

- ****

 OPEN ACCESS

A Study on the Research Trends in the Fourth Industrial Revolution in Korea Using Topic Modeling

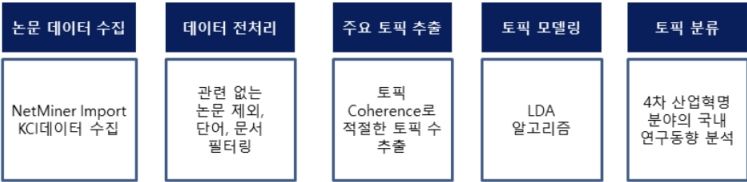
김 지 영 (Gi Young Kim)*
노 동 조 (Dong-Jo Noh)**

[illegible][illegible]

키워드: 로직모델링, 4차 산업혁명, 연구 동향, 키워드 분석, LDA 알고리즘
Topic Modeling, Fourth Industrial Revolution, Research Trends, Keyword Analysis, LDA Algorithm

- © 2023 Korea Basic Science Editors. All rights reserved. This article is distributed under the terms of the Creative Commons Attribution Non-Commercial 4.0 International License (<http://creativecommons.org/licenses/by-nc/4.0/>), which permits use, distribution and reproduction in any medium, provided the original work is properly cited, the use is for non-commercial purposes, and no modification or adaptation is made.

- 토픽모델링의 결과를 평가하는 방법 - 일관성 혼합도
- 선행연구분석 우리는 특히 텍스트마이닝 / 항만 텍스트 마이닝 이 두 개 찾은거 넣으면 될듯



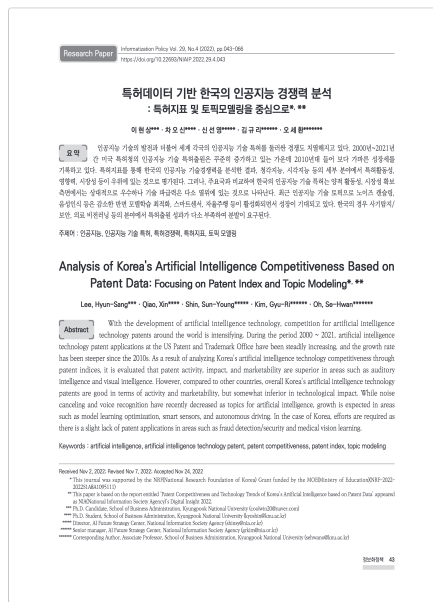
<그림 1> 연구 과정

<표 1> 데이터 전처리 과정

단 계	방 법	처 리 내 용
1단계	자료수집	• NetMiner Import로 '4차 산업혁명' 자료 수집
2단계	표준화	• 유사어 통일(제4차 산업혁명, 4차 산업혁명) • 제외어 정리(한 글자, 조사, 국가명 등) • 불용어 처리(반복적인 단어 등)
3단계	필터링	• 토픽 추출 검색 키워드 제외
4단계	관련 논문 추출	• 키워드: '4차 산업혁명' • 분석범위: 2016년 1월 ~ 2023년 8월까지 관련 논문 추출

- - 이거보니까 생각난건데, 여기의 '4차 산업혁명' 키워드 사용한 것처럼 우리도 특정 키워드가 하나 있어야 할 것 같다. 3만 5천개가 다 제대로된 애들인지 모르니까 키워드 넣고 뽑아야 할지도,,?
 - 아니면 특히 추출한 과정을 정확히 알 수 있나? 추출하는 과정이 신뢰성있을지도 모르니까
- '최적의 토픽 수를 결정하기 위하여 NetMiner의 토픽모델링 평가(Evaluation of Topic Models) 기능을 활용하였다. 토픽모델링 결과를 평가하는 토픽 일관성 지표의 c_v 점수를 사용하여 최적의 토픽 개수를 선정하였다.'
 - 최적 토픽 수 어떻게 정할지 중요할 듯.
- 애네꺼 LDA 이후 토픽 합친건 좀 별로인듯 그냥 더한게 다임.

-
- [한국의 인공지능 경쟁력 분석]



- 애네도 특허 데이터를 사용했음
- 근데 특허로 미국과 한국에서 각각 어떤 토픽의 특허 출원이 많이 일어나고 경쟁력이 어떤지 분석한 논문임.
- 인공지능 경쟁이 치열한데, 국가별 기술경쟁력 평가가 필요하지 않겠냐 이런 느낌

• 애네의 연구질문

- 최근 인공지능 기술 관련 특허출원 동향 및 추이는 어떠한가?
- 한국의 인공지능 세부 기술별 경쟁력은 어떠한가?
- 한국을 포함한 주요국의 인공지능 기술경쟁력은 어떠한가?
- 인공지능 기술의 주요 토픽 변화 양상은 무엇인가?

