

## LDA 기법을 이용한 미세먼지 이슈의 토픽모델링 분석

윤순욱\* · 김민철\*\*†

\*녹색기술센터 정책연구부 Post-doc., \*\*녹색기술센터 정책연구부 선임연구원  
(2020년 4월 3일 접수, 2020년 5월 15일 수정, 2020년 5월 20일 채택)

## Topic Modeling on Fine Dust Issues Using LDA Analysis

Yoon soonuk\* · Kim Minchul\*\*†

Green Technology Center

(Received 3 April 2020, Revised 15 May 2020, Accepted 20 May 2020)

### 요 약

본 연구에서는 최근 10년간의 미세먼지 관련 뉴스 데이터를 수집하여 LDA 분석을 통해 최적 토픽을 도출하였다. 최적 토픽으로 선별된 80개의 이슈를 미세먼지 정책의 시각에서 해석하였다. 연구결과, 기온과 같은 날씨와 관련된 정보와 미세먼지 농도가 관련되어서 이슈화되는 경향이 있었다. 다음으로 미세먼지 저감 대책의 일환으로 노후경유차 운행 제한 제도와 저감 장치 부착과 같은 이슈의 빈도수가 높았다. 국민에 대한 제도 변경 안내를 포함하여 시민과 운수업자와의 갈등도 주요한 토픽으로 나타났다. 미세먼지 문제의 해결을 위한 수소차 보급과 같은 대안도 주요 토픽으로 분석되었다. 또한 미세먼지 관련 공기청정기 등 제품 관련 주제, 취약계층을 미세먼지로부터 보호하는 정책과 관련된 주제, 연구개발을 통한 미세먼지 저감 관련 주제가 주요 화두로 제기되었다. 미세먼지 대책은 사회 이슈로 정부 정책과 밀접한 관련이 있다고 볼 수 있다. 또한 본 연구를 통해 토픽 상에서는 거시적인 정부정책 자체보다는 시민의 안전, 시혜적인 정책이나 이해관계자간의 갈등이 정부정책 변화와 연동하여 중요한 의미를 지니는 것으로 나타났다.

**주요어 :** 미세먼지, 잠재적 디리클레 할당모형(LDA), 토픽 모델링, 정부 정책, 공기청정기, 취약계층, 연구 개발

**Abstract -** In this study, the last 10 years of news data on fine dust was collected and 80 topics are selected through LDA analysis. As a result, weather-related information made up the main words for the topic, and we can see that fine dust becomes a big issue below 10 degrees Celsius. The frequency of exposure to the media and the maximum concentration of fine dust are correlated with positive. Topics related to fine dust reduction measures and the government's comprehensive measures over the past decade, topics related to products such as air purifiers related to fine dust, topics related to policies protecting vulnerable people from fine dust, and topics on fine dust reduction through R&D were found to be major topics. Measures against fine dust as a social issue can be seen to be closely related to the government's policy.

**Key words :** Fine Dust, LDA(Latent Dirichlet Allocation), Topic modeling, Government Policy, Air Purifier, Vulnerable People, R&D

### 1. 서 론

2019년 미세먼지 문제는 국가의 재난으로 선정

되었고 가장 미디어에 노출이 많이 된 키워드였다. 그러나 어떠한 이슈가 어떤 형태로 제기되었는지 그리고 이러한 이슈가 국가정책과 어떠한 관련성이 있는지에 대한 심도가 있는 연구는 부족했다. 이에 본 연구에서는 미세먼지에 대한 최근 10년의 뉴스 데이터를 수집하여 키워드를 추출하고, 이를 잠재적

†To whom corresponding should be addressed.  
Tel : +82-2-3393-3915 E-mail : eco@gtkc.re.kr

디리클레 할당모형(Latent Dirichlet Allocation, LDA) 분석을 통해 토픽을 선별하여 주요 이슈를 도출하였다. LDA 분석을 통해 도출된 미세면지 관련 이슈들이 현재 미세면지 관련 정책과의 부합성이 있는지를 분석하고, 향후 추진될 정책에 인사이트를 제공하고자 하는 것이 본 연구의 목적이다. LDA 분석과 같은 데이터마이닝은 숫자 등 구조화(structured database)되거나, 또는 문자, 문서 등처럼 비구조화(unstructured database)되어 있는 데이터베이스를 분석하여 의미있는 규칙이나 패턴을 도출하는 방법으로 정책이슈 분석에도 의미있는 방법론이 된다. 미세면지라는 국민적인 이슈에 대해 이러한 방법론을 적용해보고 사회문제 이슈화된 미세면지 키워드와 토픽을 속에 관련 정책을 분석하기로 한다.

