

# 토픽모델링을 활용한 인공지능 연구동향 분석

최 대 수\*

## 요 약

본 연구의 목적은 인공지능의 연구동향을 분석하는 것이다. 입체적인 분석을 위하여 인공지능에 대한 사회과학에서의 연구방향과 공학에서의 연구방향의 차이를 객관적으로 비교하여 제시하고자 시도하였다. 연구방법은 빅데이터 분석 방법론 중에서 토픽모델링을 활용하였으며, 분석데이터는 학술연구정보시스템에서 인공지능(AI)라는 키워드로 검색된 1000개의 영문 논문을 활용하였다. 분석결과 사회과학분야에서는 인공지능에 대하여 '인간', '영향', '미래'라는 키워드를 중심으로 형성된 그룹을 확인할 수 있었고, 공학분야에서는 '인공지능 기반의 기술개발', '시스템', '위험-보안' 등의 그룹이 형성되었다.

## Analysis of artificial intelligence research trends using topic modeling

Daesoo Choi\*

## ABSTRACT

The purpose of this study is to analyze research trends in artificial intelligence. For a three-dimensional analysis, an attempt was made to objectively compare and present the difference between the research direction of artificial intelligence in social science and engineering. For the research method, topic modeling was used among the big data analysis methodologies, and 1000 English papers searched with the keyword artificial intelligence (AI) in the academic research information system were used for the analysis data. As a result of the analysis, in the field of social science, it was possible to identify groups formed around the keywords of 'human', 'impact', and 'future' for artificial intelligence, and in the field of engineering, 'artificial intelligence-based technology development', 'system', 'Groups such as 'Risk-Security' were formed.

**Key words** : AI, Artificial Intelligence, Research Trend, Topic Modeling, LDA

접수일(2022년 11월 29일), 게재확정일(2022년 12월 31일)

\* 중부대학교/소프트웨어공학부(주저자)

## 1. 서 론

기술의 발전이 인간의 삶을 개선할 것이라는 생각은 오랫동안 지속되어 왔다. 그러나 기술의 발전으로 야기되는 부작용에 대한 우려도 꾸준히 지속되고 있다. 기술이 인간의 편의와 행복을 위해서 만들어 진 것임에는 반론의 여지가 없다. 하지만 한편으로 편의와 행복이라는 단어가 매우 주관적인 가치를 가지기 때문에, 기술이 사람을 불편하고 불행하게 만든다고 느끼는 사람도 적지 않다. 이렇게 상반되는 평가를 받는 기술 중에 대표적인 것이 인공지능이다.

인공지능은 제4차 산업혁명을 주도하는 핵심기술로 평가받고 있다. 사물인터넷을 통해 오프라인 세상의 사람과 사물들이 연결되고, 빅데이터 기술을 통해 이들이 주고받는 정보의 의미를 확인한다. 인공지능은 이러한 환경에서 학습하면서 인간의 인지능력, 학습능력, 이해능력, 추론능력 등을 갖추게 된다.

인공지능을 바라보는 긍정적 시선은 다음과 같다. 인공지능을 활용하면 편견이 없는 합리적 결정이 가능해질 것으로 보인다. 이를 통해 비이성적인 과열이 사라질 것이고, 새로운 일자리와 혁신이 증가할 수 있을 것이다. 반면, 기존의 일자리가 감소하고, 의사결정의 책임이 불분명해질 수 있으며, 인간에 대한 위협이 생길 수 있을 것이라는 부정적 생각들도 있다.

2016년 다보스 포럼에서 시작된 제4차 산업혁명에 대한 논의는 다양한 분야에서 지속되고 있다. 이러한 논의를 크게 구분한다면 기술 자체에 대한 논의와 사회-경제적 변화에 관한 논의로 구분할 수 있다. 기술 자체에 대한 논의는 공학적 차원의 논의라고 할 수 있다. 제4차 산업혁명에서 주목받은 사물인터넷, 빅데이터, 인공지능 등과 같은 혁신기술의 연구개발에 대한 부분은 관련분야 연구자들에게 매우 매력적인 주제일 것이다. 또한 이러한 혁신기술을 통하여 구현되는 서비스와 이를 통해 유발되는 사회-경제적 변화는 사회과학분야의 학자들에게 많은 영감을 줄 것으로 판단된다.

