



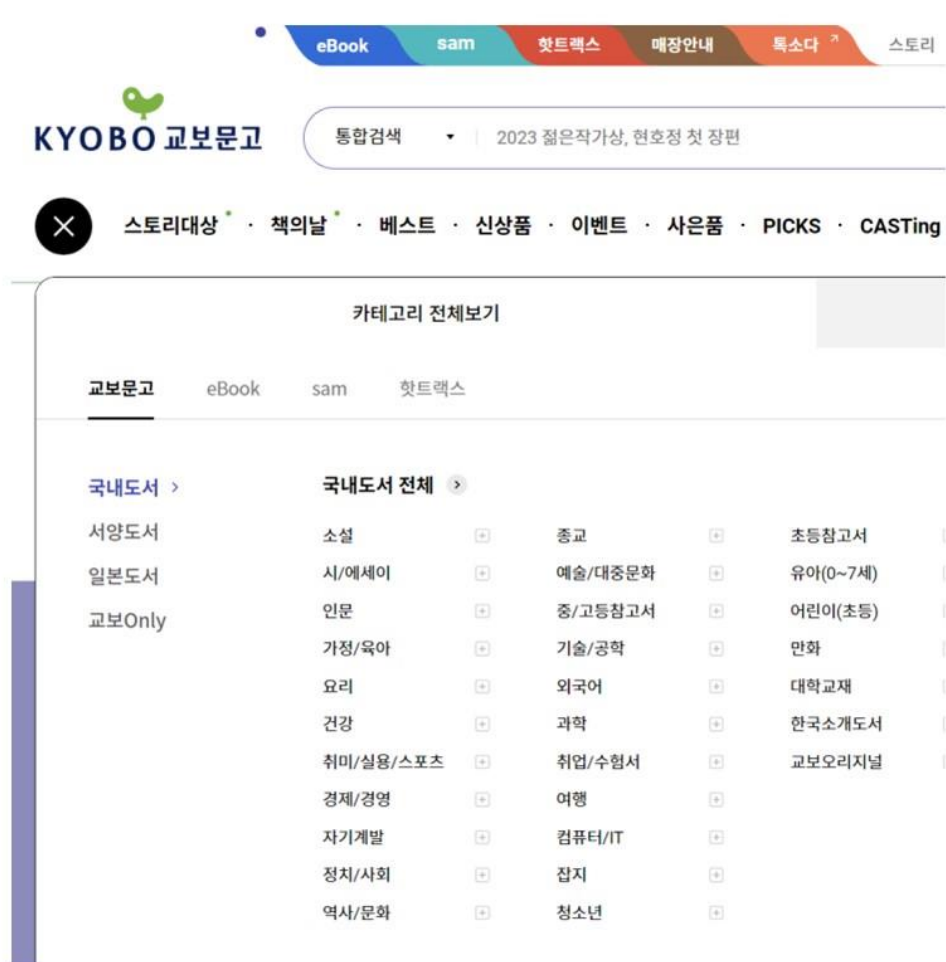
빅데이터 분석을 통한 트렌드 파악 및 사용자 맞춤 도서 추천

윤경서, 강승식
국민대학교 시빅데이터융합경영학과, 국민대학교 소프트웨어학부



학부생 경진대회 논문

02. 데이터 수집 및 전처리



데이터 수집

교보문고에서 국내도서 카테고리 중 20개를 선정하여 해당 카테고리별 베스트셀러 240권씩

제목, 저자, 내용요약, 카테고리들

Selenium을 사용한 크롤링으로 데이터 수집을 진행

= 총 4800권의 데이터 수집

전처리

re 라이브러리의 정규표현식을 사용하여 전처리 진행

03. 형태소 분석

KoNLTK	KoNLPy	Stanza
KLT2023	OKT	nouns 함수
4800권의 도서에 대해 소요시간 : 31초 약 45만개의 명사 추출	4800권의 도서에 대해 소요시간 : 1분 33초 약 54만개의 명사 추출	4800권의 도서에 대해 소요시간 : 22분 29초 약 30만개의 명사 추출

