

MARMARA UNIVERSITY – FACULTY OF ENGINEERING

CSE 4065 ASSIGNMENT #2 REPORT



Name, Surname: Hasan Şenyurt – Melisa Durmuş – Mustafa Kerem Ekinici

Student ID: 150120531 – 150119727 – 150119695

Department: Computer Engineering

IMPLEMENTATION DETAILS

Node Class:

We have node class for creating the 2D grid. Node class has some variables and two functions. Variables are score (node's score), match score, penalty scores, coordinates and Boolean variable which is for gap opening property.

One of the function is `__init__` function to define variables. Also, values of initial node of the grid is defined in this function.

```
1  class Node:
2
3      #VALUES
4      score = 0
5      match = 3
6      gapOpeningPenalty = -1
7      gapExtensionPenalty = -0.5
8      misMatchPenalty = -1
9      isOpeningGap = False
10
11  def __init__(self, xCor, yCor, nucleotid):
12      self.xCor = xCor
13      self.yCor = yCor
14      self.nucleotid = nucleotid
15      self.topEdge = 0
16      self.diagonalEdge = 0
17      self.leftEdge = 0
18
19      #Initial nucleotid values.
20      if (self.xCor == 0 and self.yCor == 0) or self.xCor!=0 and self.yCor!=0:
21          self.nucleotid = "X"
22
23      else:
24          self.nucleotid = nucleotid
```

Figure-1: Node Class

Other function is called ‘calculate_score’. This function is written to calculate node’s score by looking left, top and diagonal nodes. It compares scores which comes from left, left-diagonal and top nodes and their edges, then takes maximum value among them. That value becomes score of node that called the function.

```

1  #Function for score calculation of each node by choosing maximum score which comes from top, diagonal and left nodes, their edges.
2  def calculate_score(self,grid):
3
4      #Score of initial node is 0
5      if self.xCor == 0 and self.yCor == 0:
6          self.score = 0
7
8      #Score of first row is increasing -1 by -1.
9      elif self.xCor == 0:
10         top_x = self.xCor
11         top_y = self.yCor - 1
12         self.score = grid[top_x][top_y].score + self.gapOpeningPenalty
13         self.topEdge = self.gapOpeningPenalty
14
15     #Score of first column is increasing -1 by -1.
16     elif self.yCor == 0:
17         left_x = self.xCor - 1
18         left_y = self.yCor
19         self.score = grid[left_x][left_y].score + self.gapOpeningPenalty
20         self.leftEdge = self.gapOpeningPenalty
21
22     #Calculating scores of nodes except initial node, first row and column.
23     else:
24         #Getting coordinates of left, top and diagonal nodes.
25         left_x = self.xCor - 1; left_y = self.yCor
26         top_x = self.xCor; top_y = self.yCor - 1
27         diagonal_x = self.xCor - 1; diagonal_y = self.yCor - 1
28
29         max_left_x = 0; max_left_y = self.yCor
30         max_x = self.xCor; max_y = 0
31
32         temp_score=[]
33
34         #CHECKING MATCH!
35         if(grid[max_left_x][max_left_y].nucleotide==grid[max_x][max_y].nucleotide):
36             self.diagonalEdge = self.match
37             temp_score.append(grid[diagonal_x][diagonal_y].score+ self.match)
38
39         #CHECKING MISMATCH!
40         else:
41             temp_score.append(grid[diagonal_x][diagonal_y].score + self.misMatchPenalty)
42             self.diagonalEdge = self.misMatchPenalty
43
44         #assigning penalties to the edges which connects nodes.
45         self.topEdge = self.gapOpeningPenalty
46         self.leftEdge = self.gapOpeningPenalty
47
48         #Selecting maximum score to put the score value according to dynamic programming rules.
49         temp_score.append(grid[left_x][left_y].score+self.leftEdge)#Soldan indel
50         temp_score.append(grid[top_x][top_y].score+self.topEdge)#Üstten indel
51         self.score=max(temp_score)
52

