


```
~/workingDirectory
> run DEBUGPY_LAUNCHER_PORT=58891 /usr/local/bin/python3 /Users/dabeen/.vscode/extensions/ms-python.python-2020.3.71659/pythonFiles/lib/python/debugpy/wheels/debugpy/launcher /Users/dabeen/workingDir
ectory/bigdata/bigdata_assignment04/src/similarity0tk.py
Enter a sentence >> 반성은 자신이 저지른 죄의 대가를 조형히 견히하게 받아들이는 데서 출발한다
Enter the number of sentences to extract >> 10
다음 반성은 자신이 저지른 죄의 대가를 조형히 견히하게 받아들이는 데서 출발한다"며 "어린 나이에 참혹한 삶을 살다 하마하게 떠나간 딸을 생각하면 자신에게 주어진 형이 무겁다고 생각할지 과연 의문이 든다"
고 A씨와 B씨를 꾸짖었다. , 100.0%
정국대기 사드 47기의 추가 배치 사실에 대해 최초로 인지하게 된 과정에 대해서 묻 수석은 "26일 정의용 안보실장이 국방부 정책실장으로부터 보고를 받았으나 식언치 않은 점들이 있었다"며 "이에 이상철 안보 1차장
이 보고에 참석했던 관계자 한 명을 보고가 한한 관란 하 자신의 사무실은 따로 불러 세부적 내용을 하나하나 확인하던 중 사드 47기의 추가 배치 사실을 최초로 인지하게 된다"고 밝혔다. , 37.5%
박 대통령은 "국민연금은 전 국민이 관련이 되고 미래 세대의 복지과 소득에 영향을 큰 사안으로 관계각층의 의견수렴과 국민적 공감대가 무엇보다 중요하다"며 "지금은 지난 1년여 동안 충분한 논의를 통해서 국민
적 공감대가 형성된 공무원연금 개혁을 마무리 하는 것이 급선무이고 국민연금과 관련된 사항은 충분한 시간을 갖고 사회적 논의를 통해 신중히 결정할 사안이라고 생각한다"고 '선 공무원연금 개혁안 처리 후 국민
연금 논의' 방침을 재확인했다. , 37.5%
도둑과 인종분류는 "615 공동선언의 채택은 장장 반세기 이상 지속되어온 불신과 대결의 역사에 종지부를 찍고, 민족의 화해와 단합, 통일과 평화변명의 새 시대를 열어놓은 특대사변이었다"고 평가하고 "거기에는
북남관계를 개선하고 조국통일을 이룩하는 데서 일관하게 견지해야 할 원칙과 모든 분야에서 협력과 교류를 전면적으로 확대 발전시켜나가기 위한 가장 합리적이고 건설적인 방도들이 다 들어있다. , 37.5%
서 시장은 자신의 저서 '일하는 사람이 미래를 말한다'에서 "평화에서 벗어나 마추치거나 멀리서 보게될 때도 계념 있었다. , 37.5%
하도금 (단락, 세정구조 보충강화에 대해서는 강대도금 모두 "실로성 없는 성격식기음 대적"이라며 "대중소기업단 남용단가 합상 대응성 확보, 불공정한 하도금 대가지금 등의 불공정행위에 대해서 제재를 강화할 수
있는 정부 대책이 있어야 한다"고 주장했다. , 37.5%
한 장관은 북한에 대해서는 "만약 적이 제2평행선처럼 무모하게 도발한다면 그동안 수없이 헌명한대로 적 도발원점을 물론, 지원세력과 지휘체력까지 단호하게 응징에 도발의 대가를 뼈저리게 느끼도록 해
야 한다"고 지시했다. , 37.5%
강 부대변인은 그러면서 "1,20000 평가도면이 남 지사를 도지사로 뽑은 이유는 도정을 잘 돌보라는 것이지, 자신의 대권 욕심을 채우라는 것이 아님을 똑똑히 깨달아야 한다"며 "남 지사는 자기인식이 결여된 무책임
한 비난을 멈추고 지금이라도 결히 반성하는 자세로 도정에 충실하길 바란다. , 37.5%
홍 대표는 서 의원이 자신의 당원 자격이 없다고 주장한 데 대해서는 "나는 지난 대선 때 당의 요청대로 정계에서 신임을 받고 당은 대법원 확정판결이 난 때까지 '당원권 정지'를 정지해 현재 당원 신분을 갖고 있
다"며 "자신의 부정을 숨기기 위해 나를 억울하게 누명을 받은 사건에 대해 나에게 사과하고 반성은 하지 않고 그것을 빙기해 나의 당원권 사비를 운운하는 것은 참으로 추악무치한 반발"이라고 말했다. , 37.5%
최근에서 "기존의 경제동상 분야를 넘어 고부가가치 창출의 파트너십이 양국 간 미래협력의 방향이 되어야 한다 는 데 의견을 같이 했다"며 "이런 측면에서 이번에 양국 간 정보통신, 보건 의료, 천문우주 분야 MOU가
제결되고 양국 정책대화를 개시하게 된 점을 기쁘게 생각한다"고 말했다. , 37.5%
elapsed Time : 457.3200578689575s
```

