# PS5 DNA Sequence Alignment

This task required us to compare two ASCII strings in order to determine their edit distance, as well as provide space and time analysis on the program while it was running. Using a dynamic programming method, we were required to create a program that computed an optimum sequence alignment on two DNA sequences. The space complexity of our program might then be measured using Valgrind, a programming tool.

## Key concepts
The EDistance class represents a matrix of integers with a vector of vectors of type int. The function optDistance() would then use the min() and penalty() functions to fill the matrix from bottom to top with the smallest of three distances, then return the optimal distance between the two sequences. The alignment() method would then go backwards through the populated matrix, constructing the output alignment string depending on the best path.
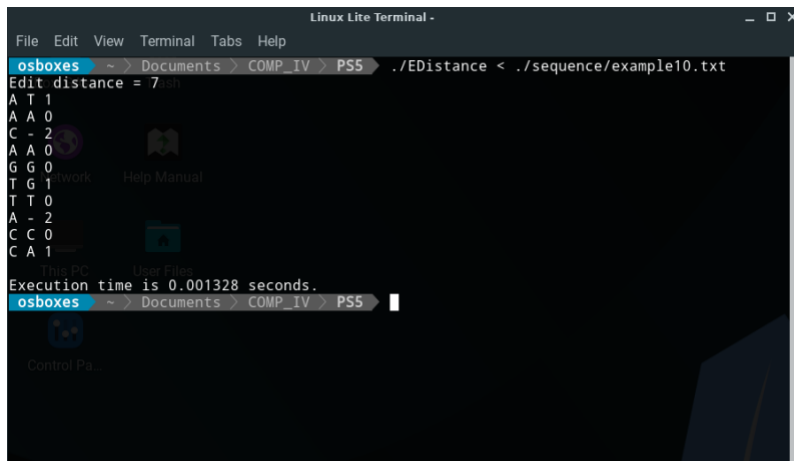
## What I Accomplished
My EDistance algorithm uses dynamic programming to fill a matrix with each computation from bottom right to top left, culminating the ideal distance at [0][0]. Tracing the matrix in reverse order from top left to bottom right also shows the alignment path. Instead of computing the same subproblem numerous times, the dynamic programming technique allows the alignment computations to be divided down into subproblems and then stored.

## What I Learned
Apart from the numerous applications of Edit Distance, I learned about the advantages of dynamic programming, which, when compared to the recursive solution for this assignment, would have had a much higher space complexity due to the number of recursive calls exceeding 2N when comparing two strings of length N. In addition to, I also gained the knowledge on how to use the Valgrind massif tool to efficiently monitor program memory use. This utility shows how much memory the heap consumed during execution. I was able to determine the projected run time and memory use of a bigger sample of strings of length N using the doubling approach.

## Output

```
 1: all:
 2:         make EDistance
 3:
 4: EDistance: main.o EDistance.o
 5:         g++ -o EDistance -g main.o EDistance.o -lsfml-system
 6:
 7: main.o: EDistance.cpp main.cpp
 8:         g++ -c -g main.cpp -ansi -pedantic -Wall -Werror
 9:
10: EDistance.o: EDistance.hpp EDistance.cpp
11:         g++ -c -g -O2 EDistance.cpp -ansi -pedantic
12:
13: clean:
14:         rm -f *.o *~ EDistance
```

```cpp
 1: #include <iostream>
 2: #include <cstring>
 3: #include <SFML/System.hpp>
 4:
 5: #include "EDistance.hpp"
 6:
 7: using namespace std;
 8:
 9: int main() {
10:
11:   // Clock
12:   sf::Clock clock;
13:   sf::Time t;
14:
15:
16:   string String1;
17:   string String2;
18:
19:   // Get input
20:   cin >> String1;
21:   cin >> String2;
22:
23:   // Initialize my class, all math done in constructor and sets relevant
24:   // member variables
25:   EDistance output(String1, String2);
26:
27:
28:   // Get desired output
29:   cout << "Edit distance = " << output.getEditDistance() << endl;
30:   cout << output.getEditString() << endl;
31:
32:   t = clock.getElapsedTime();
33:   cout << "Execution time is " << t.asSeconds() << " seconds." << endl;
34:
35:   return 0;
36: }
```

```cpp
 1: #ifndef EDistance_HPP
 2: #define EDistance_HPP
 3:
 4: #include <cstring>
 5: #include <vector>
 6:
 7: class EDistance {
 8: public:
 9:    EDistance(std::string _stringA, std::string _stringB);
10:    void printOpt() const;
11:    int getEditDistance() const;
12:    std::string getEditString() const;
13:
14: private:
15:    // Input Variables
16:    std::string stringA;
17:    std::string stringB;
18:
19:    // Constructed Variables
20:    std::vector< std::vector<int> > opt;
21:    int editDistance;
22:    std::string editString;
23:
24:    // Private Functions
25:    int optDistance();
26:    std::string alignment() const;
27:    int penalty(char a, char b);
28:    int min(int a, int b, int c);
29: };
30:
31: #endif
```

```cpp
 1: #include <iostream>
 2: #include <cstring>
 3: #include <vector>
 4: #include <sstream>
 5:
 6: #include "EDistance.hpp"
 7:
 8: using namespace std;
 9:
10: EDistance::EDistance(string _stringA, string _stringB) : stringA(_stringA
), stringB(_stringB) {
11:
12:   vector<int> temp;
13:
14:   // Populate the matrix with 0's to start
15:   for(unsigned i = 0; i < stringB.length() + 1; i++)
16:     temp.push_back(0);
17:   for(unsigned i = 0; i < stringA.length() + 1; i++)
18:     opt.push_back(temp);
19:
20:   editDistance = optDistance();
21:   editString = alignment();
22: }
23:
24: int EDistance::penalty(char a, char b) {
25:
26:   if(a == b)
27:     return 0;
28:   else
29:     return 1;
30: }
31: int EDistance::min(int a, int b, int c) {
32:
33:   if(a <= b && a <= c)
34:     return a;
35:   else if(b <= c)
36:     return b;
37:   else
38:     return c;
39: }
40:
41: int EDistance::optDistance() {
42:
43:   // Fill in the matrix with the EditDistances
44:   for(int i = opt.size() - 1; i >= 0; i--)
45:     for(int j = opt[i].size() - 1; j >= 0; j--) {
46:       if((i == opt.size() - 1) && (j == opt[i].size() - 1))
47:         opt[i][j] = 0;
48:       else if(i == opt.size() - 1)
49:         opt[i][j] = opt[i][j + 1] + 2;
50:       else if(j == opt[i].size() - 1)
51:         opt[i][j] = opt[i + 1][j] + 2;
52:       else
53:         opt[i][j] = min(opt[i+1][j+1] + penalty(stringA[i], stringB[j]),
54:                         opt[i+1][j] + 2,
55:                         opt[i][j+1] + 2);
56:     }
57:
58:   return opt[0][0];
59: }
60: string EDistance::alignment() const {
61:
62:   stringstream ss;
63:
64:   unsigned i = 0, j = 0;
```

```
 65:
 66:    while(i < opt.size() - 1 || j < opt[0].size() - 1) {
 67:      if((i < opt.size() - 1)
 68:          && (j < opt[0].size() - 1)
 69:          && (opt[i+1][j+1] <= opt[i+1][j] + 1)
 70:          && (opt[i+1][j+1] <= opt[i][j+1] + 1)) {
 71:        ss << stringA[i] << " " << stringB[j] << " " << opt[i][j] - opt[i+1
][j+1] << '\n';
 72:        i++;
 73:        j++;
 74:      }
 75:      else if(((i < opt.size() - 1) && (opt[i+1][j] <= opt[i][j+1]))
 76:              || (j == opt[0].size() - 1)) {
 77:        ss << stringA[i] << " " << "-" << " " << opt[i][j] - opt[i+1][j] <<
'\n';
 78:        i++;
 79:      }
 80:      else {
 81:        ss << "-" << " " << stringB[j] << " " << opt[i][j] - opt[i][j+1] <<
'\n';
 82:        j++;
 83:      }
 84:    }
 85:
 86:    return ss.str();
 87: }
 88:
 89: void EDistance::printOpt() const {
 90:
 91:    // Print the Matrix
 92:    for(unsigned i = 0; i < opt.size(); i++) {
 93:      for(unsigned j = 0; j < opt[i].size(); j++) {
 94:        cout.width(4);
 95:        cout << opt[i][j];
 96:      }
 97:      cout << endl;
 98:    }
 99: }
100:
101: int EDistance::getEditDistance() const {
102:
103:    return editDistance;
104: }
105:
106: string EDistance::getEditString() const {
107:
108:    return editString;
109: }
```

```
     1: /********************************************************************
     2:  *   readme
     3:  *   DNA Sequence Alignment
     4:  ********************************************************************/
     5:
     6: Name: Matthew Lorette Anaya
     7:
     8: Hours to complete assignment: 5
     9:
    10: /********************************************************************
    11:  * Explain which approach you decided to use when implementing
    12:  * (either recursive with memoization, recursive without memoization,
    13:  * dynamic programming or Hirschberg\222s algorithm). Also describe why
    14:  * you chose this approach and what its pros and cons are.
