

온톨로지 자동추출도구의 기능적 성능 평가를 위한 평가지표의 개발 및 적용

박진수

서울대학교 경영전문대학원
(jinsoo@snu.ac.kr)

조원진

서울대학교 일반대학원 경영학과
(cool97@snu.ac.kr)

노상규

서울대학교 경영전문대학원
(srho@snu.ac.kr)

.....

이제 인터넷은 시맨틱 웹(Semantic Web)의 형태로 진화 발전하고 있다. 그 결과 시맨틱 웹의 지식을 표현하는 백본인 온톨로지가 매우 중요하게 인식되고 있다. 그러나 온톨로지를 구축하는 것은 많은 시간과 자원을 필요로 하는 작업이다. 이로 인해 온톨로지 추출도구(ontology extraction tool)에 대한 개발의 필요성이 지난 십여 년간 제기되어 왔으며, 온톨로지를 자동으로 추출하거나 온톨로지 구축을 돕는 도구들이 개발되었다. 그러나 여러 온톨로지 자동추출 도구들 중에 구축하고자 하는 온톨로지의 사용자 요구사항에 적합한 자동추출도구를 선택하기 위해서는 이런 도구들에 대한 평가지표가 필요하다. 하지만, 현재 이런 도구들을 평가하기 위한 포괄적인 평가 프레임워크(comprehensive evaluation framework)가 존재하지 않는다. 본 연구에서는, 문헌연구를 수행하여 온톨로지 자동추출도구가 갖춰야 할 핵심 요소들을 찾고, 온톨로지 추출도구들을 평가하기 위한 일련의 평가지표들을 개발했다. 또한 본 연구에서 제안하는 평가지표에 따라 온톨로지 자동추출도구인 OntoLT, Text-To-Onto, TERMINAE, OntoBuilder를 평가해 보았다.

.....

논문접수일 : 2008년 10월 논문수정일 : 2008년 12월 게재확정일 : 2008년 12월 교신저자 : 조원진

1. 서론

현재 시맨틱 웹(Semantic Web)이 차세대 웹으로 주목을 받고 있다. 온톨로지는 이런 시맨틱 웹을 위한 지식표현의 백본(backbone)으로서 시맨틱 웹의 핵심 요소로 간주되고 있다. 그 결과, 다양한 분야에서 많은 연구자들이 온톨로지의 중요성을 인식하고 관련된 다양한 연구들을 해왔다. 특히 데이터 주석처리(data annotation), 데이터 통합(data integration), 지능형 응용프로그램 구축(intelligent applications building) 등의 분야에 있어 온톨로지를 적용하고자 하는 많은 노력이 있어 왔다.

국내의 온톨로지 관련 연구는 온톨로지를 이용하여 응용프로그램의 성능을 향상시키는 것에 집중되어 연구가 진행되고 있으며, 아직까지는 구축된 온톨로지의 질(quality)을 향상시키고, 온톨로지를 효과적이고 효율적으로 구축하고자 하는 것과 관련된 연구는 거의 찾아볼 수 없다. 그러나 온톨로지를 어떻게 활용할 것인가에 대한 고민을 하기 이전에 온톨로지를 어떻게 구축하고 구축된 온톨로지를 어떻게 평가할 것인지에 대해서 깊이 있게 고민해야 한다. 이와 관련해서 해결되어야 할 문제들은 다음과 같다. 첫째, 온톨로지를 구축하는 효율적인 방법의 부재로 인해 발생하는 ‘온톨로지

* 본 연구는 서울대학교 경영대학 경영연구소의 연구비 지원에 의해 수행되었음.