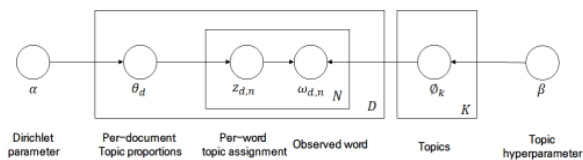
## 2. 연구 방법

### 2-1 LDA 분석 방법론

LDA는 토픽의 단어 비중과 문서의 토픽 비중이라는 두 가지 변수의 결합 확률분포에 따라 문서의 토픽을 찾는 과정이다. 두 변수 모두 양의 실수를 요소로 가지며, 모든 요소를 더한 값이 1이 되는 Dirichlet 분포를 따른다(Moon et al., 2018). 가정하는 문서생성과정을 자세히 설명하자면,  $D$ 는 말뭉치 전체문서개수,  $K$ 는 전체 토픽수(하이퍼파라미터),  $N$ 은  $d$ 번째 문서의 단어수를 의미한다. 즉, 아래 Fig. 1.의 수식과 같이 LDA 모델을 모두 정리하면  $d$ 번째 문서  $i$ 번째 단어의 토픽  $z_{d,i}$ 가  $j$ 번째에 할당될 확률을 말한다.

$$p(z_{d,i} = j | z_{-i}, w) = \frac{n_{d,k} + \alpha_k}{\sum_{i=1}^k (n_{d,i} + \alpha_i)} \times \frac{v_{k,w_{d,i}} + \beta_{w_{d,i}}}{\sum_{j=1}^V v_{k,j} + \beta_j} = AB \quad (1)$$

여기서 토픽수  $K$ 값을 지정하는 방법으로는 크게 perplexity 또는 topic coherence 점수가 있다(Seo, 2019). 본 연구에서는 2010년 D Newman이 개발한 'Coherence Score'을 이용하여 산출할 수 있다(Green Technology Center Korea, 2019).



\*출처: Park et al., 2015

Fig. 1. Schematic of LDA methodology.

### 2-2 LDA 기법을 활용한 선행연구 분석

한국보건사회연구원(2018)은 전문가 브레인스토밍 등을 통해 주요 정책 이슈 키워드를 도출하여 이를 기반으로 소셜 데이터를 대상으로 LDA 분석을 수행하였다. 도출된 보건복지 이슈와 4차 산업혁명 기술의 도입 가능성을 분석하기 위한 텍스트 수집, 자료 정제, 시각화 등을 통해 보건복지 분야에서 활성화되고 있는 산업들의 이슈와 연계를 시도하였다. 연구 결과, 블록체인 기술의 의료 관련 정부 정책, 의료 관련 플랫폼 이슈, 복지 관련 산업 등이 도출되었으며, '블록체인 기술 발전전략'의 발표시기 (2018.6)와 맞물려 이슈화된 문제들이 드러난 결과를 도출함으로써, 정책 연구 전략을 세우는 근거자료로의 활용이 가능하다고 제안하였다. Tobias R. Keller 외 (2019)는 인도의 기후변화 관련 기사를 LDA 토픽 모델링을 통해 분석하여 20년간 기후변화 범위가 확장된 것을 밝혀냈다. 최빈국의 언론 보도에 초점을 맞춘 연구 사례로, 연구 결과는 인도에서 기후변화의 범위가 20년 동안 기후변화 영향(Climate change impacts), 기후 과학(Climate science), 기후 정치(Climate politics) 및 기후변화와 사회(Climate change and society)의 주요 주제에 대해 28가지 토픽으로 분류할 수 있었다. 이를 통해 언론의 변화가 사람을 교육하고 정책을 변화시킬 수 있는 잠재력과 영향력을 가지고 있음에 대해 논의하였다. Hee Jay Kang 외 (2019)는 생화학분야의 연구 속도가 빠르게 진화하면서, 분자생물학, 합성화학 및 생물물리학 등의 방법론들이 다양화되었다는 것을 인식하고, 20년 동안 생화학분야의 연구주제를 파악하고 추세변화를 정량적으로 분석하기 위해 토픽 모델링 기법을 이용하였다. 생화학의 15가지 주요 토픽을 구분하여 시간별 연구 분야 당 연구량을 분석하였는데, 연구 결과는 생화학 산업의 동향과 매우 유사한 것으로 분석되었다. 지속적으로 데이터를 추가하여 적용하면 산업계와 학계의 실무자에게 유용한 의사 결정 도구로써 제공이 가능하다는 것을 제시하였다. Bach Xuan Tran 외 (2019)는 인공지능(AI) 기술이 의약분야에서 응용되고 있으나, 이에 대한 생산성, 과정, 주제, 연구 환경 등이 부족하다는 문제를 인지하였다. 이에 40년간의 연구 간행물을 대상으로 AI가 의학분야의 연구 환경의 변화에 미치는 영향에 대해 LDA 분석을 수행하였다. 연구 결과, 의학 분야에서 AI 적용은 임상실습(Clinical practice), 임상자료(Clinical material) 및 정책(policies)을 주요 토픽으로 분류하였으며, AI를 통해 개발도상국과 선진국 사이의 건강 관리, 의약품 제공 등 간격을 좁힐 수 있는 방법으로 기술개발의 당위성을 강조하였다.

