

# 텍스트 마이닝을 이용한 KEI 연구동향 분석

8th Bigdata Research Team Seminar : Progress Report (44p~)

2017. 07. 27

빅데이터연구팀 김도연

# 목차

- I. 연구 개요
- II. 선행연구
- III. 연구 내용
- IV. 연구 추진방법
- V. 기대효과

## 연구 개요

## 선행 연구

## 연구 내용

## 연구 추진방법

## 기대효과

## 연구 개요

- **과제명** : 텍스트 마이닝을 이용한 KEI 연구동향 분석

- **참여 연구진** :

성명	소속	주요업무	참여율(%)
강성원	빅데이터연구팀	텍스트 마이닝 분석 플랫폼 개발 및 분석	60
김도연	빅데이터연구팀	문헌 분석, 데이터 수집 및 분석	40

- **연구 기간** : 2017.1.1 ~ 2017.12.31

연구 개요

선행 연구

연구 내용

연구 추진방법

기대효과

## 초 록

### 초 록

본 연구는 자연언어분석을 이용하여 환경정책·평가연구원(이하 KEI)의 연구동향을 파악하고, KEI 연구동향이 환경연구에 대한 사회적 연구 수요와 조응하는 지 여부를 분석하였다. 연구주제 선정 범위는 연구를 수행하는 개별 연구자의 성향 및 경험에 따라 제한되기 때문에, KEI의 연구 동향이 국민적 관심과 유리될 수 있다는 우려는 지속적으로 존재해 왔다. 이러한 우려를 확인하기 위해서 본 연구는 연구동향을 나타내는 KEI 연구보고서와 연구수요를 대변하는 환경관련 언론 기사의 장기적인 추이를 비교 분석한다. 이러한 연구는 대량의 텍스트 자료 특성 추출이 필수불가결한데, 본 연구에서는 텍스트 마이닝 기법을 적용하여 이를 수행한다.

구체적으로 본 연구에서는 KEI에서 발행되는 연구보고서와 환경관련 뉴스 기사를 수집한 후 토픽 클러스터링, 연관어 분석, 키워드 네트워크 분석 등 다양한 텍스트 마이닝 기법을 이용하여 장기간의 KEI 연구동향과 사회적 연구수요 동향을 시기별로 각각 추출하여 비교분석하였다. 분석결과는 다음과 같다.

연구 개요

선행 연구

연구 내용

연구 추진방법

기대효과

## 차례

- ④ 1. 서론-
  - ① 가. 연구배경 및 목적-
  - ① 나. 연구내용 및 범위-
  - ① 다. 선행연구 동향-
  - ① 라. 본문의 구성-
- ④ 2. 텍스트 마이닝 기반 연구동향 분석 방법론-
  - ④ 가. 텍스트 마이닝 분석 기법-
    - ④ 1) LDA 분석-
    - ④ 2) 연관어 분석-
    - ④ 3) 키워드 네트워크 분석-
  - ④ 나. 연구동향분석을 위한 텍스트 마이닝 적용 가능성-
  - ④ 다. 연구 분석 절차-
- ④ 3. KEI 연구동향 분석 결과-
  - ④ 가. 분석 데이터 개요-
  - ④ 나. LDA기반 토픽 클러스터링 분석 결과-
    - ④ 1) 토픽별 KEI 연구 동향 (1993년~2016년)-
    - ④ 2) 토픽별 키워드 분석 결과-
      - ④ 가) 에너지 자원-
      - ④ 나) 파기물-
      - ④ 다) 대외협력-
      - ④ 라) 물 환경, 환경영향평가-
      - ④ 마) 기후변화-
  - ④ 다. 키워드 연관성 분석 및 네트워크 분석 결과-
    - ④ 1) 시기별 분석 결과-
      - ④ 가) 1993년~2002년-
      - ④ 나) 2003년~2007년-
      - ④ 다) 2008년~2012년-
      - ④ 라) 2013년~2016년-
- ④ 4. 환경뉴스 동향 분석 결과-
  - ④ 가. 분석 데이터 개요-
  - ④ 나. LDA기반 토픽 클러스터링 분석 결과-
    - ④ 1) 토픽별 환경뉴스 동향 (2004년~2016년)-
    - ④ 2) 토픽별 키워드 분석 결과-
      - ④ 가) 토픽1-
      - ④ 나) 토픽2-
      - ④ 다) 토픽3-
      - ④ 라) 토픽4-
      - ④ 마) 토픽5-
  - ④ 다. 키워드 연관성 분석 및 네트워크 분석 결과-
    - ④ 1) 시기별 분석 결과-
      - ④ 가) 2004년~2007년-
      - ④ 나) 2008년~2012년-
      - ④ 다) 2013년~2016년-
- ④ 5. 매체별 환경 분야 동향 비교 분석 결과-
- ④ 6. 결론 및 제언-

연구 개요

선행 연구

연구 내용

연구 추진방법

기대효과

## 연구 배경 및 목적

- KEI 연구동향이 국민적 관심에 반응하고 있는 지 여부에 대한 회의 존재
    - 개별 연구자의 시간적 제약 및 개인적 연구 성향에 의해서 연구수요 정보 파악 범위가 제한
    - 파악된 정보에 부여되는 우선순위가 개별 연구자의 선호에 영향을 받으므로 최신 정보 및 시의성 있는 연구수요 반영에 제약이 존재
  - 최근 트렌드 분석에 활발하게 사용되는 텍스트 마이닝을 통해 KEI 연구동향과 민간의 환경연구 수요 간의 관계 파악 가능
    - 텍스트 마이닝은 실시간으로 생산되는 다량의 비정형데이터 속에서 의미 있는 패턴을 발견하여 트렌드를 파악하는데 주로 활용
    - 대용량 텍스트 자료 분석이 가능하므로 KEI 연구동향과 민간의 연구수요 동향을 시기별로 트렌드를 각각 추출하여 비교 분석 가능
- ▼
- ▶ 본 연구는 텍스트 마이닝을 이용한 24년(1993~2016) KEI 연구동향 분석을 시도하여 시간적 추이 및 민간 연구수요와의 조응여부를 탐구
    - KEI 연구문헌 및 온라인 뉴스기사 분석을 병행하여 환경관련 연구공급 동향 및 연구수요 동향을 파악
    - 텍스트 마이닝 기법을 이용하여 연구수요를 파악하는 방법의 예를 제공하여 기존의 개별 연구자의 직관에 의존하는 방식을 보완하는 방법을 제공

## 선행 연구

연구 개요

선행 연구

연구 내용

연구 추진방법

기대효과

구분	구분	연구목적	연구방법	주요 연구내용
<b>주요 선행 연구</b>	1	- 과제명: 텍스트 마이닝 기법을 활용한 한국의 경제연구 동향 분석 - 연구자(년도): 송민 외(2013) - 연구목적: 텍스트 마이닝 기법 활용 외국 학술지 한국 경제분야 트렌드 분석	- 키워드 분석 - 네트워크 분석 - 토픽모델링 분석	- 외국 학술지의 한국경제 연구에 대한 연구 동향 및 지적 구조 파악
	2	- 과제명: 소셜 빅데이터를 활용한 국민 통일인식 동향 분석 - 연구자(년도): 송태민 (2015) - 연구목적: 2014년 '통일대박론' 대두 이후 통일인식 변화를 소셜빅데이터 이용 분석	- 키워드 분석 - 연관성 분석	- 소셜 미디어 통일관련 연관어 분석 - 통일관련 연관어와 통일인식간의 관계 분석
	3	- 과제명: 항공산업 미래유망분야 선정을 위한 텍스트 마이닝 기반의 트렌드 분석 - 연구자(년도): 신경식 외(2015) - 연구목적: 텍스트마이닝 트렌드 분석 활용 항공산업 미래유망분야 발굴	- 토픽 모델링 분석	- 텍스트 마이닝 기법을 적용한 항공산업 관련 논문 트렌드 분석 - 토픽 모델링 분석 활용 항공산업 미래유망부분 추출
	4	- 과제명: 빅데이터를 활용한 환경분야 정책수요 분석 - 연구자(년도): 이미숙 외(2014) - 연구목적: 매체별(뉴스, 블로그, 트위터) 환경정책수요 분석	- 감성 분석 - 연관성 분석 - 네트워크 분석	- 세부 환경분야별 소셜빅데이터 분석 - 전체문서 및 환경문서의 행복도 비교 분석
<b>본 연구</b>		- 기존의 연구는 단일 매체의 추세 파악에 집중하여 분석결과 활용이 미진 - 환경분야 문헌 텍스트 마이닝 연구는 초기 단계	- 분석 대상이 단일 매체에 국한되어 전체적인 동향 파악에 한계가 있음	- 환경분야 연구 문헌 텍스트 마이닝 기법을 적용하는 선도적 연구가 필요 - 다양한 매체의 동향을 비교 분석하여 연구 동향의 시사점을 파악하는 연구가 필요

연구 개요

선행 연구

연구 내용

연구 추진방법

기대효과

## 주요 연구 내용

### ■ KEI 연구동향 파악

- 텍스트 마이닝 기법을 활용한 연구동향 분석
  - KEI가 설립된 1993년부터 2016년까지의 KEI DB에서 제공하는 연구보고서(제목, 목차, 요약, 날짜)를 분석에 활용
  - 24년 간(1993-2016) 연구보고서 1,697건

### ■ 매체별 환경분야 이슈 비교 분석

- 연구공급 동향과 연구수요 동향 비교 분석
  - 연구공급 동향 파악 : 2004년부터 2016년(13개년)까지의 KEI DB에서 제공하는 연구보고서(제목, 날짜) 분석
  - 연구수요 동향 파악 : 2004년부터 2016년(13개년)까지의 언론매체에서 제공하는 뉴스기사(제목, 날짜) 분석
- 매체별 추출한 키워드를 시계열로 파악하고 두 시계열을 비교 분석



# Text Mining Flowchart

Plan '17	Process	Code	Description	Input	Output	Note
상반기	3월 상순 Pre-processing(1)	topic_clustering.R	- 형태소분석기 실행(KoNLP 등) - Low TF-IDF 값 제거 - 불용어처리 등 전처리 과정 (특정 단어 삭제, 특수문자 제거, 소문자로 변경 등) - Word Length는 2글자 이상 - 동의어 처리	kei.csv	DocumentTermMatrix 부록_배제키워드.hwp	- 자연언어(이명진 박사님): 한글 처리 문제(조사 등) -> 다양한 한글 전처리 방법을 통해 해결 가능함.
	3월 하순 LDA Modeling	topic_clustering.R	- LDA기반 토픽 모델링 - 토픽별 핵심 단어 출력 - 문서별 토픽번호 및 확률값 출력 - 단어별 토픽번호 및 확률값 출력	Document TermMatrix	term_topic.csv doc_prob_df.csv doc_prob_df_max.csv id_topic.csv lda_tm.csv	- 입력값: SEED = 2017, K = 5
	4월 LDAvis	topic_clustering.R	- 토픽모델링 - 2차원 시각화 및 주요 키워드 확률분포 목록 시각화	lda_tm.csv	HTML 등 웹파일	- apache-tomcat-8.5.12 사용 - 산출물 서버업로드 필요
	5월 상순 LDA Result Analysis		- 토픽별 키워드 분석 - 토픽별 연구보고서 동향 분석	id_topic.csv	id_topic_Analysis.xlsx	-1993-2016년 연구보고서 동향 분석
	5월 하순 Association Analysis(1)	Association_Analysis.R	- 지지도, 신뢰도가 0.01 이상 값 출력 - 3가지속도(지지도, 신뢰도, 향상도) 분석	1993_2002.txt 2003_2007.txt	Association.xlsx	- 연구보고서 제목 데이터 활용 - 추후으로 분석시 메트릭스가 너무 커짐 - 4개 시기별 동향 분석
	5월 하순 Network Analysis(1)	Association_Analysis.R	- 원의 크기: 언급횟이 많을수록 크기가 큼 - 원의 색깔: 매개중심성이 높을수록 색깔이 진함	2008_2012.txt 2013_2016.txt	93-02.png 03-07.png 08-12.png 13-16.png	
6월 상순 Data Collection	naver_news1.java naver_news2.java naver_news3.java	- Java jsoup를 사용하여 Web crawling - 조건: 네이버 뉴스-> 사회-> 환경 - 기간: 2004.1.1-2016.12.12 (총 137개년) - 영역: 제목, 날짜, 언론사 - 양: 193,636개		Naver_news.csv Naver_news_Analysis.xlsx	- 2004년 이전 네이버 뉴스 기사 부실	
하반기	6월 상순 Pre-processing(2)	topic_clustering.R	- 형태소분석기 실행(KoNLP 등) - Low TF-IDF 값 제거 - 불용어처리 등 전처리 과정 (특정 단어 삭제, 특수문자 제거, 소문자로 변경 등) - Word Length는 2글자 이상 - 동의어 처리	Naver_news.csv	DocumentTermMatrix 부록_제거 대상 키워드 목록.hwp	- 자연언어(이명진 박사님): 한글 처리 문제(조사 등) -> 다양한 한글 전처리 방법을 통해 해결 가능함.
	6월 상순 LDA Modeling(2)	topic_clustering.R	- LDA기반 토픽 모델링 - 토픽별 핵심 단어 출력 - 문서별 토픽번호 및 확률값 출력 - 단어별 토픽번호 및 확률값 출력	Document TermMatrix	news_term_topic.csv news_doc_prob_df.csv news_doc_prob_df_max.csv news_id_topic.csv news_lda_tm.csv	- 입력값: SEED = 2000000 K = 15
	6월 상순 LDAvis(2)	topic_clustering.R	- 토픽모델링 - 2차원 시각화 및 주요 키워드 확률분포 목록 시각화	lda_tm.csv	HTML 등 웹파일	- apache-tomcat-8.5.12 사용 - 산출물 서버업로드 필요
	6월 하순 LDA Result Analysis(2)		- 토픽별 키워드 분석 - 토픽별 연구보고서 동향 분석	news_id_topic.csv	news_id_topic_Analysis.xlsx	-2004-2016년 네이버 뉴스 기사 연도별 동향 분석
	8월 Trend Comparative Analysis					
	9월 시사전문 도출 및 정제제안					
	10월 향후 계획 수립			- 뉴스기사 댓글, 환경부 관련 페이스북, 트위터, 블로그, 유튜브 등을 통해 정제수혜자 오픈네트워킹 분석 실시		

연구 개요

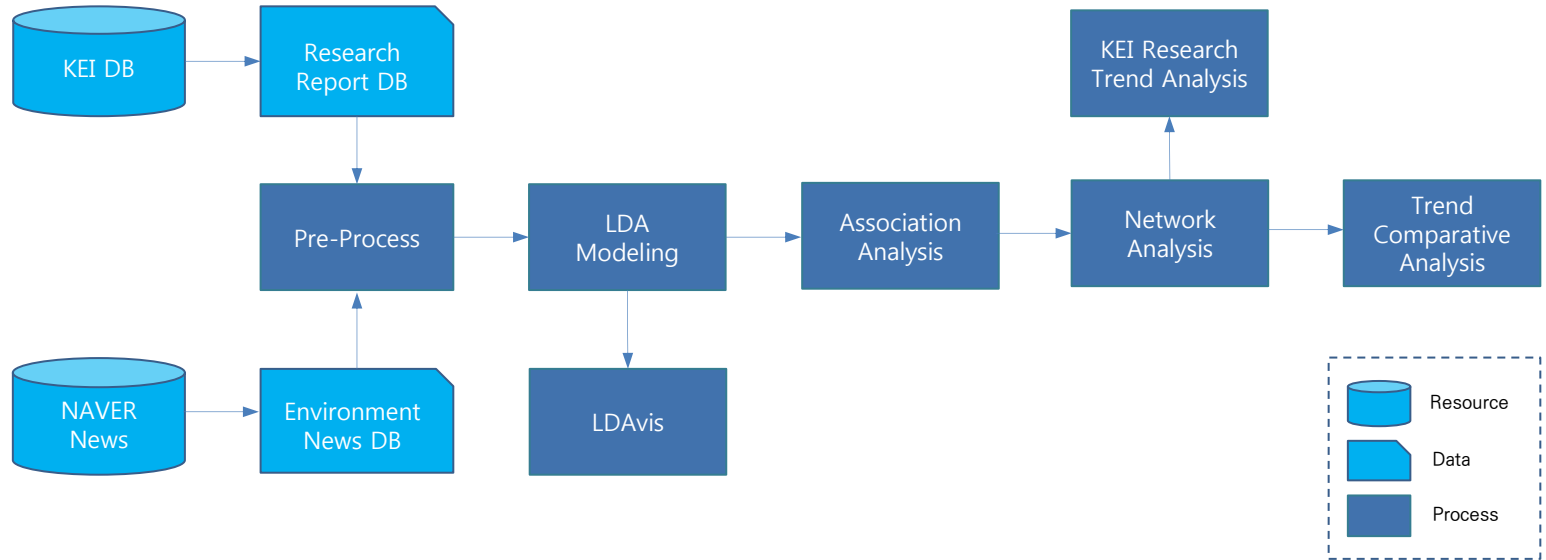
선행 연구

연구 내용

연구 추진방법

기대효과

## KEI 연구동향 분석 작업 흐름도



연구 개요

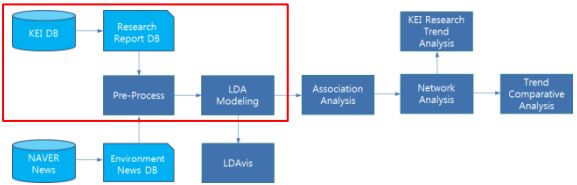
선행 연구

연구 내용

연구 추진방법

기대효과

# Text Mining Process List



Plan 2017	Process	Code	Description	Input	Output	Note
3월 상순	<b>Pre-processing(1)</b>	topic_clustering.R	<ul style="list-style-type: none"> <li>- 형태소분석기 실행(KoNLP 등)</li> <li>- Low TF-IDF 값 제거</li> <li>- 불용어처리 등 전처리 과정 (특정 단어 삭제, 특수문자 제거, 소문자로 변경 등)</li> <li>- Word Lengths는 2글자 이상</li> <li>- 동의어 처리</li> </ul>	kei.csv	DocumentTermMatrix 부록_제거 대상 키워드 목록.hwp	<ul style="list-style-type: none"> <li>- 자문의견(이명진 박사님) : 한글 처리 문제(조사 등) -&gt; 다양한 한글 전처리 방법을 통해 해결 가능함.</li> </ul>
3월 하순	<b>LDA Modeling</b>	topic_clustering.R	<ul style="list-style-type: none"> <li>- LDA기반 토픽 모델링</li> <li>- 토픽별 핵심 단어 출력</li> <li>- 문서별 토픽번호 및 확률값 출력</li> <li>- 단어별 토픽번호 및 확률값 출력</li> </ul>	Document TermMatrix	term_topic.csv doc_Prob_df.csv doc_prob_df_max.csv <b>id_topic.csv</b> lda_tm.csv	<ul style="list-style-type: none"> <li>- 입력값 : SEED = 2017, K = 5</li> </ul>
4월	<b>LDAvis</b>	topic_clustering.R	<ul style="list-style-type: none"> <li>- 토픽모델링</li> <li>- 2차원 시각화 및 주요 키워드 확률분포 목록 시각화</li> </ul>	lda_tm.csv	HTML 등 웹파일	<ul style="list-style-type: none"> <li>- apache-tomcat-8.5.12 사용</li> <li>- 산출물 서버업로드 필요</li> </ul>

연구 개요

선행 연구

연구 내용

연구 추진방법

기대효과

# Pre-processing(1)

연구 개요

선행 연구

연구 내용

연구 추진방법

기대효과

## 1. Pre-processing with R

```
#형태소 분석기 실행
system("tctstart")

#Corpus 생성
corp<-VCorpus(VectorSource(parsedData$res$content))
#특수문자 제거
corp <- tm_map(corp, removePunctuation)
#소문자로 변경
corp <- tm_map(corp, tolower)
#특정 단어 삭제
corp <- tm_map(corp, removewords,
c("전략", "연구", "평가", "마련", "조사", "관리", "보도", "분석", "구축"))
#XMET 문서 형식으로 변환
corp <- tm_map(corp, PlainTextDocument)
#Document Term Matrix 생성 (단어 Length는 2로 세팅)
dtm<-DocumentTermMatrix(corp, control=list(removeNumbers=FALSE, wordLengths=c(2,Inf)))
#한글자 단어 제외하기
colNames(dtm) = trimws(colNames(dtm))
#dtm = dtm[,nchar(colNames(dtm)) > 1]

#Sparse Terms 삭제
dtm <- removeSparseTerms(dtm, as.numeric(0.997))
#Remove low tf-idf col and row
term_tfidf <- tapply(dtm$V$row_sums(dtm)[dtm$J, dtm$J, mean) * log2(nBocs(dtm)/col_sums(dtm) > 0))
new_dtm <- dtm[,term_tfidf >= 0]
new_dtm <- new_dtm[row_sums(new_dtm)>0]
```

연구자의 판단으로 ...