이렇게 연구질문 정해놓는거 좋은거 같음.

- 토픽 모델링한 부분만 보면, 나라별 토픽의 등장확률을 집계하여 한국의 부족한 부분이 어느분야인지 분석하네.

LDA 및 BERTopic 기반 해외건설시장 뉴스 기사 토픽모델링 성능평가

백준우*, 정세환***, 지석훈***

Baik, Jooswoo*, Chung, Sehwon***, Ji, Seokho***

Evaluation of Topic Modeling Performance for Overseas Construction Market Analysis Using LDA and BERTopic on News Articles

ABSTRACT

Understanding the local conditions is a crucial factor in enhancing the success potential of overseas construction projects. This can be achieved through the analysis of news articles of the target market using topic modeling techniques. In this study, the authors aimed to analyze news articles using two topic modeling methods, namely Latent Dirichlet Allocation (LDA) and BERTopic, in order to determine the optimal approach for market condition analysis. To evaluate the alignment between the generated topics and the actual themes of the news documents, the research collected 6,273 RBC news articles, created ground truth data for individual news article topics, and finally compared this ground truth with the results of the topic modeling. The F1 score for LDA was 0.011, while BERTopic achieved a score of 0.244. These results indicate that BERTopic more accurately reflected the actual topics of news articles, making it more effective for understanding the overseas construction market.

Keywords : Overseas construction, News article, Topic modeling, Ground truth, Performance evaluation

도 록

해외건설사업 시, 현지 상황을 정확히 파악하는 것은 프로젝트 성공을 위해 매우 중요한 요소이다. 이는 프로젝트별 성공을 높이고, 리스크를 줄여 줄 수 있다. 본 연구는 Latent Dirichlet Allocation(LDA)과 BERTopic 두 토픽모델링 기법을 활용하여 뉴스 기사를 분석하고, 결과의 정확성을 검증하기 위하여 RBC 뉴스 기사를 수집하고, 실제 주제에 대한 ground truth를 생성하고, 이를 토픽모델링 도구에 적용하여 비교하였다. 그 결과 LDA의 F1 score는 0.011, BERTopic은 0.244로 나타났다. 이를 통해 BERTopic이 실제 뉴스 기사의 주제를 잘 파악하여, 해외건설사업의 주요 이슈를 파악하는 데 더욱 효과적임을 확인할 수 있었다.

주제어 : 해외건설, 뉴스 기사, 토픽모델링, Ground truth, 성능평가

1. 서 론

해외건설사업 시 현지 시장 상황을 신속하게 파악하는 것은 사업의 성공적인 수행을 위해 매우 중요하며(Javimick-Will and

Scott, 2010), 시장 분석의 뉴스 기사에는 다양한 사건들이 포함되어 있으며, 이를 분석하여 유용한 정보를 추출하고 시장 상황을 파악할 수 있다(Goldkorn et al., 2011). 그러나 뉴스 기사의 양이 매우 많아 모든 기사를 일일이 분석할 수 없다는

* 서울대학교 건축환경공학부 석사과정 (Seoul National University - master@ksee.snu.ac.kr)

*** 공저자 (Seoul National University - Seoul National University - hse@ksee.snu.ac.kr)

*** 공저자 (Seoul National University - Seoul National University - hse@ksee.snu.ac.kr)

Received August 8, 2023 / revised September 12, 2023 / accepted September 14, 2023

BERTopic

- HDBSCAN 알고리즘 적용 시, 파라미터를 조정함으로써 토픽 개수를 특정 개수로 설정하거나, 알고리즘에 의해 최적의 토픽 개수를 자동으로 찾아낼 수 있다.
- 뉴스 기사 대상으로 LDA보다 훨씬 낫다는 것을 보여주는 논문.

BERTopic

자연어 처리 기술의 발달로, 22년 BERTopic이라는 방법이 새로 등장

LDA 방법의 우수성과 유연성은 다양한 분야에서 입증되었지만 통계기반으로 토픽을 산출하기 때문에 글의 맥락을 고려하지 못한다는 단점이 있다.

하지만, BERTopic은 임베딩 방식을 통해 조금 더 정확하게 토픽이 산출될 수 있을 것으로 기대됨. 그러나 BERTopic의 경우 LDA에 비해 적용 사례가 많이 않으므로 좀 더 다양한 분야에서 성능을 검증할 필요가 있다.

따라서 본 연구에서는 LDA와 BERTopic 간의 성능을 비교하고 BERTopic을 통해 인간-로봇 상호작용 연구의 토픽을 살펴봄으로써 현재까지의 연구 동향을 파악한다.

