

## 문서로부터 계층적 개념 트리 자동 구축

김희수<sup>0</sup> 조용석<sup>1</sup> 최익규<sup>2</sup> 김민구<sup>3</sup>  
 아주대학교 정보통신전문대학원<sup>02</sup>  
 아주대학교 정보통신대학<sup>13</sup>  
 {heemanz<sup>0</sup>, laconic, ikchoi, minkoo}@ajou.ac.kr

### Automatic Construction of Concept Hierarchy from Text

Hee-soo Kim<sup>0</sup>, Yongsuk Cho<sup>1</sup>, Ikkyu Choi<sup>2</sup>, Minkoo Kim<sup>3</sup>  
 Graduate School of Information and Communication, Ajou University<sup>02</sup>  
 College of Information Technology, Ajou University<sup>13</sup>

#### 요 약

계층적 개념 트리는 개념의 전체적인 구조를 제공하여 사용자의 이해를 돕는다. 이러한 계층적 개념 트리는 특정 분야의 전문가나 지식 공학자, 혹은 개인에 의해서 제공되어 왔다. 하지만 계층적 개념 트리를 구축하기 위해서는 많은 시간과 노력이 요구된다. 따라서 계층적 개념 트리를 자동으로 구축하기 위한 시스템 필요하게 되었다. 이 논문에서는 정보 검색 및 온톨로지 연구에 있어서 계층적 개념 트리를 자동으로 구축하기 위한 기존 연구에 대해서 알아보고, 개념적 클러스터링 방법인 FCA(Formal Concept Analysis)를 사용하여 문서로부터 계층적 개념 트리를 구축하는 방법을 제안하고자 한다.

#### 1. 서 론

계층적 개념 트리(*concept hierarchy*)는 문서집합의 구조화 와 요약 을 제공 하고, 필요 한 정보 의 수월 한 접근 을 제공 하 기 위해 사용 되었 으며 [1], 또한 지 식 을 표현 하 는 데 있 어서 간결 하 면서 도 효과 적 으 로 개념 간 의 구조 를 설명 하 기 위해 사용 되었다. 이러한 개념 트리 의 대표 적 인 예 로 는, 사용 자 의 지 식 을 이용 하 여 정보 에 접근 하 는 경 로 를 제공 하 는 'Yahoo!'의 디렉터 리 서 비스, 또는 지 식 을 나타 내 거 나 WordNet과 같이 개념 간 의 is-a 관계 (*hyponym, hypernym*)를 나타 내 기 위해 사용 되는 트리 형태 의 그래프 등 이 있다. 그러나 지 금 까 지 대 부 분 의 계층 적 개념 트리 는 특정 분 야 의 전문 가 나 지 식 공 학 자 들 에 의 해 서 만 들 어 졌 다. 하지만, 계층 적 개념 트리 를 구축 하 는 작업 은 많은 시 간 과 노 려 를 필 요 로 한다. 또한, 오늘날 과 같이 급 변 하 는 환경 에 서 이미 구축 된 계층 적 개념 트리 는 지속 적 인 수정 을 통 해 새 롭 게 생 성 되 는 개념 및 현 상, 사 물 등 을 반영 하 야 할 필 요 가 있다. 따 라 서 본 논문 에 서 는 계층 적 개념 트리 의 구축 을 돕 기 위한 방 법 으 로 문서 로 부터 계층 적 개념 트리 를 구축 하 는 기존 방 법 들 에 대 해 알아 보 고, 자동 화 된 계층 적 개념 트리 의 구축 방 법 을 제 안 한 다.

본 논문 은 다음 과 같 이 구성 된 다. 2장 에 서 는 관련 연구 와 문제 점 에 대 해 서 알아 보 고, 3장 에 서 는 정형 화 된 계층 적 개념 트리 구축 방 법 에 대 해 설명 한 다. 그리고 4장 에 서 는 실험 결 과 를 보 여 주 고, 5장 에 서 는 결 론 과 추 후 과 제 에 대 해 서 언급 한 다.

