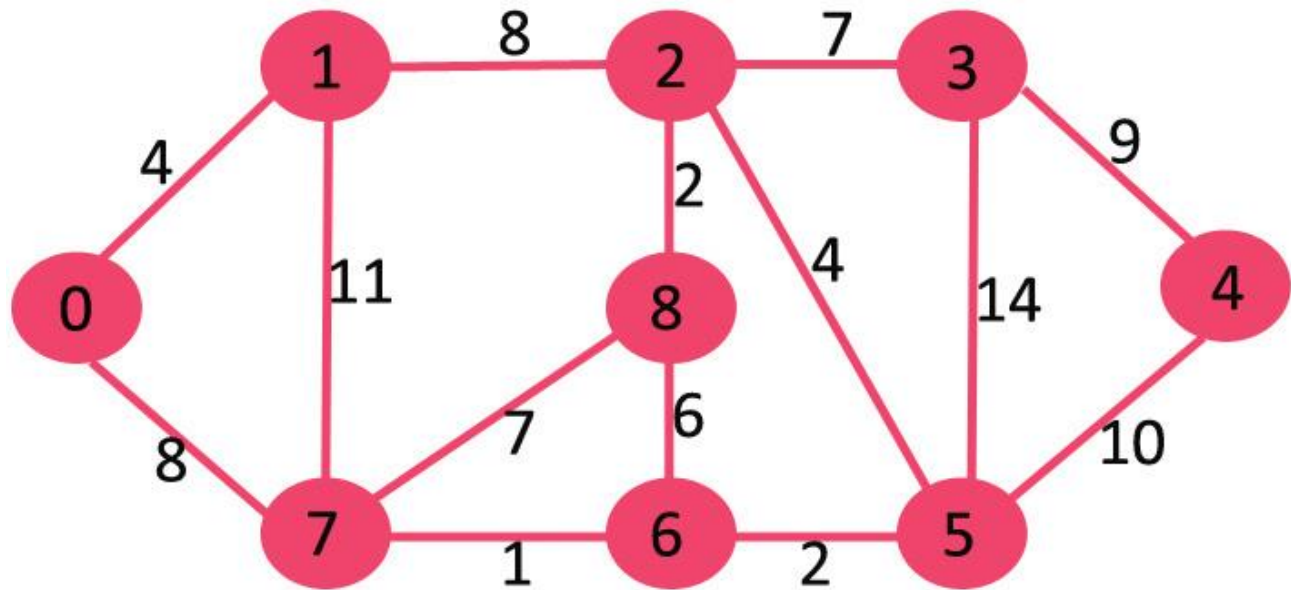
- **Definition:** An algorithm that always takes the best immediate, or local, solution while finding an answer.
- Greedy algorithms find the overall, or globally, optimal solution for some optimization problems, but may find less-than-optimal solutions for some instances of other problems.

Paul E. Black, "greedy algorithm", in Dictionary of Algorithms and Data Structures [online], Paul E. Black, ed., U.S. National Institute of Standards and Technology. 2 February 2005.  
<http://www.itl.nist.gov/div897/sqg/dads/HTML/greedyalgo.html>



