

Application 4: Applications to Genomics and Beyond

Algorithmic Thinking (Part2)

Shamsuddin Rehmani

July 13, 2019

Application 4 Description

In Project 4, you implemented dynamic programming algorithms for determining both global and local alignments of pairs of sequences. In this Application, we will demonstrate the utility of these algorithms in two domains. In the first part of the Application, we examine an interesting problem from genomics. (This is based on “Introduction to Computational Genomics”, by Nello Cristianini and Matthew W. Hahn). We will compare two sequences that have diverged from a common ancestor sequence due to mutation. (Mutation here includes base-pair substitution, which changes the sequence content, and insertion/deletion, which change the sequence lengths.) In the second part of the Application, we consider words that have spelling mistakes.

For the genomics part of the Application, you will load several protein sequences and an appropriate scoring matrix. For the spelling correction part of the Application, you will load a provided word list. To simplify these tasks, you are welcome to use this [provided code](#).

Comparing two proteins

In 1994, Walter Gehring and colleagues at the University of Basel carried out an “interesting” experiment: they were able to turn on a gene called *eyeless* in various places on the body of the fruit fly, *Drosophila melanogaster*. The result was astonishing - fruit flies developed that had whole eyes sprouting up all over their bodies. It turned out that the *eyeless* is a master regulatory gene - it controls a cascade that contains more than 2000 other genes. Turning it on anywhere in the body activates the cascade and produces a fully formed, but non-functioning, eye. Humans, as well as many other animals, have a slightly different version of the *eyeless* gene (that is, a similar, yet not identical sequence of the same gene).

This observation suggests that about 600 million years ago (the estimated time of divergence between humans and fruit flies) there was an ancestral organism that itself used some version of *eyeless*, and that throughout the evolution of humans and fruit flies this gene continued to be maintained, albeit while accumulating mutations that did not greatly affect its function. In particular, a substring of the *eyeless* protein of about 130 amino acids, known as the PAX domain, whose function is to bind specific sequences of DNA, is virtually identical between the human and fruit fly versions of *eyeless*.

In following questions, we compute the similarity between the human and fruit fly versions of the *eyeless* protein and see if we can identify the PAX domain.

Answer to Question 1

Score of local Alignment = 875

Local alignment of HumanEyelessProtein:

```
HSGVNQLGGVVFVNGRPLPDSTRQKIVELAHSGARPCDISRLQVSNGCVSKILGRYYETGSIRPRAIGGSK  
PRVATPEVVSKIAQYKRECPISFAWEIRDRLLESGVCTNDNIPSVSSINRVLRNLASEK-QQ
```

Local alignment of FruitflyEyelessProtein:

```
HSGVNQLGGVVFVGGRPLPDSTRQKIVELAHSGARPCDISRLQVSNGCVSKILGRYYETGSIRPRAIGGSK  
PRVATAEVLVSKISQYKRECPISFAWEIRDRLLENVCTNDNIPSVSSINRVLRNLAAQKEQQ
```

Answer to Question 2

- percentage of elements global alignment of local human vs consensus PAX domain that agree: 72.9%

- percentage of elements global alignment of local fruitfly vs consensus PAX that agree: 70.2%