본 연구의 첫번째 목적은 인공지능의 연구동향을 분석하여 관련분야 기술혁신의 방향을 확인하는 것이다. 인공지능 연구동향은 연구자들의 관심이 반영된 것으로 볼 수 있으며, 이를 통해 인공지능의 발전 방향과 관련연구의 세분화 경향을 확인할 수 있다. 본 연구의 두번째 목적은 인공지능에 대한 기대와 우려를 확인하는 것이다. 사회과학 분야에서는 인공지능이 사회에 미치는 영향과 활용방안 등에 대한 다양한 논의가 진행되고 있을 것으로 예상되며, 이러한 논의들은 인공지능에 대한 기대와 우려가 반영되었을 것으로 판단된다. 세번째 연구목적은 사회과학 분야에서의 주제가 공학분야의 주제가 어떻게 연결되는지를 확인하는 것이다. 사회적 수요와 공학분야의 흐름이 대부분 일치하겠지만 미스매치되는 부분도 있을 것으로 예상된다. 이러한 분석 결과는 향후 인공지능의 연구에 중요한 시사점의 근거가 될 것이다.

## 2. 이론적 배경

### 2.1 인공지능 연구동향

정명석·박성현·채병훈·이주연(2017)은 인공지능 분야 중장기 기술개발 로드맵 작성을 위하여 논문데이터를 분석하였다. 빈도분석과 키워드 네트워크 분석을 통해 연구를 진행하였으며, 시간의 흐름에 따라 이론적 연구 중심에서 실용적 연구가 많아지는 것을 확인하였다. 특히 한국의 연구 분야가 국소적이며 기술적인 부분에 집중되어 있는 것을 확인하였으며, 총체적이며 포괄적인 연구가 필요하다고 주장하였다[1].

황서이·김문기(2019)는 한국의 인공지능분야 연구동향을 분석하면서 토픽모델링과 의미연결망 분석을 활용하였다. 학술논문을 대상으로 연구 토픽을 분석한 결과, 기술과 산업분야 이외에 철학, 인문, 지적재산 분야까지 중요한 키워드를 도출하였다. 인공지능과 인간, 인공지능과 기술 영역이 주요 문제로 확인되었으며, 자율주행차가 핵심 키워드로 도출되었다. 또한 인공지능과 저작권 문제, 인공지능의 인문학적 역할에 대한 고민이 필요하

다고 주장하였다[2].

정우진·오찬희·주영준(2021)는 네트워크분석과 동적 토픽모델링을 활용하여 인공지능 분야 연구 동향을 분석하였다. 그 결과 4차 산업혁명, 빅데이터, 사물인터넷, 딥러닝, 알고리즘, 로봇, 기계학습 등의 키워드들이 한국의 인공지능 연구에서 중요한 것으로 확인되었다. 또한 중요 토픽으로 과실 책임, 국방분야 활용, 보조서비스, 기술적 실업, 게임분야 활용, 저작권, 창작품, 인문학, 기계학습, 교육, 핀테크, 정부정책, 콘텐츠 등이 도출되었다. 인공지능 분야 관련 연구가 인공지능 기술 및 활용보다는 인공지능 시대의 사회문화적 변화 관련 주제들을 논의하고 있다는 점도 확인하였다[3].

정명석·정소희·이주연(2018), 노승민(2017)은 특허데이터를 활용하여 인공지능분야 기술동향을 분석하였다[4][5].

인공지능에 대한 연구동향을 살펴보는 연구 중에서 사회과학 분야와 공학 분야를 비교하여 분석한 연구는 찾기 어렵다. 사회과학 분야에서 기술에 관심을 가지는 부분은 그 기술의 활용과 영향에 대한 관심을 대표한다고 볼 수 있다. 다시 말하면, 그 기술에 대한 사회적 관심을 반영하는 것이다. 공학 분야에서의 연구를 살펴보면, 그 기술 개발의 방향성을 알 수 있을 것이다. 따라서 사회과학 분야와 공학 분야를 비교하는 연구는 그 기술의 사회적 관심을 공학에서 반영하고 있는지 확인하는데 중요한 의미를 가질 수 있다.

