

[Day-01-Lecture-07]

관련 논문

권기성, 최운호, 김동건 (2022), "문학 작품의 거리 측정을 활용한 야담의 이본 연구", 「한국고전연구 57집」 87~120쪽, 한국고전연구학회.

- 이번에는 정렬을 단락 단위로 생성하고 연구용으로 만들어 놓은 정렬 결과를 만드는 과정입니다.
- 중간 단계는 하나 생략하고 결과 중심으로 보겠습니다.

[폴더] okssw_05 > pairs_00

- 이 폴더에 제공하는 파일은 작품집 별로 대응되는 단락을 서로 짝을 맞추어 놓은 파일들입니다.
- 이 파일들에서 어떻게 문자 단위 정렬과 대응이 이루어지는지 단계별로 확인해 보도록 하겠습니다.

[폴더] okssw_05 > pairs_01

- 이 폴더에서는 edit_dist_align.py 를 활용해서 대응되는 파일의 aligned text를 생성해 냅니다.
- backtracing & alignment는 다음 사이트의 코드를 참조했습니다.

<https://giouv.dev/2016/01/minimum-edit-distance-in-python.html>

- 이 코드를 실행하기 위해서는 다음 패키지를 설치해야 합니다.

```
1 pip install numpy
2 pip install tabulate
3 pip install argparse
```

- 사용법은 다음과 같습니다.

```
1 python --input DKIM_ZKES.txt --output DKIM_ZKES.align
```

```
1 #!python
2 #
3 #!-*-code=utf-8-*-
4 #
5 #
6 # edit_dist_align.py
7 #
8 # backtracking & alignment
9 #
10 # - https://giouv.dev/2016/01/minimum-edit-distance-in-python.html
11 #
12 #
13
```

```

14 import numpy as np
15 import tabulate as tb
16 import argparse
17
18 def wagner_fischer(word_1, word_2):
19     n = len(word_1) + 1 # counting empty string
20     m = len(word_2) + 1 # counting empty string
21
22     # init. D(istance) matrix
23
24     D = np.zeros(shape=(n, m), dtype=np.int)
25     D[:, 0] = range(n)
26     D[0, :] = range(m)
27
28     # B is the backtrack matrix. At each index, it contains a triple
29     # of booleans, used as flag.
30     # If B(i, j) = (1, 1, 0) for example,
31     # the distance computed in D(i, j) came from a deletion or a
32     # substitution. This is used to compute backtracking later.
33
34     B = np.zeros(shape=(n, m), dtype=[("del", 'b'), ("sub", 'b'), ("ins",
35 'b')])
36
37     B[1:, 0] = (1, 0, 0)
38     B[0, 1:] = (0, 0, 1)
39
40     for i, l_1 in enumerate(word_1, start=1):
41         for j, l_2 in enumerate(word_2, start=1):
42             deletion = D[i-1, j] + 1
43             insertion = D[i, j-1] + 1
44             substitution = D[i-1, j-1] + (0 if l_1 == l_2 else 2)
45
46             mo = np.min([deletion, insertion, substitution])
47
48             B[i, j] = (deletion==mo, substitution==mo, insertion==mo)
49             D[i, j] = mo
50
51     return D, B
52
53 def naive_backtrace(B_matrix):
54     i, j = B_matrix.shape[0] - 1, B_matrix.shape[1] - 1
55     backtrace_idx = [(i, j)]
56
57     while (i, j) != (0, 0):
58         if B_matrix[i, j][1]:
59             i, j = i-1, j-1
60         elif B_matrix[i, j][0]:
61             i, j = i-1, j
62         elif B_matrix[i, j][2]:
63             i, j = i, j-1
64         backtrace_idx.append((i, j))
65
66     return backtrace_idx
67
68 def align(word_1, word_2, bt):

```

```

68     aligned_word_1 = []
69     aligned_word_2 = []
70     operations = []
71
72     backtrace = bt[::-1]    # make it a forward trace
73
74     for k in range(len(backtrace) - 1):
75         i_0, j_0 = backtrace[k]
76         i_1, j_1 = backtrace[k+1]
77
78         w_1_letter = None
79         w_2_letter = None
80
81         # either substitution or no_op
82         if i_1 > i_0 and j_1 > j_0:
83             if word_1[i_0] == word_2[j_0]:
84                 w_1_letter = word_1[i_0]
85                 w_2_letter = word_2[j_0]
86                 op = " "
87                 #cost increased: substitution
88             else:
89                 w_1_letter = word_1[i_0]
90                 w_2_letter = word_2[j_0]
91                 op = "s"
92         # insertion
93         elif i_0 == i_1:
94             w_1_letter = " "
95             w_2_letter = word_2[j_0]
96             op = "i"
97         # j_0 == j_1, deletion
98         else:
99             w_1_letter = word_1[i_0]
100             w_2_letter = " "
101             op = "d"
102
103         aligned_word_1.append(w_1_letter)
104         aligned_word_2.append(w_2_letter)
105         operations.append(op)
106
107     return aligned_word_1, aligned_word_2, operations
108
109 def make_table(word_1, word_2, D, B, bt):
110     #w_1 = word_1.upper()
111     #w_2 = word_2.upper()
112     w_1 = word_1
113     w_2 = word_2
114
115     w_1 = "#" + w_1
116     w_2 = "#" + w_2
117
118     table = []
119     # table formatting in emacs, you probably don't need this line
120     table.append(["<r>" for _ in range(len(w_2)+1)])
121     table.append([""] + list(w_2))
122