    15:  ********************************************************************/
    16:
    17: Implementation of this program was done with the use of dynamic programmi
ng
    18: and a matrix. I used the algorithm on the Princeton site in order to fill
    19: said matrix. In-order to find the alignment I used backtracking top-left
to
    20: bottom right, moving from the current matrix index to the next-lowest mat
rix
    21: index. There was a certain case where if the diagonal was 1 higher than t
he
    22: downwards or rightwards option, diagonal was still the taken rout. In any

    23: case, depending on which direction I went, I either added a gap, or both
    24: letters, and incremented i and j counters to traverse back to the bottom
    25: right of the matrix.
    26:
    27:
    28:
    29:
    30: /********************************************************************
    31:  * Does your code work correctly with the endgaps7.txt test file?
    32:  *
    33:  * This example should require you to insert a gap at the beginning
    34:  * of the Y string and the end of the X string.
    35:  ********************************************************************/
    36:
    37: Kinda confused here on what this question really is. The pdf is using exa
mple10.txt
    38: And this is asking for endgaps7.txt. Seems like there is a mix up of pdfs
 between different years of this . So I'm going to use the what the HW pdf says
as there really isn't an example to compare to
    39: otherwise. Though it also says to put this all into a folder named ps3, w
hich is definitely incorrect.
    40:
    41: Input:
    42:         Ê./EDistance < ./sequence/example10.txt
    43:
    44: Expected output:
    45:
    46:         Edit distance = 7
    47:         AT1
    48:         AA0
    49:         C-2
    50:         AA0
    51:         GG0
    52:         TG1
    53:         TT0
    54:         A-2
    55:         CC0
    56:         CA1
```

```
 57:
 58: What happened:
 59:
 60:          Edit distance = 7
 61:          A T 1
 62:          A A 0
 63:          C - 2
 64:          A A 0
 65:          G G 0
 66:          T G 1
 67:          T T 0
 68:          A - 2
 69:          C C 0
 70:          C A 1
 71:
 72:          Execution time is 0.00094 seconds.
 73:
 74:
 75:
 76:
 77: /********************************************************************
 78:  * Look at your computer\222s specs in the settings.
 79:  * How much RAM does your computer have and explain what this means?
 80:  ********************************************************************/
 81:
 82: My Mac has 16gb of RAM. Random access memory gives applications a place t
o
 83: store and access data on a short-term basis. It stores the information yo
ur
 84: computer is actively using so that it can be accessed quickly.
 85:
 86:
 87: /********************************************************************
 88:  *  For this question assume M=N. Look at your code and determine
 89:  *  approximately how much memory it uses in bytes, as a function of
 90:  *  N. Give an answer of the form a * N^b for some constants a
 91:  *  and b, where b is an integer. Note chars are 2 bytes long, and
 92:  *  ints are 4 bytes long.
 93:  *
 94:  *  Provide a brief explanation.