### 3. 연구방법

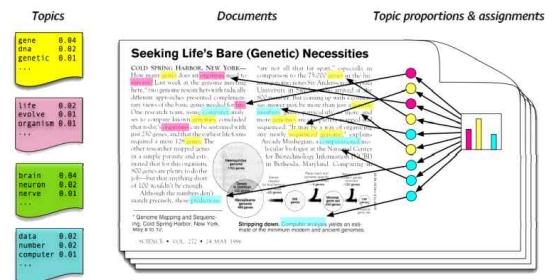
#### 3.1 데이터 분석 방법

본 연구에서 활용한 연구동향 분석 방법은 연구논문 제목에서 추출한 텍스트 마이닝이다. 텍스트 마이닝은 비정형 데이터인 텍스트 형태의 데이터를 통해 의미있는 정보를 추출하는 분석 방법이다. 빈도 분석의 결과를 검토할 수 있는 워드클라우드와 주제를 범주화 하는데 유용한 토픽모델링이 본 연구에서 활용한 분석 기법이며, 이 중에서

토픽모델링은 텍스트 마이닝의 대표적 분석 기법이다.

우선 워드클라우드의 텍스트에서 단어의 빈도를 시각화하기 위해 활용하는 방법이다. 분석의 절차는 다음과 같다. 우선 다양한 방법을 통하여 수집된 데이터를 정제하는 작업, 즉 전처리 과정을 진행한다. 의미없는 단어와 공백, 불용어 등을 제거하는 작업을 의미한다. 전처리가 끝난 데이터는 상태를 확인하기 위해 히스토그램을 그려보고, 최종적으로 원하는 워드클라우드 시각화를 진행한다. 빈도 분석은 텍스트 분석에서 가장 기본적으로 실시하는 분석 기법이며, 본 연구에서도 토픽 모델링보다 앞서 진행하였다.

연구동향을 확인하기 위하여 채택한 토픽모델링은 다양한 문서에서 정보를 추출하고, 공통된 주제를 잠재적 확률 모델로 추출하는 분석방법을 의미한다. 특히 LDA 기법은 핵심 키워드가 문서에 출현할 확률을 토픽별로 클러스터링 하여 분류하는 분석기법이다[6].



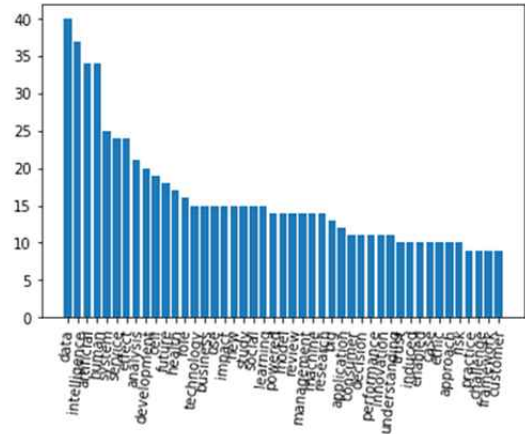
(그림 1) LDA 토픽모델링 개념(Blei, 2012)

#### 3.2 분석 대상 데이터

본 연구는 두 단계로 구분하여 진행하였다. 첫 번째 단계는 인공지능에 대한 연구가 어떤 주제로 연구되고 있는지 전반적으로 확인하기 위해 텍스트 빈도 분석을 진행하였다. 두 번째 단계에서는 연구 동향을 키워드 중심으로 분류하기 위하여 토픽 모델링을 활용하였다.