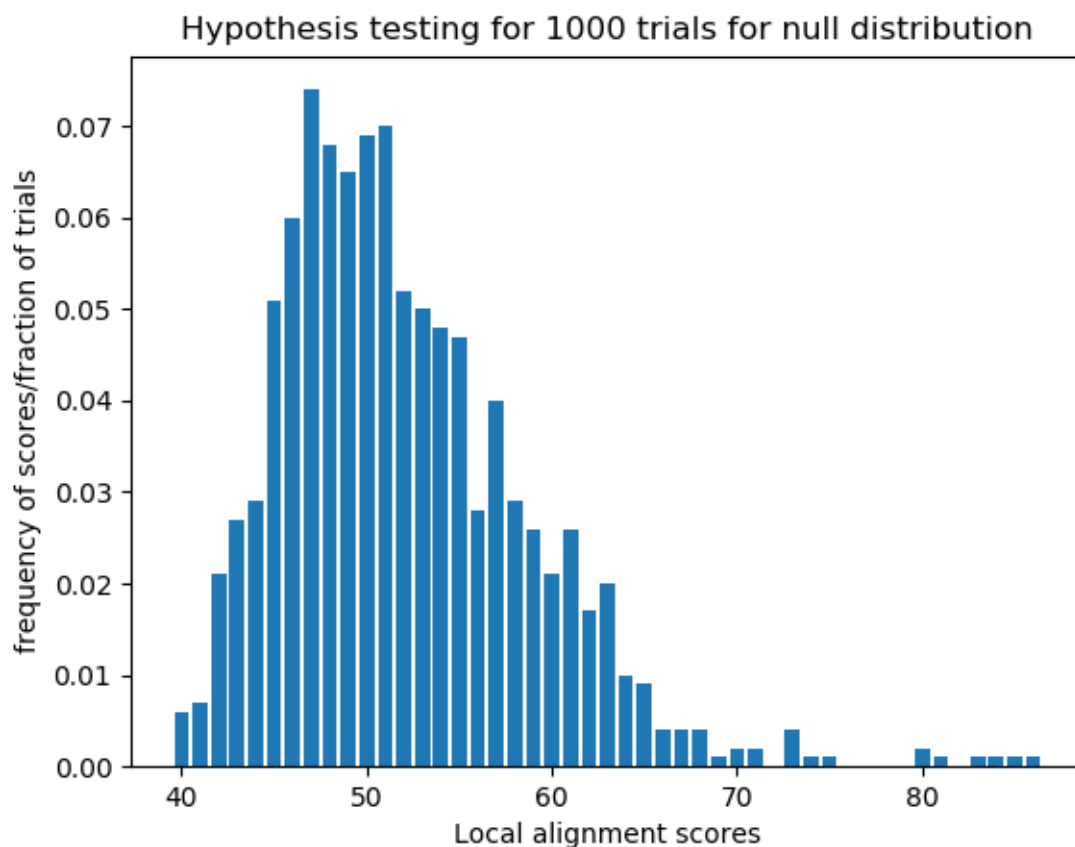
Answer to Question 3

The number of elements matching between each local alignment to the PAX consensus domain is too high to be deemed that they matched due to chance.

For local human vs PAX, 97 out of 133 elements match and the probability for all 97 of them to match due to chance would be quite small given than we have to randomly choose between 23 alphabets and that the probability of choosing an element is equal for all alphabets

The same can be said for the local fly vs PAX where 94 out of the 134 elements match

Answer to Question 4



Answer to Question 5

- mean of the distrubution: 52.0
- standard deviation of distribution: 6.72
- z-score of the distribution: 122.5

Answer to Question 6

The z-score is around 100 which means it right at the corner of the bell shaped distribution. 1σ from the mean of around 50 is about 7. And we know that 3σ covers about 99% of bell shape, i.e the likelihood of value falling within 3σ is high and the probability of value to fall outside 3σ is about 1%. The z-score we found tells us that the value of the score we found in Q1 is about 100 standard deviations away from mean (score of 875 is more than $100\sigma + \mu$). We can find the exact probability using the z test for normal distribution which gives us a probability on the order of 10^{-2000} ([source](#)) which is much smaller than the chance of winning a lottery

Answer to Question 7

- Let $x = \text{"ABC"}$ and $y = \text{" "}$. Then global alignment will return $x' = \text{"ABC"}$ and $y' = \text{"---"}$ so $|x| = 3$ and $|y| = 0$ and edit distance need 3 substitution operation to substitute " --- " to "ABC" and $score(x', y') = |x'| + |y'| - 3 = 0$. Therefore, **dash_score** = 0
- Let $x = \text{"ABC"}$ and $y = \text{"ABC"}$. Then global alignment returns the same strings and edit distance need zero operations. Hence $score(x, y) = |x| + |y| - 0 = 3 + 3 = 6$. Thus **diag_score** must be 2 for the score of x, y to be 6
diag_score = 2
- Let $x = \text{"ABC"}$ and $y = \text{"ABT"}$. Then global alignment returns the same strings with $score(x, y) = score(\text{"C"}, \text{"T"}) + 2 \times \text{diag_score}$. Hence edit distance needs substitution of "T" to "C" i.e one operation. Thus we have:
 $score(\text{"C"}, \text{"T"}) + 2 \times \text{diag_score} = |x| + |y| - 1$
off_diag_score = $3 + 3 - 1 - 2 * 2 = 1$.

Answer to Question 8

- words within 1 edit distance of string "humble":
{ "bumble", "humbléd", "tumble", "humblé", "rumble", "humblér", "humblés", "fumble", "humblý", "jumble", "mumble" }
- words within 2 edit distance of string "firefly":
{ "firefly", "tiredly", "freely", "fireclay", "direly", "finely", "firstly", "liefly", "fixedly", "refly", "firmly" }