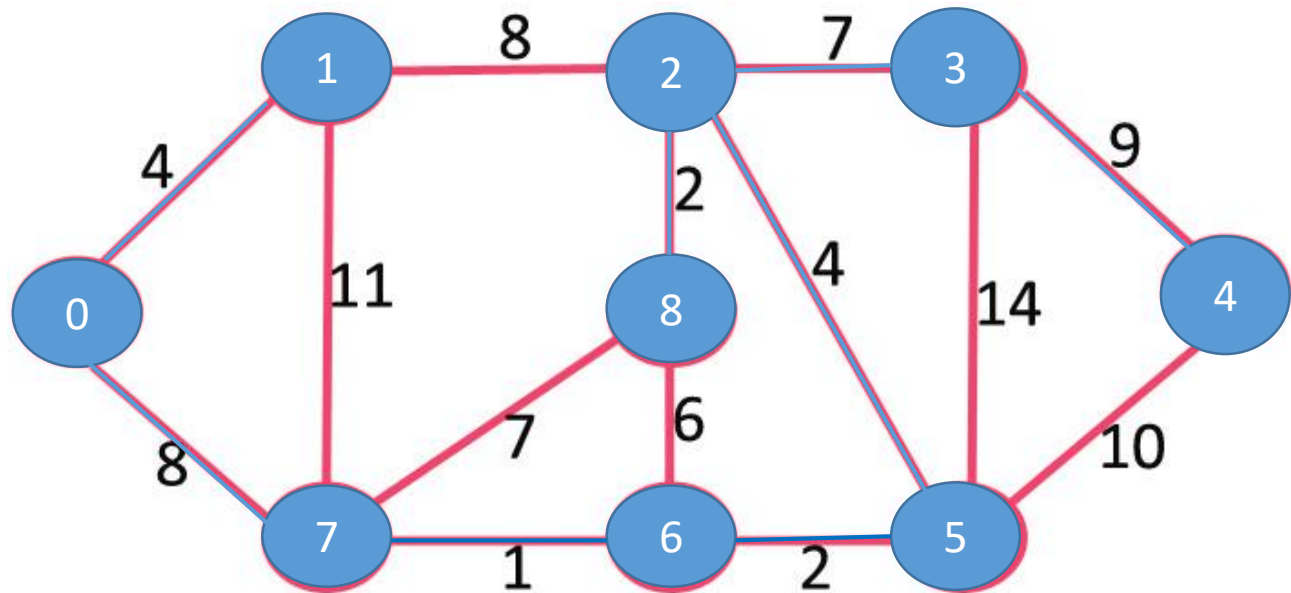
# Greedy Algorithm Examples

- Kruskal's Algorithm for Minimum Spanning Tree
  - Minimum spanning tree: a set of  $n-1$  edges that connects a graph of  $n$  vertices without any cycles and that has minimal total weight
  - Kruskal's algorithm adds the edge that connects two components with the smallest weight at each step without introducing a cycle
    - Proven to give an optimal solution



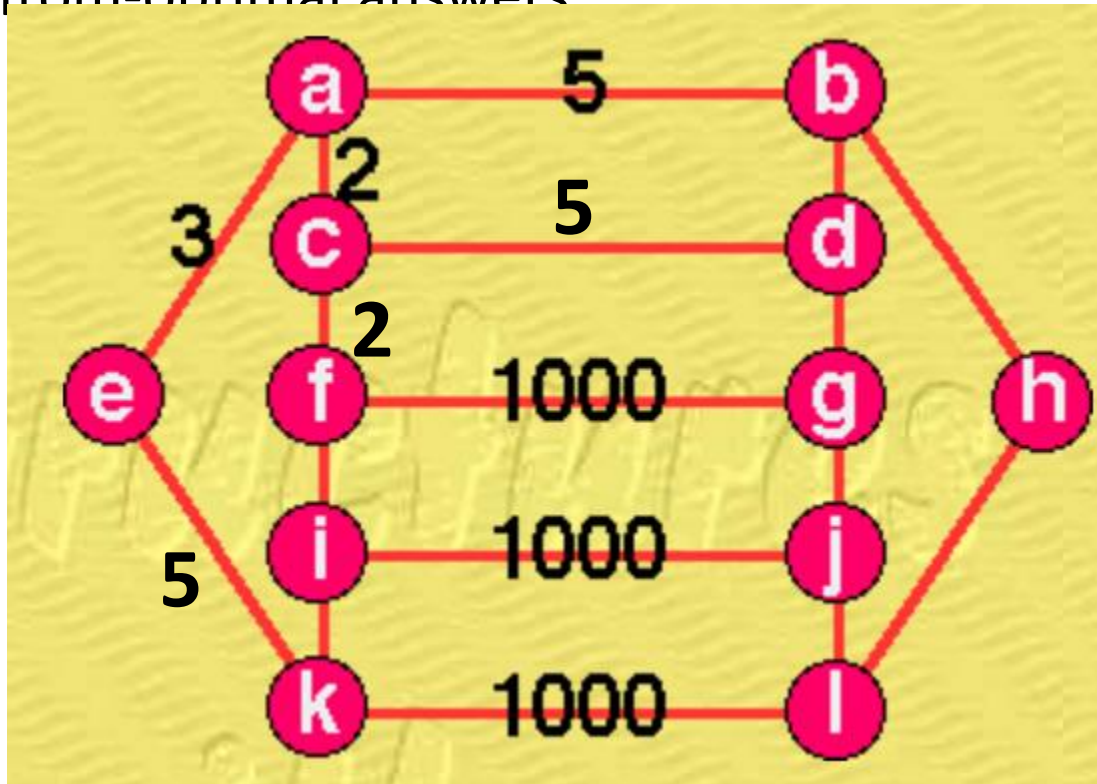
# Greedy Algorithm Examples

Weight	Src	Dest
1	7	6
2	8	2
2	6	5
4	0	1
4	2	5
6	8	6
7	2	3
7	7	8
8	0	7
8	1	2
9	3	4
10	5	4
11	1	7
14	3	5