병목(ontology bottleneck)' 문제가 있다. 이 문제를 해결하기 위해, 온톨로지 구축 프로세스를 자동화하기 위한 연구들이 수행되고 있다. 두 번째 이슈는 온톨로지를 구축할 때 사용되는 도구들의 평가와 관련된다. 온톨로지 관련 도구들은 크게 온톨로지 병합도구(ontology merging tools), 온톨로지 편집도구(ontology editing tools), 온톨로지 추출도구(ontology extraction tools)로 구분할 수 있다. 온톨로지 병합도구는 두 개 또는 그 이상의 온톨로지를 통합하여 하나의 온톨로지를 생성하는데 사용되며, 온톨로지 편집도구는 온톨로지를 구축하는 동안 온톨로지스트(ontologist)들이 도메인 지식을 얻고, 개념이나 관계를 구조화하고, 시각화하는 것을 돕는다. 그리고 온톨로지 추출도구는 자연어 처리나 기계학습 같은 기법을 적용하여 개념이나 관계를 자동생성하는 도구이다. 그러나 지금까지 온톨로지 병합도구(Lambrix and Edberg, 2003)나 온톨로지 편집도구(Murshed and Singh, 2005)를 평가하기 위한 프레임워크를 제시한 연구들은 존재하나, 온톨로지 추출도구를 위한 평가 프레임워크를 제안한 연구는 존재하지 않는다. 마지막으로, 구축된 온톨로지의 질적인 평가와 관련된 이슈가 있다. 이를 위해 다양한 방법이 제안되고 시도되고 있지만, 구축된 온톨로지가 의미적으로 그리고 실제적으로 얼마나 현실세계를 잘 반영했는지를 평가하기 위해서는 제대로 개발된 평가 지표에 따라 온톨로지가 포함하고 있는 개념과 그 관계를 평가하는 것이 가장 현실적인 접근법이라고 할 수 있다. 그러나 아직까지 이에 대한 연구는 미미하므로(Burton-Jones et al., 2005), 구축된 온톨로지의 질을 평가하기 위한 지표의 개발 및 응용에 대한 보다 적극적인 연구가 필요하다.

온톨로지 추출도구는 시맨틱 웹의 핵심인 온톨로지의 개발과 유지에 있어서 핵심적인 역할을 한

다. 온톨로지 추출도구를 이용하여 온톨로지를 자동 구축하는 경우, 다양한 도구들 중에 어떤 온톨로지 추출도구가 적합한지 결정하기 위해서 그 도구들을 객관적으로 평가할 필요가 있다. 그러나 앞에서 언급된 것처럼, 온톨로지 추출도구를 평가하는 지표들은 제안한 연구는 찾아볼 수 없다. 따라서 본 연구에서는 온톨로지 추출도구를 평가하는 프레임워크를 제안한다. 이 프레임워크는 온톨로지 추출도구의 기능적 평가지표뿐 아니라 추출된 온톨로지의 질적인 평가를 위한 지표도 포함하는 포괄적인 프레임워크이다. 이 프레임워크를 사용하여 대표적인 온톨로지 추출도구인 OntoLT, Text-To-Onto, TERMINAE, OntoBuilder라는 네 개의 도구를 실험하고 평가했다.

본 연구의 구성은 다음과 같다. 다음의 제 2장은 관련연구로서, 제 2.1절에서는 온톨로지 자동추출도구와 관련된 기존 연구들을 소개하며, 제 2.2절에서는 온톨로지 구축을 위해 사용되는 도구들을 위한 평가지표를 제시한 연구들과 함께, 구축된 온톨로지의 평가를 위한 프레임워크를 다룬 기존 연구에 대해서 살펴본다. 이어서 제 3장에서는 온톨로지 자동추출도구 중 본 연구의 실험에 사용된 도구들을 소개한다. 그리고 제 4장에서는 본 연구에서 제안하고 있는 온톨로지 추출도구를 위한 평가 프레임워크를 소개하며, 제 5장에서는 본 연구에서 제안한 프레임워크에 맞춰 실제 온톨로지 추출도구들을 평가한 결과를 제시한다. 마지막으로 제 6장에서는 본 연구를 종합 정리한다.

2. 관련 연구

2.1 온톨로지 자동 추출

컴퓨터에 의한 지능형 의미기반의 지식 및 정보

의 자동처리를 위한 방안으로써 온톨로지에 대한 관심이 높아지고 있으며, 실제 각 도메인 별로 온톨로지를 구축하고자 노력해 왔다. 그러나 온톨로지를 구축하는 것은 관련된 도메인 지식을 얻어야 하는 등 많은 시간을 요하며, 숙련된 인력을 필요로 하는 어려운 작업이다. 이에 많은 연구자들이 온톨로지를 자동생성하는데 관심을 가져왔다. 그 결과 개발된 온톨로지 추출도구들은 온톨로지스트들이 온톨로지 구축 프로세스를 보다 효율화할 수 있게 만들었다.