A Python code used to answer the Application Questions

```
1 """
2 All answers to questions in Application #4 –
3 Applications to Genomics and Beyond
4 """
5
6 import random
7 import matplotlib.pyplot as plt
8 import alg_application4.provided as alg_app4_prov
9 import alg_project4.solution as alg_proj4_sol
10
11 #####
12 # Solution Q1
13
14 # Q1: First, load the files HumanEyelessProtein and FruitflyEyelessProtein
15 # using the provided code. These files contain the amino acid sequences that
16 # form the eyeless proteins in the human and fruit fly genomes, respectively.
17 # Then load the scoring matrix PAM50 for sequences of amino acids. This scoring
18 # matrix is defined over the alphabet {A,R,N,D,C,Q,E,G,H,I,L,K,M,F,P,S,T,W,Y,V,B,Z,X,-}
19 # which represents all possible amino acids and gaps (the "dashes" in the alignment).
20 # Next, compute the local alignments of the sequences of HumanEyelessProtein and
21 # FruitflyEyelessProtein using the PAM50 scoring matrix and enter the score and
22 # local alignments for these two sequences below. Be sure to clearly distinguish
23 # which alignment is which and include any dashes ('-') that might appear in the
24 # local alignment.
25
26 # load the two amino acid sequences that form the eyeless proteins in the human
27 # and fruit fly genomes
28 human_seq = alg_app4_prov.read_protein(alg_app4_prov.HUMAN_EYELESS_URL)
29 fly_seq = alg_app4_prov.read_protein(alg_app4_prov.FRUITFLY_EYELESS_URL)
30
31 # load their scoring matrix
32 scoring_matrix = alg_app4_prov.read_scoring_matrix(alg_app4_prov.PAM50_URL)
33
34 # compute the local alignments matrix for these 2 sequences
35 local_align_matrix = alg_proj4_sol.compute_alignment_matrix(human_seq, fly_seq, scoring_matrix, False)
36
37 # compute the local alignments and the score for the two sequences
38 (score_q1, align_human, align_fly) = alg_proj4_sol.compute_local_alignment(human_seq,
39                                                                           fly_seq, scoring_matrix,
40                                                                           local_align_matrix)
41
42 print "Q1 Answers: "
43 print "Score: ", score_q1
44 print "local alignment of human: ", align_human
45 print "local alignment of fly: ", align_fly
46 print
47
48 #####
49 # Solution Q2
50 # Q2: To continue our investigation, we next consider the similarity of the two
51 # sequences in the local alignment computed in Question 1 to a third sequence.
52 # The file ConsensusPAXDomain contains a "consensus" sequence of the PAX domain;
53 # that is, the sequence of amino acids in the PAX domain in any organism. In this
54 # problem, we will compare each of the two sequences of the local alignment computed
55 # in Question 1 to this consensus sequence to determine whether they correspond to
56 # the PAX domain.
57 #
58 # Load the file ConsensusPAXDomain. For each of the two sequences of the local
59 # alignment computed in Question 1, do the following:
60 #
61 # – Delete any dashes '-' present in the sequence.
62 # – Compute the global alignment of this dash-less sequence with the ConsensusPAXDomain
63 # sequence.
64 # – Compare corresponding elements of these two globally-aligned sequences (local vs.
65 # consensus) and compute the percentage of elements in these two sequences that agree.
66
67 # load the file ConsensusPAXDomain
68 pax_domain_seq = alg_app4_prov.read_protein(alg_app4_prov.CONSENSUS_PAX_URL)
69
70 # delete the dashes from the 2 local alignments in Q1
71 hum_nodash_seq = align_human.replace('-', '')
72
73 # compute the global alignments matrix for PAX domain and local human alignment seq with no dash
74 glob_align_mat_hum_pax = alg_proj4_sol.compute_alignment_matrix(pax_domain_seq, hum_nodash_seq,
75                                                                scoring_matrix, True)
76
77 # compute the global alignment matrix for the PAX domain and local fly alignment seq with no dash
78 glob_align_mat_fly_pax = alg_proj4_sol.compute_alignment_matrix(pax_domain_seq, align_fly,
79                                                                scoring_matrix, True)
80
81 (_, align_pax_hum, align_hum_nodash) = alg_proj4_sol.compute_global_alignment(pax_domain_seq,
82                                                                              hum_nodash_seq,
83                                                                              scoring_matrix,
84                                                                              glob_align_mat_hum_pax)
85
86 (_, align_pax_fly, align_fly_nodash) = alg_proj4_sol.compute_global_alignment(pax_domain_seq,
87                                                                              align_fly,
88                                                                              scoring_matrix,
89                                                                              glob_align_mat_fly_pax)
90
91 count_pax_hum = 0.0
92 count_pax_fly = 0.0
93
94
```

```

95 for charc_i, charc_j in zip(align_pax_hum, align_hum_nodash):
96     if charc_i == charc_j:
97         count_pax_hum += 1
98
99 print "Q2 Answers: "