## 2. Pre-processing data: id\_topic.csv 생성

A	B	C	D	E	F	G	H	I	J	K
row	id	doc_topic	pContent	maxProb	Title	author	class	year	month	day
1	a1	5	환경 개선 마스터플랜 수립 지역 국내 사후 설	0.840	환경분야 공적개발원조(O)조공장	수시연구2016	06	30		
2	a2	4	추진 국내 탄소 감축 정책 해외 탄소 감축 정책	0.450	제주 탄소제로섬 추진전략이영국	수시연구2016	06	24		
3	a3	4	국내외 기술 사회 경제 시나리오 사회 경제 시	0.340	지탄소 기후변화 적응 사회참여라	수시연구2016	05	31		
4	a4	3	화학 사고 피해액 추정 제안 화학 사고 인적 상	0.805	화학사고의 경제적 손실 최소화원칙	수시연구2016	04	30		
5	a5	3	나노 폐기를 나노 물질 나노 폐기를 세계 나노	0.586	나노폐기물의 안전처리를 조지혜	수시연구2016	04	30		
6	a6	2	2015년 중남 서브 부 지역 가을 대응 2015년	0.760	가을 단계에 따른 적용형	수시연구2016	03	31		
7	a7	2	국내외 기술 국내 기술 국내 산지 정책 동향	0.352	국내 농산물 GIS기반 통합 이주제	수시연구2016	03	31		
8	a8	3	국내외 기술 최종 최종 단계 추진 토의 건강	0.542	기후변화에 따른 건강영향신용승	수시연구2016	02	28		
9	a9	2	제네바 텍스트 중재 제네바 텍스트 제네바 텍	0.502	Post-2020 신기후체제 협상이승준	수시연구2015	12	31		
10	a10	4	사물 인터넷 혁명 핵심 기술 부상 물 환경 사물	0.404	사물인터넷(IoT)을 활용한 한해진	연구	2016	10	31	
11	a11	4	나타 차별 중국 전략 아시아 인프라 투자 은행	0.998	중국의 일대일로(一帶一路)추진방	연구	2016	10	31	
12	a12	2	도시 기후 회복력 도시 기후 회복력 도시 기후	0.810	도시의 기후 회복력 확보를김동현	연구	2016	10	31	
13	a13	2	국내 지역 사회 환경 보건 문제 진단 국내 지	0.458	지역기반 환경보건정책 지원승준	연구	2016	10	31	
14	a14	2	파리 협정 핵심 파리 협정 파리 협정 적응 손	0.910	신기후체제의 기후변화 적 이승준	수시연구2016	09	30		
15	a15	4	차별친환경 차 보급 정책 동향 국내 정책 동향	0.898	대기환경비용을 고려한 친환경적	수시연구2016	09	30		
16	a16	1	경유 차 실 도로 대기 오염 물질 초과 배출 원	0.821	실제로서 경유차의 대기광공	수시연구2016	09	22		
17	a17	5	접근 국내 정보 관련 부지 시사점 건설 기	0.703	토양정보와 관련 부지의 최적복용	수시연구2016	08	31		
18	a18	2	과업 과업 실행 생물 다양 정책 여건 중간	0.896	제3차 국가생물다양성전략이현우	수시연구2016	08	30		
19	a19	4	지속 가능 발전 87년 환경 관 세계 위원회 발	0.733	국가 지속가능성 평가 등 이승준	수시연구2016	07	30		
20	a20	2	국의 지지도 공원 동향 유네스코 아시아 태	0.772	유네스코 세계지질공원 운이주제	수시연구2016	11	22		
21	a21	2	시스템 네트워크 언어대 환경 정책 에너지	0.526	시스템과 네트워크 이론을 이승준	수시연구2016	11	06		
22	a22	5	국가 지역 미래 성장 동력 미래 성장 동력	0.479	국가 및 지역 미래성장동력방상원	연구	2016	10	31	
23	a23	5	국가 지역 미래 성장 동력 미래 성장 동력	0.479	지중환경을 위한 제도 개 황상일	연구	2016	10	31	
24	a24	3	정책 정부 패러다임 주민주 승 정책 동민주	0.699	정부3.0 기반 지역기피시승김태현	연구	2016	10	31	
25	a25	3	차별친환경 차 보급 정책 동향 국내 정책	0.596	공기정보를 활용한 재해조지혜	연구	2016	10	31	
26	a26	3	국내 폐기를 활용 활용 산업 국내 폐기를	0.540	자원순환사회적 전환 추진을 이소라	수시연구2016	11	06		
27	a27	3	전기 전자 제품 활용 정책 활용 국내 일반	0.844	폐자원유용분류를 통한 전 이희선	연구	2016	10	31	
28	a28	4	물 환경 인프라 사회 수익 물 환경 인	0.605	사회적 투자수익률(SROI)이류재	연구	2016	10	31	
29	a29	2	자연 자본 여건 전망 자연 자본 특성 국내	0.393	생태계서비스 기반의 자연이현우	연구	2016	10	31	
30	a30	2	크리티컬 존 국내외 정책 동향 국외	0.527	근지표환경 일계영역(Critical)현용정	기초연구2016	12	06		
31	a31	5	국내 외 환경 재난 사후 대응 정책 국내	0.359	드론을 이용한 환경재난 시승승우	기초연구2016	12	06		
32	a32	4	건물 지속 가능 고밀 건물 환경 주다	0.648	건물부문의 환경유무하 평가승지	기초연구2016	12	06		
33	a33	3	고밀 기후 변화 노동자 대 영향 노동자	0.712	미래 고온환경 변화와 직결김동현	기초연구2016	12	06		

# Pre-processing(1)

## 〈부록〉 제거 대상 키워드 목록

채널	제거 대상 키워드
KEI 연구 보고서	10년, 1990년, 1장, 1절, 2000년, 2001년, 2002년, 2003년, 2004년, 2005년, 2006년, 2007년, 2008년, 2009년, 2010년, 2011년, 2012년, 2013년, 2014년, 2015년, 2020년, 2장, 3장, 3절, 4장, 4절, 5개년, 5장, 가능, 가다, 같다, 개념, 개발, 개선, 개요, 결과, 결론, 결정, 경우, 계수, 계획, 고려, 과제, 관련, 관리, 관점, 관하다, 구조, 구축, 그리다, 기반, 기본, 기존, 기준, 기초, 기타, 내용, 다루다, 다양, 대책, 대하다, 도출, 되다, 따른, 마련, 말다, 모형, 목록, 목적, 목차, 문헌, 미치다, 발전, 방법, 방안, 방향, 배경, 범위, 보고서, 보급, 보다, 보이다, 본론, 부록, 부문, 분석, 비교, 사업, 사용, 사항, 산정, 서다, 서론, 설정, 수립, 수행, 시기, 시스템, 업무, 업종, 여건, 연구, 영향, 요소, 요약, 우리, 운영, 위하다, 유형, 의하다, 이리하다, 이루어지다, 이용, 인하다, 작성, 적용, 적절하다, 전략, 절차, 정보, 정의, 제공, 제기, 제도, 제시, 제안, 조사, 종합, 중심, 지속가능, 지점, 차례, 참고, 처리, 체계, 체제, 초록, 추진, 측면, 통하다, 통합, 특성, 특징, 평가, 평가모형, 필요, 하다, 허용, 현황, 협약, 활용, 회의, 효과 <b>(총 154개)</b>

연구 개요

선행 연구

연구 내용

연구 추진방법

기대효과

# LDavis (Topic 1)

연구 개요

선행 연구

연구 내용

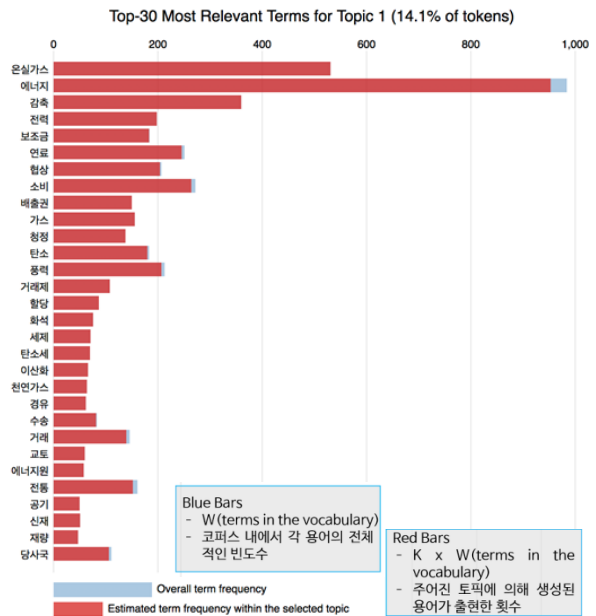
연구 추진방법

기대효과

Selected Topic: 1 Previous Topic Next Topic Clear Topic

Slide to adjust relevance metric:<sup>(2)</sup> λ = 0.05

0.0    0.2    0.4    0.6    0.8    1.0



1. saliency(term w) = frequency(w) \* [sum\_t p(t | w) \* log(p(t | w)/p(t))] for topics t; see Chuang et. al (2012)  
 2. relevance(term w | topic t) = λ \* p(w | t) + (1 - λ) \* p(w | t)/p(w); see Sievert & Shirley (2014)

- Term
- 온실가스
  - 에너지
  - 전력
  - 연료
  - 가스
  - 청정
  - 탄소
  - 환경
  - 세계
  - 탄소세
  - 이산화탄소
  - 천연가스
  - 경유
  - 공기
  - 신재생에너지

“에너지 자원”

# LDAvis (Topic 2)

연구 개요

선행 연구

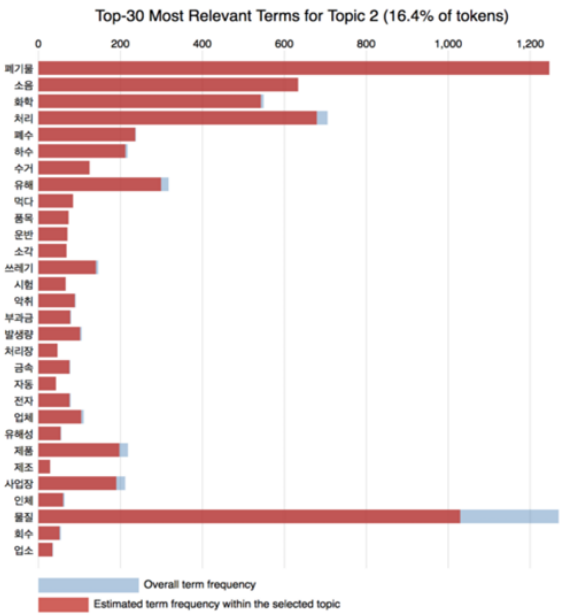
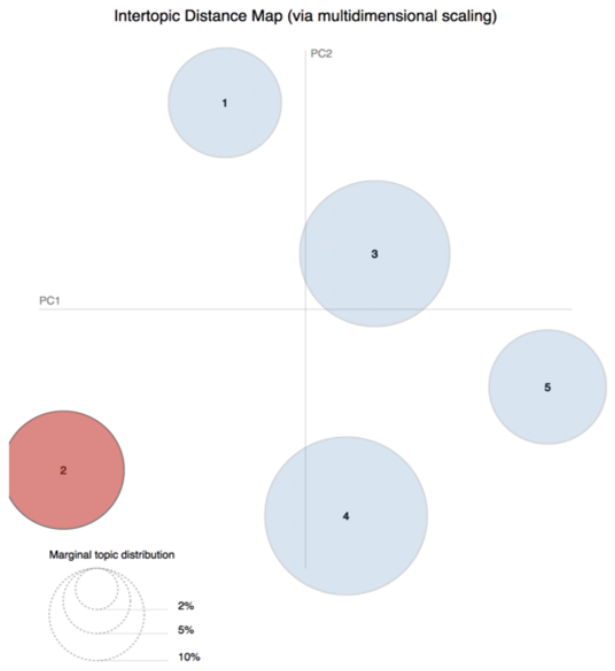
연구 내용

연구 추진방법

기대효과

Selected Topic: 2    Previous Topic    Next Topic    Clear Topic

Slide to adjust relevance metric:<sup>(2)</sup>   $\lambda = 0.05$



1. saliency(term w) = frequency(w) \* [sum<sub>t</sub> p(t | w) \* log(p(t | w)/p(t)) for topics t; see Chuang et. al (2012)  
 2. relevance(term w | topic t) =  $\lambda * p(w | t) + (1 - \lambda) * p(w | t)/p(w)$ ; see Sievert & Shirley (2014)

- Term
- 폐기물
- 소음
- 화학
- 처리
- 폐수
- 하수
- 수거
- 유해
- 막다
- 소각
- 쓰레기
- 악취
- 부담금
- 처리장
- 유해성

“폐기물”

# LDAvis (Topic 3)

연구 개요

선행 연구

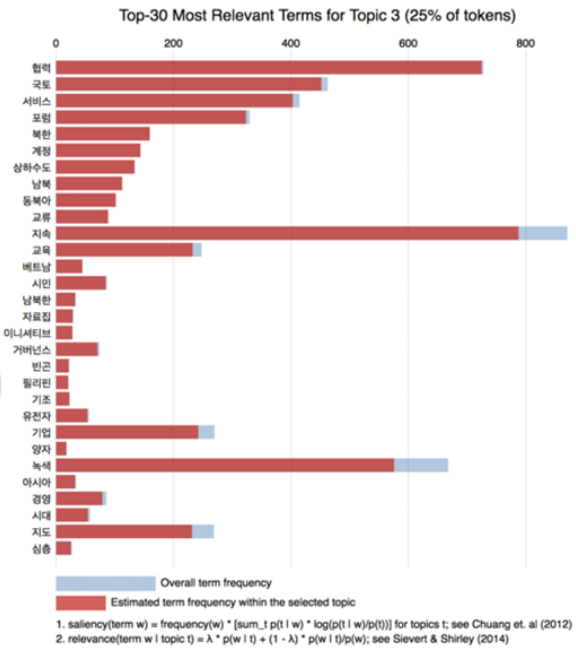
연구 내용

연구 추진방법

기대효과

Selected Topic: 3 Previous Topic Next Topic Clear Topic

Slide to adjust relevance metric:<sup>(2)</sup>  
λ = 0.04 0.0 0.2 0.4 0.6 0.8 1.0



- Term
- 협력
  - 포럼
  - 북한
  - 상하수도
  - 남북
  - 동북아
  - 교류
  - 지속
  - 베트남
  - 시민
  - 남북한
  - 이니셔티브
  - 거버넌스
  - 필리핀
  - 아시아

“대의 협력”



# LDavis (Topic 4)

연구 개요

선행 연구

연구 내용

연구 추진방법

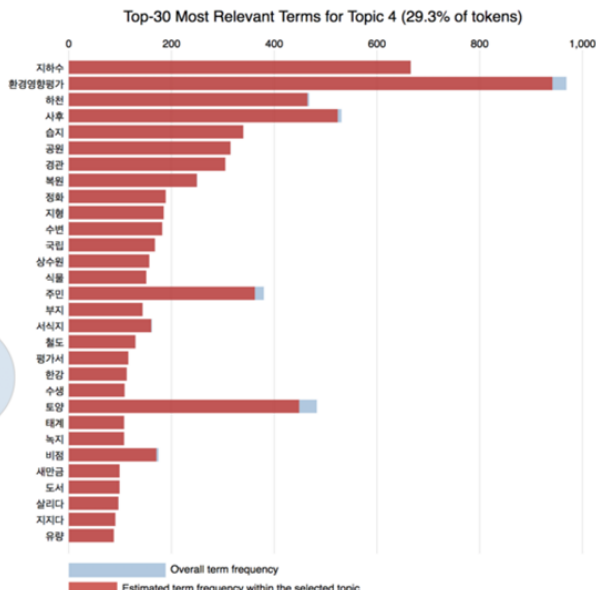
기대효과

Selected Topic: 4 Previous Topic Next Topic Clear Topic



Slide to adjust relevance metric:<sup>(2)</sup> λ = 0.05

0.0 0.2 0.4 0.6 0.8 1.0



- | Term   |
|--------|
| 지하수    |
| 환경영향평가 |
| 하천     |
| 습지     |
| 정화     |
| 지형     |
| 수변     |
| 상수원    |
| 부지     |
| 서식지    |
| 한강     |
| 수생     |
| 도양     |
| 녹지     |
| 새만금    |

“물환경”  
+  
“환경영향평가”

1. saliency(term w) = frequency(w) \* [sum\_t p(t | w) \* log(p(t | w)/p(t))] for topics t; see Chuang et. al (2012)  
 2. relevance(term w | topic t) = λ \* p(w | t) + (1 - λ) \* p(w | t)/p(w); see Sievert & Shirley (2014)

# LDAvis (Topic 5)

연구 개요

선행 연구

연구 내용

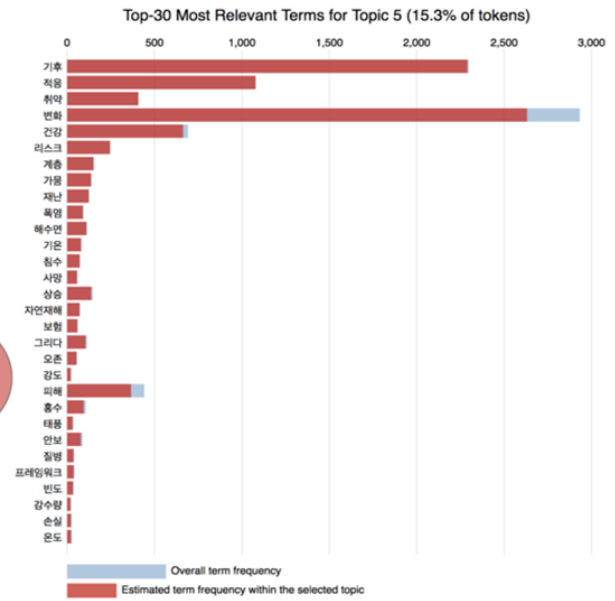
연구 추진방법

기대효과

Selected Topic: 5 Previous Topic Next Topic Clear Topic



Slide to adjust relevance metric:<sup>(2)</sup>  
 $\lambda = 0.05$  0.0 0.2 0.4 0.6 0.8 1.0



1.  $saliency(term\ w) = frequency(w) * [\sum_{t=1}^T p(t|w) * \log(p(t|w)/p(t))]$  for topics  $t$ ; see Chuang et. al (2012)  
 2.  $relevance(term\ w\ l\ topic\ t) = \lambda * p(w|t) + (1 - \lambda) * p(w|l)/p(w)$ ; see Sievert & Shirley (2014)

- Term
- 기후
  - 변화
  - 가뭄
  - 재난
  - 폭염
  - 해수면
  - 기온
  - 침수
  - 사망
  - 자연재해
  - 오존
  - 홍수
  - 태풍
  - 강수량
  - 온도

“기후변화”

# LDAvis (Topic 전체)

연구 개요

선행 연구

연구 내용

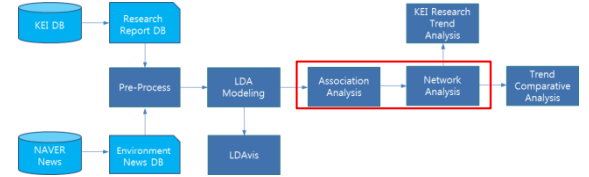
연구 추진방법

기대효과



No.	Topic 1	Topic 2	Topic 3	Topic 4	Topic 5
Title	에너지 자원	폐기물	대외협력	물 환경, 환경영향평가	기후변화
1	온실가스	폐기물	협력	지하수	기후
2	에너지	소음	포럼	환경영향평가	변화
3	전력	화학	북한	하천	가뭄
4	연료	처리	상하수도	습지	재난
5	가스	폐수	남북	정화	폭염
6	청정	하수	동북아	지형	해수면
7	탄소	수거	교류	수변	기온
8	풍력	유해	지속	상수원	침수
9	세제	막다	베트남	부지	사망
10	탄소세	소각	시민	서식지	자연재해
11	이산화탄소	쓰레기	남북한	한강	오존
12	천연가스	약취	이니셔티브	수생	홍수
13	경유	부담금	거버넌스	토양	태풍
14	공기	처리장	필리핀	녹지	강수량
15	신재생에너지	유해성	아시아	새만금	온도
...					