## 2-3 데이터 수집

LDA 분석 대상 언론 기사는 한국언론진흥재단이 제공하는 뉴스 빅데이터 분석 서비스인 빅카인즈 (BIGKinds, <http://www.kinds.or.kr>)의 뉴스 통합 데이터베이스를 이용하였다. 검색어는 ‘미세먼지’를 활용하였으며, 최근 10년 (2010.1.1.~2019.10.31.)의 언론사 기사를 수집하였다. 검색결과 중 미세먼지와 관련없는 ‘스포츠’, ‘문화’ 섹션에 분류된 기사를 배제하였다. 각 연도별 최종 수집된 뉴스 DB는 152,990건이었다. 데이터는 엑셀파일 형태로 수집되며, 기사 데이터가 포함하는 정보는 ‘뉴스 식별자’, ‘일자’, ‘언론사’, ‘기고자’, ‘제목’, ‘분류’, ‘인물’, ‘위치’, ‘특성추출’, ‘본문’이 있는데, 그 중 특성추출 부분을 분석에 활용하였다.

빅카인즈를 이용한 뉴스추출 방식은 신문과 방송사 뉴스의 반정형·비정형 상태의 뉴스도 정형화되었고 무엇보다 최근 10년간 국내 이슈로는 가장 많은 데이터가 축적되어 있기에 연구의 질적 향상에 도움이 된다. 또한 코딩 대신 자연어 처리(natural language processing: NLP) 방식을 활용하기 때문에 코더 간 신뢰도를 측정할 필요가 없다는 장점이 있다(이은별 외, 2017). 언론 기사는 대중에게 가장 많이 노출되는 설득 수단이며, 대중이 ‘무엇에 대해’ 생각할지 전달하고 그 의제가 주요하게 다뤄지도록 하는데 기여한다. 특히, 국민들의 미세먼지에 대한 관심도가 증가하면서, 미세먼지와 관련된 정책이 마련되고 있다. 언론 기사를 통해 미세먼지와 연계되어 자주 언급되는 이슈가 무엇인지, 유관 정책, 산업, 과학기술 분야의 토픽이 무엇인지를 언론 기사를 대상으로 정성, 정량적인 분석을 하고자 하며, 이러한 분석을 통해 미세먼지 정책과 이슈와의 연계성 등을 분석하고자 한다.