100 print ("percentage of elements global alignment of local human vs PAX that agree: " +
101         str(round(count_pax_hum/len(align_pax_hum) * 100, 2)) + "%")
102
103 for charc_i, charc_j in zip(align_pax_fly, align_fly_nodash):
104     if charc_i == charc_j:
105         count_pax_fly += 1
106
107 print ("percentage of elements global alignment of local fly vs PAX that agree: " +
108         str(round(count_pax_fly/len(align_pax_fly) * 100, 2)) + "%")
109 print
110
111 #####
112 # Solution Q3
113 # Q3: Examine your answers to Questions 1 and 2. Is it likely that the level of
114 # similarity exhibited by the answers could have been due to chance? In particular,
115 # if you were comparing two random sequences of amino acids of length similar to
116 # that of HumanEyelessProtein and FruitflyEyelessProtein, would the level of agreement
117 # in these answers be likely? To help you in your analysis, there are 23 amino acids
118 # with symbols in the string "ACBEDGFIHKMLNQPSRTWVXYZ". Include a short justification
119 # for your answer.
120 #
121 # Ans: The number of elements matching between each local alignment to the PAX
122 # consensus domain is too high to be deemed that they matched due to chance.
123 #
124 # For local human vs PAX, 97 out of 133 elements match and the probability
125 # for all 97 of them to match due to chance would be quite small given that we have
126 # to randomly choose between 23 alphabets and that the probability of choosing
127 # an element is equal for all alphabets
128 #
129 # The same can be said for the local fly vs PAX where 94 out of the 134 elements
130 # match
131
132 print "Q3 Answers:"
133 print "length of global alignment of human vs pax: ", len(align_pax_hum)
134 print "number of elements matching in this global alignment: ", int(count_pax_hum)
135 print "length of global alignment of fly vs pax: ", len(align_pax_fly)
136 print "number of elements matching in this global alignment: ", int(count_pax_fly)
137 print
138
139 #####
140 # Solution Q4
141 #
142 # Q4: Write a function generate_null_distribution(seq_x, seq_y, scoring_matrix, num_trials)
143 # that takes as input two sequences seq_x and seq_y, a scoring matrix scoring_matrix, and a
144 # number of trials num_trials. This function should return a dictionary scoring_distribution
145 # that represents an un-normalized distribution generated by performing the following process
146 # num_trials times:
147 #
148 # - Generate a random permutation rand_y of the sequence seq_y using random.shuffle().
149 # - Compute the maximum value score for the local alignment of seq_x and rand_y using the score
150 #   matrix scoring_matrix.
151 # - Increment the entry score in the dictionary scoring_distribution by one.
152 #
153 # Use the function generate_null_distribution to create a distribution with 1000 trials using
154 # the protein sequences HumanEyelessProtein and FruitflyEyelessProtein (using the PAM50 scoring matrix).
155 # Important: Use HumanEyelessProtein as the first parameter seq_x (which stays fixed) and
156 # FruitflyEyelessProtein as the second parameter seq_y (which is randomly shuffled) when calling
157 # generate_null_distribution. Switching the order of these two parameters will lead to a slightly
158 # different answers for question 5 that may lie outside the accepted ranges for correct answers.
159 #
160 # Next, create a bar plot of the normalized version of this distribution using plt.bar in
161 # matplotlib (or your favorite plotting tool). (You will probably find CodeSkulptor too slow to
162 # do the required number of trials.) The horizontal axis should be the scores and the vertical axis
163 # should be the fraction of total trials corresponding to each score. As usual, choose reasonable
164 # labels for the axes and title. Note: You may wish to save the distribution that you compute in
165 # this Question for later use in Question 5.
166
167 def generate_null_distribution(seq_x, seq_y, scoring_matrix, num_trials):
168     """
169     Takes as input two sequences seq_x and seq_y, a scoring matrix scoring_matrix, and a
170     number of trials num_trials. This function should return a dictionary scoring_distribution
171     that represents an un-normalized distribution generated by performing the following process
172     num_trials times:
173
174     - Generate a random permutation rand_y of the sequence seq_y using random.shuffle().
175     - Compute the maximum value score for the local alignment of seq_x and rand_y using the score
176       matrix scoring_matrix.
177     - Increment the entry score in the dictionary scoring_distribution by one.
178
179     Arguments:
180         seq_x {string} — a sequence of alphabets representing amino acids
181         seq_y {string} — a sequence of alphabets representing amino acids
182         scoring_matrix {dict of dict} — the scoring matrix for the all alphabet plus "-" combination
183         num_trials {integer} — the number of trial
184
185     Returns:
186         dict — a dictionary representing the scoring distribution
187     """
188
189     # initialize the scoring distribution
190     scoring_distribution = {}

