

[Day-02-Lecture-01]

관련 논문

최운호, 김동건 (2009), "군집분석 기법을 이용한 텍스트의 계통 분석-수궁가 '고고천변' 대목을 대상으로", [인문논총 62] pp. 203~229, 서울대학교 인문학연구원.

최운호, 김동건 (2012), "「춘향가」 서두 단락의 어휘 사용 유사도를 이용한 판본 계통 분류 연구", [한국정보기술학회논문지 10(4)] pp. 111~117, 한국정보기술학회.

[폴더] spr0d > spr0d_001

- spr000_002_옛글자.hwp (pdf)
 - 춘향가 서두단락이 있는 이 파일을 읽어봅시다.
 - 이 파일의 서두 단락에 사용된 어휘들을 비교하기 위해서는 다음과 같은 문제를 해결해야 합니다.
 - 표기법 정규화
 - 형태분석 > 어휘형태 추출
 - 어휘형태의 사용 양상에 따른 비교/거리 척도(distance metric) 결정
- spr001_001.txt
 - 판본별로 형태 분석을 모아놓은 결과이다.
- spr001_002.txt
 - 어휘 형태만 추출
- spr001_003.R
 - Jaccard Similarity 계산 샘플

```
1 #
2 # spr001_003.R
3 #
4
5 rm(list=ls())
6 CBSH = c('강산정기/N', '군산만학부형문/N', '나/V', '남원부/N', '녹주/N',
7 '동/N', '산수정기/N', '생겨나/V', '생기/V', '생하/V', '서/N', '서시/N', '수
8 려/N', '쌍각산/N', '악야/N', '어리/V', '왕소군/N', '저라산/N', '적/N', '적성
9 강/N', '전라좌도/N', '절대가인/N', '종출/N', '지리산/N', '춘향/N', '타/V')
10 CJJB = c('간성지장/N', '계계승승/N', '금고옥촉/N', '기운/N', '남/V', '덕/N',
11 '버금/N', '산하/N', '성자성손/N', '속종대왕/N', '시절/N', '없/V', '요순/N',
12 '용양호위/N', '우탕/N', '의관문물/N', '있/V', '조정/N', '존비/N', '좌우보
13 필/N', '주석지신/N', '즉위/N', '피/V', '향곡/N', '흐르/V')
14 CJKS = c('강산정기/N', '군산만학부형문/N', '금강활아미수/N', '나/V', '남원
15 부/N', '녹주/N', '동/N', '산수정기/N', '생기/V', '생장/N', '서/N', '서시/N',
16 '설도/N', '수려/N', '쌍각산/N', '악야계/N', '어리/V', '왕소군/N', '저라산/N',
17 '적성강/N', '절대가인/N', '제/N', '종출/N', '지리산/N', '춘향/N', '타/V', '태
18 어나/V', '호남좌도/N', '환출/N')
```

```

9  CJSH = c('곳곳이/B', '관왕묘/N', '나/V', '남녀/N', '남북강성/N', '남원부/N',
'당당하/V', '대방국/N', '동/N', '만고충신/N', '모시/V', '북통운암/N', '산수정
기/N', '서/N', '수/N', '승지/N', '아니/B', '어리/V', '옛날/N', '일색/N',
'있/V', '적성강/N', '지리산/N', '충렬/N', '하/V', '호남좌도/N')