온톨로지 추출을 위해 사용되는 기법들은 기계학습, 지식획득, 자연어처리, 그리고 정보검색과 같은 다양한 분야의 방법들을 적용하고 있다. 지난 십여 년간, 지식획득 프로세스의 자동화를 위해 몇 가지 접근법들이 제안되어 왔으며(Srikant and Agrawal, 1995; Faure and Nedellec, 1998; Maedche and Staab, 2000; Staab et al., 2001), 특히 자연어처리와 기계학습 기법이 매우 활발히 적용되고 있다.

온톨로지 자동 구축을 위한 대부분의 연구들이 텍스트로부터 온톨로지를 추출하는 것을 지향해 왔으며 많은 시스템들이 개발되어 왔다. 대표적인 시스템으로 OntoLearn(Velardi et al., 2002), Text-To-Onto(Maedche and Volz, 2001), ASIUM(Faure and Nedellec, 1999), OntoLT(Buitelaar et al., 2004), TERMINAE(Biebow and Szulman, 1999), OntoBuilder(Gal and Modica, 2004), SOAT(Wu and Hsu, 2002), SVETLAN(Chaelandar and Grau, 2000), Mo'K Workbench(Bisson et al., 2000) 등이 있다.

반면 국내 연구의 경우, 온톨로지 자동추출도구의 개발과 관련된 연구를 거의 찾아볼 수가 없다. 온톨로지 생성과 관련된 국내 연구의 대부분이 아직은 수작업의 효율 및 정확도를 높이는 방안에 치

중되어 있으며(윤현주 et al., 2004; 송도규, 2005; 구미숙 et al., 2006) 온톨로지 추출과 관련된 연구도 대부분이 시소러스에서 온톨로지를 추출하는 것에 제한되어 있다(정도현과 김태수, 2003). 따라서 다양한 소스 데이터를 입력 받아서 개념과 그 관계를 자동 추출하는 도구에 대한 연구가 필요하다고 판단된다. 온톨로지 자동생성을 위한 방법론을 제안한 국내의 논문을 살펴보면 다음과 같다. 송도규(2005)는 세종전자사전이라는 정형의 전자사전에서 자동 프로그램을 이용하여 쉽고 빠르게 OWL의 대용량 온톨로지를 작성할 수 있는 방법론을 제시했다. 세종전자사전의 체언과 용언을 각각 OWL의 클래스와 속성으로 분류한 후, 도메인 별로 나누어 각 도메인 별로 온톨로지를 구축할 수 있으며, 반의어, 동의어 관계어 등의 데이터를 적절히 활용할 수도 있음을 보였다. 구미숙 등(2006)은 연관규칙 알고리즘을 이용하여 온톨로지 반자동생성 기법을 제안했다. 인터넷 상에 존재하는 관광정보 사이트를 통해서 찾아낸 관광정보들을 XML 문서로 변환한 후, 각 문서에 존재하는 태그를 추출하여 연관규칙 알고리즘을 적용하였다. 그 결과 도출된 빈발패턴 중에서 서로 관련 있는 개념의 쌍을 추출한 후 온톨로지 자동생성의 기반을 마련하여 온톨로지를 구축하였다. 윤현주 등(2004)은 비구조화된 텍스트로 이루어진 문서에서 원하는 정보를 추출하기 위해 사람이 문서를 분석하고 규칙을 생성하여 그 생성된 규칙에 의해 문서를 구조화하고 자동으로 온톨로지를 구축하는 방법을 제안했다.

2.2 평가 프레임워크

2.2.1 온톨로지 구축 도구에 대한 평가 프레임워크

앞서 언급했듯이, 온톨로지 구축을 위해 사용되

는 도구들로는 크게 온톨로지 병합도구, 온톨로지 편집도구, 그리고 온톨로지 추출도구로 구분할 수 있다. 온톨로지에 대한 관심이 높아지면서 온톨로지 병합도구나 온톨로지 편집도구를 평가하기 위한 프레임워크들이 제안되어 왔다. 반면 온톨로지 추출도구를 위한 평가 프레임워크를 제안한 연구는 거의 찾아볼 수 없다. 따라서 본 연구에서는 온톨로지 추출도구를 위한 평가지표를 개발하고자 한다. 그러나 온톨로지 병합도구나 편집도구를 위한 평가지표에 대한 연구 결과들 역시 온톨로지 추출도구를 위한 평가지표를 개발하는 경우 중요한 참조연구가 될 수 있으므로, 본 절에서는 온톨로지 병합도구와 편집도구의 평가와 관련된 기존 연구들에 대해서 살펴보도록 하겠다.