```

```

191
192     for dummy_idx in range(num_trials):
193         # generate a random permutation for the seq_y
194         rand_y = list(seq_y)
195         random.shuffle(rand_y)
196         rand_y = "".join(rand_y)
197
198         # compute the local alignments matrix for seq_x and rand_y
199         local_align_matrix = alg_proj4_sol.compute_alignment_matrix(seq_x, rand_y, scoring_matrix, False)
200
201         # compute the local alignment score for the two sequences
202         (score, _, _) = alg_proj4_sol.compute_local_alignment(seq_x, rand_y, scoring_matrix, local_align_matrix)
203
204         # increment the entry score in the scoring_distribution
205         if scoring_distribution.has_key(score):
206             scoring_distribution[score] += 1
207         else:
208             scoring_distribution[score] = 1
209
210
211     return scoring_distribution
212
213 # get the scoring distribution for 1000 trial for HumanEyelessProtein and FruitflyEyelessProtein
214 scoring_dist = generate_null_distribution(human_seq, fly_seq, scoring_matrix, 1000)
215
216 sum_of_vals = float(sum(scoring_dist.values()))
217 norm_values = [value/sum_of_vals for value in scoring_dist.values()]
218 plt.bar(scoring_dist.keys(), norm_values)
219 plt.title('Hypothesis testing for 1000 trials for null distribution')
220 plt.xlabel('Local alignment scores')
221 plt.ylabel('frequency of scores/fraction of trials')
222 plt.xlim(None, None)
223 plt.ylim(None, None)
224 plt.show()
225
226 #####
227 # Solution Q5
228 #
229 # Q5: What are the mean and standard deviation for the distribution that you
230 # computed in Question 4?
231 # What is the z-score for the local alignment for the human eyeless protein vs.
232 # the fruitfly eyeless protein based on these values?
233
234 # calculate the mean of the distribution
235 num_trials = 1000.0
236 mean_mu = sum([key * value for key, value in scoring_dist.items()])/num_trials
237
238 # calculate the standard deviation
239 sigma = sum([((key - mean_mu) ** 2) * value for key, value in scoring_dist.items()])/num_trials
240 sigma = sigma ** 0.5
241
242 # use mean and standard deviation to calculate the z score
243 z_score = (score_q1 - mean_mu)/sigma
244
245 print "Q5 Answers"
246 print "mean of the distrubution: ", mean_mu
247 print "standard deviation of distribution: ", sigma
248 print "z-score of the distribution: ", z_score
249 print
250
251 #####
252 # Solution Q6
253 #
254 # Q6: Based on your answers to Questions 4 and 5, is the score resulting from the local
255 # alignment of the HumanEyelessProtein and the FruitflyEyelessProtein due to chance? As
256 # a concrete question, which is more likely: the similarity between the human eyeless
257 # protein and the fruitfly eyeless protein being due to chance or winning the jackpot
258 # in an extremely large lottery? Provide a short explanation for your answers
259 #
260 # Ans: the z-score is around 100 which means it right at the corner of the bell
261 # shaped. 1 standard deviation from the mean of around 50 is about 7. And
262 # we know that 3 standard deviations(sigma) covers about 99% of bell shape, i.e the
263 # likelihood of value falling within 3 sigma is high and the probability of
264 # value to fall outside 3 sigma is about 1%. The z-score we found tells us that
265 # the value of the score we found in Q1 is about 100 standard away from mean
266 # (score of 875 is more than 100 * sigma + mean)
267 # We can find the exact probability using the z test for normal distribution
268 # which gives us a probability on the order of 10-2000
269 # (https://www.wolframalpha.com/input/?i=Probability+of+100+standard+deviations)
270 # which is much smaller than the chance of winning a lottery
271
272 #####
273 # Solution Q7
274 #
275 # Not surprisingly, similarity between pairs of sequences and edit distances between pairs
276 # of strings are related. In particular, the edit distance for two strings x and y can be
277 # expressed in terms of the lengths of the two strings and their corresponding similarity
278 # score as follows: |x| + |y| - score(x, y) where score(x, y) is the score returned by the
279 # global alignment of these two strings using a very simple scoring matrix that can be
280 # computed using build_scoring_matrix.
281 # Determine the values for diag_score, off_diag_score, and dash_score such that the score
282 # from the resulting global alignment yields the edit distance when substituted into the
283 # formula above
284 #
285 # Ans:
286 # 1-Let x = "ABC" and y = "". Then global alignment will return x' = "ABC" and y' = "----"