# Text Mining Process List



Plan 2017	Process	Code	Description	Input	Output	Note
5월 상순	LDA Result Analysis		<ul style="list-style-type: none"> <li>- 토픽별 키워드 분석</li> <li>- 토픽별 연구보고서 동향 분석</li> </ul>	id_topic.csv	id_topic_Analysis.xlsx	-1993-2016년 연구보고서 연도별 동향 분석
5월 하순	Association Analysis(1)	Association_Analysis.R	<ul style="list-style-type: none"> <li>- 지지도, 신뢰도가 0.01 이상 값 출력</li> <li>- 3가지측도(지지도, 신뢰도, 향상도) 분석</li> </ul>	1993_2002.txt 2003_2007.txt	Association.xlsx	<ul style="list-style-type: none"> <li>- 연구보고서 제목 데이터 활용</li> <li>- 초록으로 분석시 매트릭스가 너무 커짐</li> <li>- 4개 시기별 동향 분석</li> </ul>
	Network Analysis(1)	Association_Analysis.R	<ul style="list-style-type: none"> <li>- 원의 크기 : 언급량이 많을수록 크기가 큼</li> <li>- 원의 색깔 : 매개중심성이 높을수록 색깔이 진함</li> </ul>	2008_2012.txt 2013_2016.txt	93-02.png 03-07.png 08-12.png 13-16.png	

연구 개요

선행 연구

연구 내용

연구 추진방법

기대효과

## 토픽별 KEI 연구보고서 동향 분석

연구 개요

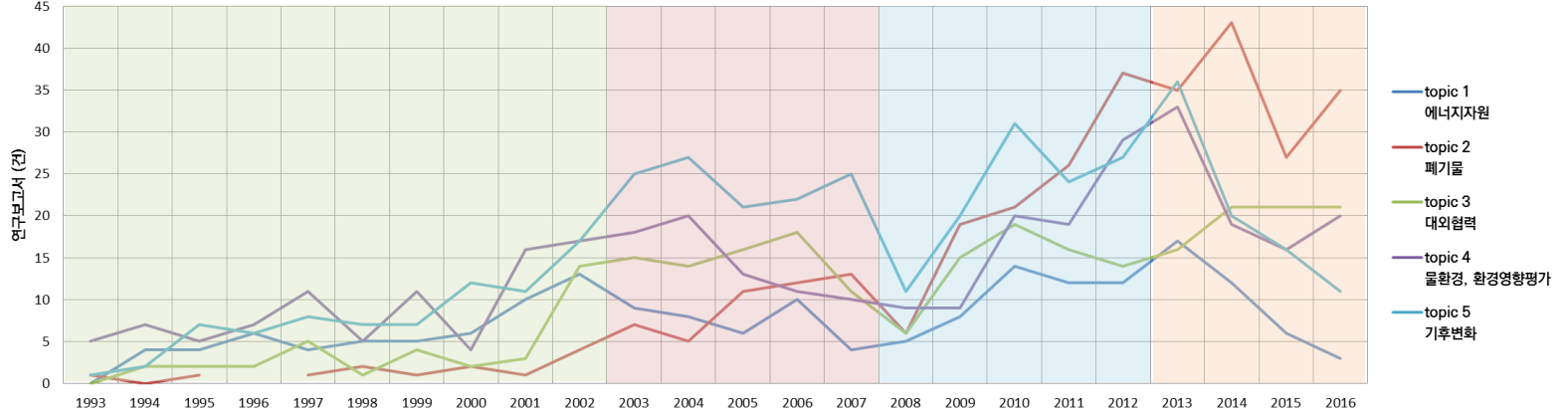
선행 연구

연구 내용

연구 추진방법

기대효과

토픽별 KEI 연구보고서 동향



- 1993~ 2002년도: 전반적으로 토픽별 연구추세가 비슷함.
- 2003~ 2007년도: 기후변화 관련 연구가 활발하게 진행  
물 환경/환경영향평가, 에너지자원 관련 연구는 감소하는 추세를 보임.
- 2008~ 2012년도: 폐기물, 물 환경/환경영향평가 연구가 급증함.
- 2013~ 2016년도: 폐기물 관련 연구가 활발하게 진행  
2015년을 기점으로 연구의 양이 적어짐.

Doc Topic	Title	1993~ 2002	2003~ 2007	2008~ 2012	2013~ 2016	총합
1	에너지자원	57	37	51	38	183
2	폐기물	13	48	109	140	310
3	대외협력	35	74	70	79	258
4	물 환경, 환경영향평가	88	72	86	88	334
5	기후변화	78	120	113	83	394
NA(영문, 한문)		5	27	11	0	43
총합		276	378	440	428	1,522

# 1. 키워드 연관성 및 네트워크 분석(1993-2002년)

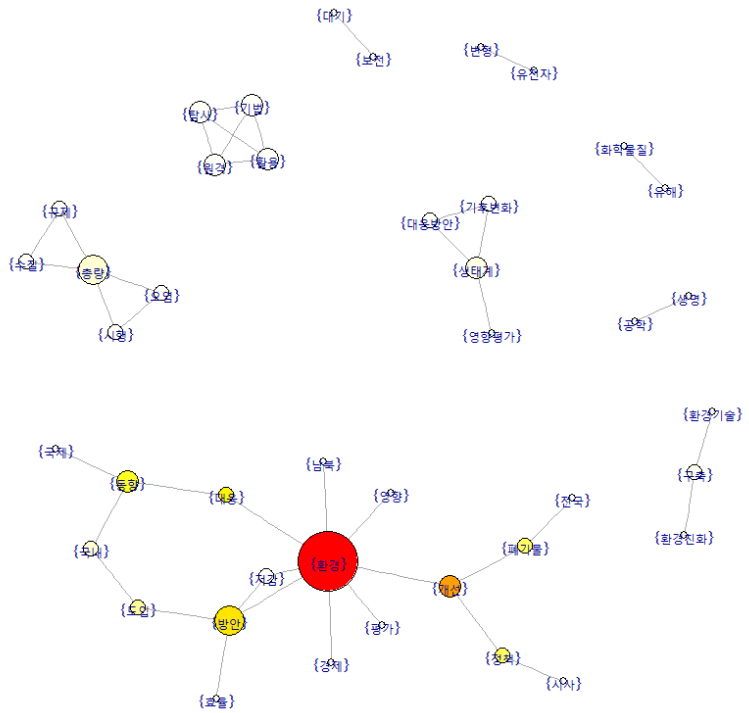
연구 개요

선행 연구

연구 내용

연구 추진방법

기대효과



no	lhs		rhs	support	confidence	lift
1	유전자	=>	변형	0.0123	1.0000	81.3333
2	변형	=>	유전자	0.0123	1.0000	81.3333
3	기후변화	=>	생태계	0.0123	0.7500	45.7500
4	생태계	=>	기후변화	0.0123	0.7500	45.7500
5	기후변화	=>	대응방안	0.0123	0.7500	36.6000
6	대응방안	=>	기후변화	0.0123	0.6000	36.6000
7	영향평가	=>	생태계	0.0123	1.0000	61.0000
8	생태계	=>	영향평가	0.0123	0.7500	61.0000
9	남북	=>	환경	0.0123	0.7500	4.8158
10	환경	=>	남북	0.0123	0.0789	4.8158
11	생태계	=>	대응방안	0.0123	0.7500	36.6000
12	대응방안	=>	생태계	0.0123	0.6000	36.6000
13	규제	=>	수질	0.0123	0.5000	20.3333
14	수질	=>	규제	0.0123	0.5000	20.3333
15	환경친화	=>	구축	0.0164	0.4000	7.5077
16	구축	=>	환경친화	0.0164	0.3077	7.5077

- \* 연관성 분석 평가지표
1. 지지도(support) =  $P(X \cap Y)$
  2. 신뢰도(confidence) =  $P(X \cap Y) / P(X)$
  3. 향상도(lift) =  $P(X \cap Y) / P(X) * P(Y) \Rightarrow$  lift=1(독립), lift<1(음의 연관성), lift>1(양의 연관성)

• 수질오염총량제 시행, 원격탐사기법 활용, 환경친화 기술, 유전자 변형, 전국 폐기물 개선 등의 연구가 활발했음.

\* 키워드 네트워크 분석 결과  
 1. 원의 크기 : 언급량이 높을수록 크다.  
 2. 원의 색깔 : 매개중심성이 높을수록 진하다.(하얀색<노란색<주황색<빨간색)

## 2. 키워드 연관성 및 네트워크 분석(2003-2007년)

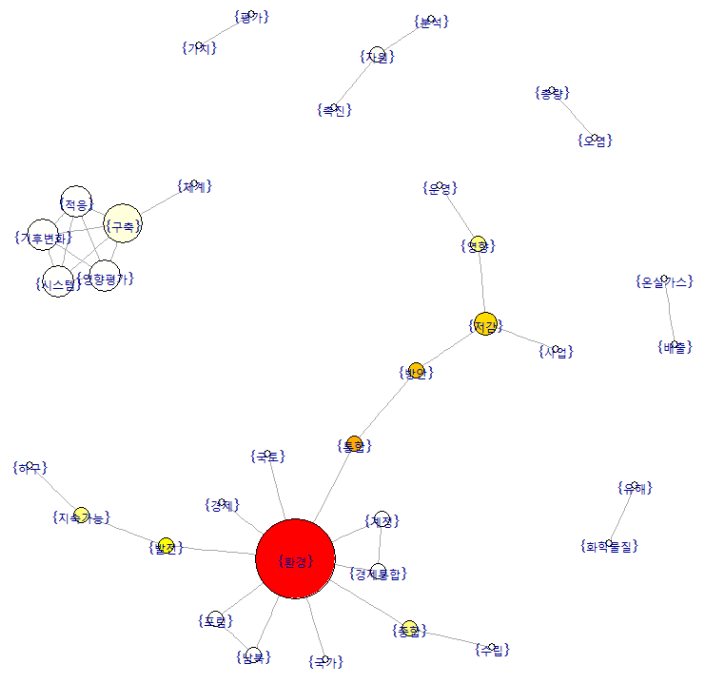
연구 개요

선행 연구

연구 내용

연구 추진방법

기대효과



no	lhs		rhs	support	confidence	lift
1	남북	=>	포럼	0.0117	1.0000	68.2000
2	포럼	=>	남북	0.0117	0.8000	68.2000
3	경제통합	=>	환경	0.0147	1.0000	4.5467
4	환경	=>	경제통합	0.0147	0.0667	4.5467
5	영향평가	=>	기후변화	0.0147	0.8333	23.6806
6	기후변화	=>	영향평가	0.0147	0.4167	23.6806
7	총량	=>	오염	0.0117	0.5714	27.8367
8	오염	=>	총량	0.0117	0.5714	27.8367
9	화학물질	=>	유해	0.0117	0.6667	37.8889
10	유해	=>	화학물질	0.0117	0.6667	37.8889
11	자원	=>	분석	0.0117	0.5000	12.1786
12	분석	=>	자원	0.0117	0.2857	12.1786
13	시스템	=>	기후변화	0.0117	0.4444	12.6296
14	기후변화	=>	시스템	0.0117	0.3333	12.6296
15	경제	=>	환경	0.0147	0.5556	2.5259
16	환경	=>	경제	0.0147	0.0667	2.5259

- \* 연관성 분석 평가지표
1. 지지도(support) =  $P(X \cap Y)$
  2. 신뢰도(confidence) =  $P(X \cap Y) / P(X)$
  3. 향상도(lift) =  $P(X \cap Y) / P(X) * P(Y) \Rightarrow$  lift=1(독립), lift<1(음의 연관성), lift>1(양의 연관성)

- 기후변화 영향평가 및 적응시스템 구축, 온실가스 배출, 환경경제통합 계정 키워드가 새롭게 등장함.
- 전구간에 이어 유해화학물질, 남북 키워드는 계속 등장함.

\* 키워드 네트워크 분석 결과  
 1. 원의 크기 : 언급량이 높을수록 크다.  
 2. 원의 색깔 : 매개중심성이 높을수록 진하다.(하얀색<노란색<주황색<빨간색)

### 3. 키워드 연관성 및 네트워크 분석(2008-2012년)

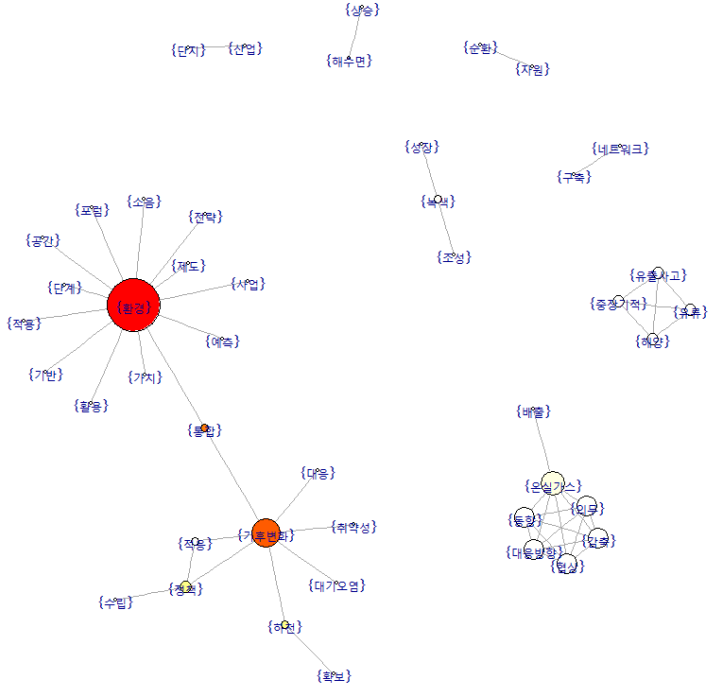
연구 개요

선행 연구

연구 내용

연구 추진방법

기대효과



no	lhs		rhs	support	confidence	lift
1	상승	=>	해수면	0.0101	1.0000	99.5000
2	해수면	=>	상승	0.0101	1.0000	99.5000
3	순환	=>	자원	0.0101	1.0000	66.3333
4	자원	=>	순환	0.0101	0.6667	66.3333
5	대응방향	=>	감축	0.0101	1.0000	39.8000
6	감축	=>	대응방향	0.0101	0.4000	39.8000
7	대응방향	=>	온실가스	0.0101	1.0000	22.1111
8	온실가스	=>	대응방향	0.0101	0.2222	22.1111
9	의무	=>	감축	0.0101	1.0000	39.8000
10	감축	=>	의무	0.0101	0.4000	39.8000
11	의무	=>	온실가스	0.0101	1.0000	22.1111
12	온실가스	=>	의무	0.0101	0.2222	22.1111
13	협상	=>	온실가스	0.0101	1.0000	22.1111
14	온실가스	=>	협상	0.0101	0.2222	22.1111
15	중장기적	=>	유출사고	0.0151	1.0000	56.8571
16	유출사고	=>	중장기적	0.0151	0.8571	56.8571

\* 키워드 네트워크 분석 결과  
 1. 원의 크기 : 언급량이 높을수록 크다.  
 2. 원의 색깔 : 매개중심성이 높을수록 진하다.(하얀색<노란색<주황색<빨간색)

- \* 연관성 분석 평가지표
1. 지지도(support) =  $P(X \cap Y)$
  2. 신뢰도(confidence) =  $P(X \cap Y) / P(X)$
  3. 향상도(lift) =  $P(X \cap Y) / P(X) * P(Y) \Rightarrow$  lift=1(독립), lift<1(음의 연관성), lift>1(양의 연관성)

- 기후변화, 온실가스 키워드의 매개중심성이 높아짐.
- 해양 유류 유출사고, 녹색성장 조성, 해수면 상승, 소음 키워드가 새롭게 등장함.



## 4. 키워드 연관성 및 네트워크 분석(2013-2016년)

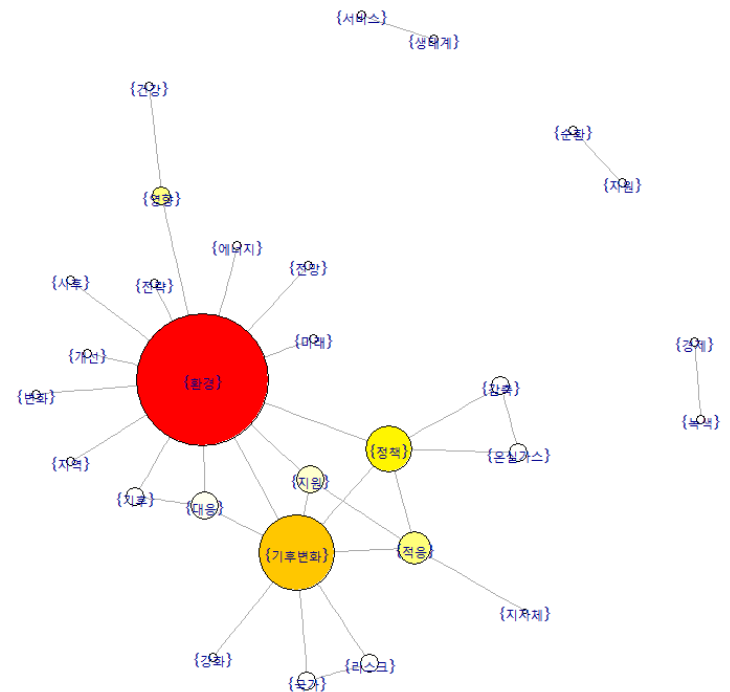
연구 개요

선행 연구

연구 내용

연구 추진방법

기대효과



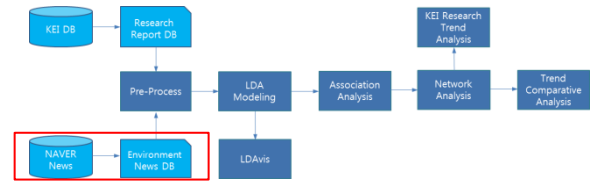
\* 키워드 네트워크 분석 결과  
 1. 원의 크기 : 언급량이 높을수록 크다.  
 2. 원의 색깔 : 매개중심성이 높을수록 진하다.(하얀색<노란색<주황색<빨간색)

no	lhs		rhs	support	confidence	lift
1	지자체	=>	적용	0.0125	0.7143	8.4244
2	적용	=>	지자체	0.0125	0.1471	8.4244
3	감축	=>	온실가스	0.0175	0.8750	43.8594
4	온실가스	=>	감축	0.0175	0.8750	43.8594
5	감축	=>	정책	0.0125	0.6250	5.8285
6	정책	=>	감축	0.0125	0.1163	5.8285
7	온실가스	=>	정책	0.0125	0.6250	5.8285
8	정책	=>	온실가스	0.0125	0.1163	5.8285
9	녹색	=>	경제	0.0175	0.7000	20.0500
10	경제	=>	녹색	0.0175	0.5000	20.0500
11	서비스	=>	생태계	0.0150	0.6000	16.0400
12	생태계	=>	서비스	0.0150	0.4000	16.0400
13	리스크	=>	국가	0.0125	0.5000	10.0250
14	국가	=>	리스크	0.0125	0.2500	10.0250
15	리스크	=>	기후변화	0.0125	0.5000	3.3417
16	기후변화	=>	리스크	0.0125	0.0833	3.3417

- \* 연관성 분석 평가지표
1. 지지도(support) =  $P(X \cap Y)$
  2. 신뢰도(confidence) =  $P(X \cap Y) / P(X)$
  3. 향상도(lift) =  $P(X \cap Y) / P(X) * P(Y)$  => lift=1(독립), lift<1(음의 연관성), lift>1(양의 연관성)

- 전구간에 이어 기후변화 키워드의 매개중심성이 높아짐.
- 환경 키워드와 관련하여 건강, 미래, 전망, 에너지 키워드가 새롭게 등장함.

## Text Mining Process List



Plan 2017	Process	Code	Description	Input	Output	Note
6월 상순	<b>Data Collection</b>	naver_news1.java naver_news2.java naver_news3.java	<ul style="list-style-type: none"> <li>- Java jsoup을 사용하여 Web crawling</li> <li>- 조건: 네이버 뉴스 -&gt; 사회-&gt; 환경</li> <li>- 기간: 2004.1.1~2016.12.12 (총 13 개년)</li> <li>- 영역: 제목, 날짜, 언론사</li> <li>- 양: 193,636개</li> </ul>		Naver_news.csv Naver_news_Analysis.xlsx	- 2004년 이전 네이버 뉴스 기사 부실

연구 개요

선행 연구

연구 내용

연구 추진방법

기대효과

# Data Collection

연구 개요

선행 연구

연구 내용

연구 추진방법

기대효과

- 네이버 뉴스 > 사회 > 환경 관련 기사 전체 수집  
파싱 도구 : JAVA Jsoup



05.24 (수) 청주 19°C 주요뉴스 > 文대통령, 일자리상황판 설치...경제정책, 일자리로 완성...