\* 본 논문 은 과학 기술 부 의 국가 지정 연구 실 사업 (과 제 명: 차 세 대 인터 넷 을 위한 지능 형 온톨로지 자동 생 성 시 스템 개발, 과 제 번 호: M10302000087-03J0000-04400) 지원 으 로 수행 되 었 음

#### 2. 관련 연구

문서 로 부터 계층 적 개념 트리 를 구축 하 는 방 법 으 로 는 개념 의 정 의 에 따 라 크게 두 가 지 로 분류 할 수 있다. 첫 번째 방 법 은 개념 을 단 어 들 의 집 합 으 로 간 주 하 여 계층 적 개념 트리 를 구축 하 는 방 법 으 로 계층 적 군 집 화 (*clustering*)를 이용 하 는 것 이다. 하지만 계층 적 군 집 화 를 사용 한 방 법 은 형성 된 클러스터 에 대한 명 명 문제 등 과 같이 사용 자 에 게 직 관 적 인 정보 를 제공 하 기 어 렵 고 [1], 클러스터링 의 결 과, 클러스터 가 편 중 되는 현 상 을 보 이 는 문제 를 가 지 고 있다. 다른 방 법 은 문서 에 존재 하 는 어휘 를 개념 으 로 간 주 하 고 통계 적 방 법 을 이용 하 여 계층 적 개념 트리 를 구축 하 는 것 이다 [2][3][4]. 이러한 방 법 들 은 일반 적 으 로 검색 된 문서 집합 들 로 부터 정보 의 수월 한 접근 을 제공 하 기 위해 사용 되 지 만, 개념 간 의 관계 가 엄밀 한 의미 의 is-a 관계 가 아 닌 연구 자 가 제 안 한 방 법 에 의 해 부여 된 개념 간 의 관계 에 의 해 구축 된 것 이다. 이에 반 하 여 개념 적 클러스터링 기법 인 FCA (*Formal Concept Analysis*)는 개념 에 대한 속 성 을 정의 하 고 속 성 의 포함 관계 를 사용 하 여 계층 적 개념 트리 를 생 성 하 는 수학적 방 법 을 제시 한 다 [5]. Gu Tao 는 온톨로지 의 계층 적 개념 트리 의 구축 을 돕 기 위해 FCA 를 사용 하 는 방 법 을 제 안 하였 다 [6]. 또한 FCA 를 사용 하 여 문서 로 부터 온톨로지 의 개념 트리 를 자동 으 로 구축 하 기 위해 문장 에 서 의 동사 와 목적 어 를 추출 하 여, 목적 어 를 개념 으 로, 동사 를 개념 의 속 성 으 로 간 주 하 여 FCA 에 적용 하 는 방 법 이 연구 되 었 다 [7]. 그러나, 이 방 법 은 개념 의 특 징 을 설명 하 기 에 부적 합 한 동사, 특히 일반 적 으 로 자주 사용 되는 동사 가 개념 의 속 성 으 로 추가 됨 으 로 써, 합리 적 인 is-a 관계 를 형성 하 는 것 을 방해 한 다.

이 논문에서는 FCA를 사용하여 계층적 개념 트리를 구성하는 방법의 타당성과 발생하는 문제에 대한 해결 방법에 대한 실험 결과를 보여준다.