```

Figure-2: Score Calculation Function of Node Class

Backtracking Function:

Outside of node class, we have backtracking function to find correct path in grid to make sequence alignment. Main idea of backtracking with gap openings and extension scores is that preferring indel penalties to mismatches in specific situation such as putting indel and making mismatch are available at same point because gap extension penalty is less than mismatch, so if there will be more than one indel, penalty will be less than mismatches. If blocks are designed for the purpose above. We checked indels first to prefer them to mismatches. After checking indels, matches and mismatches are checked. (node's score – node's edge != surrounding node's scores).

```
1 #backtracking function which backtracks from end of the grid to starting point of grid according to dynamic pairwise alignment rules.
2 def backTracking(grid,seq1,seq2):
3
4     x = len(seq1)
5     y = len(seq2)
6
7     seq1_align = ""
8     seq2_align = ""
9
10    gap_extension_scores = 0 # this extension scores will be added to total scores.
11
12    while x+y != 0:
13
14        ...
15        Main idea of backtracking with gap openings and extension scores is that preferring indel penalties to mismatches in
16        specific situation such as putting indel and making mismatch are available at same point because
17        gap extension penalty is less than mismatch, so if there will be more than one indel, penalty will be less than mismatches.
18        ...
19
20    #Checking top edge of the node (INDEL)
21    if y != 0 and grid[x][y].score - grid[x][y].topEdge == grid[x][y-1].score:
22        seq1_align += '-'
23        seq2_align += grid[0][y].nucleotid
24
25        #gap opened
26        grid[x][y].isOpeningGap = True
27
28        #checking if gap opened before, if yes than it is gap extension.
29        if y != len(seq2):
30            if grid[x][y+1].isOpeningGap == True:
31                gap_extension_scores += 0.5
32            y -= 1
33
34    #Checking left edge of the node (INDEL)
35    elif x != 0 and grid[x][y].score - grid[x][y].leftEdge == grid[x-1][y].score:
36        seq1_align += grid[x][0].nucleotid
37        seq2_align += '-'
38
39        #gap opened
40        grid[x][y].isOpeningGap = True
41
42        #checking if gap opened before, if yes than it is gap extension.
43        if x != len(seq1):
44            if grid[x+1][y].isOpeningGap == True:
45                gap_extension_scores += 0.5
46            x -= 1
47
48    #Checking if there is a mismatch or match.
49    elif (x != 0 or y != 0) and grid[x][y].score - grid[x][y].diagonalEdge == grid[x-1][y-1].score: #çapraz
50        seq1_align += grid[x][0].nucleotid
51        seq2_align += grid[0][y].nucleotid
52        x -= 1
53        y -= 1
54
55    #print the results
56    print(seq1_align[::-1])
57    print(seq2_align[::-1])
58    print("Score:", grid[-1][-1].score + gap_extension_scores)
```

Figure-3: Backtracking Function

Main Function:

In main function, grid is created and scores of nodes are assigned. Scores of nodes which are on first row and column of grid increase -1 by -1. Also, score of initial node which has (0,0) coordinate in grid is 0. Test inputs are defined as well in main function. After all definitions and assignments, backtracking function is called for each test input. Results are printed on the console.

```
1  def main():
2      test_inputs=["test1.seq", "test2.seq", "test3.seq", "test4.seq", "test5.seq"]
3
4      for file in test_inputs:
5          with open("Test Inputs/"+file, "r") as f:
6              seq1 = f.readline()[::-1]
7              seq2 = f.readline()[::-1]
8
9              grid = [ [Node] * (len(seq2)+1) for i in range(len(seq1)+1)]
10             grid[0][0] = Node(0, 0, 'X')
11
12             #Calculating initial score of second sequence (-1 indels)
13             for i in range(1,len(seq2)+1):
14                 node = Node(0,i,seq2[i-1])
15                 grid[0][i] = node
16                 node.calculate_score(grid)
17
18             #Calculating initial score of first sequence (-1 indels)
19             for i in range(1,len(seq1)+1):
20                 node = Node(i,0,seq1[i-1])
21                 grid[i][0] = node
22                 node.calculate_score(grid)
23
24             #putting other nodes to the grid and calculates the scores
25             for i in range(1,len(seq1)+1):
26                 for y in range(1, len(seq2) + 1):
27                     node = Node(i, y, "X")
28                     grid[i][y] = node
29                     node.calculate_score(grid)
30
31             #printing the results
32             print("")
33             print("Result of " + file)
34             backTracking(grid,seq1,seq2)
35
36
37  main()
```

