

In [1]:

```
1
2 import numpy as np
3 from numba import njit
4 import edlib
5
6
7 def get_rc(s):
8     map_dict = {'A': 'T', 'T': 'A', 'G': 'C', 'C': 'G', 'N': 'N'}
9     l = []
10    for c in s:
11        l.append(map_dict[c])
12    l = l[::-1]
13    return ''.join(l)
14 def rc(s):
15     map_dict = {'A': 'T', 'T': 'A', 'G': 'C', 'C': 'G', 'N': 'N'}
16     l = []
17     for c in s:
18         l.append(map_dict[c])
19     l = l[::-1]
20     return ''.join(l)
21
22 def seq2hashtable_multi_test(refseq, testseq, kmersize=15, shift = 1):
23     rc_testseq = get_rc(testseq)
24     testseq_len = len(testseq)
25     local_lookuptable = dict()
26     skiphash = hash('N'*kmersize)
27     for iloc in range(0, len(refseq) - kmersize + 1, 1):
28         hashedkmer = hash(refseq[iloc:iloc+kmersize])
29         if(skiphash == hashedkmer):
30             continue
31         if(hashedkmer in local_lookuptable):
32
33             local_lookuptable[hashedkmer].append(iloc)
34         else:
35             local_lookuptable[hashedkmer] = [iloc]
36     iloc = -1
37     readend = testseq_len-kmersize+1
38     one_mapinfo = []
39     preiloc = 0
40     while(True):
41
42         iloc += shift
43         if(iloc >= readend):
44             break
45
46         #if(hash(testseq[iloc: iloc + kmersize]) == hash(rc_testseq[-(iloc + kmersize): -iloc])):
47             #continue
48
49         hashedkmer = hash(testseq[iloc: iloc + kmersize])
50         if(hashedkmer in local_lookuptable):
51
52             for refloc in local_lookuptable[hashedkmer]:
53
54                 one_mapinfo.append((iloc, refloc, 1, kmersize))
55
56
57         hashedkmer = hash(rc_testseq[-(iloc + kmersize): -iloc])
58         if(hashedkmer in local_lookuptable):
59             for refloc in local_lookuptable[hashedkmer]:
60                 one_mapinfo.append((iloc, refloc, -1, kmersize))
61         preiloc = iloc
62
63
64
65
66     return np.array(one_mapinfo)
67
68 def get_points(tuples_str):
69     data = []
70     num = 0
71     for c in tuples_str:
72         if(ord('0') <= c <= ord('9')):
73             num = num * 10 + c - ord('0')
74         elif(ord(',') == c):
75             data.append(num)
76             num = 0
77     if(num != 0):
78         data.append(num)
79     return data
80
81 def calculate_distance(ref, query, ref_st, ref_en, query_st, query_en):
82     A = ref[ref_st: ref_en]
83     a = query[query_st: query_en]
84     _a = rc(query[query_st: query_en])
85     return min(edlib.align(A, a)['editDistance'], edlib.align(A, _a)['editDistance'])
86
87 def get_first(x):
88     return x[0]
```

```
89
90
91 def calculate_value(tuples_str, ref, query):
92
93     slicepoints = np.array(get_points(tuples_str.encode()))
94     if(len(slicepoints) > 0 and len(slicepoints) % 4 == 0):
95         editdistance = 0
96         aligned = 0
97         preend = 0
98         points = np.array(slicepoints).reshape((-1, 4)).tolist()
99         points.sort(key=get_first)
100         for onetuple in points:
101             query_st, query_en, ref_st, ref_en = onetuple
102             if(preend > query_st):
103                 return 0
104             if(query_en - query_st < 30):
105                 continue
106             preend = query_en
107             if((calculate_distance(ref, query, ref_st, ref_en, query_st, query_en)/len(query[query_st:query_en])) > 0.1):
108                 continue
109             editdistance += calculate_distance(ref, query, ref_st, ref_en, query_st, query_en)
110             aligned += len(query[query_st:query_en])
111         return max(aligned - editdistance, 0)
112     else:
113         return 0
114
```

```
In [2]: 1 #实验1
2 ref = 'TATTATAGTCTTCATTCTGTGTATTAGATTACTAAAGCATATTACTTCTGTCTAAATGAAATTTTGTATCCTTTAGTCAGCATCTTCCCTTCTCCATCCACTTTTCTCTCCAGCCTCTGGTAATCAACATTCTACCACAATTCTGTGAGTTCTACTTTTTTAGATTCCCATGTAAGTAAGATCATGCTGTATTTGTCTTTGTGTGCCTGGCTTATCTGAGTTAAGC
3 query = 'TATTATAGTCTTCATTCTGTGTATTAGATTACTAAAGCATATTACTTCTGTCTAAATGAAATTTTGTATCCTTTAGTCAGCATCTTCCCTTCTCCATCCACTTTTCTCTCCAGCCTCTGGTAATCAACATTCTACCACAATTCTGTGAGTTCTACTTTTTTAGATTCCCATGTAAGTAAGATCATGCTGTATTTGTCTTTGTGTGCCTGGCTTATCTGAGTTAA
4
```

```
In [3]: 1 data = seq2hashtable_multi_test(ref, query, kmersize=9, shift = 1)
2 data.shape
```

Out[3]: (63746, 4)

```
In [4]: 1 #实验2
2 ref = 'TGATTTAGAACGGACTTAGCAGACATTGAAACTCGAGGGGTATAGCAATAGATGCCAAAAAGGTAAGCGCCATAAGCGTGGTTCTACGAGCCAGGTGCTCATGCCTAAGTTCTGCGCCTTCGCTGTCACTTGGAATACTGTAATGGATCATGCCTAGGTTATGCGCCTTCGGGGTCACTTCAACATACTGTAATGGATCATGCCTAGGTTTTGCGTGTTTCGCTGTCA
3 query = 'TGATTTAGAACGGACTTAGCAGACATTGAAACTCGAGGGGTATAGCAATAGATGCCAAAAAGGTAAGCGCCATAAGCGTGGTTCTACGAGCCAGGTGCTCATGCCTAAGTTCTGCGCCTTCGCTGTCACTTGGGAAATACTGTAATGGATCATCGGTAGGTTATGCGCCTTCGGGGTCACTTCAACATACTGTAACGGATCGTGCCTAGGTTTTGCGTATTCGCTGT
4
```

```
In [5]: 1 data = seq2hashtable_multi_test(ref, query, kmersize=9, shift = 1)
2 data.shape
```

Out[5]: (2347, 4)

```
In [ ]: 1
1 在这里设计你的算法
```

```
In [ ]: 1 #Design a algorithm
2 def function(data):
3     return [[q_st, q_en, r_st, r_en], ...]
```

```
In [ ]: 1
1 Result
```

```
In [ ]: 1 tuples_str = str(youfunction(data))
```

```
In [8]: 1 #Score
2 calculate_value(tuples_str, ref, query)
```

Out[8]: 2090

```
In [ ]: 1
```