- 사회
- 사건/사고
- 교육
- 노동
- 언론
- 환경 >
- 인권/복지
- 식품/의료
- 지역
- 인물
- 사회 일반
- 속보

### 환경

**국립수목원, 희귀식물 77% 보전... 국제기준 초과 달성**  
(포전=연합뉴스) 김소연 기자 · 산림청 국립수목원이 국내 희귀식물의 77.2%를 보전함으로써 국제기구의 권고 기준을 조기예 초 ... [연합뉴스](#) 2017-05-23 14:53

**대전충남 '봄기름 심각'... 여름에도 비 적을듯**  
(대전=연합뉴스) 김소연 기자 · 심각한 봄 가뭄이 이어지고 있는 대전·충남지역에 여름에도 뽕님보다 적은 비가 내릴 것으로 예보 ... [연합뉴스](#) 2017-05-23 14:47

**"폭염 막아라" 건물 온도 낮추는 '쿨루프' 조성 구슬땀**  
(부산=연합뉴스) 김재홍 기자 · 부산시가 지난해 시범 추진한 '쿨루프' 조성 사업을 올해 크게 확대해 폭염에 ... [연합뉴스](#) 2017-05-23 14:35

**산림청, 헬기 이용한 산림병해충 항공 방제 나서**  
- 5~7월 소나무재선충병 출몰이 매개충 활동 시기 고려 - 경남과 제주, 경기 등 전국 41개 시군구 7236ha 방제 [대 ... [이데일리](#) 2017-05-23 14:25

**<날씨 이야기>5월24일 수요일(음력 4월29일)**  
전국이 흐리고 비가 오다가 아침에 서쪽 지역부터 그치기 시작해 오후에 대부분 그치겠다. 아침 최저기온은 12도에서 18도, 낮 ... [문화일보](#) 2017-05-23 14:21

id	title	year	month	day	time	source	page
a1	[날씨]주말 대체로 흐리고 미세먼지 주의	2016	01	01	10:36:00 PM	아시아경제	0
a2	고창 동물저수지, 거대한 거북이 모양의 가창오리 군무	2016	01	01	8:19:00 PM	뉴스1	1
a3	고창 동물저수지, 하늘을 뒤덮는 가창오리 군무	2016	01	01	7:28:00 PM	뉴스1	2
a4	고창 동물저수지, 가창오리의 화려한 군무	2016	01	01	7:28:00 PM	뉴스1	3
a5	고창 동물저수지, 노을 속 가창오리 군무	2016	01	01	7:24:00 PM	뉴스1	4
a6	고창 동물저수지 가창오리 군무	2016	01	01	7:23:00 PM	뉴스1	5
a7	고창 동물저수지, 가창오리 군무	2016	01	01	7:23:00 PM	뉴스1	6
a8	파타고니아코리아 '쓰레기 없는 바다'	2016	01	01	4:15:00 PM	뉴스1	7
a9	양양 알바다에 펼쳐진 환경 캠페인	2016	01	01	4:15:00 PM	뉴스1	8
a10	'쓰레기 없는 바다' 메시지 전하는 서퍼들	2016	01	01	4:15:00 PM	뉴스1	9
a11	고속도로, 차량 흐름 대체로 원활...정체 구간은?	2016	01	01	3:47:00 PM	한국경제	10
a12	군산 탁류길 해돋이 문화제	2016	01	01	2:08:00 PM	뉴스1	11
a13	'군산새만금 해돋이 행사'	2016	01	01	2:08:00 PM	뉴스1	12
a14	'행복한 한해가 되기를 바랍니다'	2016	01	01	2:08:00 PM	뉴스1	13
a15	미세먼지와 맞는 새해 첫 날, 수도권·중부내륙 골목...오후계 개선	2016	01	01	11:54:00 AM	에럴드경제	14
a16	2016년 새해맞이 인파	2016	01	01	11:09:00 AM	뉴스1	15
a17	소백산 제2연화봉 새해 첫 일출 장관	2016	01	01	10:26:00 AM	뉴스1	16
a18	2016 새해 해돋이 인파	2016	01	01	9:11:00 AM	뉴스1	17
a19	2016 새해 기분 좋은 출발	2016	01	01	9:06:00 AM	뉴스1	18
a20	2016 새해 희망찬 출발	2016	01	01	9:06:00 AM	뉴스1	19
a21	새해 소원 비는 해돋이 관광객들	2016	01	01	9:02:00 AM	뉴스1	0
a22	2016 해돋이 인파로 가득찬 영일대해수욕장	2016	01	01	8:54:00 AM	뉴스1	1
a23	2016 새해를 반기는 시민들	2016	01	01	8:53:00 AM	뉴스1	2
a24	'2016 원정개 솟아라'	2016	01	01	8:53:00 AM	뉴스1	3
a25	전주 황원산서 옛돼지매 출몰...1마리 사살	2016	01	02	10:02:00 PM	연합뉴스	0
a26	중부 지방, 극심한 가뭄에 물 확보 '안간힘'	2016	01	02	8:23:00 PM	MBC 뉴스	1
a27	<날씨> 더 포근한 일요일...낮 최고 7~16도(3일)	2016	01	02	8:00:00 PM	연합뉴스	2
a28	3일 가뭄 가뭄 많아...미세먼지 주의	2016	01	02	5:43:00 PM	뉴스1	3
a29	'주말 날씨' 미세먼지 농도 1 노약자 외출 자제	2016	01	02	3:28:00 PM	데일리안	4
a30	중국 베이징정성 진도 6.4 규모 지진, 한국에는 영향 없어	2016	01	02	3:13:00 PM	세계일보	5
a31	희석빛 도심	2016	01	02	10:56:00 AM	뉴스1	6
a32	'다가오는 미세먼지'	2016	01	02	10:56:00 AM	뉴스1	7
a33	'하늘은 흐려도 우리는 즐겁게!'	2016	01	02	10:56:00 AM	뉴스1	8
a34	'미세 먼지도 같이 닦아 볼까'	2016	01	02	10:56:00 AM	뉴스1	9
a35	'미세먼지를 닦자!'	2016	01	02	10:56:00 AM	뉴스1	10
a36	'먼지를 닦자'	2016	01	02	10:56:00 AM	뉴스1	11
a37	'아무 것도 안보이네'	2016	01	02	10:55:00 AM	뉴스1	12
a38	'아무 것도 볼 수 없네'	2016	01	02	10:55:00 AM	뉴스1	13
a39	'새해에도 찾아온 미세먼지'	2016	01	02	10:55:00 AM	뉴스1	14
a40	해수담수화 돌파구 찾을까...6일 첫 대화협약체 모임	2016	01	02	8:41:00 AM	연합뉴스	15
a41	빛나간 원숭이 사냥, 사람의 끝은 유기	2016	01	02	8:00:00 AM	연합뉴스	16
a42	2일 날씨 전국 골목, 미세먼지 일부 지역 제외 '보통' 예상	2016	01	02	7:16:00 AM	세계일보	17
a43	전북 주요 하천 8곳 생태독성 '이상 무'	2016	01	02	7:00:00 AM	연합뉴스	18
a44	새해 첫 주요일, 전국 흐리고 포근...'안개 주의'	2016	01	02	6:06:00 AM	뉴스1	19
a45	[금주 뉴스 포토8]2015 헬조선 국어이	2016	01	02	6:00:00 AM	뉴스1	0
a46	예민권 '기내 호텔식당 반입' 놓고 승객폭질 성형	2016	01	02	2:13:00 AM	연합뉴스	1
a47	새해 첫출근길 귀를 많아...미세먼지 농도 '나쁨'	2016	01	03	8:28:00 PM	에럴드경제	0

연구 개요

선행 연구

연구 내용

연구 추진방법

기대효과

## Data Collection

- 네이버 뉴스 기사 데이터 산출 범위

구분	내용
채널	네이버 뉴스
산출 조건	네이버 뉴스 -> 사회 분야 -> 환경 분야
산출 기간	2004-01-01 00:00:00 ~ 2016-12-12 23:59:59 (총 13개년)
산출 영역	제목, 날짜(년, 월, 일, 시간), 언론사
산출 유형	지면기사, 보도자료
언론사	EBN, EPA연합뉴스, JTBC, KBS 뉴스, MBC IMTV, MBC 뉴스, MBN, OSEN, SBS, SBS CNBC, SBS funE, SBS 뉴스, TV리포트, TV조선, Y-STAR, YTN, YTN 현장생중계, ZDNet Korea, 강원일보, 경향신문, 광주드림, 국민일보, 국정브리핑, 내일신문, 노컷뉴스, 뉴스1, 뉴시스, 대전일보, 데일리 서프라이즈, 데일리 안, 동아일보, 디지털데일리, 디지털타임스, 라디오코리아, 레이디경향, 마이데일리, 매경이코노미, 매일경제, 매일신문, 머니S, 머니투데이, 문화일보, 미디어오늘, 부산일보, 블로터, 서울경제, 서울신문, 세계일보, 소년한국일보, 스타뉴스, 스포츠경향, 스포츠동아, 스포츠서울, 스포츠서울닷컴, 스포츠조선, 스포츠한국, 시사 IN, 시사저널, 신동아, 아시아경제, 아이뉴스24, 업코리아, 엑스포츠뉴스, 연합뉴스, 연합뉴스 TV, 오마이TV, 오마이뉴스, 이데일리, 이코노미21, 이코노믹리뷰, 인터뷰365, 일간스포츠(OLD), 일다, 전자신문, 제주일보, 조선비즈, 조선일보, 조세일보, 주간경향, 주간동아, 주간한국, 중앙SUNDAY, 중앙일보, 참세상, 참세상 vod, 컬처뉴스, 코메디닷컴, 쿠키뉴스, 파이낸셜뉴스, 팝뉴스, 프라임경제, 프레시안, 프로메테우스, 한겨레, 한겨레21, 한국경제, 한국경제TV, 한국일보, 헤럴드POP, 헤럴드경제, 헬스조선 (총 101개)
산출 양	193,636개

## Data Collection

연구 개요

선행 연구

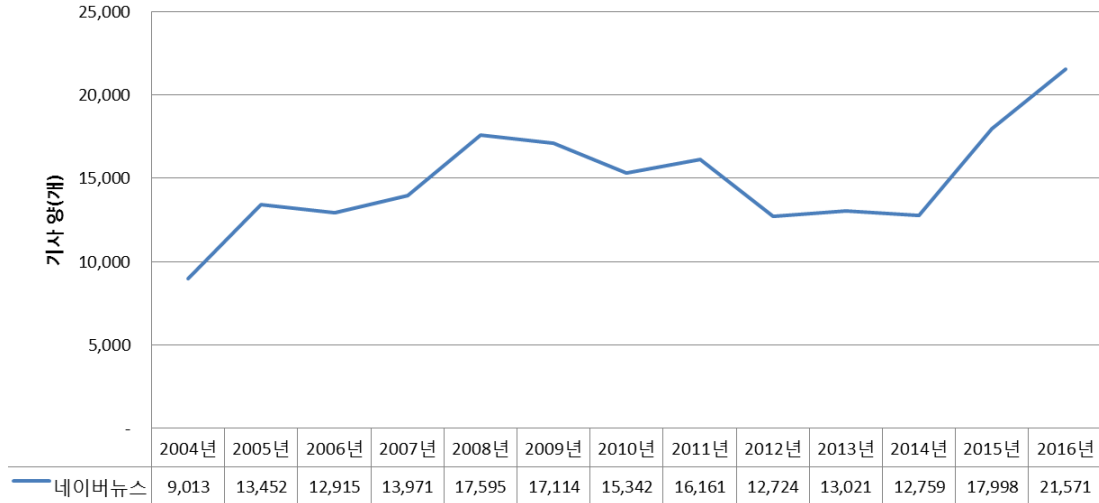
연구 내용

연구 추진방법

기대효과

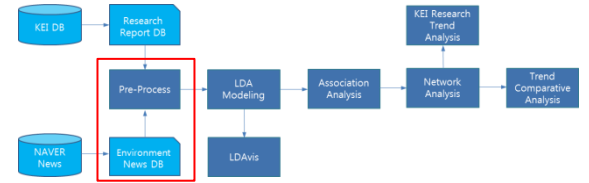
- 환경분야 네이버 뉴스 기사 데이터 기초 분석

〈연도별 환경분야 네이버뉴스 산출 양 추이〉





# Text Mining Process List



Plan 2017	Process	Code	Description	Input	Output	Note
6월 상순	Pre-processing(2)	topic_clustering.R	<ul style="list-style-type: none"> <li>- 형태소분석기 실행(KoNLP 등)</li> <li>- Low TF-IDF 값 제거</li> <li>- 불용어처리 등 전처리 과정 (특정 단어 삭제, 특수문자 제거, 소문자로 변경 등)</li> <li>- Word Lengths는 2글자 이상</li> <li>- 동의어 처리</li> </ul>	Naver_news.csv	DocumentTermMatrix 부록_제거 대상 키워드 목록.hwp	<ul style="list-style-type: none"> <li>- 자문의견(이명진 박사님) : 한글 처리 문제(조사 등) -&gt; 다양한 한글 전처리 방법을 통해 해결 가능함.</li> </ul>

연구 개요

선행 연구

연구 내용

연구 추진방법

기대효과

## Pre-processing(2)

### 〈부록〉 제거 대상 키워드 목록

채널	제거 대상 키워드
네이버 뉴스기사	<p>날씨, 전국, 환경, 오후, 내일, 뉴스, 오늘, 주말, 기상, 관리, 곳곳, 아침, 지역, 영향, 사업, 올해, 국내, 일부, 규모, 종합, 개발, 주변, 조사, 우리, 대상, 활동, 연구, 회의, 생활, 마련, 처리, 가능, 첫날, 최고, 전면, 확인, 수준, 작업, 발령, 경향, 필요, 도입, 활용, 발표, 준비, 효과, 협력, 때문, 잇따, 시스템, 제품, 사회, 하계, 주요, 현상, 마지막, 이후, 대표, 시대, 개월, 단계, 계획, 위해, 가지, 구간, 언제, 통합, 운영, 개체, 차례, 아래, 프로그램, 구역, 기록, 등록, 보고, 연속, 이전, 하기, 재차, 이름, 반기, 들이, 양식, 부분, 누구, 목표, 구조, 기관, 이야기, 중심, 재개, 가득, 설명, 평년기온, 건립, 다양, 가운데, 업무, 다음, 모습, 공간, 하나, 기간, 완화, 초록, 행위, 구경, 공식, 주춤, 구상, 시행, 유의, 일반, 동안, 사전, 시내, 저녁, 낮, 오전, 과정, 최종, 진입, 작전, 자동, 연간, 제도, 특집, 현실, 구름많고, 구축, 방식, 본부, 생각, 선언, 중요, 포함, 사례, 일보, 중순, 노력, 개화, 표지, 쌀쌀, 월일, 유명, 기획, 광역, 그림, 기대, 구성, 관찰, 가로, 수립, 이젠, 전략, 사건, 제외, 추정, 하루, 이틀, 삼일, 월요일, 화요일, 수요일, 목요일, 금요일, 토요일, 일요일, 1월, 2월, 3월, 4월, 5월, 6월, 7월, 8월, 9월, 10월, 11월, 12월, 나흘째, 삼일째, 이틀째, 각종, 이곳, 저곳, 사흘째, 사흘, 나흘, 진짜, 모두, 분야, 표시, 특유, 정기, 단기, 구석, 근본, 기본, 기초, 파악, 수도, 모양, 특정, 데이터, 질문, 채택, 정도, 일교차, 일교, 소식, 전달, 전원, 정신, 직접, 발달, 제기, 선택, 개념, 내부, 봄기운, 이하, 업종, 역할, 어제, 다음주, 이번주, 이번, 금주, 내년, 작년, 세계 (총 233개)</p>

연구 개요

선행 연구

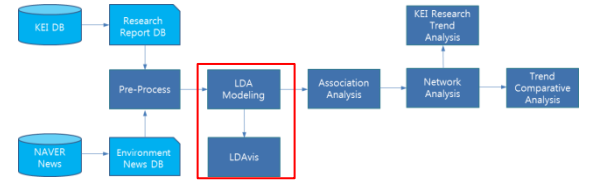
연구 내용

연구 추진방법

기대효과



# Text Mining Process List



Plan 2017	Process	Code	Description	Input	Output	Note
6월 상순	<b>LDA Modeling(2)</b>	topic_clustering.R	<ul style="list-style-type: none"> <li>- LDA기반 토픽 모델링</li> <li>- 토픽별 핵심 단어 출력</li> <li>- 문서별 토픽번호 및 확률값 출력</li> <li>- 단어별 토픽번호 및 확률값 출력</li> </ul>	Document TermMatrix	news_term_topic.csv news_doc_Prob_df.csv news_doc_prob_df_max.csv news_id_topic.csv news_lda_tm.csv	<ul style="list-style-type: none"> <li>- 입력값 : SEED = 2000000 K = 10</li> </ul>
6월 상순	<b>LDavis(2)</b>	topic_clustering.R	<ul style="list-style-type: none"> <li>- 토픽모델링</li> <li>- 2차원 시각화 및 주요 키워드 확률분포 목록 시각화</li> </ul>	lda_tm.csv	HTML 등 웹파일	<ul style="list-style-type: none"> <li>- apache-tomcat-8.5.12 사용</li> <li>- 산출물 서버업로드 필요</li> </ul>

연구 개요

선행 연구

연구 내용

연구 추진방법

기대효과

# LDAvis (Topic 1)

연구 개요

선행 연구

연구 내용

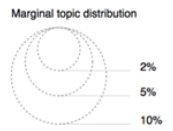
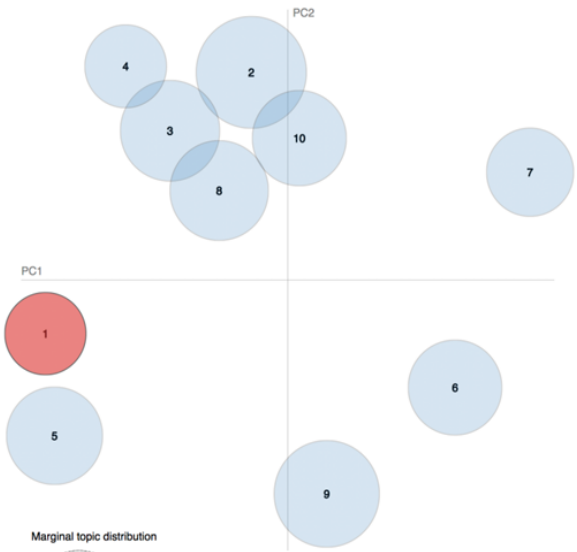
연구 추진방법

기대효과

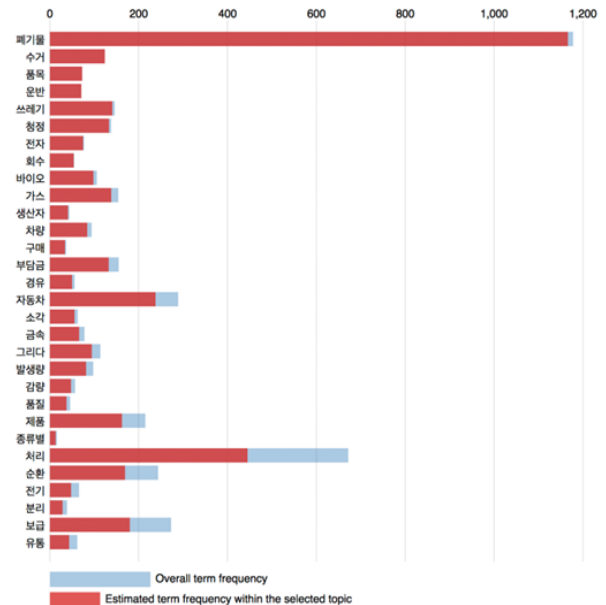
Selected Topic: 1 Previous Topic Next Topic Clear Topic

Slide to adjust relevance metric:<sup>(2)</sup>  
 $\lambda = 0.05$  0.0 0.2 0.4 0.6 0.8 1.0

Intertopic Distance Map (via multidimensional scaling)



Top-30 Most Relevant Terms for Topic 1 (7.3% of tokens)



1. saliency(term w) = frequency(w) \* [sum\_t p(t | w) \* log(p(t | w)/p(t))] for topics t; see Chuang et. al (2012)  
 2. relevance(term w | topic t) =  $\lambda * p(w | t) + (1 - \lambda) * p(w | t)/p(w)$ ; see Sievert & Shirley (2014)

- Term
- 폐기물
  - 수거
  - 운반
  - 쓰레기
  - 청정
  - 가스
  - 차량
  - 부담금
  - 경유
  - 자동차
  - 소각
  - 발생량
  - 처리
  - 순환
  - 진기

“폐기물”

# LDAvis (Topic 2)

연구 개요

선행 연구

연구 내용

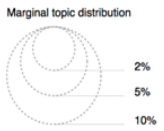
연구 추진방법

기대효과

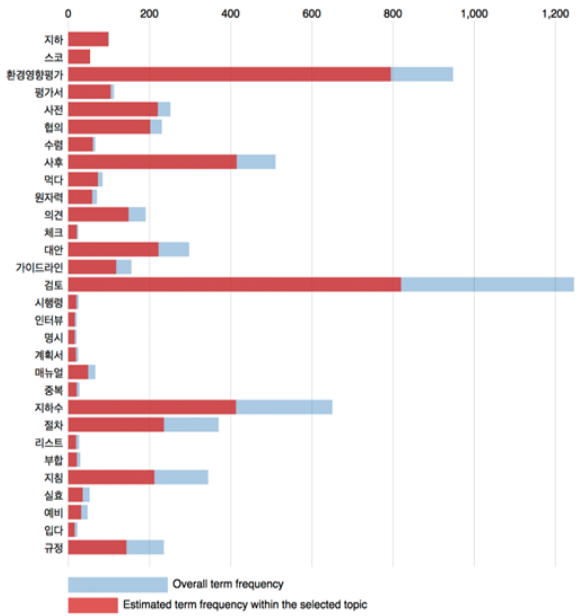
Selected Topic: 2 Previous Topic Next Topic Clear Topic

Slide to adjust relevance metric:<sup>(2)</sup>  
 $\lambda = 0.05$  0.0 0.2 0.4 0.6 0.8 1.0

Intertopic Distance Map (via multidimensional scaling)



Top-30 Most Relevant Terms for Topic 2 (13.4% of tokens)



1. saliency(term w) = frequency(w) \* [sum\_1 p(t | w) \* log(p(t | w)/p(t))] for topics t; see Chuang et. al (2012)  
 2. relevance(term w | topic t) =  $\lambda * p(w | t) + (1 - \lambda) * p(w | t)/p(w)$ ; see Sievert & Shirley (2014)

- Term
- 환경영향평가
- 평가서
- 사전
- 협의
- 수렴
- 사후
- 의견
- 체크
- 가이드라인
- 검토
- 매뉴얼
- 절차
- 리스트
- 지침
- 예비

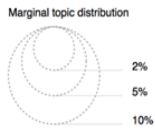
“환경영향평가”