3. 계층적 개념 트리의 생성

어휘의 통계적 방법을 통한 계층적 개념 트리의 구축은 단지 문서 집합에서 그 어휘의 빈도수에 기반하기 때문에 이러한 방법에 의해 추출된 개념간의 관계가 is-a 관계를 의미하지 않을 수 있다. 인공지능의 과거 연구들은 지식을 표현하는 방법으로 개념을 다음과 같이 정의했다. 개념은 속성들의 집합으로 정의되며, 각 개념과 관련된 속성 (attribute)은 저마다의 제약(role restrictions)을 갖는다[8]. 이러한 정의는 계층적 개념 트리를 구축하기 위한 중요한 정보를 제공한다. 계층적 개념 트리를 구축하기 위해 개념이 가지고 있는 속성이 중요하게 작용한다는 것이다. 예를 들면, '동물'과 '사람', '개'의 개념을 분류하는 기준은 각각의 개념이 가지고 있는 속성에 기반한다. '동물'의 속성은 모두 '사람'과 '개'의 개념에 포함되며, '사람'과 '개'는 '동물'이 가지는 속성에 '젓을 먹인다'라는 속성을 추가적으로 갖는다. 또한 '사람'과 '개'를 구별하는 속성은 '말한다'와 '젓는다'로 생각할 수 있다. FCA는 개념이 가지고 있는 속성을 기반으로 계층적 개념 트리를 구축하기 위해 사용된다.

3.1 계층적 개념 트리 구축 시스템의 구조

시스템의 전체 구조는 다음과 같다.

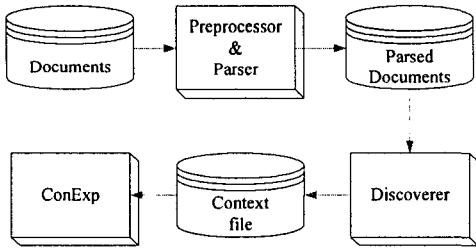


그림 1 시스템의 구조도

먼저, 전처리 과정으로 문서들을 문장 단위로 나누고 Parser를 통해 문법적인 태깅을 한다. 그 다음, Discoverer를 사용하여 문장에서의 목적어-동사를 찾아 formal context를 나타내는 context 파일을 생성한다. Formal context에 대해서는 3.3절에서 설명한다. 마지막으로 ConExp(Concept Explorer)를 사용하여 context 파일로부터 계층적 개념 트리를 구축한다.

시스템의 구성 요소 중 Parser와 ConExp는 각각 Apple Pie Parser (<http://nlp.cs.nyu.edu/app/>)와 ConExp(<http://www.mathematik.tu-darmstadt.de/ags/ag1/fag1/Software/ConExp/>)를 사용했다.

3.2 개념의 추출

문서로부터 개념을 추출하기에 앞서, 개념을 속성들의 집합으로 정의한다. 개념과 속성은 각각 문장에서 목적어와 동사의 형태로 문서에 출현한다는 가정하에 개념을 정의한다[7]. 하지만, 동사는 의미를 전달하기 위해 다양한 목적어를 취할 수 있다. 바꾸어 말하면, 문서에서 다양한 특정 동사는 목적어와 일대다의 관계를 맺으며, 이러한 관계는

FCA에 의해 구축된 계층적 개념 트리의 정확도를 떨어뜨리는 요인으로 작용한다. 이러한 문제를 해결하기 위해 개념에 포함되는 속성의 필터링이 필요하다. 본 연구에서는 개념과 속성간의 포함관계를 판단하기 위한 척도로서 KL (Kullback-Leibler) divergence score를 사용하며, 이에 대한 정의는 다음과 같다.

$$KL\ contribution_c(a) = P_c(a) \log_2 \frac{P_c(a)}{P_c(a)}$$

$$P_c(a) = \frac{freq(c \wedge a)}{\sum_{a' \in C} freq(c \wedge a')}$$

$$P_c(a) = \frac{\sum_{c \in C} freq(c \wedge a)}{\sum_{c \in C} \sum_{a' \in C} freq(c \wedge a')}$$

C는 전체 개념의 집합을 나타낸다. 개념 c에 포함된 속성 a의 중요도는 KL divergence score으로 나타낼 수 있으며, 속성이 개념을 설명하기 위해 자주 사용되면, 그 속성은 개념을 표현하는 중요한 정보로 사용될 수 있음을 의미한다.