## 2-4 이슈 도출 및 분석

최종적으로 구축한 기후변화 적응 분야별 언론 기사 DB를 활용하여 키워드 빈도 기준의 워드클라우드 분석을 R프로그램으로 수행하여 6대 적응 분야별 50대 주요 키워드를 시각화하여 보여주었다. 그리고 파이썬(python)을 통하여 토픽 분석 기법 중 하나인 LDA로부터 적응 분야별 주요 이슈를 분석·발굴하고, 주요 이슈에 대한 미래수요 분석을 수행하였다. 도출된 토픽을 구성하는 키워드들의 비중, 빈도를 확인하고, 토픽에 적합한 기사의 원문을 통해 이슈를 분류하였다. 지역 키워드와 같이 뉴스기사의 특성상 빈도수가 높은 이슈를 제외한 다른 이슈들의 키워드를 조합하여 기사를 검색하여 이슈가 의미하는 바를 찾아냈다. 이를 통해 미세먼지 이슈와 관련 산업, 정책 동향을 제시하였다.

## 3. 연구 결과

### 3-1 LDA 분석을 통한 토픽 도출 결과

수집된 뉴스기사 데이터를 토대로 미세먼지 관련 최적의 토픽수를 산출하기 위한 perplexity 지수는

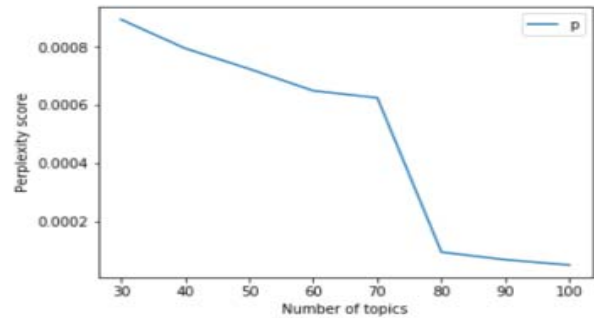
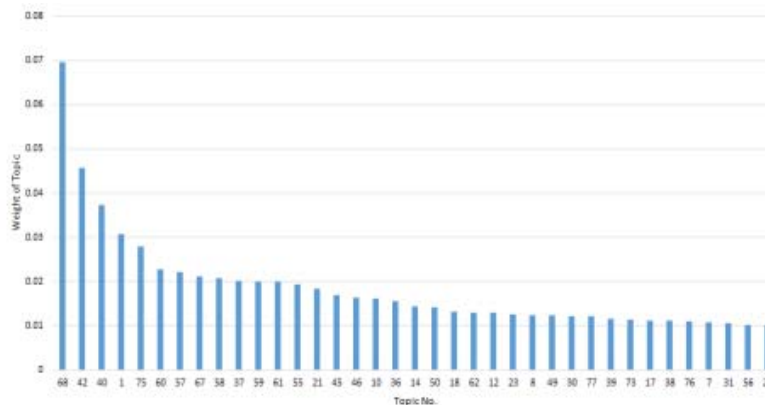


Fig. 2. Optimal number of topics.