# LDAvis (Topic 3)

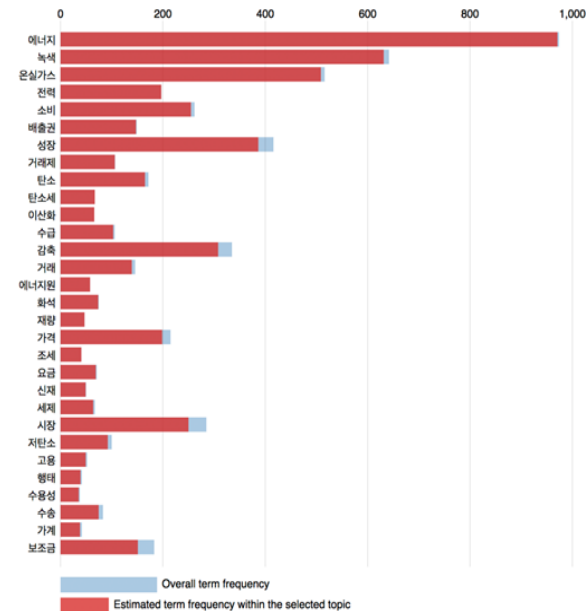
Selected Topic: 3    Previous Topic    Next Topic    Clear Topic

Slide to adjust relevance metric:<sup>(2)</sup>    
 $\lambda = 0.05$     0.0    0.2    0.4    0.6    0.8    1.0

Intertopic Distance Map (via multidimensional scaling)



Top-30 Most Relevant Terms for Topic 3 (10.9% of tokens)



1.  $saliency(term\ w) = frequency(w) * [\sum_t p(t|w) * \log(p(t|w)/p(t))]$  for topics  $t$ ; see Chuang et. al (2012)  
 2.  $relevance(term\ w\ i\ topic\ t) = \lambda * p(w\ i\ t) + (1 - \lambda) * p(w\ i\ t)/p(w)$ ; see Sievert & Shirley (2014)

- Term
- 에너지
- 녹색
- 온실가스
- 전력
- 배출권
- 거래제
- 이산화탄소
- 탄소세
- 감축
- 에너지원
- 화석
- 신재생에너지
- 세계
- 저탄소
- 보조금

“에너지 자원”

연구 개요

선행 연구

연구 내용

연구 추진방법

기대효과

# LDAvis (Topic 4)

연구 개요

선행 연구

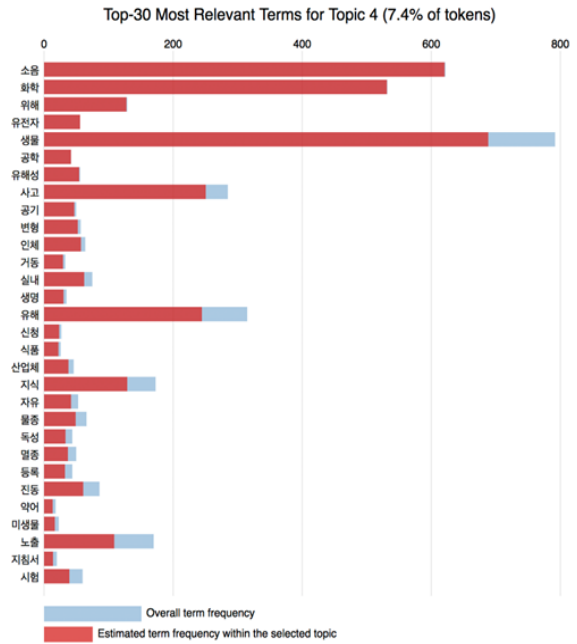
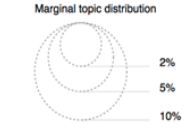
연구 내용

연구 추진방법

기대효과

Selected Topic: 4 Previous Topic Next Topic Clear Topic

Slide to adjust relevance metric:<sup>(2)</sup>  
 $\lambda = 0.05$  0.0 0.2 0.4 0.6 0.8 1.0



1. saliency(term w) = frequency(w) \* [sum\_t p(t | w) \* log(p(t | w)/p(t))]; for topics t; see Chuang et. al (2012)  
 2. relevance(term w | topic t) =  $\lambda * p(w | t) + (1 - \lambda) * p(w | t) / p(w)$ ; see Sievert & Shirley (2014)



# LDAvis (Topic 5)

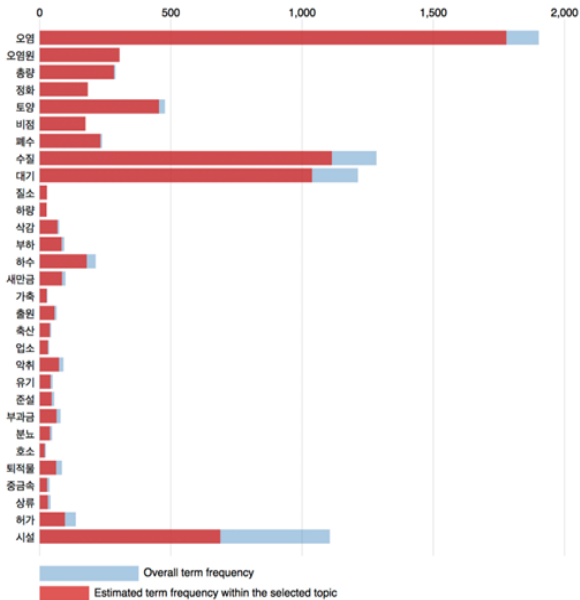
Selected Topic: 5 Previous Topic Next Topic Clear Topic

Slide to adjust relevance metric:<sup>(2)</sup>  
 $\lambda = 0.05$  0.0 0.2 0.4 0.6 0.8 1.0

Intertopic Distance Map (via multidimensional scaling)



Top-30 Most Relevant Terms for Topic 5 (10.2% of tokens)



1. saliency(term w) = frequency(w) \* [sum\_t p(t | w) \* log(p(t | w)/p(t))] for topics t; see Chuang et. al (2012)  
 2. relevance(term w | topic t) =  $\lambda * p(w | t) + (1 - \lambda) * p(w | t)/p(w)$ ; see Sievert & Shirley (2014)

- Term
- 오염
- 오염원
- 정화
- 폐수
- 수질
- 대기
- 하수
- 새만금
- 가축
- 축산
- 약취
- 분노
- 퇴적물
- 중금속
- 상류

“수질오염”

연구 개요

선행 연구

연구 내용

연구 추진방법

기대효과

# LDAvis (Topic 6)

연구 개요

선행 연구

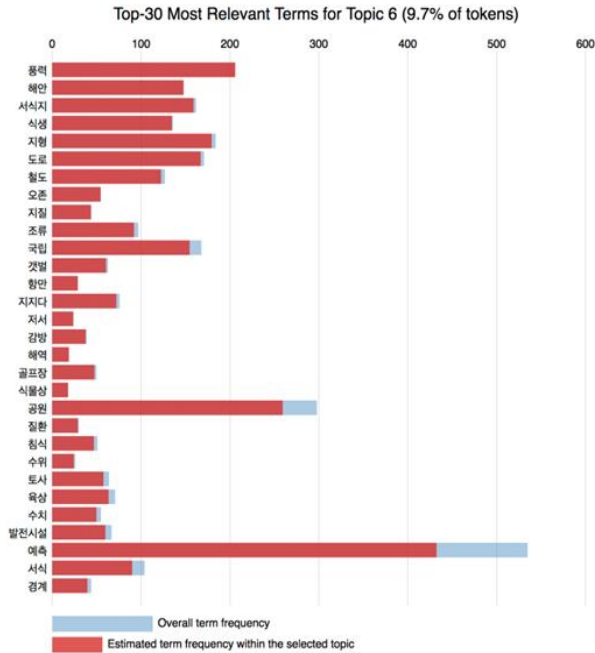
연구 내용

연구 추진방법

기대효과

Selected Topic: 6    Previous Topic    Next Topic    Clear Topic

Slide to adjust relevance metric:<sup>(2)</sup>  
λ = 0.05    0.0    0.2    0.4    0.6    0.8    1.0



1. saliency(term w) = frequency(w) \* [sum\_t p(t | w) \* log(p(t | w)/p(t))] for topics t; see Chuang et. al (2012)  
2. relevance(term w | topic t) = λ \* p(w | t) + (1 - λ) \* p(w | t)/p(w); see Sievert & Shirley (2014)



# LDAvis (Topic 7)

연구 개요

선행 연구

연구 내용

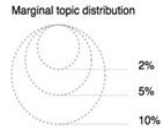
연구 추진방법

기대효과

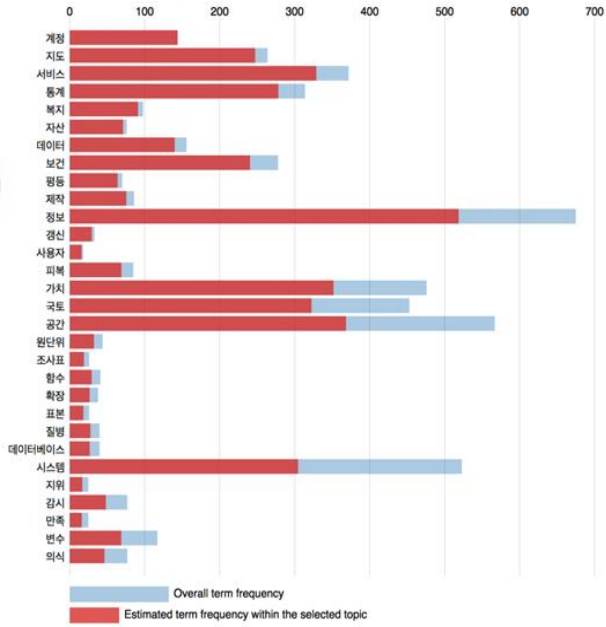
Selected Topic: 7   Previous Topic   Next Topic   Clear Topic

Slide to adjust relevance metric:<sup>(2)</sup>   $\lambda = 0.05$

Intertopic Distance Map (via multidimensional scaling)



Top-30 Most Relevant Terms for Topic 7 (8.4% of tokens)



1. saliency(term w) = frequency(w) \* [sum\_t p(t | w) \* log(p(t | w)/p(t))]; for topics t; see Chuang et. al (2012)  
 2. relevance(term w | topic t) =  $\lambda * p(w | t) + (1 - \lambda) * p(w | t)/p(w)$ ; see Sievert & Shirley (2014)

- Term
- 서비스
- 통계
- 복지
- 데이터
- 보건
- 정보
- 피복지도
- 국토
- 공간
- 조사표
- 항수
- 표본
- 질병
- 데이터베이스
- 변수

“보건”  
+  
“데이터”



# LDAvis (Topic 8)

연구 개요

선행 연구

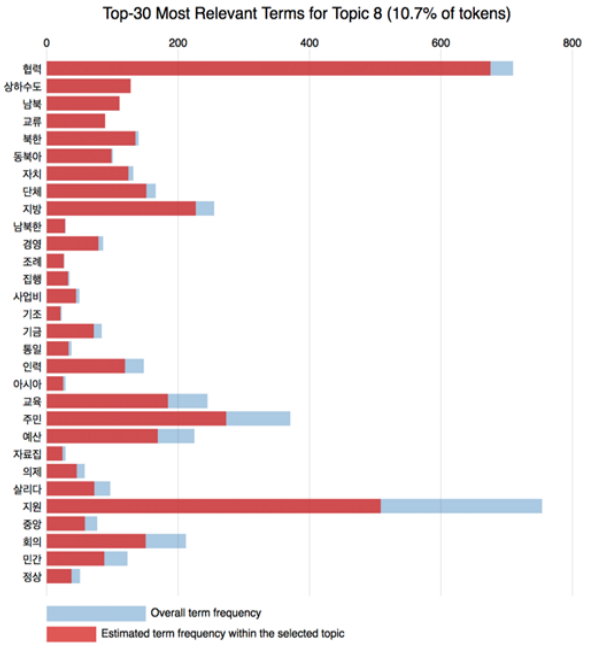
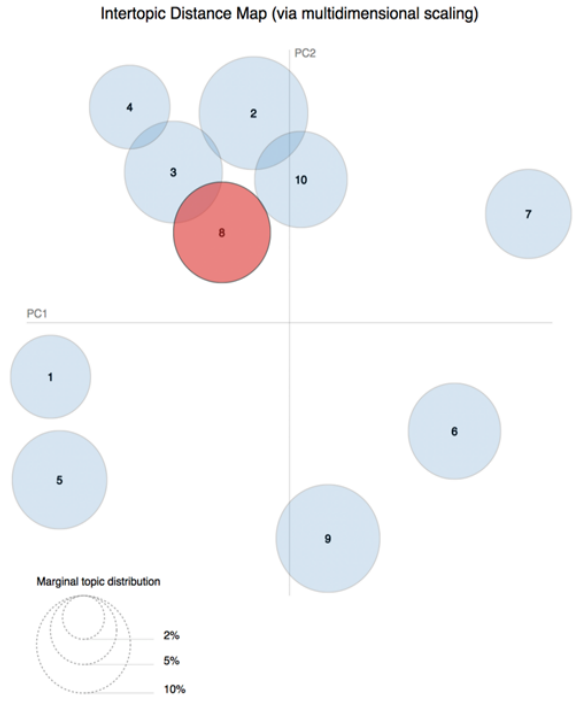
연구 내용

연구 추진방법

기대효과

Selected Topic: 8 Previous Topic Next Topic Clear Topic

Slide to adjust relevance metric:<sup>(2)</sup> 0.0 0.2 0.4 0.6 0.8 1.0  
 $\lambda = 0.05$



1. saliency(term w) = frequency(w) \* [sum\_t p(t | w) \* log(p(t | w)/p(t))] for topics t; see Chuang et. al (2012)  
 2. relevance(term w | topic t) =  $\lambda * p(w | t) + (1 - \lambda) * p(w | t)/p(w)$ ; see Sievert & Shirley (2014)



# LDAvis (Topic 9)

연구 개요

선행 연구

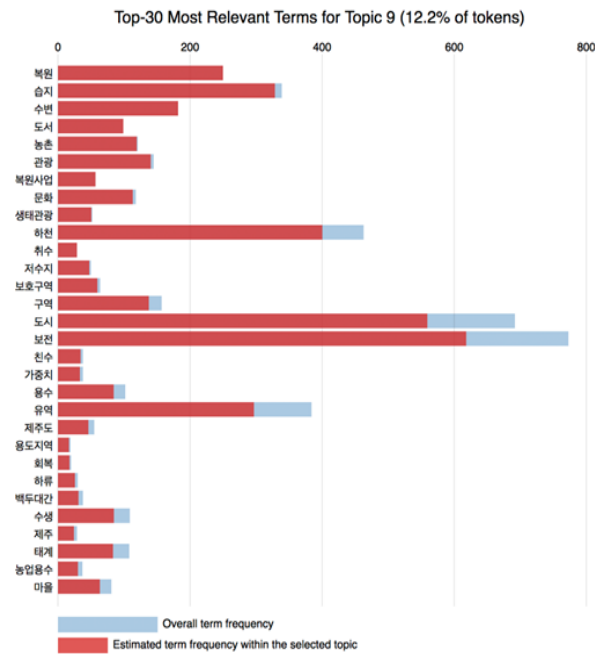
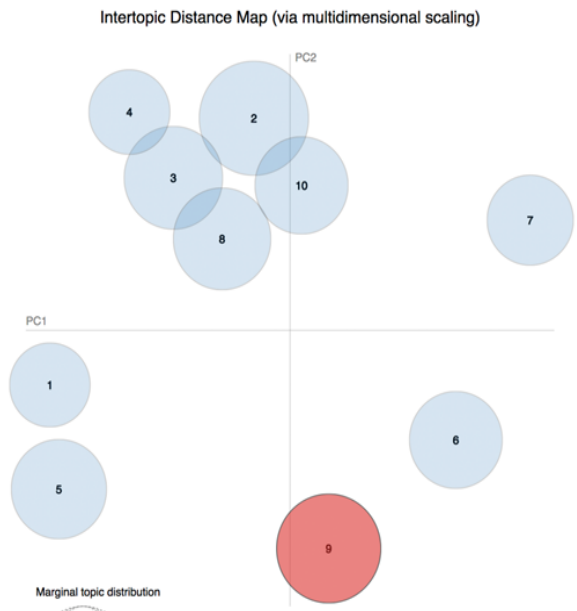
연구 내용

연구 추진방법

기대효과

Selected Topic: 9    Previous Topic    Next Topic    Clear Topic

Slide to adjust relevance metric:<sup>(2)</sup>  
λ = 0.05    0.0    0.2    0.4    0.6    0.8    1.0



- Term
- 복원
  - 습지
  - 수변
  - 복원사업
  - 하천
  - 취수
  - 저수지
  - 친수
  - 용수
  - 유역
  - 하류
  - 백두대간
  - 수생
  - 제주
  - 농업용수

“물환경”

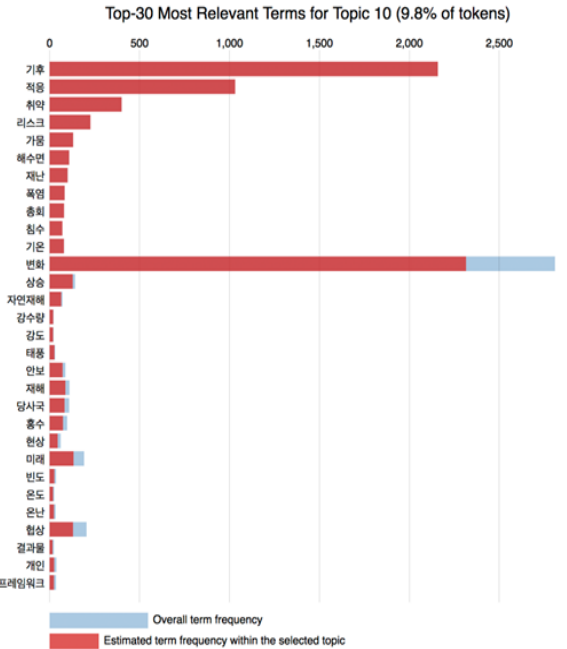
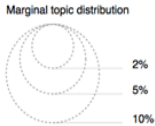
1. saliency(term w) = frequency(w) \* [sum\_t p(t | w) \* log(p(t | w)/p(t))]; for topics t; see Chuang et. al (2012)  
2. relevance(term w | topic t) = λ \* p(w | t) + (1 - λ) \* p(w | t)/p(w); see Sievert & Shirley (2014)

# LDAvis (Topic 10)

Selected Topic: 10    Previous Topic    Next Topic    Clear Topic

Slide to adjust relevance metric:<sup>(2)</sup>   $\lambda = 0.05$

Intertopic Distance Map (via multidimensional scaling)



1. saliency(term w) = frequency(w) \* [sum\_i p(i | w) \* log(p(i | w)/p(i))]; for topics t; see Chuang et. al (2012)  
 2. relevance(term w | topic t) =  $\lambda * p(w | t) + (1 - \lambda) * p(w | t)/p(w)$ ; see Slovert & Shirley (2014)

- Term
- 기후
- 적응
- 취약
- 리스크
- 가뭄
- 재난
- 폭염
- 침수
- 기온
- 변화
- 상승
- 자연재해
- 태풍
- 홍수
- 온난

“기후변화”

연구 개요

선행 연구

연구 내용

연구 추진방법

기대효과

## LDAvis (Topic 전체)

연구 개요

선행 연구

연구 내용

연구 추진방법

기대효과



No.	Topic 1	Topic 2	Topic 3	Topic 4	Topic 5	Topic 6	Topic 7	Topic 8	Topic 9	Topic 10
Title	폐기물	환경영향평가	에너지 자원	유전자 변형, 소음	오염원	해양, 풍력	보건, 데이터	대외협력	물 환경	기후 변화
1	폐기물	환경영향평가	에너지	소음	오염	풍력	서비스	협력	복원	기후
2	수거	평가서	녹색	화학	오염원	해안	통계	상하수도	습지	적응
3	운반	사전	온실가스	유전자	총량	서식지	복지	남북	수변	취약
4	쓰레기	협의	전력	생물	정화	식생	데이터	교류	복원사업	리스크
5	청정	수렴	배출권	공학	도양	지형	보건	북한	하천	가뭄
6	가스	사후	거래제	유해성	폐수	지질	정보	동북아	취수	재난
7	차량	의견	이산화탄소	변형	대기	조류	피복지도	남북한	저수지	폭염
8	부담금	체크	탄소세	인체	새만금	갯벌	국토	경영	친수	침수
9	경유	가이드라인	감축	생명	축산	항만	공간	통일	용수	기온
10	자동차	검토	에너지원	식품	약취	해역	조사표	아시아	유역	변화
11	소각	매뉴얼	화석	독성	부과금	침식	함수	의제	하류	상승
12	발생량	절차	신재생에너지	멸종	분뇨	수위	표본	중앙	백두대간	자연재해
13	처리	리스트	세제	진동	호소	토사	질병	회의	수생	태풍
14	순환	지침	저탄소	미생물	퇴적물	육상	데이터베이스	민간	제주	홍수
15	전기	예비	보조금	시험	중금속	발전시설	변수	정상	농업용수	온난

...