구현 방법

- 형태소 분석 방식에 의한 유사도 측정

아래와 같은 순서로 형태소 분석에 의한 유사도 측정 알고리즘을 구현하였다.

1. 텍스트 파일을 한 문장씩 리스트에 저장해 반환한다.
2. 한 문장씩 형태소를 분석한다.
3. collections.Counter를 이용하여 각 형태소에 대한 빈도를 측정한다.
4. 형태소 갯수가 적은 문장을 기준으로 공통 형태소 갯수를 측정한다.
5. 입력받은 문장의 형태소 갯수로 공통 형태소 갯수를 나누고 100을 곱하여 공통 형태소 유사도를 계산하여 sentenceDict에 저장한다.
6. 유사도가 높은 순서대로 정렬 한 후 사용자 입력으로 받은 n 만큼 문장을 출력한다.

- 1을 구현한 openFile

```
def openFile(filePath):
    lines = None
    with open(filePath) as f:
        try:
            lines = f.read().splitlines()
        except:
            f.close()
    return lines
```

- 2를 구현한 코드

```
for line in sentenceList:
    if line == '':
        continue
    sentencePos = han.pos(line)
```

- 3을 구현한 코드

```
sentencePosCount = Counter(sentencePos)
```

- 4를 구현한 코드

```
common = 0
for morpheme in inputPosCount:
    if morpheme in sentencePosCount:
        common += min(inputPosCount[morpheme], sentencePosCount[morpheme])
```

- 5를 구현한 코드

```
similarity = 100 * common / inputLen
sentenceDict[line] = similarity
```

- 6을 구현한 코드

```
hanResult = konlpyHannanum(sen, sentenceList)
hanResult = sorted(hanResult.items(), key=lambda x: x[1], reverse=True)

for i in range(Num):
    print(f'{hanResult[i][0]} : {hanResult[i][1]}%')
```

문제점

- 공통사항 : sample.txt (15000줄)
- 실행시간 : Komoran > Okt > Hannanum

1. Komoran KoNLPy 패키지의 경우 다른 형태소분석기보다 실행시간이 적게 걸렸으나 다르게 같은 입력 문장이더라도 105%의 유사도를 반환하였다.

```
oktResult = [('반성은', 'Noun'), ('자신', 'Noun'), ('이', 'Josa'), ('저지르', 'Verb'), ('죄', 'Noun'), ('의', 'Josa'), ('대가', 'Noun')
komoResult = [('반성', 'NNG'), ('은', 'JX'), ('자신', 'NNG'), ('이', 'JKS'), ('저지르', 'VV'), ('ㄴ', 'ETM'), ('죄', 'NNG'), ('의', 'J')
hanResult = [('반성', 'N'), ('은', 'J'), ('자신', 'N'), ('이', 'J'), ('저지르', 'P'), ('ㄴ', 'E'), ('죄', 'N'), ('의', 'J'), ('대가', 'N')]
```

2. 느린 실행시간