```

```

123     max_n_len = len(str(np.max(D)))
124
125     for i, l_1 in enumerate(w_1):
126         row = [l_1]
127         for j, l_2 in enumerate(w_2):
128             v, d, h = B[i, j]
129             direction = ("↑" if v else "") + ("↖" if d else "") + ("←" if
h else "")
130             dist = str(D[i,j])
131
132             cell_str = "{direction} {star}{dist}
{star}".format(direction=direction, star=" *"[((i,j) in bt)], dist=dist)
133             row.append(cell_str)
134             table.append(row)
135
136     return table
137
138
139 def get_lv_table(w1, w2):
140     word_1 = w1
141     word_2 = w2
142
143     D, B = wagner_fischer(word_1, word_2)
144     bt = naive_backtrace(B)
145
146     edit_distance_table = make_table(word_1, word_2, D, B, bt)
147
148     print("Edit Distance with Backtrace:")
149     print("#+ATTR_HTML: :border 2:rules all : frame border :style text-
align: right")
150     print(tb.tabulate(edit_distance_table, stralign="right",
tablefmt="orgtbl"))
151
152 def get_alignment(w1, w2):
153     word_1 = w1
154     word_2 = w2
155
156     D, B = wagner_fischer(word_1, word_2)
157     bt = naive_backtrace(B)
158
159     alignment_table = align(word_1, word_2, bt)
160
161     #print("\nAlignment")
162     #print(tb.tabulate(alignment_table, tablefmt="orgtbl"))
163
164     return alignment_table
165
166
167 def make_alignment(w1, w2):
168     word_1 = w1
169     word_2 = w2
170
171     D, B = wagner_fischer(word_1, word_2)
172     bt = naive_backtrace(B)
173

```

```

174     edit_distance_table = make_table(word_1, word_2, D, B, bt)
175     alignment_table = align(word_1, word_2, bt)
176
177     print("Edit Distance with Backtrace:")
178     print("#+ATTR_HTML: :border 2:rules all : frame border :style text-align: right")
179     print(tb.tabulate(edit_distance_table, stralign="right",
180                       tablefmt="orgtbl"))
181
182     print("\nAlignment")
183     print(tb.tabulate(alignment_table, tablefmt="orgtbl"))
184
185 def make_process(str_f_in, str_f_out):
186
187     f_out = open(str_f_out, "w", encoding="utf-8")
188
189
190     n_linecnt = 0
191     with open(str_f_in, "r", encoding="utf-8") as f_in:
192         for bline in f_in:
193             bline = bline.strip("\n")
194             cur_data = bline.split("\t")
195             str_src = cur_data[5]
196             str_target = cur_data[11]
197             bscr = cur_data[0]
198             btarget = cur_data[6]
199
200
201             print("#seqpair" + "\t" + bline, file=f_out)
202             aligned_result = get_alignment(str_src, str_target)
203             print("#seqsrc" + "\t" + "\t".join(aligned_result[0]),
204                   file=f_out)
205             print("#seqtar" + "\t" + "\t".join(aligned_result[1]),
206                   file=f_out)
207             print("#seqopr" + "\t" + "\t".join(aligned_result[2]),
208                   file=f_out)
209
210
211             n_linecnt += 1
212             print("current line: {0}".format(n_linecnt))
213
214
215     f_out.close()
216
217 if __name__ == "__main__":
218
219     # argument parsing
220     argpar = argparse.ArgumentParser(description="edit_dist_align.py --input infile_name --output outfile_name")
221
222     argpar.add_argument('--input', required=True, help="Input File Name to be aligned")
223     argpar.add_argument('--output', required=True, help="Output File Name")

```

```
222  
223     args = argpar.parse_args()  
224  
225     make_process(args.input, args.output)
```

[폴더] okssw_05 > pairs_02 > html_arc

- 이 폴더에서 실제 정렬된 파일들을 검토해 보도록 합시다.