Murshed와 Singh(2005)는 온톨로지 편집도구들을 평가하기 위한 지표들을 제안했다. 그들의 프레임워크는 기능성(functionality), 재사용성(reusability), 데이터 저장용량(data storage), 복잡성(complexity) 외에 다수의 지표들을 포함하고 있다. 이 지표들에 기초하여, Protégé, Metis, OntoEdit을 실제로 평가했으며, 그 결과 Protégé는 대용량의 온톨로지를 구축하는 경우에 적합하고, Metis는 중간 사이즈의 온톨로지를 구축하기에 적합하다는 결론을 내렸으며, 마지막으로 OntoEdit은 비즈니스 모델링에 적합하다고 판단했다. Lambrix와 Edberg(2003)는 가장 잘 알려진 온톨로지 병합도구들이라고 할 수 있는 PROMPT와 Chimaera가 생명정보학(Bioinformatics) 관점에서 현존하는 생명정보 분야 온톨로지를 병합하는데 적합한지를 평가했다. 테스트 온톨로지로서 Gene 온톨로지와 Single 온톨로지를 사용하였으며, 그들의 평가지표는 이용가능성(availability), 안정성(stability), 표현언어(representation language), 기능성(functionality), 사용자 인터페이스(user in-

terface)를 포함하고 있다. 제안된 평가지표에 기초하여 평가한 결과, PROMPT와 Chimaera 모두 현존하는 생명정보 분야의 온톨로지를 병합하는데 적합하다고 결론 내렸다.

2.2.2 온톨로지의 질적 평가

구축된 온톨로지의 질을 평가하려는 다양한 시도가 이루어져왔다. 이런 다양한 노력은 크게 네 개의 범주로 분류될 수 있다(Brank et al., 2006). 첫 번째 접근법은 금본위제(gold standard)로서, 구축된 온톨로지를 잘 구축되었다고 판단되는 미리 정해놓은 온톨로지와 비교를 통해서 이루어진다. 두 번째는 응용프로그램 기반 평가 방법이다. 구축된 온톨로지를 실제로 적용될 응용프로그램에 적용해봄으로써, 온톨로지가 제대로 기능을 하는지 판단해 보는 것이다. 예를 들어, 의료문헌정보 검색시스템의 성능을 향상시키기 위해 온톨로지를 구축한 경우, 그 온톨로지를 검색시스템에 도입하여 검색시스템의 성능을 어느 정도 향상시켰는지를 평가함으로써 온톨로지 자체의 질을 평가할 수 있다. 세 번째 접근법은 데이터 주도형 방법이다. 온톨로지와 그것이 참조하는 문제영역의 데이터 집합 사이의 적합성 및 관련도를 평가함으로써 온톨로지의 질을 평가한다. 마지막 방법은 전문가나 온톨로지스트 등의 사람에 의존하는 접근법이다. 이 접근법은 온톨로지가 미리 정의된 일련의 지표나, 표준, 필요사항 등을 얼마나 잘 충족하는지를 실제 사람이 평가함으로써 행해진다. 본 연구에서 수행된 온톨로지의 질을 평가하는 시도는 미리 정의된 지표들을 가지고 사람에 의해 수행되었기 때문에 마지막 접근법에 해당된다.

Burton-Jones 등(2005)은 구축된 온톨로지의 질을 평가하기 위한 프레임워크를 제안했다. 그들

은 기호학 이론에 기초하여 구문적(syntactic), 의미적(semantic), 실질적(pragmatic), 그리고 사회적(social) 측면에서 온톨로지의 질을 평가했다. 그리고 그들은 실제로 지표들을 조작화(operation-
alize)하여 “Ontology Auditor”라는 프로토타입을 구현해서 기 구축된 온톨로지들을 평가했다. 평가 결과, 그들이 제안한 지표들은 실행 가능하며, 온톨로지들 사이의 질적인 측면에 있어서도 광범위한 차이가 있음을 강조했다.