 95:  *
 96:  *  What is the largest N that your program can handle if it is
 97:  *  limited to 8GB (billion bytes) of memory?
 98:  ********************************************************************/
 99:
100: N^2 is the area of the matrix, the number of integer slots that need to b
e
101: filled in. 4 is the size of an integer in bytes.
102:
103: a = 4
104: b = 2
105: largest N = ~44,721
106:
107: Explanation:
108:          4 * 44,721^2 = 7,999,871,364 just shy of 8gb.
109:
110: /********************************************************************
111:  * Run valgrind if you can and attach the output file to your submission.

112:  * If you cannot run it, explain why, and list all errors you\222re seein
g.
113:  * If you can run it successfully, does the memory usage nearly match tha
t
114:  * found in the question above?
115:  * Explain why or why not.
```

```
 116: /*********************************************************************
 117:
 118: ----------------------------------------------------------------------
------
 119:   n          time(i)          total(B)   useful-heap(B) extra-heap(B)    sta
cks(B)
 120: ----------------------------------------------------------------------
-------
 121:  67  6,808,807,498    3,085,122,584    3,084,577,402         545,182
     0
 122:  68  6,870,407,882    3,146,681,624    3,146,125,562         556,062
     0
 123:  69 30,570,148,197    3,201,904,240    3,201,338,395         565,845
     0
 124:
 125: It does not, its actually quite different and I'm not entirely sure as to
why. Not
 126: Sure if I'm reading valgrind output wrong or my equation is.
 127:
 128:
 129: /*********************************************************************
 130:  *  For each data file, fill in the edit distance computed by your
 131:  *  program and the amount of time it takes to compute it.
 132:  *
 133:  *  If you get segmentation fault when allocating memory for the last
 134:  *  two test cases (N=20000 and N=28284), note this, and skip filling
 135:  *  out the last rows of the table.
 136:  *********************************************************************/
 137:
 138: data file         distance      time (seconds)
 139: ---------------------------------------------
 140: ecoli2500.txt      118                  0.125216
 141: ecoli5000.txt      160                  0.334861
 142: ecoli7000.txt      194                  0.521017
 143: ecoli10000.txt     223            1.4272
 144: ecoli20000.txt     3135           74.6052
 145: ecoli28284.txt     8394          177.645
 146:
 147: /*********************************************************************
*
 148:  *  Here are sample outputs from a run on a different machine for
 149:  *  comparison.
 150:  *********************************************************************
/
 151:
 152: data file         distance      time (seconds)
 153: ---------------------------------------------
 154: ecoli2500.txt      118          0.171
 155: ecoli5000.txt      160          0.529
 156: ecoli7000.txt      194          0.990
 157: ecoli10000.txt     223          1.972
 158: ecoli20000.txt     3135         7.730
 159:
 160:
 161:
 162: /*********************************************************************
 163:  *  For this question assume M=N (which is true for the sample files
 164:  *  above). By applying the doubling method to the data points that you
 165:  *  obtained, estimate the running time of your program in seconds as a
 166:  *  polynomial function a * N^b of N, where b is an integer.
 167:  *  (If your data seems not to work, describe what went wrong and use
 168:  *  the sample data instead.)
 169:  *
 170:  *  Provide a brief justification/explanation of how you applied the
 171:  *  doubling method.
```

```
172:  *
173:  *  What is the largest N your program can handle if it is limited to 1
174:  *  day of computation? Assume you have as much main memory as you need.
175:  **************************************************************/
176: a =
177: b =
178: largest N =
179:
180: /
181: ***************************************************************
182:  *  Did you use the lambda expression in your assignment? If yes, where
183:  *  (describe a method or provide a lines numbers)
184:  **************************************************************/
185: No
186:
187:
188:
189: ***************************************************************
190:  *  List whatever help (if any) you received from the course TAs,
191:  *  instructor, classmates, or anyone else.
192:  **************************************************************/
193:
194: N/a
195:
196: /**************************************************************
197:  *  Describe any serious problems you encountered.
198:  **************************************************************/
199:
200:
201:
202: /**************************************************************
203:  *  List any other comments here.
204:  **************************************************************/
205:
```