각 단계의 분석을 위하여 활용한 데이터는 한

국학술정보원(KERIS)에서 운영하는 학술연구정보시스템(www.riss.kr)의 해외전자정보서비스에서 확보하였다. 영어로 작성된 연구논문을 대상으로 진행하였으며, 다보스 포럼에서 4차산업혁명이 주요이슈로 언급되었던 2016년부터 2022년까지 출간된 논문의 제목을 텍스트 데이터로 정리하여 활용하였다. SCOPUS, SCIE, SCI, SSCI, AHCI의 목록에 해당되는 저널에 출간된 논문을 활용하였으며, ‘인공지능(AI)’을 검색어로 활용하여 정확도 순으로 추출한 논문을 대상으로 진행하였다. 추출한 논문의 수는 공학분야 500편과 사회과학분야 500편으로 총 1000편이다.

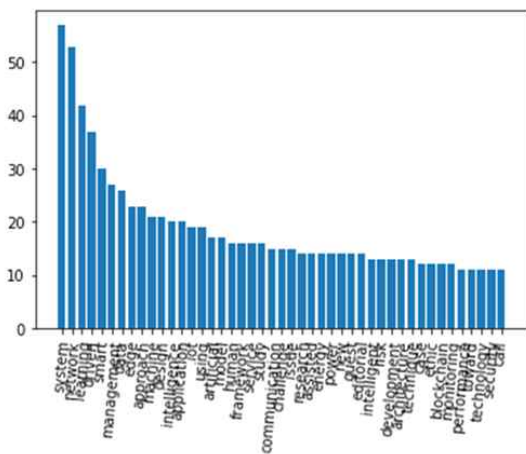


(그림 3) 사회과학분야 키워드 빈도

## 4. 연구결과

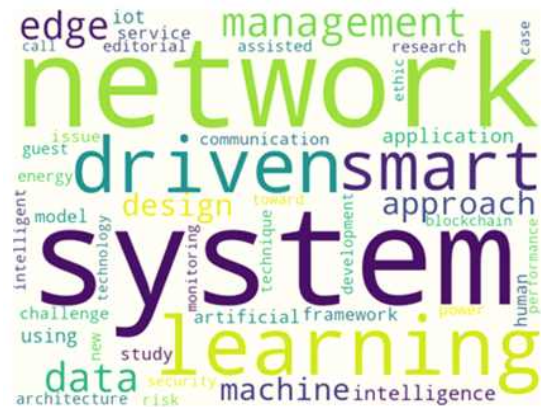
### 4.1 키워드 빈도 분석 결과

빈도 분석을 통해 확인한 결과를 우선 히스토그램으로 표현하였다. 검색어에 해당되는 AI, Artificial, Intelligence 등의 단어를 제외하면, 공학분야의 논문들에서는 system, network, learning, driven, smart 등 단어들이 빈도가 높은 것으로 나타났다으며, 사회과학 논문들에서는 data, human, system, service, effect 등이 많은 것으로 나타났다. 관련 히스토그램은 그림2, 그림3에서 확인할 수 있다.



(그림 2) 공학분야 키워드 빈도

그림4, 그림5를 통해 워드클라우드를 활용하여 시각화된 모습을 확인할 수 있다. 공학분야에서 빈도가 높은 단어는 기술의 개념과 개발에 관련된 용어로 볼 수 있고, 사회과학 분야에서 높은 빈도의 단어들은 기술의 활용과 영향에 대한 내용을 포함하고 있는 것으로 보인다. 어느 정도 경향성을 확인할 수 있고, 일부분은 예상했던 결과를 보이고 있지만, 구체적으로 공학분야와 사회과학 분야 연구들의 관계를 도출하고, 비교 분석하기에는 한계가 있는 것을 확인할 수 있다.



(그림 4) 공학분야 워드클라우드



(그림 5) 사회과학분야 워드클라우드

### 4.2 토픽모델링 결과

인공지능의 주요 연구 분야를 그룹핑하기 위하여 공학분야 500개 논문과 사회과학분야 500개 연구를 활용하여 LDA 분석을 진행하였다. 겹치는 부분을 최소화하기 위하여 5개의 토픽을 설정하였으며 하나의 토픽에 15개의 단어를 추출하는 방식으로 진행하였다. 분야별 분석결과는 표1과 표2에서 확인할 수 있다.

