

관련 논문

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

Levenshtein Distance

- Minimum Edit Distance
- Dynamic Programming
 - 최적의 해를 찾기 위한 동적 계획법(dynamic programming)을 사용하려고 합니다.
 - 자연어처리(nlp)에서는 tabular parsing 또는 CYK algorithm이라고도 소개되어 있습니다.
 - 동적 프로그래밍은 하나의 큰 문제를 여러 개의 작은 문제로 나누어서 그 결과를 저장하여 다음 계산에 활용하는 방식으로 계산 부담을 줄여 줍니다.
 - Levenshtein Distance 또는 Minimum Edit Distance는 Source String을 Target String으로 변환할 때 Insertion/Deletion/Substitution 연산을 적용하여 가장 적은 비용이 드는 경로를 찾는 방법을 사용합니다.
- okssw_02 폴더의 파이썬 프로그래밍을 사용하여서 두 스트링의 거리를 LD (Levenshtein Distance)를 사용하여서 적용해 봅시다.
- 분절음(segmental sound), 문자(character) 단위의 비교가 아닌 feature를 사용한다든지 함으로써 Insertion/Deletion의 Level을 feature 단위로 다루려면 distance metric을 직접 적용한 자신만의 모듈을 만들어야 합니다.
- 우리 수업에서는 이렇게 자세히 들어가지 않으니 difflib를 사용하도록 하겠습니다.

```
1 | $> pip install difflib
```

```
1  #! lv_test.py
2  #
3  #!-*-coding: utf-8 -*-
4  #
5
6  from difflib import ndiff
7
8  def lv_align(str1, str2):
9      result = ""
10     pos, removed = 0, 0
11
12     for x in ndiff(str1, str2):
13         if pos < len(str1) and str1[pos] == x[2]:
14             pos += 1
15             result += x[2]
16             if x[0] == "-":
```

```

17         removed += 1
18         continue
19     else:
20         if removed > 0:
21             removed -= 1
22         else:
23             result += "-"
24     return result
25
26 def lv_align_max(str1, str2):
27     res01 = lv_align(str1, str2)
28     res02 = lv_align(str2, str1)
29
30     len01 = len(res01)
31     len02 = len(res02)
32
33     if (len01 > len02):
34         return res01
35
36     return res02
37
38
39 def iterative_levenshtein(s, t):
40
41     rows = len(s) + 1
42     cols = len(t) + 1
43     dist = [ [0 for x in range(cols)] for x in range(rows)]
44
45     # source prefixed can be transformed into empty strings
46     # by deletions:
47
48     for i in range(1, rows):
49         dist[i][0] = i
50
51     # target prefixes can be created from an empty source string
52     # by inserting the characters
53
54     for i in range(1, cols):
55         dist[0][i] = i
56
57     for col in range(1, cols):
58         for row in range(1, rows):
59             if s[row-1] == t[col-1]:
60                 cost = 0
61             else:
62                 cost = 1
63             dist[row][col] = min(dist[row-1][col] + 1,      # deletion
64                                 dist[row][col-1] + 1,      # insertion
65                                 dist[row-1][col-1] + cost)  # substitution
66
67     return dist[row][col]
68
69
70 if __name__ == "__main__":
71

```

```
72
73     str01 = "korea"
74     str02 = "koreea"
75     print("Distance: ", iterative_levenshtein(str01, str02), sep="\t")
76     res01 = lv_align("korea", "koreea")
77     res02 = lv_align("koreea", "korea")
78     print(len(res01), res01)
79     print(len(res02), res02)
80     print("'" * 80)
81
82
83     """
84     str01 = "橘屬木從登聲"
85     str02 = "橘屬从木登聲"
86     print("Distance: ", iterative_levenshtein(str01, str02), sep="\t")
87     res01 = lv_align(str01, str02)
88     res02 = lv_align(str02, str01)
89     print(len(res01), res01)
90     print(len(res02), res02)
91     print("'" * 80)
92     """
93
94
95     """
96     str01 = "與之相狎恩情之篤如山如海"
97     str02 = "與之相狎恩情如山女如海"
98     print("Distance: ", iterative_levenshtein(str01, str02), sep="\t")
99     res01 = lv_align(str01, str02)
100    res02 = lv_align(str02, str01)
101    print(len(res01), res01)
102    print(len(res02), res02)
103    """
```