\* Skipping the Topic Ratio of 0.01 or less

Fig. 3. Weight distribution by topics



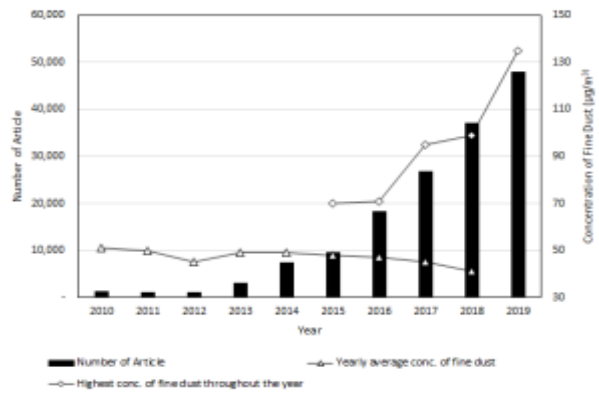




75번 토픽은 주요 단어들의 구성을 통해 미세먼지 저감 대책으로의 노후경유차 운행제한 제도와 저감 장치 부착과 같은 조치로서 대국민 제도 안내 이상으로 이를 둘러싼 시민, 운수업자 등들의 갈등도 많이 소개되었다. 67번 토픽은 2018년 공기청정기의 렌탈 사업 시장 점유율이 2위로 높아졌다는 이슈에서 도출된 토픽으로 공기청정기 소비형태 변화가 미세먼지와 관련된 민간경제와 산업에도 영향을 미치고 국민의 심리도 반영한다는 것을 파악할 수 있었다. 58번 토픽은 2016년 미세먼지 종합대책 발표를 통해 경유차와 석탄화력발전 감축 정책이 이슈에서 도출된 토픽이었다. 2019년에는 미세먼지 특별법이 시행되고 국무총리 소속으로 미세먼지 개선기획단도 설치되었다. 그러나 법제 변화보다는 시민들의 생활에 더욱 관련성이 높은 정부대책과 규제와 관련된 토픽이 부각되었다. 61번 토픽은 경기도 어린이집 공기청정기 보급사업 관련 이슈로부터 도출되었다. 지자체들은 유치원 공기청정기 설치 예산 지원을 통해 미세먼지로부터 어린이, 노인 등 취약계층을 보호하고 환경권을 보장하는 노력을 하고 있음을 알 수 있다. 55번 토픽은 2016년 수소차 출시와 맞물려 수소차가 미세먼지 저감하는 역할을 하는 이슈로부터 도출되었다. 정부의 '수소차 보급 및 시장 활성화 계획'과 2019년 1월 수소경제 활성화 로드맵이 수소차와 관련된 이슈였지만 수소 저장장치의 안전관리 이슈와 충전소 보급 확대의 어려움과 같은 현실적인 문제도 제기되었다. 45번 토픽은 극심한 미세먼지로 인해 미세먼지 저감 과학기술 개발이 적극적으로 이루어지고 있는 이슈로 도출되었다. 드론을 이용한 먼지포집, 화학물질 살포, 정전기 등 다양한 미세먼지 저감 방법의 연구 개발에 주목했다. 법제상 2019년 3월 26일, 미세먼지의 배출량 정보를 수집하고 분석하는 미세먼지센터의 설치가 강행규정화 되었고, 미세먼지에 대한 국제 협력에 있어서 과학적인 데이터 확보와 연구 개발은 주요한 이슈로 나타났다.

### 3-3 미디어 분석 결과와 미세먼지 현상과의 상관관계

기사 데이터를 수집한 Kinds 사이트를 통해 미세먼지 키워드의 연도별 기사 수를 조사하였다. 이와 더불어 실제 미세먼지 연평균 농도와 연중 최고치는 다음 그림과 같다. 미세먼지 연평균 농도는 2013년 이후 점차 감소세를 나타내고 있으나, 연중 농도 최고치는 2015년 이후 가파르게 상승하고 있다. 또한, 미세먼지 고농도 발생일 수 역시 2015년 5일에서 2017년 16일로 급증하여, 2019년도 역시 16일로 지속되고 있다. 이러한 미세먼지 현상이 반영되어 미디어에서 미세먼지라는 키워드를 다루는 횟수가 증가했다고도 볼 수 있다.



\*\* No record on yearly highest concentration before 2015

Fig. 6. Yearly number and concentration of article related fine dust

## 4. 결론

본 연구는 사회문제 이슈화된 미세먼지를 키워드로 뉴스 데이터베이스를 수집하여 텍스트 마이닝 기법인 LDA 분석을 이용하여 정량화하였다. 분석 결과, 날씨관련 정보가 토픽의 주요 단어를 구성하고 있었고 10도씨 이하 낮은 기온에서 미세먼지가 더욱 이슈화 됨을 알 수 있다. 미디어에 노출되는 빈도와 미세먼지 농도 최고치는 정(+)의 상관관계가 있고, 토픽화되는 빈도가 높은 정책에 주목할 수 있다. 지난 10년간 미세먼지 저감 대책 및 정부의 종합대책과 관련한 토픽, 미세먼지와 관련한 공기청정기와 같은 상품이나 소비트렌드와 관련된 토픽, 미세먼지로부터 취약계층을 보호하는 정책에 대한 토픽, 수소차 및 R&D를 통한 미세먼지 저감에 대한 토픽들이 주요한 토픽으로 나타났다. 다만 이 이면에 최근 3년간 미세먼지 관련 기사가 급증하였고 대중적인 정책이슈는 소개가 되었지만 장기적인 관점의 저감 대책과 R&D 대책은 기사로 접하기 어려운 실정이라고 판단된다.

