

# CONCOR 중심성을 이용한 건설기술 키워드 가중치 분석

정성윤\*

\*한국건설기술연구원 미래스마트건설연구본부  
e-mail:syjeong@kict.re.kr

## Construction Technology Keyword Weighting Analysis Using the CONCOR Centrality

Seong-Yun Jeong\*

\*Dept. of Future & Smart Construction Research,  
Korea Institute of Civil Engineering and Building Technology

### 요약

본 논문은 건설기술정보시스템에서 서비스되는 건설기술정보에 내재한 유의미한 키워드가 건설기술 네트워크에서 얼마나 영향력을 갖는지를 확인하고자 하였다. 이를 위해 불용어 전처리 작업, LDA 분석 및 토픽 모델링 등의 과정을 통해 건설기술에 영향력이 있는 키워드를 추출하였다. 추출한 키워드 간의 구조적 동위성 기반의 중심성과 상관계수를 추정하여 묶은 4개의 블록을 생성하였다. 생성한 블록을 가지고 가중치가 높은 50개의 키워드를 표시한 네트워크와 키워드의 가중치를 추정하였다.

## 1. 서론

건설기술정보시스템은 건설기술과 관련한 서지정보와 원문 자료를 서비스하는 시스템이다. 서지정보는 해당 건설기술정보에 대한 의미를 함축적으로 내포하고 있다. 이러한 내포된 정보에서 영향력이 높은 키워드를 뽑을 수 있다고 판단하였다. 이를 위해 건설기술정보시스템에서 서비스하는 30,045건의 건설기술정보를 수집하였다[1]. 수집한 정보는 분석에서 필요하지 않은 불용어 처리 등 전처리 작업을 진행하였다. 잠재 디리클레 할당(Latent Dirichlet Allocation) 분석 기법과 토픽 모델링을 통해 유사한 의미를 갖는 키워드를 하나의 그룹으로 군집화한 토픽을 생성하였다. 토픽에 포함된 키워드가 네트워크에서 얼마나 영향력을 갖는지를 추정한다. 이때 많이 사용하는 중심성 알고리즘으로는 단계(Degree), 근접(Closeness), 매개(Between), 아이겐벡터(Eigenvector) 등이 대표적이다. 본 연구는 이들 중심성 이외에 구조적 동위성(structural equivalence)을 측정하여 군집에 내포된 특성 또는 의미를 추정하는 CONCOR 중심성 알고리즘을 건설기술 정보에 내포된 유의미한 정보를 분석하는 데 활용할 수 있는지를 확인하고자 하였다. 이를 위해서 전처리 과정을 통해 얻은 어휘를 기초 데이터로 하여 CONCOR 중심성에 적용하여 가중치가 상위 50위 안에 포함된 키워드를 중심으로 한 네트워크를 표시하였다. 다음으로, 4개의 블록에 포함된 가중치가 높

은 키워드를 추정하였다. 블록에서 상위를 차지한 키워드는 주로 도로와 시설물의 건설공사와 관련한 키워드가 네트워크의 중심에서 영향력을 갖는 것으로 유추할 수 있었다.

## 2. 이론적 고찰

텍스트 마이닝(Text Mining)은 특정 콘텐츠를 구성하는 어휘의 품사와 이웃 어휘 간의 선·후행 관계, 어휘의 발생빈도 등을 통계적으로 분석하여 유의미적인 키워드를 유추하는 이론을 말한다. 잠재 디리클레 할당 알고리즘을 이용하여 유사한 의미를 갖는 키워드를 묶은 토픽 모델링 과정과 토픽 모델링을 통해 키워드 간의 연관관계를 나타내는 네트워크 이론(Network Theory)을 사용한다. 이때 키워드 간의 네트워크는 꼭짓점과 연결선 등으로 이용하여 키워드 간의 연결 관계를 표현한다. 보통 네트워크 중심성 알고리즘으로 중심성 계산식에 따라 단계, 근접, 매개, 아이겐벡터를 가장 많이 사용하는데 네트워크에 있는 키워드의 가중치를 직접 산출한다[2]. 반면에, CONCOR 중심성은 키워드 간 중심성과 공동 출현 행렬을 반복적인 피어슨 상관계수를 수렴하여 유사한 의미를 갖는 키워드를 서로 묶어준 군집을 찾는다[3]. CONCOR 중심성은 기존에 키워드 간의 연결 개수, 중개 역할 등 키워드의 영향력을 기반으로 하지만 CONCOR 중심성은 유사한 의미를 갖는 군집을 이용한다는 점이 차이가 있다.