<표 1> 공학분야 LDA 분석결과

번호	주요 단어	토픽 레이블
0	0.060*ai" + 0.022*based" + 0.009*drive n" + 0.008*network" + 0.007*learning" + 0.006*approach" + 0.005*system" + 0.004*g" + 0.004*issue" + 0.004*data" + 0.004*using" + 0.004*information" + 0.003*machine" + 0.003*model" + 0.003*w ireless"	driven
1	0.118*ai" + 0.017*system" + 0.014*ena bled" + 0.011*learning" + 0.010*based" + 0.008*network" + 0.006*design" + 0.006*machine" + 0.006*edge" + 0.005*ma nagement" + 0.005*application" + 0.005* human" + 0.005*intelligence" + 0.004*s mart" + 0.004*approach"	system
2	0.066*ai" + 0.029*based" + 0.014*netw ork" + 0.011*learning" + 0.008*risk" + 0.006*service" + 0.005*blockchain" + 0.005*framework" + 0.005*management" + 0.005*system" + 0.005*intelligence" + 0.005*research" + 0.005*g" + 0.004*en abled" + 0.004*prediction"	risk
3	0.097*ai" + 0.017*based" + 0.013*data"	data

	+ 0.012*driven" + 0.009*enabled" + 0.009*network" + 0.007*g" + 0.006*manag ement" + 0.005*ethic" + 0.005*model" + 0.005*issue" + 0.005*special" + 0.005*l earning" + 0.005*using" + 0.004*system"	
4	0.066*ai" + 0.026*based" + 0.014*syste m" + 0.013*smart" + 0.009*network" + 0.007*study" + 0.006*monitoring" + 0.006*g" + 0.006*edge" + 0.006*manufac turing" + 0.006*approach" + 0.005*focu sed" + 0.005*design" + 0.005*detection" + 0.005*section"	smart

<표 2> 사회과학분야 LDA 분석결과

번호	주요 단어	토픽 레이블
0	0.026*ai" + 0.008*based" + 0.007*huma n" + 0.006*service" + 0.005*using" + 0.004*data" + 0.004*analysis" + 0.004*sy stem" + 0.003*management" + 0.003*ar tificial" + 0.003*intelligence" + 0.003*po wered" + 0.003*research" + 0.002*use" + 0.002*innovation"	human
1	0.068*ai" + 0.008*data" + 0.008*intellig ence" + 0.007*based" + 0.006*artificial" + 0.005*new" + 0.004*human" + 0.004* analysis" + 0.004*impact" + 0.004*devel opment" + 0.004*big" + 0.004*using" + 0.004*model" + 0.003*system" + 0.003* application"	data
2	0.056*ai" + 0.010*based" + 0.007*data" + 0.005*effect" + 0.004*technology" + 0.004*using" + 0.004*service" + 0.004*b" + 0.004*consumer" + 0.003*research" + 0.003*enabled" + 0.003*intelligence" + 0.003*social" + 0.003*artificial" + 0.003*u"	effect
3	0.041*ai" + 0.011*based" + 0.006*huma n" + 0.005*cell" + 0.005*data" + 0.004* service" + 0.004*intelligence" + 0.004*a nalysis" + 0.004*effect" + 0.003*role" + 0.003*social" + 0.003*using" + 0.003*a rtificial" + 0.003*study" + 0.003*system"	cell
4	0.085*ai" + 0.008*artificial" + 0.007*int elligence" + 0.007*future" + 0.006*base d" + 0.006*system" + 0.006*human" + 0.005*development" + 0.005*learning" + 0.005*machine" + 0.004*data" + 0.004* health" + 0.004*business" + 0.004*revie w" + 0.004*using"	future

토픽모델링 결과에서 토픽레이블을 결정하는 부분은 종합적 해석이 포함되어야 하므로 연구자의 해석이 중요하다. 각 토픽마다 같은 단어가 반

복하여 나타나기 때문에 다음과 같은 방법을 적용하였다. 주요 단어 중에서 비중이 높은 단어를 후보로 선정하고 나머지 단어들과의 관계를 검토하여 최종 선정하였다.