## References

1. Green Technology Center Korea, 2019. A Study on the Development of Domestic Climatic Adaptation Industry in the 4th Industry Revolution Era (in Korean). Seoul
2. Jeong JW, Lee JM, Choi SY. 2018. Analysis of news regarding the disabled labor using text mining techniques(in Korean). Reinterpretation of Disability, pp. 48-100
3. Kang HJ, Kim C, Kang K. 2019. Analysis of the

- Trends in Biochemical Research Using Latent Dirichlet Allocation (LDA). Processes Vol. 7, No. 6, [www.doi:10.3390/pr7060379](http://www.doi:10.3390/pr7060379)
4. Keller TR, Hase V, Thaker J, Mahl D, Schafer M. 2020. News Media Coverage of Climate Change in India 1997–2016: Using Automated Content Analysis to Assess Themes and Topics. Environmental Communication Vol. 14, No. 2, [www.doi:10.1080/17524032.2020.1716033](http://www.doi:10.1080/17524032.2020.1716033)
  5. Kim JH, Cho JH. 2019. Investigation of Effects of Individuals Social Viewing of Fine Dust Information Obtained through Social Media on Behavioral Intentions of Disease Prevention : Application of Health Beliefs Model(in Korean). Korean Journal of Broadcasting and Telecommunication Studies Vol. 33, No. 4, pp.37-65
  6. Kim MC, Yoon SU, Kim HM. 2020. A Study on the Hydrogen Economic Law for the Realization of Hydrogen Society in Korea(in Korean). Soongsil Law Review 46, pp.1-30
  7. Kim MC. 2019. Legislation of climate change adaptation has become a global trend(in Korean), KACCC ADAPTATION
  8. Kim MC. 2019. Proactive Legislative Evaluation of Hydrogen Economy Legislation in Response to Climate Change, Korea Legislation Research Institute(in Korean), Legislative Evaluation Issue Paper 19-14-①. 2-40
  9. Kim YW, Lee HS, Jang YJ, Lee HJ. 2015. How Does Media Construct Particulate Matter Risks? : A News Frame and Source Analysis on Particulate Matter Risks(in Korean). Korean Journal of Journalism & Communication Studies Vol. 59, No. 2, pp.121-154
  10. Korea Institute for Health & Social Affairs. 2018. Social big data trend analysis based on health and welfare issues in 2018(in Korean). Sejong.
  11. Lee EB, John JN, Baek JS. 2017. A Study of Multicultural Space in Seoul : Analysing the Coverage of Foreign Communities with News Big Data Analytics BigKinds for 27 Years(in Korean). Journal of Media Economics & Culture Vol. 15, No. 2, pp.7-43
  12. Moon MR, Kim MC, Kim JW. 2019 A Study on the Fine Dust-related Bills in the National Assembly – Based on the Revision of the Special Act on Fine Dust(in Korean). Hannam Journal of Law&Technology Vol. 25, No. 4, pp.87-115
  13. Moon SH, Chung SH, Chi SH. 2018. Topic Modeling of News Article about International Construction Market Using Latent Dirichlet Allocation (in Korean). Journal of the Korean Society of Civil Engineers Vol. 38, No. 4, pp.595-599
  14. Park HJ. et al.. 2015. Prediction of correct answer rate for English scholastic ability test using text mining. Industrial Engineering & Management Systems in Proceedings, pp. 2,277-2,288
  15. Seo Dae-ho. 2019. Catch! Text Mining with Python(in Korean), BJpublic Soo-Sang Lee. 2018. Network Analysis Methods Applications and Limitations(in Korean), Chung-Ram
  16. Tran BX, Nghiem S, Sahin O, Vu T, Ha GH, Vu GT, Pham HQ, Do HT, Latkin CA, Tam W, Ho C, Ho R. 2019. Modeling Research Topics for Artificial Intelligence Applications in Medicine: Latent Dirichlet Allocation Application Study. J Med Internet Res [www.doi:10.2196/15511](http://www.doi:10.2196/15511)
  17. [https://biz.chosun.com/site/data/html\\_dir/2018/04/09/2018040900053.html](https://biz.chosun.com/site/data/html_dir/2018/04/09/2018040900053.html)
  18. <http://biz.newdaily.co.kr/site/data/html/2016/06/03/2016060310058.html>
  19. <https://www.yna.co.kr/view/AKR20170525088700061>
  20. <https://www.yna.co.kr/view/AKR2016081118900003>
  21. <https://www.sedaily.com/NewsView/1L1EHL6AHJ>