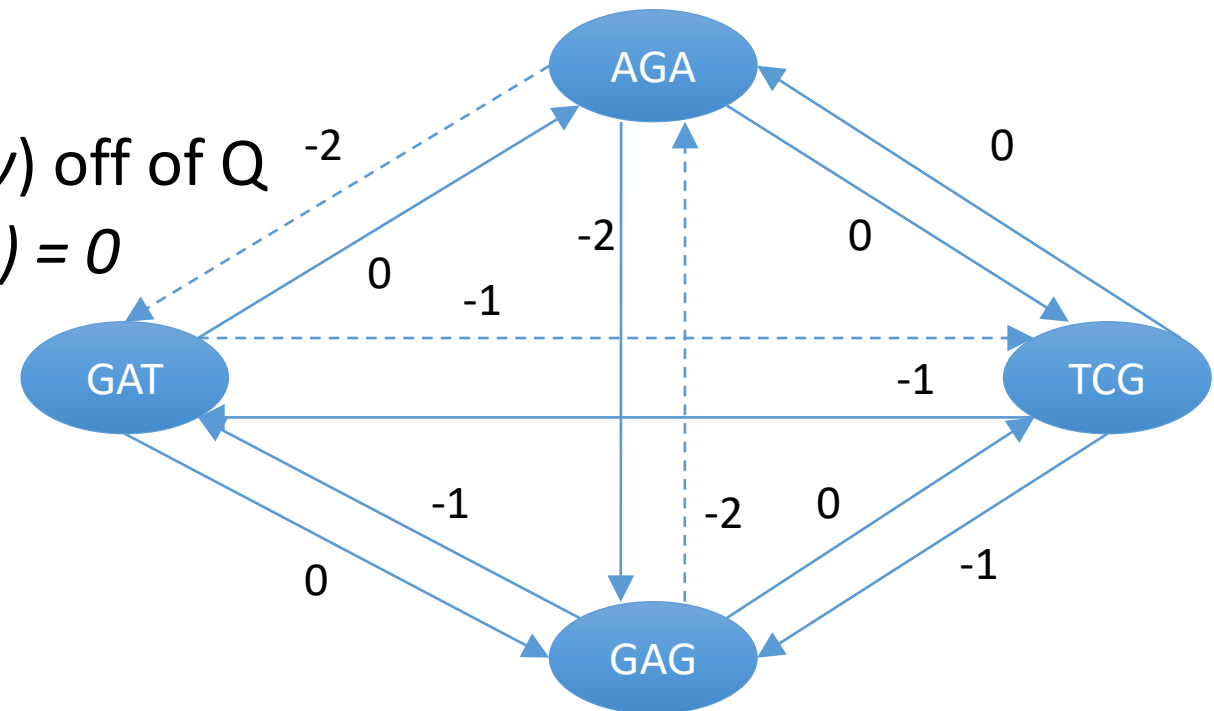
# Greedy Algorithm Examples

- Traveling Salesman Problem
  - Greedy algorithm chooses to visit closest vertex at each step
  - Can give far-from-optimal answers



# The Greedy Algorithm

- Let  $G$  be a graph with fragments as vertices, and no edges to start
- Create a queue,  $Q$ , of overlap edges, with edges in order of increasing weight
- While  $G$  is disconnected
  - Pop the next possible edge  $e = (u,v)$  off of  $Q$
  - If  $outdegree(u) = 0$  and  $indegree(v) = 0$  and  $e$  does not create a cycle
    - Add  $e$  to  $G$



# The Greedy Algorithm

- While  $G$  is disconnected
  - Pop the next possible edge  $e = (u,v)$  off of  $Q$
  - If  $\text{outdegree}(u) = 0$  and  $\text{indegree}(v) = 0$  and  $e$  does not create a cycle
    - Add  $e$  to  $G$

- GAG  $\rightarrow$  AGA -2
- AGA  $\rightarrow$  GAG -2
- AGA  $\rightarrow$  GAT -2
- GAG  $\rightarrow$  GAT -1
- TCG  $\rightarrow$  GAT -1
- TCG  $\rightarrow$  GAG -1
- GAT  $\rightarrow$  TCG -1

Path: GAGATCG  
Path length: -5  
String length: 7