연구 개요

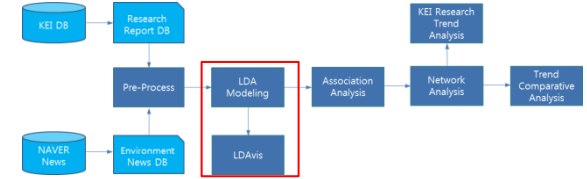
선행 연구

연구 내용

연구 추진방법

기대효과

## Text Mining Process List



Plan 2017	Process	Code	Description	Input	Output	Note
6월 하순	LDA Result Analysis(2)		<ul style="list-style-type: none"> <li>- 토픽별 키워드 분석</li> <li>- 토픽별 네이버 환경뉴스 동향 분석</li> </ul>	news_id_topic.csv	news_id_topic_Analysis.xlsx	<ul style="list-style-type: none"> <li>- 2004~2016년 네이버 뉴스 기사 연도별 동향 분석</li> <li>- 2004~2016년 네이버 뉴스와 KEI 연구보고서 비교 분석</li> </ul>

중 간 자 문 회 의 (2017.06.29)

## KEI 연구보고서 LDAvis (2004~2016)

연구 개요

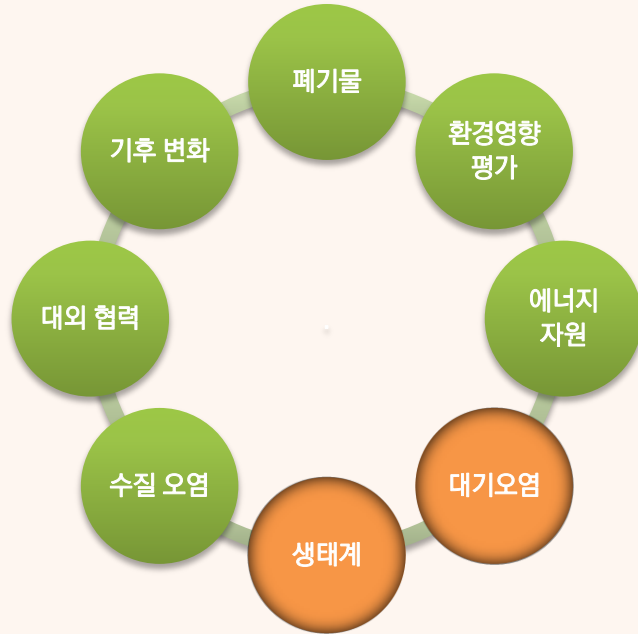
선행 연구

연구 내용

연구 추진방법

기대효과

### KEI 연구보고서



No.	Topic 1	Topic 2	Topic 3	Topic 4	Topic 5	Topic 6	Topic 7	Topic 8
Title	폐기물	환경영향 평가	에너지 자원	대기오염	생태계	수질오염	대외협력	기후변화
1	폐기물	영향	에너지	대기	생태	지하수	협력	적응
2	처리	제도	전력	화학	생물	수질	협상	폭염
3	시설	사후	온실가스	먼지	습지	비점	포럼	해수면
4	배출	검토	원료	오염	서식지	오염원	아세안	침수
5	하수	개선	석탄	초미세	자연환경	새만금	동북아	범람
6	쓰레기	정책	화석	천식	식물	용담댐	의정서	태풍
7	발생량	주민	재생	황사	서식	녹조	교토	기상이변
8	폐수	지역	천연가스	방사능	외래	취수	베트남	자외선
9	총량제	갈등	신재생	미세먼지	야생동물	농업용수	국제	진동
10	약취	설문	화력	방사성	멸종	하량	아시아	자연재해

## 매체별 LDAvis 결과 비교(2004-2016)

연구 개요

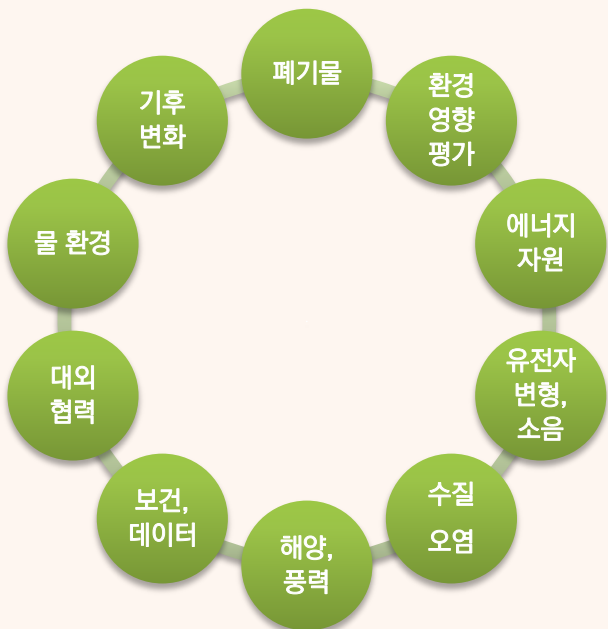
선행 연구

연구 내용

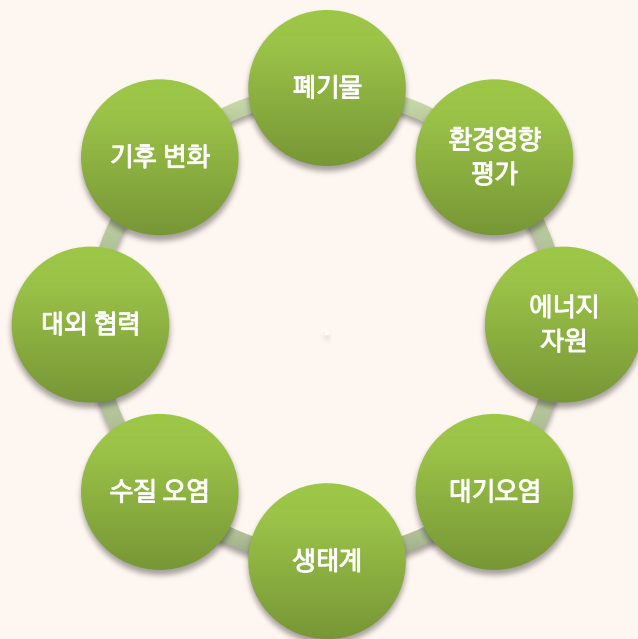
연구 추진방법

기대효과

NAVER 환경뉴스



KEI 연구보고서



## 매체별 LDAvis 결과 비교(2004-2016)

연구 개요

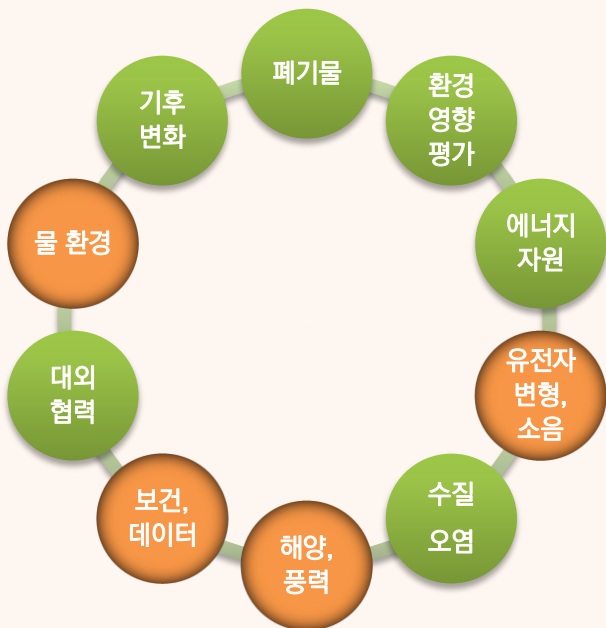
선행 연구

연구 내용

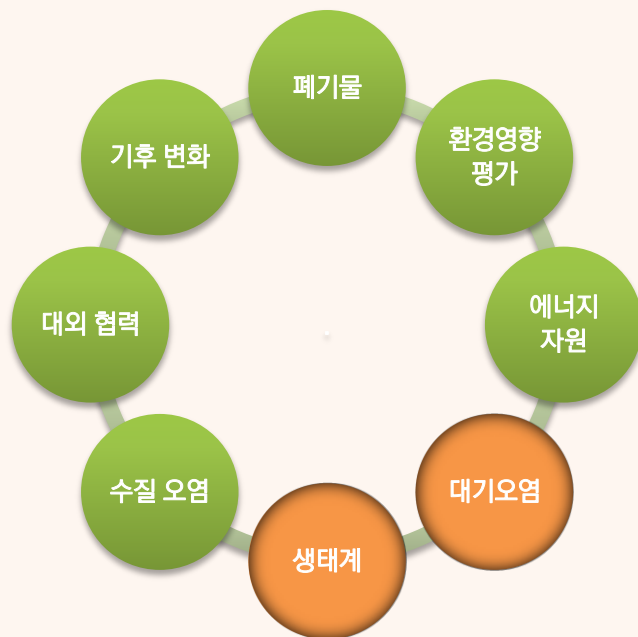
연구 추진방법

기대효과

NAVER 환경뉴스

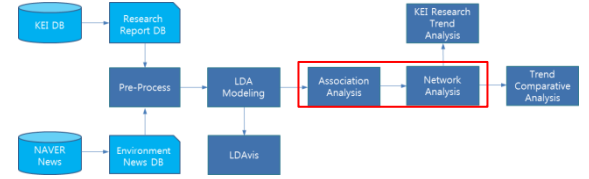


KEI 연구보고서





## Text Mining Process List



Plan 2017	Process	Code	Description	Input	Output	Note
7월	<b>Association Analysis(2)</b>	Association_Analysis.R	<ul style="list-style-type: none"> <li>- 지지도, 신뢰도가 0.01 이상 값 출력</li> <li>- 3가지측도(지지도, 신뢰도, 향상도) 분석</li> </ul>	2004_2007.txt	Association.xlsx	<ul style="list-style-type: none"> <li>- 네이버뉴스 제목 데이터 활용</li> <li>- 본문으로 분석시 매트릭스가 너무 커짐</li> <li>- 3개 시기별 동향 분석</li> </ul>
	<b>Network Analysis(2)</b>	Association_Analysis.R	<ul style="list-style-type: none"> <li>- 원의 크기 : 언급량이 많을수록 크기가 큼</li> <li>- 원의 색깔 : 매개중심성이 높을수록 색깔이 진함</li> </ul>	2008_2012.txt 2013_2016.txt	04-07.png 08-12.png 13-16.png	

연구 개요

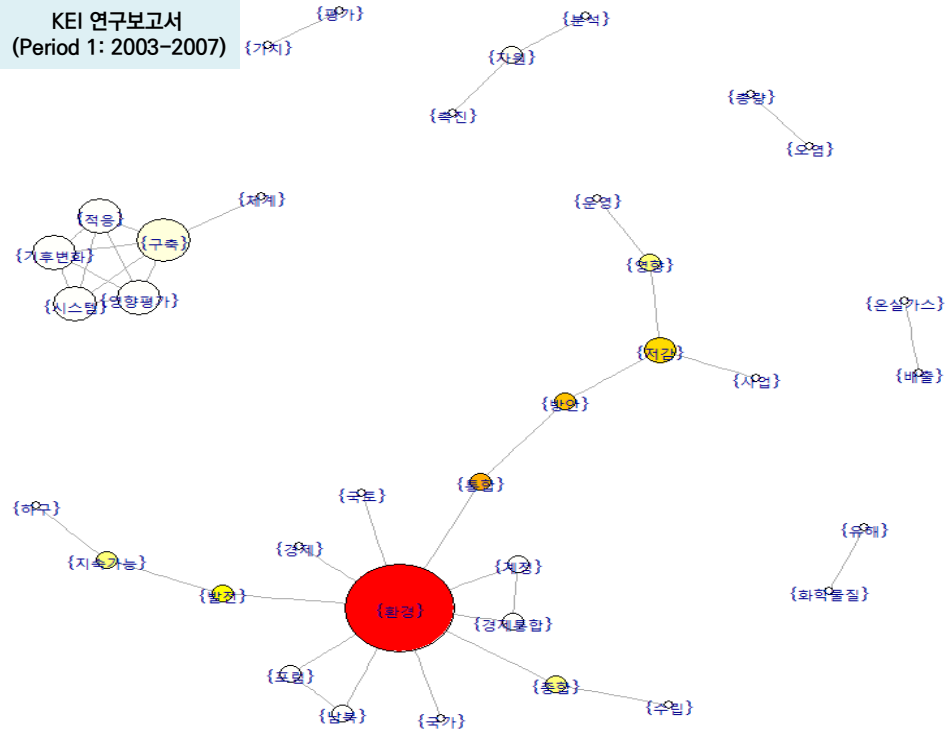
선행 연구

연구 내용

연구 추진방법

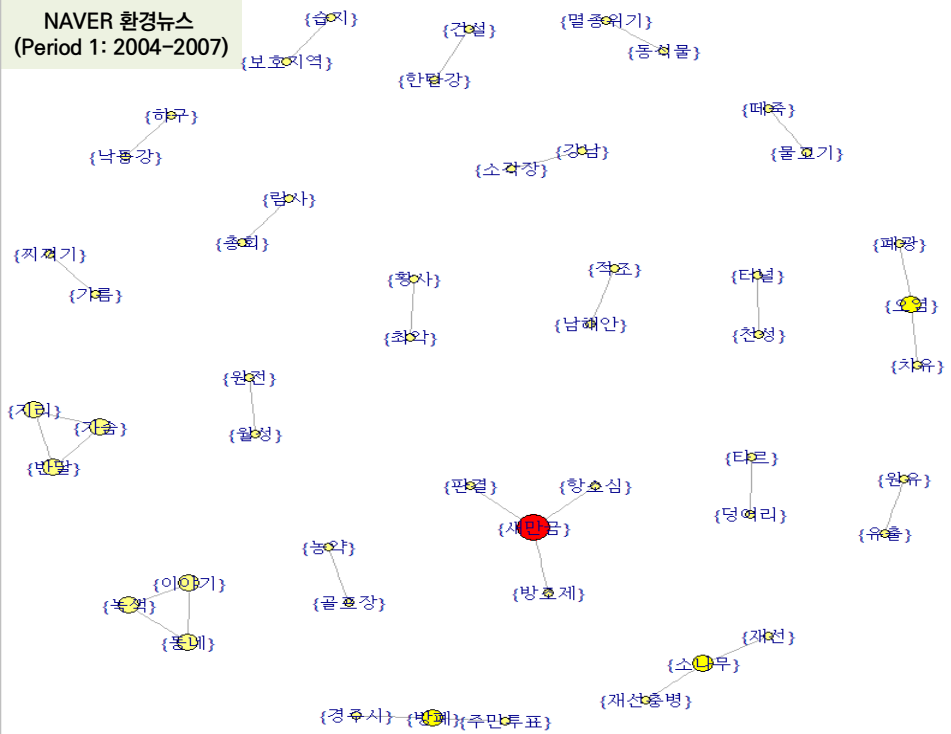
기대효과

### KEI 연구보고서 (Period 1: 2003-2007)

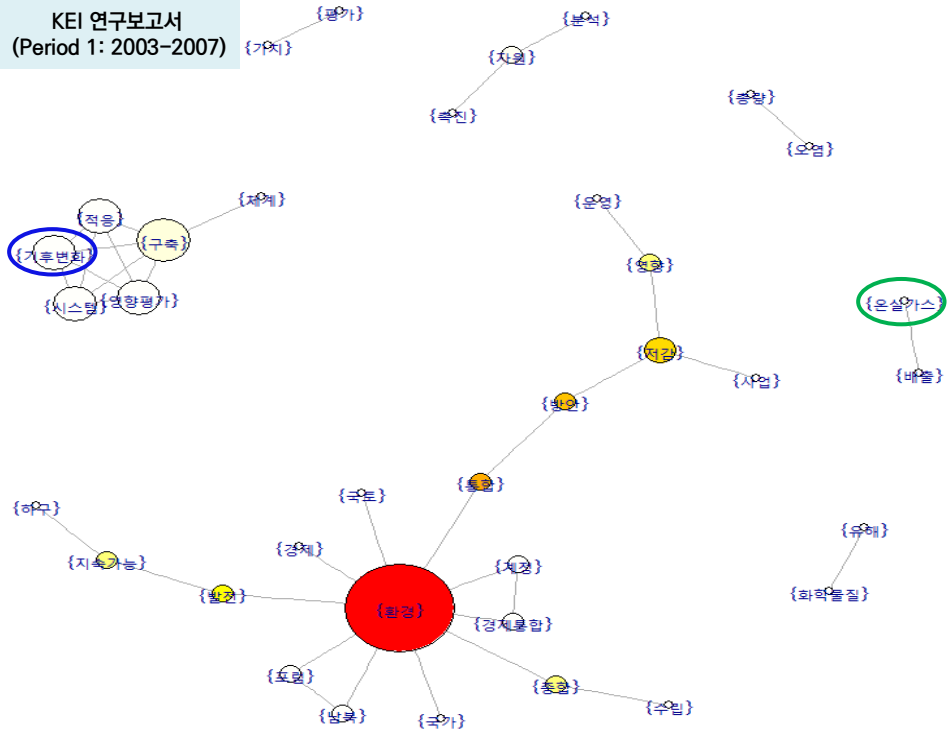


	A		B	신뢰도	지지도	항상도
1	하구	=>	지속가능	0.015	0.833	21.859
2	지속가능	=>	하구	0.015	0.385	21.859
3	경제통합	=>	환경	0.015	1.000	4.547
4	환경	=>	경제통합	0.015	0.067	4.547
5	경제	=>	환경	0.015	0.556	2.526
6	환경	=>	경제	0.015	0.067	2.526
7	남북	=>	포럼	0.012	1.000	68.200
8	포럼	=>	남북	0.012	0.800	68.200
9	자원	=>	분석	0.012	0.500	12.179
10	분석	=>	자원	0.012	0.286	12.179
11	시스템	=>	기후변화	0.012	0.444	12.630
12	기후변화	=>	시스템	0.012	0.333	12.630

### NAVER 환경뉴스 (Period 1: 2004-2007)



	A		B	신뢰도	지지도	항상도
1	동식물	=>	멸종위기	0.002	0.620	63.038
2	멸종위기	=>	동식물	0.002	0.235	63.038
3	주민투표	=>	방폐	0.002	0.671	32.029
4	방폐	=>	주민투표	0.002	0.098	32.029
5	방조제	=>	새만금	0.002	0.870	49.356
6	새만금	=>	방조제	0.002	0.108	49.356
7	남해안	=>	적조	0.002	0.489	60.799
8	적조	=>	남해안	0.002	0.231	60.799
9	월성	=>	원전	0.002	0.884	80.033
10	원전	=>	월성	0.002	0.139	80.033
11	경주시	=>	방폐	0.001	0.821	39.163
12	방폐	=>	경주시	0.001	0.062	39.163



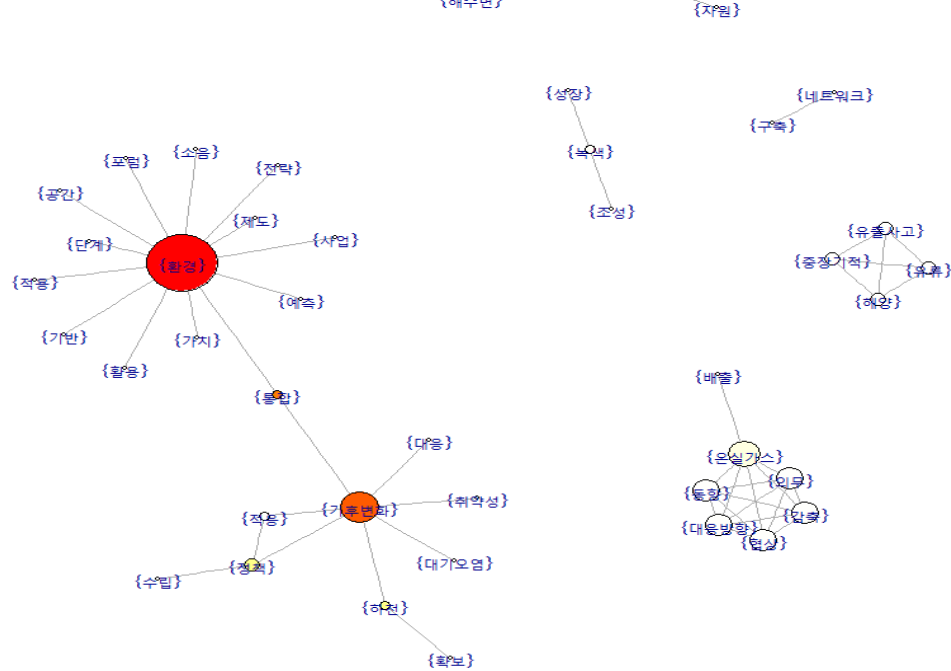
1. 기후변화 2. 온실가스 3. 태안 기름 유출 사고 4. 녹색성장 5. 환경오염 대책

- A. 기후변화 영향평가 및 적응시스템 구축, 온실가스 배출, 환경경제통합계정, 유해화학물질, 남북도림 키위드가 등장함.
- B. 기후변화 키워드를 중심으로 영향평가, 적응시스템 구축의 키워드가 나타남.