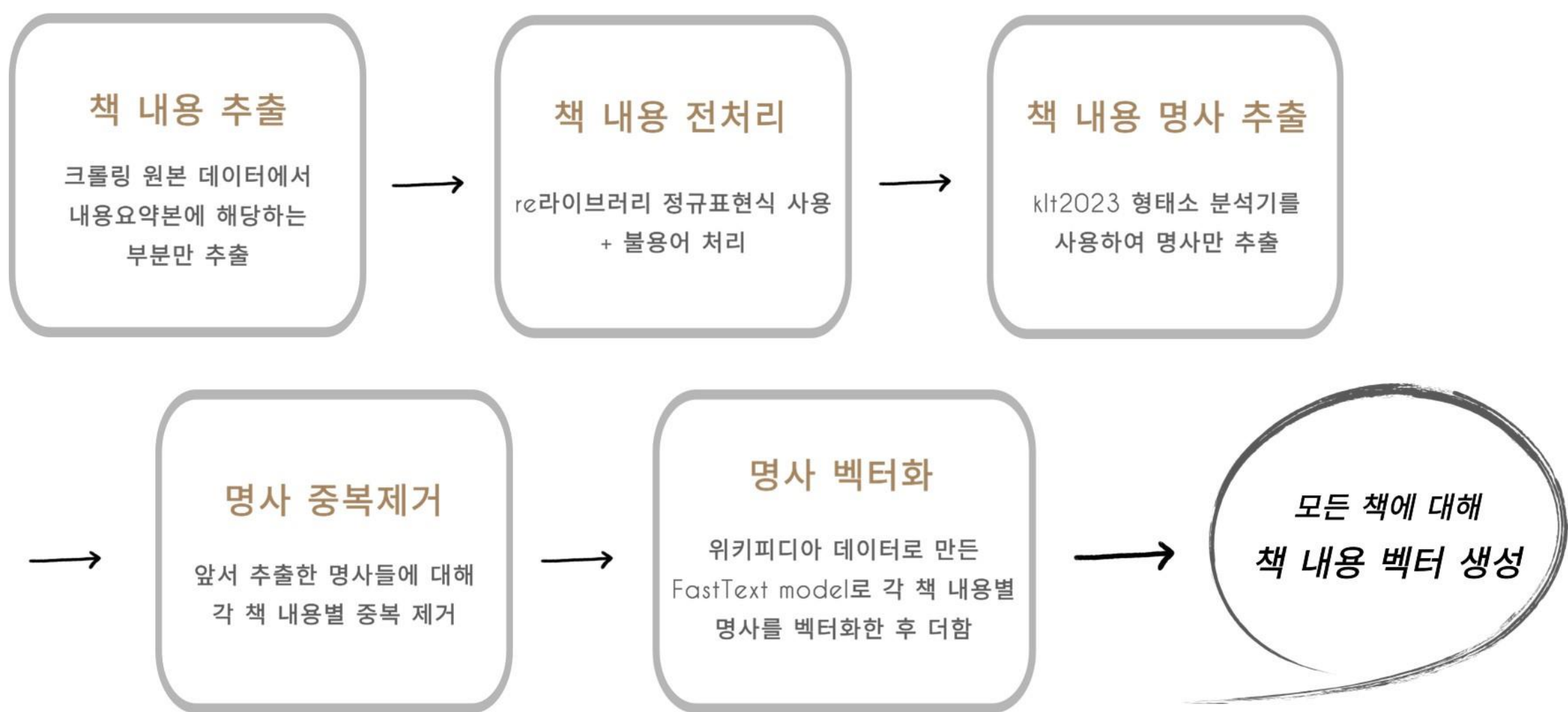
KoNLTK로 명사만
추출한
위키피디아 데이터
ko_wiki_KMA_uff8.txt