Figure-4: Main Function

TEST RESULTS

Result of test1.seq

C--GAGACCGA-CG-AAGAGGTTTGCCCC-CAAC-CAGGTTCC-CTGATCACGTAACCT--ACCGGCCAAAAGGAC-T-GGCCTTA-CTAAGGCCT-TT--GTC-TACT-GC-G-G--GT-C--CGG-G-GGC-CGTT-GG-T-T-TCGGCAGAAC-
ACTTC

CCTG-GACCGAGCTTAA-A--TT-G-CTAGCAATACAGATGCCGCT--TC-C-T---TGGGGAG-GGTGTGTAGGATGTAGG--TTAACGAATGCAAGTTCGGGGTA-TCGCAGAGTCGTGCTACGGCGTGGCAC-TTAGGGTCTCTCGGGAAAAAG
AGTAG

Score: 186.5

Result of test2.seq

GTGT-GGTTGCTTGATCACTCCGTGTACATGTGACAACCGAACGAGTTGATCCAGCTTTGTTAAGTCAGCTTCGAATG-CGGTAGCTCTCAAA-TATGAT-ATGA-C-TCTTGGGGTAGATGCTG-GG--G-ACCTATTGCGCCC--AAAGCGATAT
TCGGGGC-A-CCGGTTTAGGGTACCCTATCAA-GGCGAT-AC-TTCGATGCAGTGTG-ATGCCGG-A-G-GTGTG-ATCAGCATGTGAGA-GC-T

G-GTAGGTT-CGTGAAGCACTCC-TGGAC-TCTGACCAC--AAC-A-TA-ATCAAGC--GC-A-G--AG-TT-GAATGACCAAGCGCTCAAGCT-T--TCACGAACGTCTCA--TA--TCC-GCGGTCGTACAAAC-GCGCTCTTAAAC-A-AT
-CGT--CTATCCA-TCCGGGG-A--TACCAATGG-G-TGACATTA-A--CTG-G-GCAGGCCGGCACGCG-GACGCGACAG-AT-TGAAATGGAT

Score: 337.0

Result of test3.seq

CGGG-GAAAGAC-GGA-AT--GCATCGACCA-TCGGACAAT-GCCTCAC-TGGAAGCGTG-CT-GATTTTTGCGCAA-CGAGCCTCGGA-CCTCCCG-C-TCAAACCTACGA-AA--AT-GACT-CCACCAGCACTGAACCA-AC--TG--G--CCT
CCAG-TGA-AGATA--G-TATC-T--A--CA--A-AT-C-GT-TT-G--C-CGGGA-GAAACATTT-GTATGGTAG-TCAGCGTTC-TGCACGTC-ACG-TAATCGTTCAGTAGTTGTGGGGTATACCAGGTCGTAATGAGAT-GTTAGTA-TGAATC
GTTTATAACG-CTTGTCATAGGAGTT-C---C-GAATAATCGTCACT-AGGTCAA-T-G---GCCCCCTACT-TGTAATAA-CTA-CGT-CTGATTGAGAAACAGATCGTTAGTCATCGGTTAATAGTCCCGAAGATAACGGGTTCTTGCGT--TT
-TTGCGAATACTACTATC-TCATGGCGATGG--AG-C---GT---T-TGG---TCCCAT-CCAGCC-G-CGCGAGTATCACT-TGTTGCGCTGCA-C--CTG-TCCGACCTTTCGATGG--GGAGCTC----CTTTCATTGGTT--GTAG-GTACT
-AAGGGTGAGCAATGTCACTGACCCAA-GGAGCC-GTGTGTAATTTCCACTTGCTCAGAAAAGC-C-TC-GACAAT--C-T--G-GGA-CC-GACACCTGAGTGA-TGGCTT-ACA-GACCAGGGGAGGGG-G-GCAGGTTCCCTGGCCACAGAAT
GGCACGCCCTGAGGAGGCACCGGCCAA-CGTCGC-TGG