3.3 FCA를 사용한 계층적 개념 트리의 구축

FCA는 주로 데이터의 분석을 위해 사용되며, 이러한 데이터는 개념의 형식적인 추상화(formal abstractions)의 단위로 구조화된다. FCA를 이해하기 위한 정의는 다음과 같다.

정의 1. Formal context (G, M, I)는 두 개의 집합 G와 M, 그리고 G와 M사이의 관계 I로 구성된다. G와 M의 원소들은 각각 그 context의 objects와 attributes를 의미한다. 또한 object g가 attribute m과 관계가 있을 때, glm 또는 (g, m) ∈ I로 나타내며, g는 m을 갖는다는 것을 의미한다.

A ⊆ G, B ⊆ M을 만족하는 집합 A와 B에 대하여,  
 A' = {m ∈ M | (g, m) ∈ I, for all g ∈ A},  
 B' = {g ∈ G | (g, m) ∈ I, for all m ∈ B}  
 A'와 B'를 각각 정의하면, 다음과 같이 formal concept과 상위-하위개념을 정의할 수 있다.

정의 2. Formal context (G, M, I)에서 A ⊆ G, B ⊆ M일 때, A' = B, B' = A를 만족하는 (A, B)를 formal concept이라고 한다.

정의 3. (A<sub>1</sub>, B<sub>1</sub>)와 (A<sub>2</sub>, B<sub>2</sub>)가 formal context의 formal concept일 때, concept간의 (A<sub>1</sub>, B<sub>1</sub>) ≤ (A<sub>2</sub>, B<sub>2</sub>)가 성립하기 위한 필요충분 조건은 A<sub>1</sub> ⊆ A<sub>2</sub> 또는 B<sub>1</sub> ⊇ B<sub>2</sub>이다. (≤는 concept의 포함관계를 나타낸다.)

간단한 예로 '생물', '동물', '사람', '개'의 개념이 있고, 이것을 formal context로 나타내면 표 1과 같다고 가정하자.

objects	Attributes			
	살아있다	움직인다	운다	젓다
생물	0			
동물	0	0		
고양이	0	0	0	
개	0	0		0

표 1. Formal Context의 예

위 예에서 네 개의 objects는 정의 2를 만족하는 formal concept이다. 그리고 정의 3에 의해 그림 1과 같은 계층적

개념 트리를 얻을 수 있으며, 속성이 개념을 분류하는 기준으로 사용되는 것을 볼 수 있다.

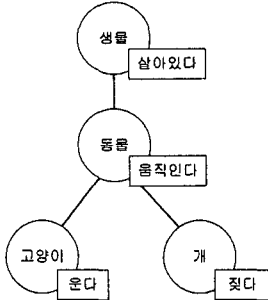


그림 2. FCA로 구성된 계층적 개념 트리의 예

본 논문에서는 FCA를 수행하기 위해 ConExp 프로그램을 사용했다.

4. 실험 결과

이 논문에서 제한된 계층적 개념 트리의 구축을 위해 사용된 문서집합은 TREC(Text REtrieval Conference)에서 제공하는 Text Research Collection Volume 2의 Associated Press (1988)를 사용하였으며, AP와 관련된 개념을 선택하기 위해 'Yahoo!' 디렉터리 서비스의 'Business & Economy(BE)'와 'Society & Culture(SC)', 'Social Science(SS)', 'Government(GOV)'에 포함된 카테고리 제목으로부터 각각 253, 273, 220, 163개의 단어를 선택했다. 또한, KL-divergence 0.25를 임계값으로 개념을 속성의 집합으로 정의했다. 마지막으로 FCA를 사용하여, 네 개의 개념 집합에서 각각의 계층적 개념 트리를 구성하였다. 그림 3은 구축된 개념 트리에서 의미 있는 is-a 관계를 간추린 것이다. 그림 3에서 개념은 '개념:속성'의 형식으로 표현되며, 표 2는 실험 결과에 대한 통계적 결과를 보여준다.