전처리

- re라이브러리 정규표현식
- 불용어 처리

약 3006만개의 명사로
Word2Vec model
& FastText model
생성

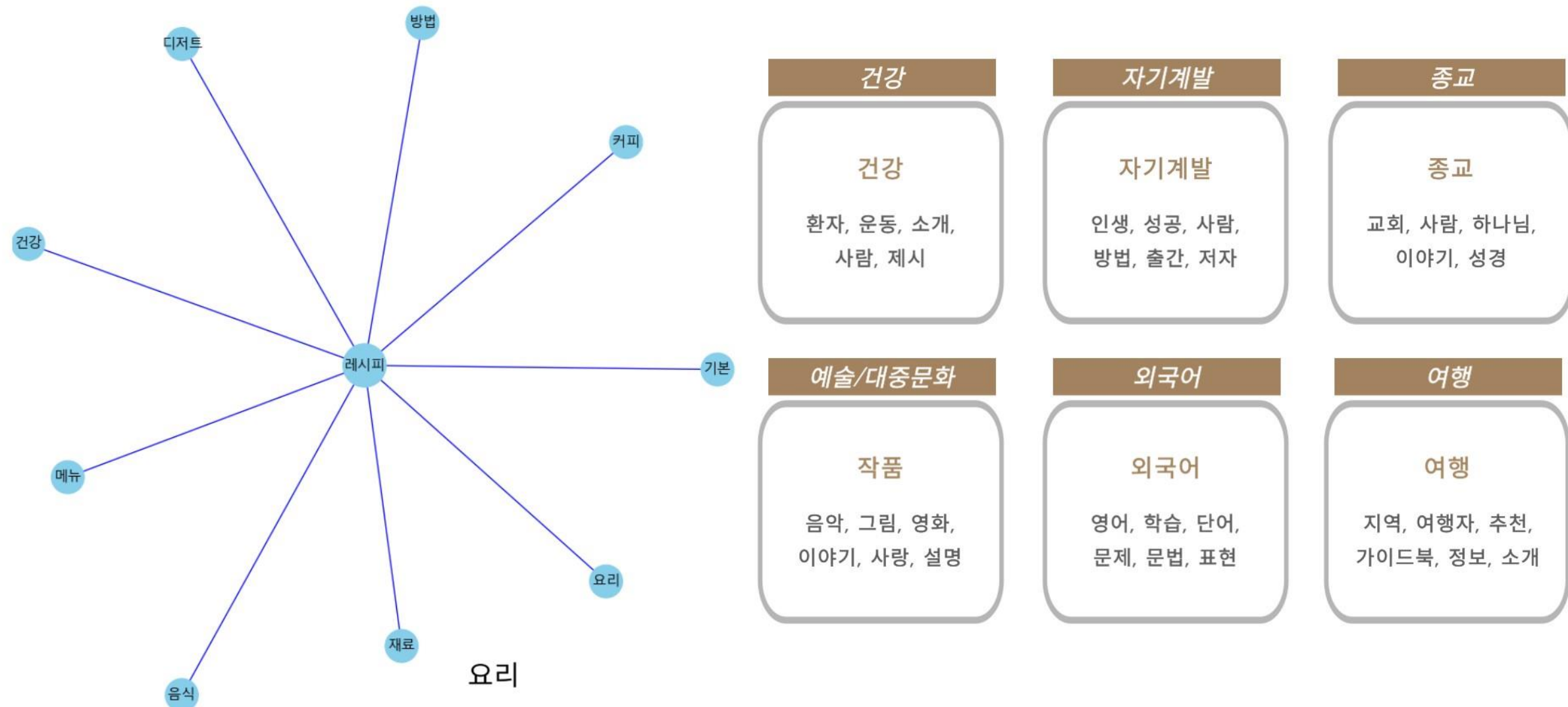
04. 책 내용 벡터 생성



05. 워드클라우드 CountVectorizer



06. 키워드 시각화 plot4kcc_keyword



07. 임베딩 모델 활용 도서 추천

책 내용으로 추천

사용자가 문장을 입력하면 해당 문장과 책 내용 벡터의 유사도를 구해 유사도가 가장 높은 책 상위n권을 추천해준다.