```

```

287 # so  $|x| = 3$  and  $|y| = 0$  and edit distance need 3 substitution operation to substitute "----"
288 # to "ABC" and  $\text{score}(x', y') = |x'| + |y'| - 3 = 0$ . Therefore,  $\text{dash\_score} = 0$ 
289 #  $\text{dash\_score} = 0$ 
290
291 # 2—Let  $x = \text{"ABC"}$  and  $y = \text{"ABC"}$ . Then global alignment returns the same strings and edit
292 # distance need zero operations. Hence  $\text{score}(x, y) = |x| + |y| - 0 = 3 + 3 = 6$ . Thus
293 #  $\text{diag\_score}$  must be 2 for the score of  $x, y$  to be 6
294 #  $\text{diag\_score} = 2$ 
295
296 # 3—Let  $x = \text{"ABC"}$  and  $y = \text{"ABT"}$ . Then global alignment returns the same strings with
297 # total score of score of string "C" and "T" + 2 *  $\text{diag\_score}$ . Hence edit distance needs
298 # substitution of "T" to "C" i.e one operation. Thus we have:
299 #  $\text{score}(\text{"C"}, \text{T}) + 2 * \text{diag\_score} = |x| + |y| - 1$ 
300 #  $\text{off\_diag\_score} = 3 + 3 - 1 - 2 * 2 = 1$ 
301 #  $\text{off\_diag\_score} = 1$ 
302
303 #####
304 # Solution Q8
305 #
306 # Q8: For this final question, we will implement a simple spelling correction function that
307 # uses edit distance to determine whether a given string is the misspelling of a word.
308 # To begin, load this list of 79339 words. Then, write a function
309 # check_spelling(checked_word, dist, word_list) that iterates through word_list and returns
310 # the set of all words that are within edit distance dist of the string checked_word.
311 # Use your function check_spelling to compute the set of words within an edit distance of
312 # one from the string "humble" and the set of words within an edit distance of two from the
313 # string "firefly".
314 #
315
316 def check_spelling(checked_word, dist, word_list):
317     """[summary]
318
319     Arguments:
320         checked_word {[type]} — [description]
321         dist {[type]} — [description]
322         word_list {[type]} — [description]
323
324     Returns:
325         {[type]} — [description]
326     """
327
328     # initialize the set that will contain a list of words within edit distance dist
329     # from word_list
330     within_dist = set()
331
332     # initialize a set of alphabets
333     alphabets = set(['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm',
334                     'n', 'o', 'p', 'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z'])
335
336     # compute the scoring matrix for of scores calculated in Q7
337     score_mat = alg_proj4_sol.build_scoring_matrix(alphabets, 2, 1, 0)
338
339     for word in word_list:
340         # compute the global alignment matrix
341         align_mat = alg_proj4_sol.compute_alignment_matrix(word, checked_word, score_mat, True)
342
343         # compute the global alignment score of the word
344         (score, _, _) = alg_proj4_sol.compute_global_alignment(word, checked_word, score_mat, align_mat)
345
346         # add the word to the within_dist set if it edit distance is within edit distance dist
347         if (len(word) + len(checked_word) - score) <= dist:
348             within_dist.add(word)
349
350     return within_dist
351
352 # load the word list
353 word_list = alg_app4_prov.read_words(alg_app4_prov.WORD_LIST_URL)
354
355 # compute the set of word that are within edit distance 1 of string "humble"
356 with_wrds_humble = check_spelling("humble", 1, word_list)
357
358 # compute the set of word that are within edit distance 2 of string "firefly"
359 with_wrds_firefly = check_spelling("firefly", 2, word_list)
360
361 print "Q8 Answers:"
362 print "words within 1 edit distance of string 'humble': \n", with_wrds_humble
363 print "words within 2 edit distance of string 'firefly': \n", with_wrds_firefly