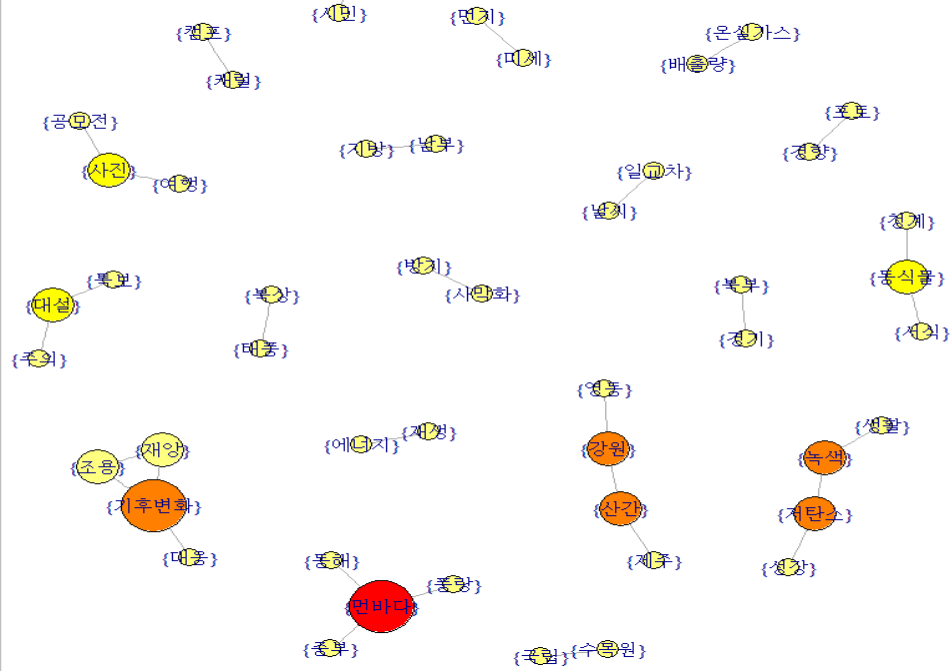
- A. 동식물-멸종위기, 지리산-반달가슴, 습지-보호지역, 람사르-총회  
 - 2004년 멸종위기 동식물 보호강화를 위해 멸종위기동식물의 범위를 대폭 확대함.  
 - 2006년 환경부에서 멸종위기 동식물 54종 증식복원사업 추진함.
- B. 새만금-판결-항소심-방조제  
 - 2006년 3월 새만금 사업은 긴 법정다툼 끝에 대법원 확정 판결이 나고 방조제 공사를 추진함.
- C. 원유-유출, 타르-덩여리, 기름-짜짜기, 물고기-폐죽음  
 - 2007년 12월 태안 기름 유출 사고 발생함.
- D. 남해안-적조, 경주지-방폐장-주민투표, 월성-원전, 황사-최악, 낙동강-하구 등의 키워드가 등장함.

KEI 연구보고서  
(Period 2: 2008-2012)



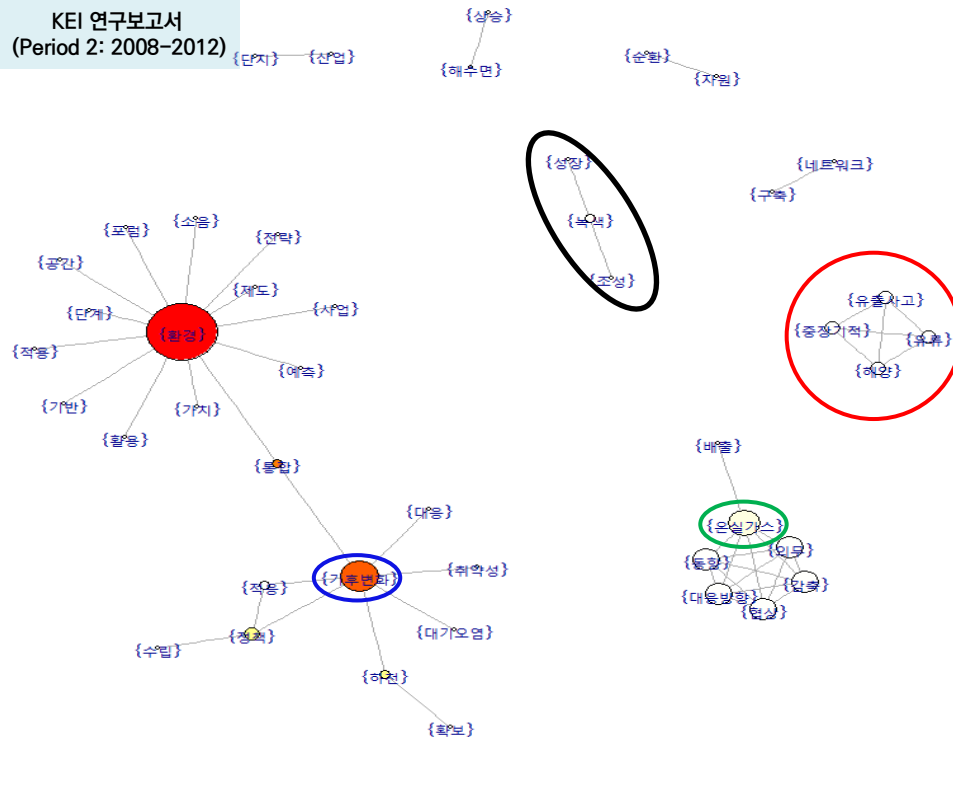
	A		B	신뢰도	지지도	항상도
1	중장기적	=>	유출사고	0.015	1.000	56.857
2	유출사고	=>	중장기적	0.015	0.857	56.857
3	중장기적	=>	유류	0.015	1.000	56.857
4	유류	=>	중장기적	0.015	0.857	56.857
5	대응방향	=>	동향	0.010	1.000	56.857
6	동향	=>	대응방향	0.010	0.571	56.857
7	대응방향	=>	감축	0.010	1.000	39.800
8	감축	=>	대응방향	0.010	0.400	39.800
9	대응방향	=>	온실가스	0.010	1.000	22.111
10	온실가스	=>	대응방향	0.010	0.222	22.111
11	협상	=>	온실가스	0.010	1.000	22.111
12	온실가스	=>	협상	0.010	0.222	22.111

NAVER 환경뉴스  
(Period 2: 2008-2012)



	A		B	신뢰도	지지도	항상도
1	대응	=>	기후변화	0.002	0.584	43.980
2	기후변화	=>	대응	0.002	0.143	43.980
3	대설	=>	주의	0.002	0.409	27.919
4	주의	=>	대설	0.002	0.126	27.919
5	재앙	=>	기후변화	0.002	0.566	42.564
6	기후변화	=>	재앙	0.002	0.121	42.564
7	저탄소	=>	녹색	0.002	0.612	39.608
8	녹색	=>	저탄소	0.002	0.100	39.608
9	복상	=>	태풍	0.001	0.673	62.289
10	태풍	=>	복상	0.001	0.138	62.289
11	대설	=>	특보	0.001	0.317	39.795
12	특보	=>	대설	0.001	0.179	39.795

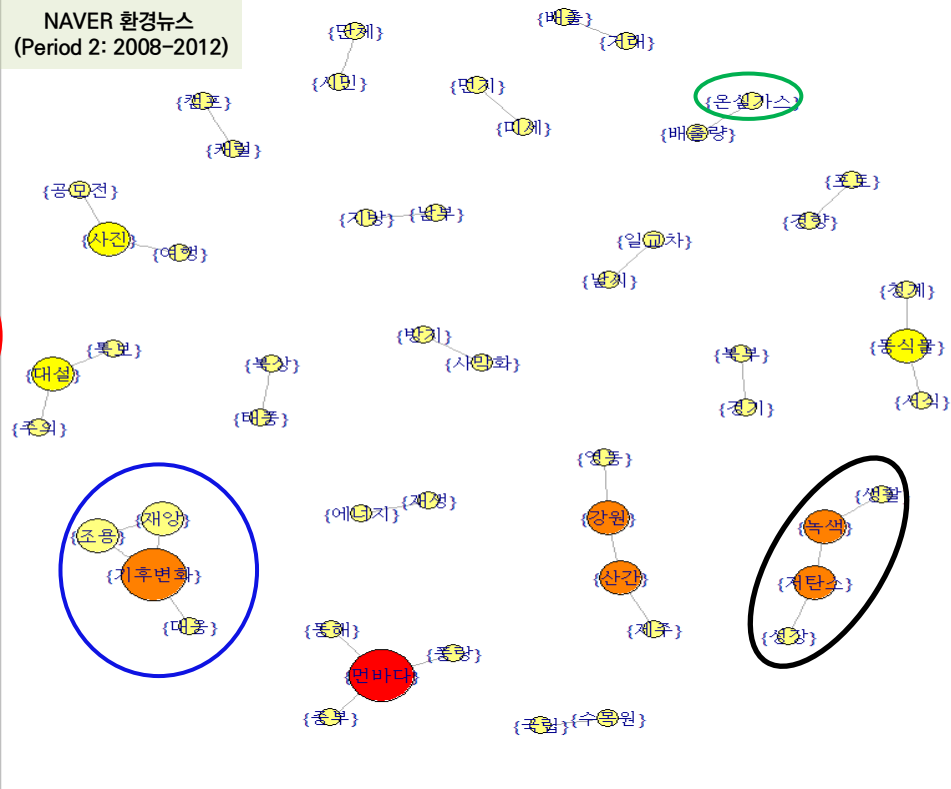
KEI 연구보고서  
(Period 2: 2008-2012)



1. 기후변화 2. 온실가스 3. 태안 기름 유출 사고 4. 녹색성장 5. 환경오염 대책

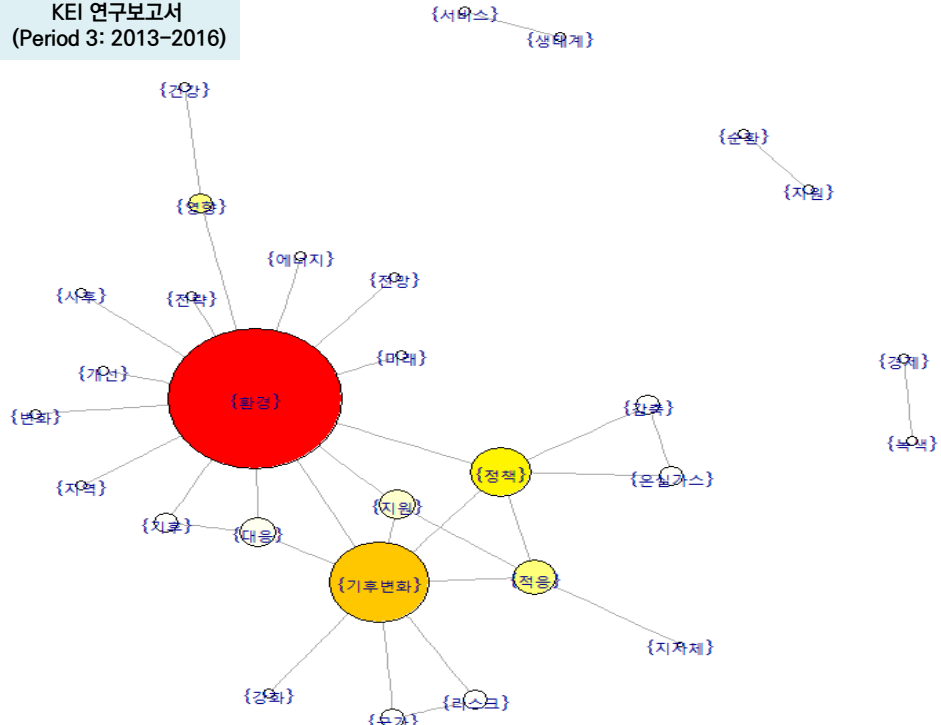
- A. 전 구간 대비 기후변화, 온실가스 키워드의 매개중심성이 높아짐.
- B. 해양 유류 유출사고, 녹색성장 조성, 해수면 상승, 소음 키워드가 새롭게 등장함.
- C. period1 NAVER 환경뉴스에서 등장한 원유-유출, 타르-덩어리, 기름-찌꺼기, 물고기-떼죽음 키워드가 의미하는 '태안 기름 유출 사고'의 중장기적 영향분석 및 제도개선 방안 마련과 관련한 연구가 진행되었음을 확인할 수 있음.

NAVER 환경뉴스  
(Period 2: 2008-2012)



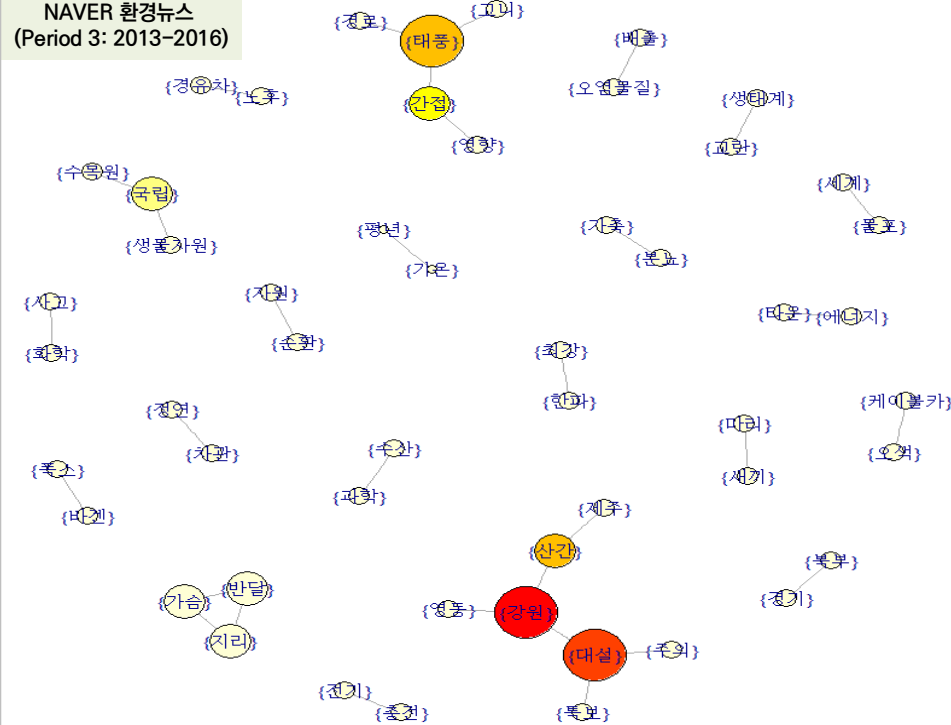
- A. 기후변화, 온실가스, 녹색성장 키워드가 새롭게 등장하고, KEI 연구보고서와 마찬가지로 중요하게 다뤄짐.
- B. 미세먼지, 재생에너지, 대설-추위-특보, 사막화-방지, 태풍-복상, 날씨-일교차 등의 키워드가 새롭게 등장함.
- C. period1에 이어 '동식물 보호'에 대한 기사가 많았음을 동식물-청계-서식 키워드를 통해 알 수 있음.

# KEI 연구보고서 (Period 3: 2013-2016)



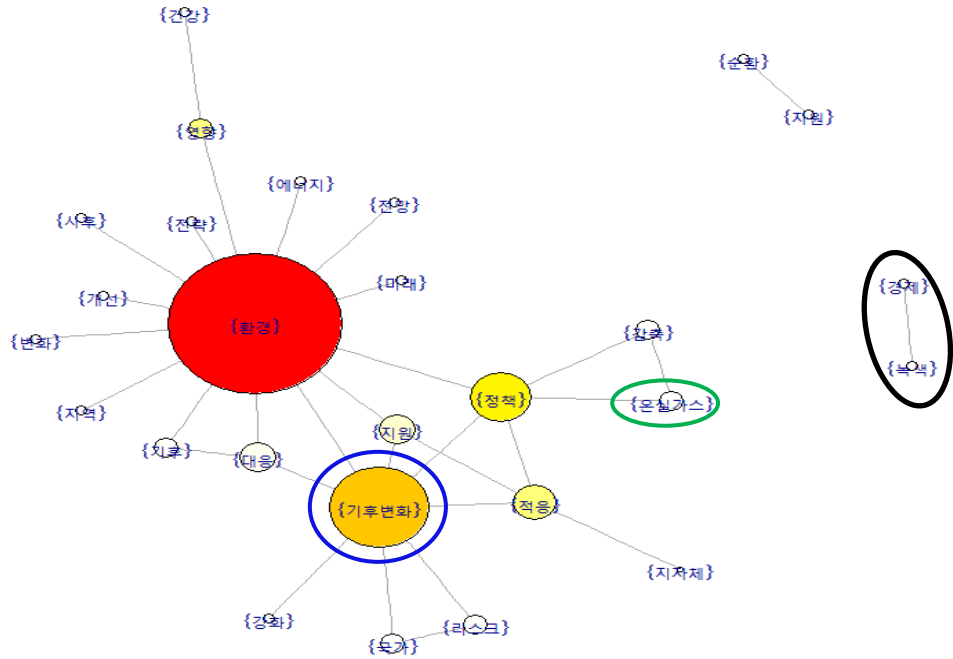
	A		B	신뢰도	지지도	향상도
1	정책	=>	기후변화	0.030	0.279	1.865
2	기후변화	=>	정책	0.030	0.200	1.865
3	대응	=>	기후변화	0.025	0.357	2.387
4	기후변화	=>	대응	0.025	0.167	2.387
5	전략	=>	환경	0.020	0.500	1.728
6	환경	=>	전략	0.020	0.069	1.728
7	기후변화	=>	환경	0.015	0.100	0.346
8	환경	=>	기후변화	0.015	0.052	0.346
9	리스크	=>	기후변화	0.013	0.500	3.342
10	기후변화	=>	리스크	0.013	0.083	3.342
11	미래	=>	환경	0.013	0.455	1.571
12	환경	=>	미래	0.013	0.043	1.571

# NAVER 환경뉴스 (Period 3: 2013-2016)



	A		B	신뢰도	지지도	향상도
1	간접	=>	태풍	0.002	0.960	51.279
2	태풍	=>	간접	0.002	0.081	51.279
3	대설	=>	특보	0.001	0.380	40.323
4	특보	=>	대설	0.001	0.128	40.323
5	최강	=>	한파	0.001	0.688	62.212
6	한파	=>	최강	0.001	0.105	62.212
7	고니	=>	태풍	0.001	0.675	36.055
8	태풍	=>	고니	0.001	0.061	36.055
9	반달	=>	지리	0.001	0.685	153.626
10	지리	=>	반달	0.001	0.237	153.626
11	가슴	=>	지리	0.001	0.667	149.529
12	지리	=>	가슴	0.001	0.237	149.529

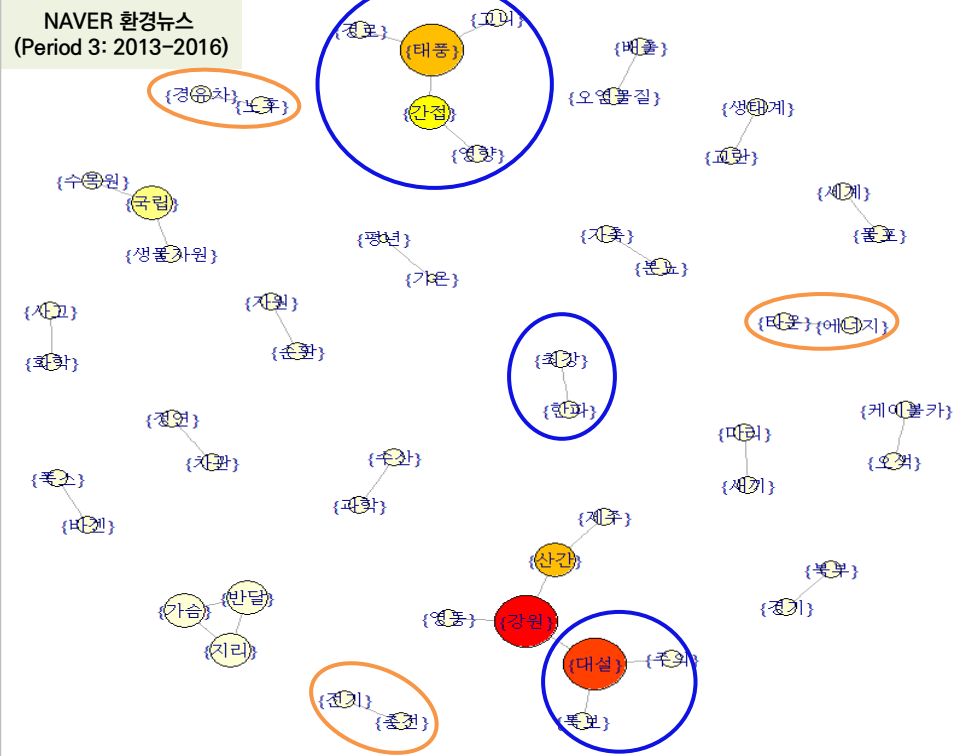
KEI 연구보고서  
(Period 3: 2013~2016)



1. 기후변화 2. 온실가스 3. 대안 기름 유출 사고 4. 녹색성장 5. 환경오염 대책

- A. period2에 이어 **기후변화** 키워드의 매개중심성이 높아짐
- B. 환경 키워드를 중심으로 **건강, 미래, 전망, 에너지** 키워드가 새롭게 등장함.
- C. period2에 '녹색' 키워드의 연관어였던 '성장' 키워드는 사라지고, '**경제**'가 연관어로 나타남.
- D. period2에 분산되어 있던 키워드 네트워크들이 **환경, 기후변화, 정책** 키워드를 중심으로 합쳐짐.