마음을 편안하게
해주는 위로의 글

Word2Vec

[[('행복을 담아줄게', '나란'), 0.3582575959687917],
[('있는 그대로', '김지훈'), 0.3402578870821218],
[('보이지 않는 곳에서 애쓰고 있는 나에게', '최대호'), 0.33851055987823103],
[('당신과 아침에 싸우면 밤에는 입맞춤 겁니다.', '유래혁'), 0.3367892482086977],
[('엄마의 말 공부 일력 365', '이일숙'), 0.3345570815563324]]

FastText

[[('당신이 좋아지면, 밤이 깊어지면', '이희연'), 0.7639079225989486],
[('나로서 충분히 괜찮은 사람', '김재식'), 0.761079103235505],
[('당신과 아침에 싸우면 밤에는 입맞춤 겁니다.', '유래혁'), 0.7595394689581971],
[('나는 당신이 행복했으면 좋겠습니다.', '박찬휘'), 0.7582670581220978],
[('그대 뻐가하는 것이 아니라 의욕하는 것이다', '오명선'), 0.7533357641240881]]

단어 조합으로 추천

사용자가 positive word와 negative word를 입력하면 해당 단어들의 조합으로 유사 단어를 추출한 뒤, 유사 단어들의 벡터 합과 책 내용 벡터간의 유사도를 구해 유사도가 가장 높은 상위 1권의 책을 추천해준다.

Word2Vec

positive_word (여러 단어를 쓸 경우, 띄어쓰기로 구분해주세요.): 모험 전쟁 성공
negative_word (여러 단어를 쓸 경우, 띄어쓰기로 구분해주세요.): 공주 연애
[[('이미 시작된 전쟁', '이철'), '정치/사회'], 0.27360490246192454]]

FastText

positive_word (여러 단어를 쓸 경우, 띄어쓰기로 구분해주세요.): 모험 전쟁 성공
negative_word (여러 단어를 쓸 경우, 띄어쓰기로 구분해주세요.): 공주 연애
[[('로마인 이야기 2: 한니발 전쟁', '시오노 나나미'), '역사/문화'], 0.5823006581838892]]

07. 임베딩 모델 활용 도서 추천

카테고리로 추천

사용자가 단어를 입력하면 해당 단어와 각 카테고리에 해당하는 단어간의 유사도를 구한 뒤, 가장 유사도가 높은 카테고리의 책들의 내용벡터와 입력 단어 벡터와의 유사도를 구해 유사도가 높은 상위 n권을 추천해준다.

성공

상위 5개

(0.04292066,
'대중문화'),
[[('SAVE THE CAT!: 흥행하는 영화 시나리오의 8가지 법칙', '블레이크 스나이더'), 0.44993079561155186],
[('우리는 왜 일명중을 사랑하는가', '조위'), 0.34890712733006674],
[('시나리오 어떻게 쓸 것인가 세트', '로버트 맥키'), 0.343034047968187],
[('데이비드 호크니', '마르코 리빙스턴'), 0.327195837798465],
[('연기하지 않는 연기', '해럴드 거스킨'), 0.31482199930558946]]]