〈표 1〉 LDA와 BERTopic의 비교

Topic Model	Advantage	Disadvantage
LDA	<ul style="list-style-type: none">- 비교적 간단하고 계산 비용이 저렴한 알고리즘임- 워드 임베딩 기반 접근법에 비해 적은 수의 주제를 생성함- 해석이 용이함- 다수의 연구에서 검증된 방법임	<ul style="list-style-type: none">- 단어주머니 방식을 사용하여 문맥을 고려하지 않음- 대용량의 데이터셋에서의 성능이 보장되지 않음- 너무 일반적인 주제나 단어를 추출하는 경향이 있음
BERTopic	<ul style="list-style-type: none">- 워드 임베딩 방식을 사용하여 문맥을 고려하여 주제를 추출함- 대용량의 데이터셋에 적합한 알고리즘임- 단어의 희소성을 고려한 방법임	<ul style="list-style-type: none">- 비교적 적은 용량의 데이터셋에서의 검증이 이루어지지 않음- 해석이 용이하지 않음

이후 결과, (둘의 성능 비교, 토픽이 잘 형성되었는지 확인하기 위해 일관성과 혼잡도를 비교하는 방식.) 보편적인 토픽 모델의 평가 방식인가봄.(주제 추출)

우수하다고 나왔음. (이런 논문이 좀 많음.)

BERTopic을 활용한 인간-로봇 상호작용 동향 연구

최근까지 LDA가 가장 일반적으로 활용되는 토픽 모델링 방법론이었으나, 머신러닝이 토픽모델링에 응용되기 시작하면서, 기존 방법의 단점을 보완하여 정확도를 높이려는 시도들이 이루어지고 있다. 대표적으로 2021년, BERT기반 임베딩 처리와 클래스 기반 TF-IDF등을 활용하여 일관된 토픽을 생성하는 BERTopic 기술이 제안되었는데, 기존 통계기반 토픽모델링 방법론들에 대비하여 높은 주제 일관성과 다양성을 보이는 것으로 확인되었다.

BERTopic을 활용한 텍스트마이닝 기반 인공지능 반도체 기술 및 연구동향 분석 -특허와 논문 빅데이터를 중심으로

회의할 내용.

최근에 경영과학회, 산업공학회 토픽모델링 관련 주제로 강장히 핫한게 BERTopic

기존의 LDA(Latent Dirichlet Allocation)와 같은 주제 모델링 기법에 비해 더 정확하고 세분화된 주제 분류를 가능하게 합니다.

1. 항만 특허 주제로 어떤 연구 질문을 선정할지. (= 연구목적을 분명히 하자.)
 1. 항만 특허 기술 관련 특허출원 동향 및 추이는 어떠한가?
 2. 항만 기술의 토픽 변화 양상?
 3. 아니면, 한국의 항만 특허도 수집해서 미국과 한국의 경쟁력을 비교해볼 수 도 있을듯
 1. 위 논문중에 KCI 논문 크롤링하는 프로그램을 사용해서 금방 똑딱 크롤링 하는 것 같음.
2. BERTopic의 등장. LDA보다 우월하고, 최근 토픽모델링 연구 이걸로 많이 하는 듯.
 - 항만 디지털화 특허 출원 동향 및 추이
 - HOT 토픽, COLD 토픽
 - GTM 활용 공백기술 뭔지?
 - 항만 디지털화 기술의 토픽 변화 양상
 -

연구 방법론에 대한 스케치 / 개념도

항만 관련 특허 데이터(약 3만5천여개)

- 목차

국가	출원번호	출원일자	공개번호	공개일자	공고번호	공고일자	등록번호	등록일자	발명의명칭	IPC분류	출원인	발명자/고안자	대리인	요약	청구항	UPC	CPC 분류
----	------	------	------	------	------	------	------	------	-------	-------	-----	---------	-----	----	-----	-----	-----------

1. 특허 데이터의 요약 및 청구항에 우리의 주제가 되는 '디지털화'관련 특허를 찾기 위해서
2. 키워드 선별기준 (- 찾아보고)
3. 3만 5천개의 특허에서 선별한 키워드를 통해 유효 특허를 분류한다.
4. 데이터 전처리 - 소문자화, 특수문자, 불용어, 숫자 등을 제거하는 작업.
5. 이후 LDA, BERTopic 각각 해보고,
6. 이후 둘 중 토픽 모델링 성능 평가
7. BERTopic 우수하다
8. 여기서 나온 토픽 사용해서
9. 시계열 분석 or GTM? or 을 통해 결론