NAVER 환경뉴스  
(Period 3: 2013~2016)



- A. 기후변화 키워드는 사라지고 **태풍, 대설** 키워드의 매개중심성이 높아짐.
- B. 환경오염 대책을 위한 키워드가 등장함.
  - **노후 경유차**: 노후 경유차 폐차 후 두 달 안에 새차 구입 시 개별소비세를 줄여주는 제도
  - **전기 충전**: 환경 친화적 전기자동차의 개발 및 보급촉진을 위한 다양한 지원제도
  - **에너지 대안**: 환경부는 친환경 에너지타운을 전국으로 확산을 위해 노력
- C. period2에 이어 '**동식물 보호**'에 대한 기사가 많았음을 국립-수목원-생물자원, 지리-반달-기슴 키워드를 통해 알 수 있음.

연구 개요

선행 연구

연구 내용

연구 추진방법

기대효과

## Trend Comparative Analysis

### 1. 기후변화

KEI 연구보고서에서 '기후변화' 키워드는 시간이 지날수록 매개중심성(Betweenness Centrality, Cb)이 높아지고, 언급량이 많아지는 주요 키워드로 판단된다. 초기에는 기후변화 영향평가 및 적응시스템 구축관련 연구가 진행되었고, 이후 기후변화 대응을 위한 연구들(하천공간 확보방안 연구, 기후변화 적응정책 수립 연구)이 진행되었다. 또한, 기후변화와 대기오염을 함께 연구하는 경향을 보였다. 최근에는 이러한 기후변화 적응정책을 지원하기 위한 연구가 진행되었으며, 기후변화에 따른 국가 리스크 연구도 진행되었다.

NAVER 환경뉴스에서 '기후변화' 키워드는 Period 2(2008-2012)에 '재앙' 키워드와 함께 많이 등장하기 시작하였다. 이후 기후변화 키워드보다는 '태풍', '최강 한파', '대설 주의' 등의 기후변화 관련 키워드들이 많이 언급되었다. 이를 통해 '기후변화'의 이슈는 앞으로 계속 중요하게 다뤄질 문제이기 때문에 향후 KEI는 기후변화의 세부적인 현상들(가뭄, 온난화, 태풍, 홍수 등)에 대한 연구가 개별적으로 진행되어야 한다고 판단된다.

따라서 본 연구에서는 가장 중요하다고 판단되는 '기후변화' 키워드를 word2vec(skip-gram)분석을 통해 구체적으로 살펴보고자 한다.



연구 개요

선행 연구

연구 내용

연구 추진방법

기대효과

## Trend Comparative Analysis

2. 온실가스
3. 태안 기름 유출 사고
4. 녹색성장
5. 환경오염 대책

-> 향후 결과 작성 예정

## 기후변화 관련 키워드 Word2Vec(Skip-gram) 분석

연구 개요

선행 연구

연구 내용

연구 추진방법

기대효과

1. 가뭄	식수난	역부족	극심한	발작물	끌어오다
	2.29E-01	2.52E-01	2.76E-01	2.76E-01	2.97E-01
목마르다	해갈	타들다	저수율	단비	급수
3.07E-01	3.09E-01	3.18E-01	3.24E-01	3.65E-01	3.65E-01
저수량	물그릇	일조량	트라우마	생육	부족
3.77E-01	3.78E-01	3.83E-01	3.87E-01	3.90E-01	3.95E-01
기우제	용수	장강	전전금금	타들어가다	황하
3.97E-01	4.04E-01	4.10E-01	4.14E-01	4.26E-01	4.28E-01
한시름	해소	강수량	농업용수	안희정	홍수
4.28E-01	4.30E-01	4.33E-01	4.35E-01	4.39E-01	4.47E-01
이변	마른장마	최악	소양강댐	벼농사	모자라다
4.47E-01	4.55E-01	4.56E-01	4.57E-01	4.57E-01	4.60E-01
지독	광동	다목	상수	말라가	기근
4.61E-01	4.63E-01	4.63E-01	4.64E-01	4.66E-01	4.74E-01
강우량	엘니뇨	폭등	서북	밑바닥	부영양화
4.76E-01	4.80E-01	4.81E-01	4.81E-01	4.82E-01	4.83E-01
마르다	대지	농사	42년	수관	배춧값
4.84E-01	4.87E-01	4.90E-01	4.92E-01	4.93E-01	4.93E-01
남강댐	겹치다	당호	도움	호전	절수
4.94E-01	4.97E-01	5.01E-01	5.02E-01	5.03E-01	5.03E-01
저지대	흙탕물	4분	지천	저수	뱅해
5.04E-01	5.06E-01	5.06E-01	5.07E-01	5.07E-01	5.07E-01
초비상	자연재해	대비책	후난	방심	댐
5.07E-01	5.08E-01	5.08E-01	5.08E-01	5.11E-01	5.13E-01

2. 온난화	온난	2080년	표면	지중해	2100년
	1.58E-01	2.52E-01	2.77E-01	2.91E-01	2.94E-01
툰드라	해빙	아열대	2040년	난대림	세기말
2.97E-01	2.98E-01	3.04E-01	3.05E-01	3.21E-01	3.22E-01
기후대	지표면	급변	팽창	방귀	북반구
3.24E-01	3.28E-01	3.33E-01	3.34E-01	3.35E-01	3.41E-01
생장	플랑크톤	엽	해수면	금세기	심상찮다
3.42E-01	3.42E-01	3.43E-01	3.43E-01	3.50E-01	3.51E-01
지구	상승	산호초	그린란드	임팩트	영구동토
3.52E-01	3.52E-01	3.53E-01	3.54E-01	3.60E-01	3.64E-01
빙하	마그마	말매미	빠르다	이변	경고
3.66E-01	3.70E-01	3.75E-01	3.76E-01	3.79E-01	3.80E-01
급상승	빙상	북극	2000년	난민	장강
3.85E-01	3.90E-01	3.91E-01	3.91E-01	3.91E-01	3.92E-01
칭장	2050년	징후	온기	대박	트림
3.93E-01	3.94E-01	3.96E-01	3.97E-01	4.02E-01	4.04E-01
티베트	식량	난류성	급격	투발루	피시
4.04E-01	4.04E-01	4.05E-01	4.06E-01	4.08E-01	4.08E-01
북극곰	양미리	수확량	빙봉	아열대성	황제펭귄
4.09E-01	4.09E-01	4.10E-01	4.12E-01	4.13E-01	4.14E-01
70년	돌변	이대로	과학자	녹아내리다	시한폭탄
4.15E-01	4.15E-01	4.15E-01	4.16E-01	4.16E-01	4.17E-01
몰디브	90년	복어	재배지	패턴	코알라
4.17E-01	4.18E-01	4.19E-01	4.20E-01	4.21E-01	4.21E-01

## 기후변화 관련 키워드 Word2Vec(Skip-gram) 분석

연구 개요

선행 연구

연구 내용

연구 추진방법

기대효과

3. 태풍	크리	나스	간접	복상	할롱
	1.35E-01	1.38E-01	1.43E-01	1.55E-01	1.58E-01
오키나와	서진	12호	윙	나크	사니
1.65E-01	1.69E-01	1.70E-01	1.73E-01	1.74E-01	1.78E-01
16호	풍	11호	풍	10호	롤라
1.79E-01	1.83E-01	1.84E-01	1.88E-01	1.89E-01	1.95E-01
경로	15호	낭	중형	이파	마니
1.99E-01	2.00E-01	2.03E-01	2.04E-01	2.08E-01	2.16E-01
차바	저기압	유동	파탁	빠져나가다	파스
2.18E-01	2.18E-01	2.24E-01	2.26E-01	2.26E-01	2.29E-01
17호	피토	린	로사	카스	랍
2.30E-01	2.33E-01	2.39E-01	2.51E-01	2.54E-01	2.55E-01
우사기	도라지	필리핀	모라꽃	나리	룩
2.55E-01	2.59E-01	2.60E-01	2.61E-01	2.63E-01	2.63E-01
저압부	상륙	카눈	고다	강타	메아리
2.63E-01	2.64E-01	2.65E-01	2.72E-01	2.78E-01	2.87E-01
18호	파	벤	푸	진로	초속
2.88E-01	2.91E-01	2.92E-01	2.97E-01	3.01E-01	3.02E-01
폭풍우	동진	변질	13호	온대	열도
3.10E-01	3.14E-01	3.16E-01	3.16E-01	3.20E-01	3.24E-01
하이옌	예상	21호	폭우	비온다	남동쪽
3.25E-01	3.29E-01	3.29E-01	3.29E-01	3.36E-01	3.36E-01
흙	비슷하다	볼라	14호	규슈	위력
3.38E-01	3.42E-01	3.45E-01	3.50E-01	3.52E-01	3.54E-01

4. 폭설	춘설	적설량	눈	미시령	적설
	2.74E-01	2.89E-01	3.13E-01	3.35E-01	3.56E-01
춘삼월	밤새	진부령	빙판길	냉해	고갯길
3.59E-01	3.64E-01	3.75E-01	3.75E-01	3.79E-01	3.81E-01
결빙	제설	함박눈	한파	첫눈	혹한
3.86E-01	3.91E-01	3.94E-01	3.96E-01	3.96E-01	3.98E-01
시애틀	기습	속출	빙판	대설	폭탄
3.99E-01	4.00E-01	4.02E-01	4.02E-01	4.03E-01	4.03E-01
진눈깨비	<b>교통대란</b>	대설주의보	첫서리	갑작스럽다	흄비
4.09E-01	4.12E-01	4.18E-01	4.20E-01	4.20E-01	4.22E-01
동사	엄습	폭우	급강하	강추위	강원산
4.23E-01	4.24E-01	4.25E-01	4.29E-01	4.33E-01	4.34E-01
장대비	김화	축대	기우제	비닐하우스	<b>교통사고</b>
4.43E-01	4.43E-01	4.45E-01	4.46E-01	4.50E-01	4.51E-01
고립	최전방	봄눈	내리다	때늦다	첫얼음
4.54E-01	4.55E-01	4.56E-01	4.57E-01	4.60E-01	4.61E-01
난기류	전야	마스	얼다	서리	추돌
4.61E-01	4.64E-01	4.64E-01	4.67E-01	4.70E-01	4.70E-01
역부족	용평	퍼붓다	엄화칼슘	얼기	간밤
4.70E-01	4.72E-01	4.72E-01	4.74E-01	4.75E-01	4.76E-01
덕장	보슬비	집중호우	맹추위	날뛰기	남서부
4.76E-01	4.76E-01	4.76E-01	4.77E-01	4.77E-01	4.78E-01
정겹다	초겨울	비야	하늘길	대관령	몰아치다
4.78E-01	4.78E-01	4.79E-01	4.79E-01	4.79E-01	4.80E-01

## 기후변화 관련 키워드 Word2Vec(Skip-gram) 분석

연구 개요

선행 연구

연구 내용

연구 추진방법

기대효과

5. 폭우	집중호우	장대비	침수	큰비	속출	차바
	0.195054	0.223173	0.231008	0.300321	0.302235	0.306135
장맛비	카눈	밤새	아수라장	축대	윙	나스
0.306364	0.306911	0.307038	0.307994	0.316263	0.317301	0.317346
나크	만조	이파	우박	태풍	이재민	할퀴다
0.318857	0.319282	0.32342	0.323558	0.329145	0.330582	0.331176
폭탄	크리	저기압	폭풍우	풍	갑작스럽다	풍
0.332907	0.334293	0.335054	0.338794	0.340517	0.342885	0.345542
간밤	국지	낙뢰	빠져나가다	범람	12호	초속
0.34609	0.350649	0.351028	0.352279	0.352389	0.352854	0.354327
하교	파탁	메아리	호우	서진	사니	저지대
0.358288	0.359075	0.362549	0.366078	0.367529	0.370019	0.370566
마니	휩쓸다	강우	기압골	할롱	모라꽃	10호
0.371405	0.37564	0.377975	0.37819	0.37842	0.379908	0.381248
로사	파스	오키나와	롤라	발기다	강타	물벼락
0.38386	0.384685	0.387816	0.389042	0.389309	0.390812	0.391989
하이엔	예상	서늘하다	천호동	간접	나리	16호
0.392596	0.394707	0.395308	0.395629	0.396531	0.397616	0.399087
15호	교통대란	정전	찾아들다	푸	소낙성	마른장마
0.399589	0.39974	0.399782	0.400582	0.401197	0.40392	0.405028
퍼붓다	비	17호	돌풍	황해도	카스	휩쓸리다
0.405908	0.406354	0.40647	0.407229	0.40882	0.409197	0.411872
18호	쏟아지다	높이	흩비	비온다	전도	규슈
0.413255	0.413439	0.416529	0.416646	0.417325	0.418635	0.419606

6. 한파	강추위	맹추위	추위	영하	초겨울
	2.04E-01	2.32E-01	2.52E-01	2.53E-01	2.66E-01
꽃샘추위	춡다	동장군	김화	대설	엄습
2.77E-01	2.80E-01	2.84E-01	3.01E-01	3.02E-01	3.10E-01
혹한	적설	12도	칼바람	13도	건조
3.22E-01	3.29E-01	3.30E-01	3.31E-01	3.32E-01	3.35E-01
7도	울겨울	대설주의보	꽃샘	매섭다	강풍
3.37E-01	3.42E-01	3.44E-01	3.44E-01	3.45E-01	3.46E-01
세밀	수은주	8도	운종일	6도	14도
3.47E-01	3.50E-01	3.52E-01	3.53E-01	3.54E-01	3.59E-01
전야	때늦다	소한	15도	최전방	전날
3.59E-01	3.59E-01	3.65E-01	3.66E-01	3.67E-01	3.69E-01
백운	5도	삼한사온	최강	18도	입춘
3.70E-01	3.78E-01	3.79E-01	3.79E-01	3.81E-01	3.82E-01
모스크바	0도	3도	마스	첫얼음	다소
3.82E-01	3.83E-01	3.83E-01	3.86E-01	3.87E-01	3.88E-01
4도	수험	첫서리	9도	2도	소집
3.89E-01	3.90E-01	3.92E-01	3.93E-01	3.93E-01	3.95E-01
폭설	폴리다	적설량	1도	불조심	맹위
3.96E-01	3.96E-01	3.96E-01	3.96E-01	3.98E-01	3.98E-01
누그러지다	동파	11도	폭염	미시령	수그러지다
3.99E-01	3.99E-01	4.03E-01	4.03E-01	4.04E-01	4.05E-01
진부령	빙판길	38도	눈발	성탄	수능
4.05E-01	4.05E-01	4.06E-01	4.06E-01	4.06E-01	4.07E-01

## 기후변화 관련 키워드 Word2Vec(Skip-gram) 분석

연구 개요

선행 연구

연구 내용

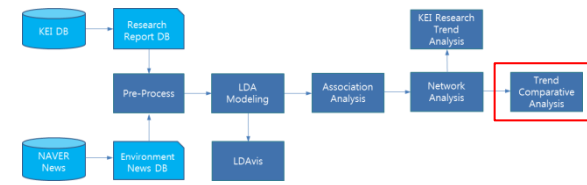
연구 추진방법

기대효과

7. 해수면	상승	2100년	2040년	지중해
	2.16E-01	2.40E-01	2.66E-01	2.97E-01
금세기	수온	저지대	지표면	2080년
3.00E-01	3.16E-01	3.20E-01	3.22E-01	3.36E-01
팽창	온난화	투발루	닥치다	급상승
3.38E-01	3.43E-01	3.49E-01	3.57E-01	3.58E-01
표면	해빙	심상찮다	수확량	방글라데시
3.59E-01	3.64E-01	3.82E-01	3.83E-01	3.85E-01
39년	직면	툰드라	세기말	시한폭탄
3.88E-01	3.89E-01	3.91E-01	3.91E-01	3.95E-01
마그마	평균	2050년	침몰	난민
3.97E-01	3.97E-01	3.99E-01	4.01E-01	4.02E-01
지구	과학자	공황	섬나라	자연재해
4.04E-01	4.05E-01	4.06E-01	4.06E-01	4.06E-01
임팩트	일조량	연평균	식탁	산호초
4.07E-01	4.07E-01	4.12E-01	4.15E-01	4.16E-01
난대림	온난	타격	다습	연령
4.17E-01	4.18E-01	4.18E-01	4.18E-01	4.19E-01
초흔	양미리	세기	빈도	피시
4.20E-01	4.23E-01	4.23E-01	4.23E-01	4.26E-01
이래	경고	기후대	50년	43년
4.27E-01	4.31E-01	4.32E-01	4.33E-01	4.34E-01
둔화	돌변	갈치	90년	가라앉다
4.34E-01	4.34E-01	4.37E-01	4.38E-01	4.38E-01

8. 홍수	가물막이	사방댐	후난	보	제방
	2.81E-01	3.14E-01	3.31E-01	3.39E-01	3.74E-01
홍수조절	무너지다	탁수	극하다	천보	싼샤
3.82E-01	3.87E-01	3.91E-01	3.94E-01	3.95E-01	3.97E-01
소방방	유속	홍수피	대홍수	붕괴	댐
3.99E-01	4.04E-01	4.08E-01	4.10E-01	4.11E-01	4.13E-01
광동	임하	휴탕물	황하	조절	저수량
4.16E-01	4.19E-01	4.21E-01	4.21E-01	4.22E-01	4.24E-01
미보	두만강	다목	침하	물그릇	산사태
4.25E-01	4.25E-01	4.28E-01	4.29E-01	4.29E-01	4.30E-01
이변	대강	치수	군남	범람	팔당댐
4.31E-01	4.38E-01	4.38E-01	4.40E-01	4.41E-01	4.46E-01
가뭄	비만	임진강	저수조	결여	항공업
4.47E-01	4.51E-01	4.51E-01	4.51E-01	4.52E-01	4.53E-01
쓰촨	기근	부작용	재해	침수	합
4.54E-01	4.55E-01	4.55E-01	4.55E-01	4.56E-01	4.58E-01
필승	준설	급하다	수위	속도전	강경
4.58E-01	4.59E-01	4.61E-01	4.61E-01	4.62E-01	4.63E-01
역행	땀질	강우	수자원	탄강	물난리
4.64E-01	4.64E-01	4.65E-01	4.65E-01	4.69E-01	4.69E-01
저류지	보강	예당저수지	누수	역부족	황강
4.69E-01	4.69E-01	4.70E-01	4.70E-01	4.71E-01	4.72E-01
강바닥	지천	곡보	저수율	도암	퇴적
4.73E-01	4.74E-01	4.75E-01	4.75E-01	4.77E-01	4.79E-01

## Text Mining Process List



Plan 2017	Process	Code	Description	Input	Output	Note
8월	<b>Trend Comparative Analysis</b>					
9월	시사점 도출 및 정책제언					
10월	<b>향후 계획 수립</b>		<ul style="list-style-type: none"> <li>- 기후변화 세부현상들(가뭄, 온난화, 태풍, 폭설, 폭우, 한파, 해수면상승 등)에 대한 연구</li> <li>- 분야별(정치, 경제, 사회, 세계, it/과학, 오피니언) 기후변화 세부현상들에 대한 인식 분석</li> </ul>			

연구 개요

선행 연구

연구 내용

연구 추진방법

기대효과

연구 개요

선행 연구

연구 내용

연구 추진방법

기대효과

## 학술적 기대효과

- 장기간에 걸친 KEI 연구 동향을 정리하여 추후 환경연구 기획에 필요한 정보를 원내외 연구진에게 제공
- 환경분야 텍스트 마이닝 분석기반 플랫폼 개발의 기초 구성
  - 환경관련 키워드 빈도 분석, 연관성 분석, 토픽 클러스터링 등 다양한 텍스트 마이닝 분석 기법 집적 가능
  - 추후 이들 기법을 자동으로 처리하는 플랫폼을 구축하는 기초로 활용 가능

연구 개요

선행 연구

연구 내용

연구 추진방법

기대효과

## 후속 연구

- 매체별 환경문제 인식 성향 분석을 소셜미디어, 전통미디어, 전문사이트(학술논문), 공공기관 발간문건 등으로 확대하여 연구동향과 사회적 인식간의 관계파악 범위를 확대
- KEI 제공 발간물 데이터 시각화 서비스를 구축하여 사용자의 이용 편이 증진
  - 대량의 KEI 발간물 데이터에 대한 정보를 사용자가 효율적으로 파악할 수 있도록 정보 전달력을 제고
- 환경연구 트렌드 분석을 활용하여 미래 환경연구 수요 예측에 반영
  - 기존의 정량적 전망을 활용한 미래 환경문제 예측을 반영하는 연구수요 예측과 매체 분석을 통해 수요자 선호를 반영하는 연구수요 예측을 병행



연구 개요

선행 연구

연구 내용

연구 추진방법

기대효과

## 정책 개발

- 환경정책 수요자의 선호를 정책개발에 활용하여 “환경서비스 품질수준 제고<sup>1)</sup>” 도모 가능
    - 매체별 환경분야 연구동향과 사회적 요구를 비교분석한 결과를 근거로 수요자의 선호를 파악하여 정책 개발 기초 자료로 활용
- 1) 국정과제 95. 생활환경 취약지역 개선 및 환경질 개선의 과제개요

**Thank you.**