모험

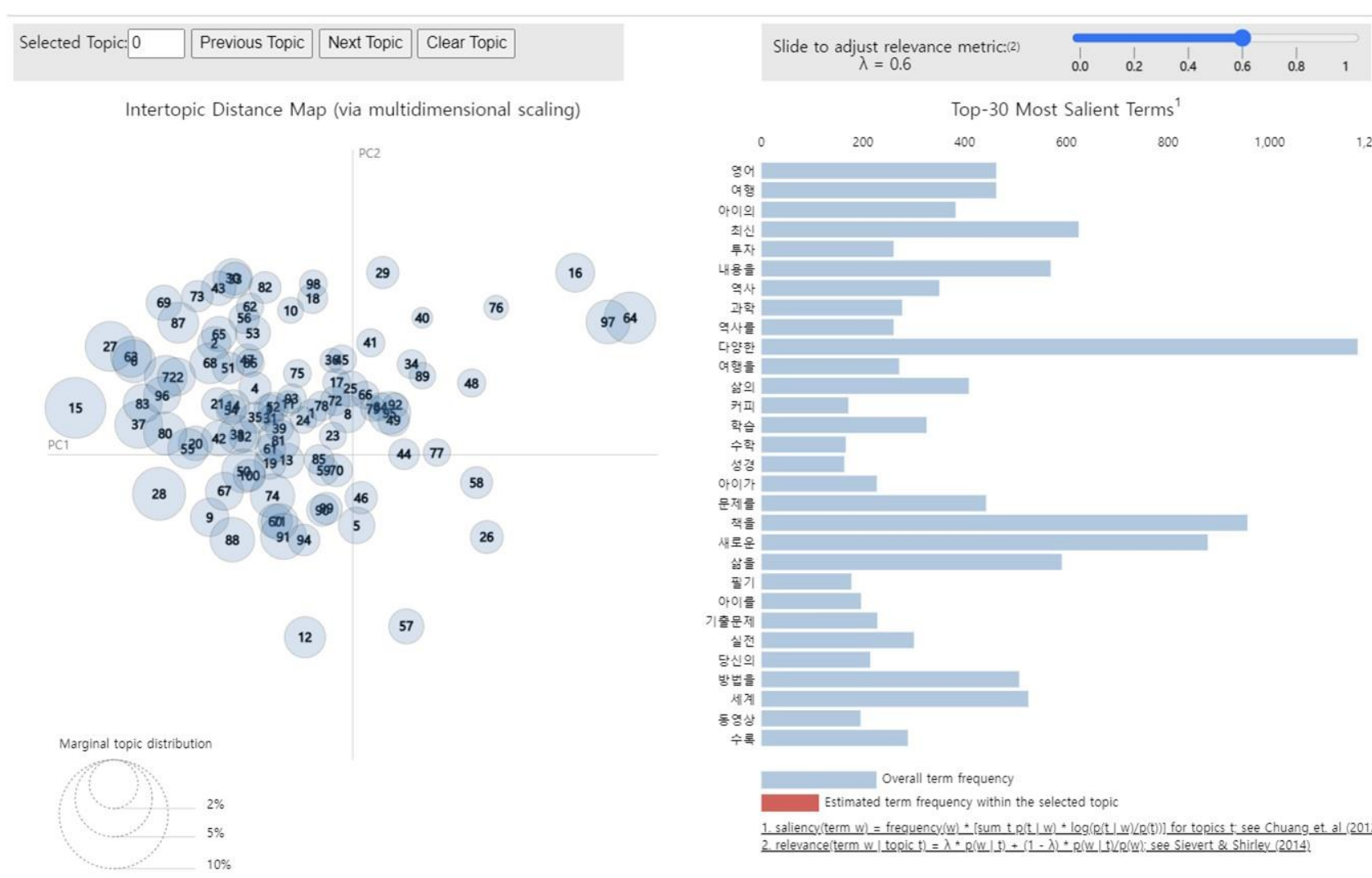
상위 5개

(0.1325179,
'소설'),
[[('해리 포터 시리즈 1~4권 세트(해리포터 20주년 개정판)', 'J. K. 롤링'), 0.6852860696465063],
[('해리 포터와 마법사의 돌 1(해리포터 20주년 개정판)', 'J. K. 롤링'), 0.6755744803449041],
[('해리 포터와 마법사의 돌 2(해리포터 20주년 개정판)', 'J. K. 롤링'), 0.6755744803449041],
[('해리 포터와 마법사의 돌(해리포터 20주년 개정판)', 'J. K. 롤링'), 0.6755744803449041],
[('해리 포터와 아즈카반의 죄수 1(해리포터 20주년 개정판)', 'J. K. 롤링'), 0.665523252850332]]]