Society & Culture	Social Science
:introduce └recreation:say └exhibit:review :included └hippie:call └skinhead:see,look :affect └gay:play,support └lesbian: :enjoy └romance:attract └wildflower:tell	:say └sociology:teach └economics:impose └anthropology:get └archaeology:take └psychology:avoid,receive :use └directory:make,consider └reference:become philosophy:attack └marxism:avoid,consider,ban,modify :ban └marxism:avoid,consider,attack,modify └propaganda:dismiss
Business & Economy	Government
:wear └badge:flash └jewelry:steal	pact:demand └conscription:sought

그림 3 계층적 개념 트리 구축 결과

실험 결과는 'use' 등과 같이 개념 집합에 명시되지 않은 많은 잠재적 개념이 상위 개념으로 나타나는 것을 볼 수 있으며, 개념을 정의하기 위해 사용된 속성의 수에 비해

	BE	SC	SS	GOV
No. of concepts	80	92	38	22
No. of attributes	65	69	31	21
Total no. of is-a relations	51	63	28	8
No. of exact is-a relations	2	8	10	1
Precision of relations (%)	3.92	12.69	35.71	12.5

표 2 계층적 개념 트리 구축의 평가

개념의 수가 비교적 많음을 볼 수 있으며, 이것은 속성을 효과적으로 구별하지 못함을 의미한다. 또한, 특수한 개념 간에 존재하는 is-a관계가 비교적 정확한 결과를 보인다.

5. 결론 및 추후 과제

본 논문에서는 문서로부터 계층적 개념 트리를 구축하기 위해 FCA를 사용하는 방법을 제시했다. 실험 결과는 문서에 등장하는 개념간의 is-a 관계를 발견할 수 있으나, 그것의 정확도가 현저히 떨어짐을 보여준다. 이것의 이유는 동사는 취하는 목적어와 전치사의 종류에 따라 다양한 의미를 가지는데 반하여 제한한 방법에서는 동일한 동사를 같은 의미의 속성으로 간주하기 때문이다. 계층적 개념 트리 구축 시스템의 정확도를 높이기 위한 추후 과제로써 동사의 의미를 구별하기 위한 방법을 도입하고, 개념-속성간의 관계를 정의하는 다양한 방법에 대한 실험과 상위 개념으로 나타나는 잠재적 개념의 발견에 대한 연구를 진행할 것이다.

6 참고 문헌

- [1] Dawn Lawrie and W. Bruce Croft, "Discovering and Comparing Topic Hierarchies, Proceedings of RIAD2000 conference, pp. 314-330, 2000.
- [2] Mark Sanderson and Bruce Croft, Deriving Concept Hierarchies from Text, Proceedings of the 22<sup>th</sup> Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, pp. 206-213, 1999.
- [3] Dawn Lawrie, W. Bruce Croft, and Arnold Rosenberg, "Finding Topic Words for Hierarchical Summarization", Proceedings of the 24<sup>th</sup> Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, pp. 349-357, 2001.
- [4] Dawn J. Lawrie and W. Bruce Croft, "Generating Hierarchical Summaries for Web Searches", Proceedings of the 26<sup>th</sup> Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, pp. 457-458, 2003.
- [5] Uta Priss, "Formal Concept Analysis", <http://www.fcachome.org.uk/>
- [6] Gu Tao, "Using Formal Concept Analysis (FCA) for Ontology Structuring and Building", <http://www.ntu.edu.sg/home5/pg04247176/index.htm>
- [7] Philipp Cimiano, Steffen Staab and Julien Tane, "Automatic Acquisition of Taxonomies from Text: FCA meets NLP", Proceedings of the GI Workshop Lehren - Lernen - Wissen - Adaptivität (LLWA), 2003.
- [8] Natalya F. Noy and Deborah L. McGuinness, "Ontology Development 101: A Guide to Creating Your First Ontology", Stanford Knowledge Systems Laboratory Technical Report KSL-01-05, 2001.