10 CKSH = c('강산정기/N', '군산만학부형문/N', '금강활아미수/N', '나/V', '남원
부/N', '녹주/N', '동/N', '산수정기/N', '생기/V', '생장/N', '서/N', '서시/N',
'설도/N', '수려/N', '쌍각산/N', '약야계/N', '어리/V', '왕소군/N', '저라산/N',
'적성강/N', '절대가인/N', '제/N', '종출/N', '지리산/N', '춘향/N', '타/V', '태
어나/V', '호남좌도/N', '환출/N')
11 CKYR = c('숙종대왕/N', '즉위/N', '초/N')
12 CKYS = c('강산정기/N', '군산만학부형문/N', '금강활아미수/N', '나/V', '남원
부/N', '동/N', '문군/N', '산수정기/N', '생겨나/V', '생기/V', '서/N', '설
도/N', '생겨나/N', '어리/V', '영웅열사/N', '왕소군/N', '우리나라/N', '적성
강/N', '절대가인/N', '제/N', '지리산/N', '춘향/N', '타/V', '호남좌도/N', '환
생/N')
13 CLSU = c('고요직설/N', '나/V', '때/N', '법/N', '숙종/N', '시절/N', '아동
국/N', '여상/N', '요순시절/N', '은주/N', '이윤/N', '있/V', '자고로/B', '진실
로/B', '충신/N', '현성지국/N', '홍모우순풍/N')
14 CPBS = c('간성지장/N', '계계승승/N', '금고옥촉/N', '기운/N', '남/V', '덕
화/N', '산하/N', '성자성손/N', '숙종대왕/N', '시절/N', '없/V', '요순/N', '용
왕후위/N', '있/V', '조정/N', '존비/N', '즉위/N', '피/V', '향곡/N', '흐르/V')
15
16 # #####
17 # Sample 01: sim 0 == dist 1
18 # #####
19
20 intersect(CBSH, CJJB)
21 union(CBSH, CJJB)
22 length(intersect(CBSH, CJJB))
23 length(union(CBSH, CJJB))
24 jac_sim = length(intersect(CBSH, CJJB)) / length(union(CBSH, CJJB))
25 jac_sim
26
27
28 # #####
29 # Sample 02: sim 1 == dist 0
30 # #####
31 intersect(CJKS, CKSH)
32 length(intersect(CJKS, CKSH))
33 union(CJKS, CKSH)
34 length(union(CJKS, CKSH))
35 jac_sim = length(intersect(CJKS, CKSH)) / length(union(CJKS, CKSH))
36
37 # #####
38 # Sample 03: sim = 0.5428571
39 # #####
40 intersect(CJKS, CKYS)
41 length(intersect(CJKS, CKYS))
42 union(CJKS, CKYS)
43 length(union(CJKS, CKYS))
44 jac_sim = length(intersect(CJKS, CKYS)) / length(union(CJKS, CKYS))
45 jac_sim
46 jac_dist = 1 - jac_sim
47 jac_dist