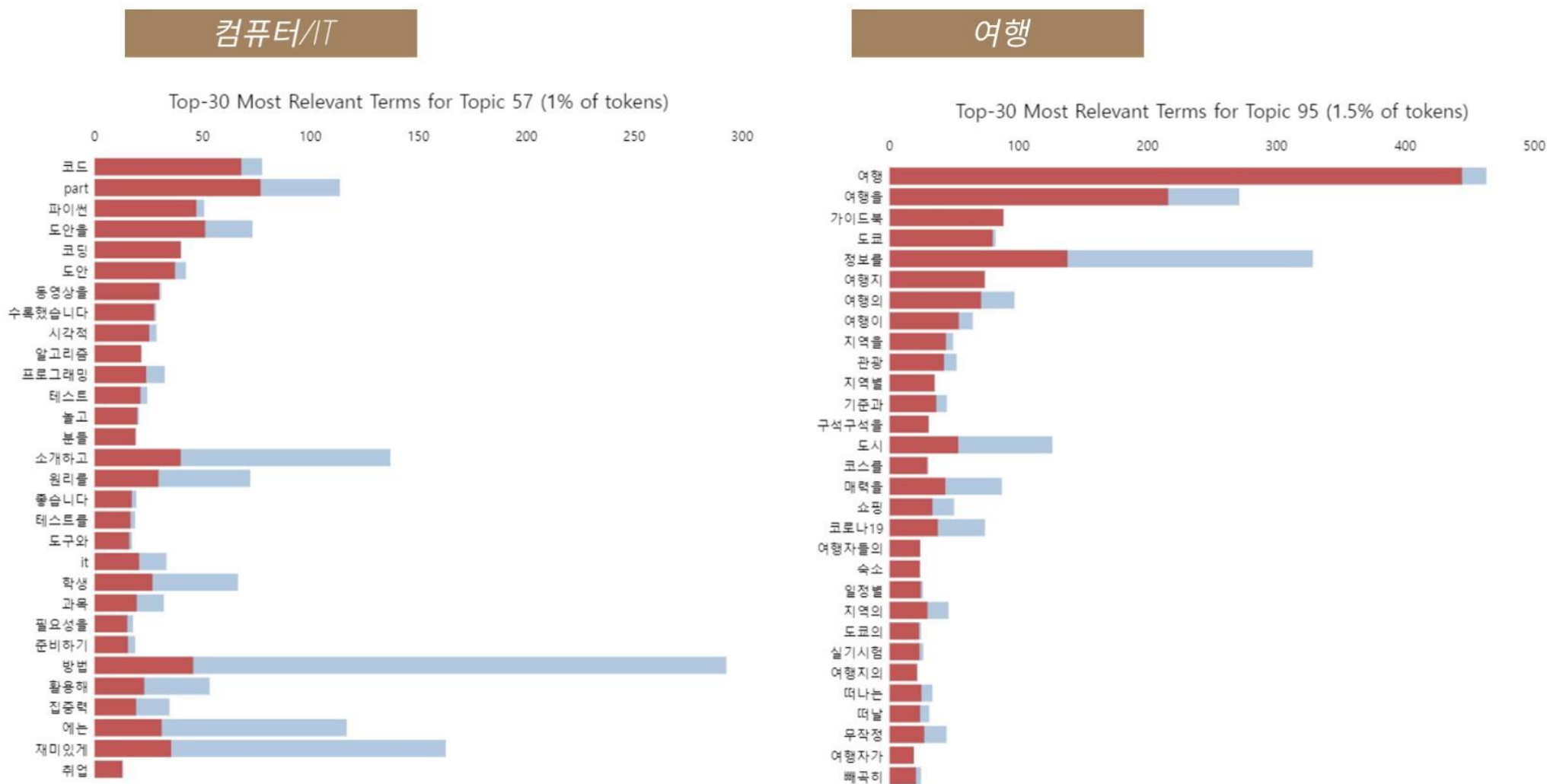
08. LDA 주제 분석

$\lambda = 0.6$

1에 가까울수록 토픽에서 가장 많이 나온 단어 기준
0에 가까울수록 전체적인 비율에 비해서 해당토픽에서 자주 나오는 횟수 기준



08. LDA 주제 분석



결론

- 사용자에게 맞춤형 도서 추천을 위한 도서 데이터 분석 연구
- 도서 데이터를 수집하여 카테고리별 트렌드 파악
- 도서 데이터의 카테고리별 특징 파악
- 사용자 취향과 일치하는 도서 추천 연구 수행
- 워드임베딩 기법을 이용하여 도서 추천의 효율성 향상
- 도서 데이터 분석을 통해 트렌드 및 키워드 파악
- 독자들의 요구사항에 적합한 도서 추천 방법 연구