AGTGAGAA-GACCGGATATATGCAACGAACAATCGCAAAATAGC-TC-CATGTACGC-CTGTCCCG-TAATACCGCATGCGAGCCTCGAAGCC-CCCGACGTCAAACCT-CAACAAGTATAGGGTTCCACC-GTA-TGAGCCACAGAATTCAGTTCCCT
CC-GCTGATAG-TCCCGGT-TCCTCCAATCAGCAGATTCTGACTTAGGTCAGTGAAGAAACCTGGCGTATTGT-GATGAAATTTCTGTG---GTGGACGCTAC-CGT-C--TAG--GTAGAG---CCA---C-TAGTGC-ATTGT-AC-ACTG--TC
TTTT-TC-CGGCTA-TAATCGGAGTTTCAGGCTGACT--T-GCCACCCAGGTCAAATAGTATGTC---ACGGTGTA-TGATCTCTCGGGCTGAT--AGCAAGGGG-CG--G-CGTCGGGCACT---CCTTGAAT-TGACCA-TT-TGGCGTGCCTT
GTTGCC--T-C--CT-TCCTCAAG-C-AGGGTTAGGCTTCAGTAAGTGTGGGGGTGGC-TGCC-GAAAGACG-G-GC-TCGGGGTGTT-GCC--CAACGTCTGGTC-GAC-TTCCCATATCTGG-GCACAAGACGT--AT-GGACCCGT-GCG--CT
TAATGGTG--CATT-TCACGCGAGA-AACGGAGGCAG-GTGCCAT--CCGCG-GC--GAAAGATCATCAGACA-TGACATAACGAGGAACCCGAT-CCGGAGTGAATA-CGTCACATG-CCAGGGG-GCGGTGTGCAG-TTCTC-GACCGGA-AAC
GCC-CACCC--ACGACGTACCG-CCAAAGCGTCGCCTGG

Score: 1114.5

Result of test4.seq

AGGC--CGAAAACGTCGCGAAT-T-GACCTGGCGACGCCGCCGAACGGGACCTCCGTAGT-GTGG--GA-GGTCATCAATCTCGTTCGCTAGCGGTGAC-ACCAATCACTATAAGTCTGTGATGAC-
---CTTC-AA---GTCA--AATATAGATCCTGGC--CGCT-CC-A-CGGG-C-T---TAAGTCGTTCTCCGAAGGT-A-CGATCTGGTTGGAT-GCTTCCGTCTA--AA-CA--AGAAGA-TA--AT--CG

Score: 160.5

Result of test5.seq

CGG--GTAGTTAACCTA-CAG-CATAGAGTCGCGAGATAAAGTGCAGGA-GTCTTTCGCGGCAGATTCTGACCTCAACCAGTGTCTACTTTC-TGGCA-TCACGAATC-TGCCCATAGGTCCGTGAGT-CCATATGA
AGGAAGTAGTTAGCC-TAACAGGCATAGAGTCGCGACATAT-GTGAAG-ATGTCACTTCG-GT-A--TTCAAACCTCATGCA--T-C-A-TTGCTTG-AGTCG---TCCTGGAGCATAG-TCCCTGAGTGCCATATGA

Score: 247.5

Figure-5: Results of Five Test Input

Results can be seen above. Our program works properly according to dynamic pairwise sequence alignment algorithm.