```

[폴더] spr0d > spr0d_002

- spr002_001.txt
 - header가 있는 Similarity Score 파일
- spr002_002.txt
 - header를 삭제한 Similarity Score 파
- spr002_003.py
 - okss_004_007.py 를 변경해서 만든 파일.
 - Jac. Sim.는 유사도 척도이기 때문에 이 값을 distance로 변경하는 코드로 수정한다.
 - [Line 45]에서 Sim.를 Dist.로 변경하는 코드로 수정.

```
1  #!python
2  #
3  #!-*-coding=utf-8-*-
4  #
5  # okss_004_007.py > spr002_003.py
6  #
7
8  from tqdm import tqdm
9
10 if __name__ == "__main__":
11
12     with open("spr002_002.txt", "r", encoding="utf-8") as f_in:
13         m_lines = [l.strip() for l in f_in.readlines()]
14
15         dic_book_code = {}
16         # 작품 ID 구하기
17         print("Preprocessing book codes.")
18         for line in tqdm(m_lines):
19
20             cur_elts = line.split("\t")
21             if cur_elts[0] in dic_book_code.keys():
22                 dic_book_code[cur_elts[0]] += 1
23             else:
24                 dic_book_code[cur_elts[0]] = 1
25
26             if cur_elts[1] in dic_book_code.keys():
27                 dic_book_code[cur_elts[1]] += 1
28             else:
29                 dic_book_code[cur_elts[1]] = 1
30
31         tbl_size = len(dic_book_code.keys())
32         tbl_dist = [ [0] * tbl_size for _ in range(tbl_size)]
33         list_keys = list(dic_book_code.keys())
34
35         for idx_i in range(len(list_keys)):
36             dic_book_code[list_keys[idx_i]] = idx_i
37
```

```

38
39     # Table Index 안에 값 채워 넣기.
40     print("\nFilling the distance values into the matrix.")
41     for line in tqdm(m_lines):
42         cur_elts = line.split("\t")
43         elt_01 = cur_elts[0]
44         elt_02 = cur_elts[1]
45         elt_dist = 1 - float(cur_elts[2])
46         idx_01 = dic_book_code[elt_01]
47         idx_02 = dic_book_code[elt_02]
48         tbl_dist[idx_01][idx_02] = elt_dist
49         tbl_dist[idx_02][idx_01] = elt_dist
50
51     # Table 출력하기
52     print("\nPrinting the table to the output file.")
53
54     with open("spr002_003.txt", "w", encoding="utf-8") as f_out:
55         # Header 출력
56         str_header = '\t'.join(list_keys)
57         print("", str_header, sep="\t", file=f_out)
58         for idx_i in tqdm(range(len(tbl_dist))):
59             list_line = list(map(str, tbl_dist[idx_i]))
60             str_key = list_keys[idx_i]
61             str_line = "\t".join(list_line)
62             print(str_key, str_line, sep="\t", file=f_out)
63

```

R Code (MDS, HClust)

```

1  #
2  #
3
4  #install.packages("vegan", "amap", "scatterplot3d")
5
6  require(ape)
7  require(MASS)
8  require(graphics)
9  require(vegan)
10 require(amap)
11 require(scatterplot3d)
12
13 options(digits=22)
14
15 setwd('d:/current_work/kwonks_drill/sprod/sprod_002')
16
17 kdkjac = read.table('spr002_003.txt')
18 kdkjac.d = as.dist(kdkjac)
19
20 kdkjac.mds1 = cmdscale(kdkjac.d, k=1)
21 kdkjac.mds2 = cmdscale(kdkjac.d, k=2)

```

```

22 kdkjac.mds3 = cmdscale(kdkjac.d, k=3)
23 kdkjac.mds4 = cmdscale(kdkjac.d, k=4)
24 kdkjac.mds5 = cmdscale(kdkjac.d, k=5)
25
26 plot(kdkjac.mds2[, 1], kdkjac.mds2[, 2], "p")
27
28 #####
29 # HC (Cophenetic Cor.)
30 #####
31
32 kdk_hc.average = hclust(kdkjac.d, method="average")
33
34 #####
35 # Hierarchical clustering (average method)
36 #####
37 x11(72, 48)
38 plot(kdk_hc.average, hang=-1, cex=.8, main="Hierarchical clustering
(average)\nJaccard similarity (Distance) Index", sub="", xlab="work IDs",
ylab="Distance (0~1)")
39
40
41 #####
42 # Clustering 16 clusters
43 #####
44
45 x11(72, 48)
46 plot(kdk_hc.average, hang=-1, cex=.8, main="Cut into 16 clusters", sub="",
xlab="work IDs", ylab="Distance (0~1)")
47 rect.hclust(kdk_hc.average, k=16, border="red")
48
49 #####
50 #
51 # Unrooted Tree
52 #
53 #####
54
55 kdk_hc.tree = as.phylo.hclust(kdk_hc.average)
56 x11(64, 64)
57 plot(kdk_hc.tree, type="u", cex=0.8, font=3)
58
59 #####
60 # Left-Justified Clade
61 #####
62 kdk_hc.tree = as.phylo.hclust(kdk_hc.average)
63 x11(64, 64)
64 plot(kdk_hc.tree, type="c", use.edge.length=FALSE, direction="l", cex=.8)

```