```

B All functions for project 4 used in the application

```

1 """
2 All the algorithms for Project#4 — Computing
3 Alignments of Sequences
4 """
5
6 def build_scoring_matrix(alphabet, diag_score, off_diag_score, dash_score):
7     """
8     The function computes a scoring matrix(a dictionary of dictionaries)
9     whose entries are indexed by pairs of characters in alphabet plus '-'
10     The score for any entry indexed by one or more dashes is dash_score.
11     The score for the remaining diagonal entries is diag_score. Finally,
12     the score for the remaining off-diagonal entries is off_diag_score.
13
14     Arguments:

```

```

15     alphabet {set} — a set of characters
16     diag_score {integer} — the diagonal score for the scoring matrix
17     off_diag_score {integer} — the off diagonal score
18     dash_score {[type]} — the score that includes atleast one dash
19
20 Returns:
21     """
22     a dictionary of dictionary — the scoring matrix
23     """
24     # initialize the scoring matrix
25     scoring_matrix = {}
26     # add the dash character to the set of alphabets
27     all_chars = alphabet.union("-")
28
29     # create a dictionary of dictionary for scores
30     for char_i in all_chars:
31         scoring_matrix[char_i] = {}
32         for char_j in all_chars:
33             if (char_i == "-" or char_j == "-"):
34                 scoring_matrix[char_i][char_j] = dash_score
35             elif (char_i == char_j):
36                 scoring_matrix[char_i][char_j] = diag_score
37             else:
38                 scoring_matrix[char_i][char_j] = off_diag_score
39
40     return scoring_matrix
41
42 def compute_alignment_matrix(seq_x, seq_y, scoring_matrix, global_flag):
43     """
44     Takes as input two sequences seq_x and seq_y whose elements share a common
45     alphabet with the scoring matrix scoring_matrix. The function computes and
46     returns the alignment matrix for seq_x and seq_y. If global_flag is True,
47     the global alignment matrix is computed. If global_flag is False, local
48     alignment matrix is computed
49
50     For global alignment matrix S, each entry S[i][j] contains the maximum score
51     over every possible global alignment of the pair of sequences seq_x[0...i-1] and
52     seq_y[0...j-1]
53
54     For local alignment matrix S, each entry entry S[i][j] contains the maximum score
55     over every possible alignment of the pair of sequences seq_x[0...i-1] and
56     seq_y[0...j-1] except for the case when S[i][j] < 0, S[i][j] is set to zero
57
58     Arguments:
59         seq_x {string} — a string of alphabets whose elements share a common alphabet
60                         with scoring matrices
61         seq_y {string} — a string of alphabets whose elements share a common alphabet
62                         with scoring matrices
63         scoring_matrix {a dictionary of dictionaries} — the scoring matrix with each
64                                                         alphabet plus '-' combination
65                                                         scores
66         global_flag {boolean} — flag to choose between global and local alignment
67
68     Returns:
69         """
70         list of lists — an alignment matrix for seq_x and seq_y
71         """
72         # let align_matrix = S. Then initialize S[0][0] = 0
73         align_matrix = [[0]]
74
75         # set the alignment matrix's column 0 to appropriate score either based on local or global
76         # alignment selected
77         for idx in range(1, len(seq_x) + 1):
78             align_matrix.append([align_matrix[idx - 1][0] + scoring_matrix[seq_x[idx - 1]]["-"] ])
79             if not global_flag and align_matrix[idx][0] < 0:
80                 align_matrix[idx][0] = 0
81
82         # set the alignment matrix's row 0 to appropriate score either based on local or global
83         # alignment selected
84         for idx in range(1, len(seq_y) + 1):
85             align_matrix[0].append(align_matrix[0][idx - 1] + scoring_matrix["-"][seq_y[idx - 1]])
86             if not global_flag and align_matrix[0][idx] < 0:
87                 align_matrix[0][idx] = 0
88
89         # set the remaining alignment matrix(S) values based on previous 3 values of S[i-1][j-1], S[i-1][j],
90         # and S[i][j-1]
91         for idx_x in range(1, len(seq_x) + 1):
92             for idx_y in range(1, len(seq_y) + 1):
93                 val1 = align_matrix[idx_x - 1][idx_y - 1] + scoring_matrix[seq_x[idx_x - 1]][seq_y[idx_y - 1]]
94                 val2 = align_matrix[idx_x - 1][idx_y] + scoring_matrix[seq_x[idx_x - 1]]["-"]
95                 val3 = align_matrix[idx_x][idx_y - 1] + scoring_matrix["-"][seq_y[idx_y - 1]]
96                 align_matrix[idx_x].append(max(set([val1, val2, val3])))
97                 if not global_flag and align_matrix[idx_x][idx_y] < 0:
98                     align_matrix[idx_x][idx_y] = 0
99
100     return align_matrix
101
102 def compute_global_alignment(seq_x, seq_y, scoring_matrix, alignment_matrix):
103     """
104     Takes as input two sequences seq_x and seq_y whose elements share a common alphabet
105     with the scoring matrix scoring_matrix. This function computes a global alignment of
106     seq_x and seq_y using the global alignment matrix alignment_matrix. The function returns
107     a tuple of the form (score, align_x, align_y) where score is the score of the global
108     alignment align_x and align_y. Note that align_x and align_y should have the same
109     length and may include the padding character '-'.
110
111     Arguments:

```



```

111     seq_x {string} — a string of alphabets
112     seq_y {string} — a string of alphabets
113     scoring_matrix {dictionary of dictionaries} — the scoring matrix
114     alignment_matrix {list of lists} — the global alignment matrix
115
116 Returns:
117     tuple — returns (score, align_x, align_y) where score is the score of the global
118             alignment align_x and align_y.
119 """
120
121 # initialize the indicies to the lengths of the sequences
122 idx_x = len(seq_x)
123 idx_y = len(seq_y)
124 align_x = ""
125 align_y = ""
126 score = 0
127
128 while (idx_x != 0 and idx_y != 0):
129     if (alignment_matrix[idx_x][idx_y] == (alignment_matrix[idx_x - 1][idx_y - 1]
130     + scoring_matrix[seq_x[idx_x - 1]][seq_y[idx_y - 1]])):
131         align_x = seq_x[idx_x - 1] + align_x
132         align_y = seq_y[idx_y - 1] + align_y
133         score += scoring_matrix[seq_x[idx_x - 1]][seq_y[idx_y - 1]]
134         idx_x -= 1
135         idx_y -= 1
136
137     else:
138         if (alignment_matrix[idx_x][idx_y] == (alignment_matrix[idx_x - 1][idx_y]
139         + scoring_matrix[seq_x[idx_x - 1]]["-"])):
140             align_x = seq_x[idx_x - 1] + align_x
141             align_y = "-" + align_y
142             score += scoring_matrix[seq_x[idx_x - 1]]["-"]
143             idx_x -= 1
144
145         else:
146             align_x = "-" + align_x
147             align_y = seq_y[idx_y - 1] + align_y
148             score += scoring_matrix["-"][seq_y[idx_y - 1]]
149             idx_y -= 1
150
151 while idx_x != 0:
152     align_x = seq_x[idx_x - 1] + align_x
153     align_y = "-" + align_y
154     score += scoring_matrix[seq_x[idx_x - 1]]["-"]
155     idx_x -= 1
156
157 while idx_y != 0:
158     align_x = "-" + align_x
159     align_y = seq_y[idx_y - 1] + align_y
160     score += scoring_matrix["-"][seq_y[idx_y - 1]]
161     idx_y -= 1
162
163
164 return (score, align_x, align_y)
165
166 def compute_local_alignment(seq_x, seq_y, scoring_matrix, alignment_matrix):
167     """
168     Takes as input two sequences seq_x and seq_y whose elements share a common alphabet
169     with the scoring matrix scoring_matrix. This function computes a local alignment of
170     seq_x and seq_y using the local alignment matrix alignment_matrix. The function returns
171     a tuple of the form (score, align_x, align_y) where score is the score of the local
172     alignment align_x and align_y. Note that align_x and align_y should have the same
173     length and may include the padding character '-'.
174
175 Arguments:
176     seq_x {string} — a string of alphabets
177     seq_y {string} — a string of alphabets
178     scoring_matrix {dictionary of dictionaries} — the scoring matrix
179     alignment_matrix {list of lists} — the global alignment matrix
180
181 Returns:
182     tuple — returns (score, align_x, align_y) where score is the score of the global
183             alignment align_x and align_y.
184 """
185
186 # initialize the variables
187 max_value = float("-inf")
188 idx_x = -1
189 idx_y = -1
190 align_x = ""
191 align_y = ""
192 score = 0
193
194 # find the location (row, col) of the maximum value in alignment_matrix
195 for row in range(len(alignment_matrix)):
196     for col in range(len(alignment_matrix[row])):
197         if alignment_matrix[row][col] > max_value:
198             max_value = alignment_matrix[row][col]
199             idx_x = row
200             idx_y = col
201
202 while alignment_matrix[idx_x][idx_y] != 0 and idx_x != 0 and idx_y != 0:
203     if (alignment_matrix[idx_x][idx_y] == (alignment_matrix[idx_x - 1][idx_y - 1]
204     + scoring_matrix[seq_x[idx_x - 1]][seq_y[idx_y - 1]])):
205         align_x = seq_x[idx_x - 1] + align_x
206         align_y = seq_y[idx_y - 1] + align_y

```

```

207         score += scoring_matrix[seq_x[idx_x - 1]][seq_y[idx_y - 1]]
208         idx_x -= 1
209         idx_y -= 1
210
211     else:
212         if (alignment_matrix[idx_x][idx_y] == (alignment_matrix[idx_x - 1][idx_y]
213         + scoring_matrix[seq_x[idx_x - 1]]["-"])):
214             align_x = seq_x[idx_x - 1] + align_x
215             align_y = "-" + align_y
216             score += scoring_matrix[seq_x[idx_x - 1]]["-"]
217             idx_x -= 1
218
219         else:
220             align_x = "-" + align_x
221             align_y = seq_y[idx_y - 1] + align_y
222             score += scoring_matrix["-"][seq_y[idx_y - 1]]
223             idx_y -= 1
224
225     return (score, align_x, align_y)

```