공학분야에서 도출된 토픽레이블을 살펴보면 다음과 같다. 'driven'이라는 키워드는 기반 기술에 대한 연구로 해석하였다. 인공지능 기반의 다양한 기술에 대한 연구들이 여기에 포함된다고 볼 수 있다. 'system'에 대한 연구도 하나의 그룹을 형성하였고, 'risk'라는 레이블에는 보안과 관련된 연구들이 포함되었을 것으로 예상된다. 그 밖에도 'data', 'smart'라는 레이블로 그룹이 형성된 것을 확인할 수 있다.

사회과학분야의 토픽레이블은 확실히 차이를 보인다. 그 중에서도 'human'으로 명명된 그룹은 인공지능과 인간의 관계, 인간을 위한 서비스 등이 포함되었을 것으로 예상할 수 있으며, 'effect', 'future' 그룹에서 인공지능이 미치는 영향에 대한 관심을 확인할 수 있다.

## 5. 결 론

본 연구의 목적은 인공지능의 연구동향을 분석하여 그 흐름을 확인하는 것이다. 보다 구체적으로 표현하면 이러한 검증은 통하여 인공지능에 대한 기대와 우려, 즉 사회과학에서의 관심을 확인하고자 하였으며, 공학분야에서의 관심과 어떻게 다른지 확인하는데 초점을 두었다.

분석결과, 사회과학에서의 인공지능에 대한 연구는 '인간'에 대한 관심을 중심으로 큰 그룹을 이루고 있는 것을 확인하였다. 구체적인 연구의 주제를 세밀하게 언급하는 부분에는 한계가 있겠지만, 인공지능이 인간에 인지-사고 체계를 닮아있는 만큼 이에 대한 관심과 인간의 삶에 미치는 영향에 대한 부분이 다양하게 연구되고 있다는 사실을 확인할 수 있다.

공학분야에서는 인공지능 기반의 기술개발, 시스템, 데이터 등에 대한 연구 그룹을 확인할 수 있었고, 보안-위협에 대한 내용이 하나의 그룹을

형성하고 있는 부분이 인상적이었다.

인공지능에 대하여 사회과학분야에서 다루는 주제와 공학분야에서 다루는 주제를 매칭하여 분석하기에는 데이터의 특성과 방법론에서 한계를 확인할 수 있었다. 또한 시계열 데이터가 아닌 일정한 시각에 수집한 데이터를 통해 분석을 진행하였기 때문에 연구의 큰 흐름을 확인하기는 어려움이 있었다. 그럼에도 불구하고, 사회과학분야와 공학분야의 인공지능에 대한 관심의 차이를 객관적 검증방법으로 제시하였다는데 본 연구의 의의가 있다고 할 것이다.

## 참고문헌

- [1] 정명석, 박성현, 채병훈, 이주연, "논문데이터 분석을 통한 인공지능 분야 주요 연구 동향 분석", 디지털융복합연구, 15(5), pp.225-233, 2017.
- [2] 황서이, 김문기. "국내 인공지능분야 연구동향 분석:토픽모델링과 의미연결망분석을 중심으로", 디지털콘텐츠학회논문지, 20(9), pp.1847-1855, 2019.
- [3] 정우진, 오찬희, 주영준, "네트워크 분석과 동적 토픽모델링을 활용한 국내 인공지능 분야 연구동향 분석", 한국문헌정보학회지, 55(4), pp. 141-157, 2021.
- [4] 정명서, 정소희, 이주연. "국내외 특허데이터 기반의 인공지능분야 기술동향 분석", 디지털융복합연구, 16(6), pp.187-195, 2018.
- [5] 노승민, "특허분석을 통한 인공지능 기술분야의 연구동향", 디지털콘텐츠학회논문지, 18(2), pp.423-428, 2017.
- [6] Blei, D. M. "Probabilistic topic models", Communication of the ACM, 55(4), pp.77-84, 2012.

— [ 저자 소개 ] —



최 대 수 (Daesoo Choi)  
2021년 2월 고려대학교 문학사  
2011년 8월 서강대학교 경제학석사  
2022년 8월 고려사이버대 정보학석사  
2018년 8월 동국대학교 기술창업학박사  
2014년 4월~현재 중부대학교  
소프트웨어공학부 교수  
email : daesoo100@gmail.com