지금까지 온톨로지 추출도구에 대한 연구와 온톨로지 관련 도구들을 평가하기 위한 프레임워크들을 제안한 연구들, 그리고 구축된 온톨로지의 질적인 측면을 평가하고자 한 연구들에 대해 살펴보았다. 살펴본 바와 같이, 온톨로지 자체의 질적인 측면을 평가하거나 온톨로지 병합도구나 편집도구들의 기능적 특성을 평가하기 위한 프레임워크에 대해서는 연구가 진행되어 있다. 그러나 앞서 언급했듯이 온톨로지 추출도구에 대한 평가지표는 존재하지 않는다. 온톨로지 추출도구들은 병합도구들과 편집도구들과는 다른 특성들을 가지므로 온톨로지 추출도구를 평가하기 위한 새로운 프레임워크의 개발이 절실하다.

3. 온톨로지 자동추출도구

3.1 OntoLT¹⁾

OntoLT는 독일의 인공지능 연구 센터인 DFKI에서 개발된 도구로서, 온톨로지 편집도구인 Protégé에 플러그인하여 사용할 수 있다(Buitelaar et al., 2004). 이 도구는 미리 정의된 매핑 규칙에 따라 언어학적 주석처리(linguistic annotation)가 된 문서집합으로부터 개념과 그 관계를 자동적으

로 추출할 수 있다. 그러나 언어학적 주석처리 기능은 OntoLT에 포함되어 있지 않으므로, OntoLT를 사용하여 온톨로지를 추출하기 위해서는 입력 데이터가 미리 언어학적으로 주석처리 되어 있어야 한다.

OntoLT를 사용하여 온톨로지를 추출하는 과정은 다음과 같다. OntoLT는 사용자가 매핑 규칙을 정의할 수 있는 선행조건 언어(precondition language)를 제공하는데, 이 선행조건들은 언어학적 주석처리에 대해 XPASS 표현들로서 구현된다. 만약 선행조건이 만족되면, 매핑 규칙들이 발견된 후보들을 어떤 방식으로 기술할지에 대한 운영자(operator)를 활성화시킨다. 이 도구에는 많은 매핑 규칙들이 포함되어 있으며, 사용자가 추가적으로 규칙을 정의할 수도 있다(Buitelaar et al., 2004).

3.2 Text-To-Onto²⁾

Text-To-Onto는 독일의 Karlsruhe 대학의 AIFB Institute에서 개발되었다(Maedche and Volz, 2001). 이 도구는 초기 온톨로지(initial ontology)로부터 도메인 온톨로지를 구축하기 위한 환경을 가지고 있다. Text-To-Onto는 입력 데이터로써, 일반 텍스트, 반구조화된 텍스트, 사전들, 기존의 온톨로지, 데이터베이스 등을 사용할 수 있다. 초기 추출 과정의 결과는 온톨로지가 사용될 도메인과 관련 없는 개념들도 포함하고 있는 도메인 온톨로지이다. 도메인과 관련 없는 개념들은 도메인 온톨로지의 어휘에 더 적합하도록 색출하여 제거된다. 따라서 최종에 남는 결과는 특정 도메인과 관련 있는 개념들만을 포함하고 있는 도메인 온톨로지가 된다. 온톨로지리스트들에 의해 감독되는 전체 과정은 순환적인 프로세스로 진행되며, 정제되

1) <http://olp.dfki.de/OntoLT/OntoLT.htm/>

2) <http://ontoserver.aifb.unikarlsruhe.de/texttoonto/>

어 완전한 도메인 온톨로지를 구축하게 된다.

3.3 TERMINAE³⁾

TERMINAE는 프랑스의 Paris-Nord 대학의 정보 연구소에서 개발되었다(Biébow and Szulman, 1999). 이 도구는 언어학적 도구(linguistic tools)와 지식공학 도구(knowledge engineering tools)를 통합했다. 언어학적 도구는 입력 데이터에서 용어의 발생을 분석하여 용어적 형태(terminological form)를 추출한다. 이후 온톨로지스트들이 그 용어의 의미를 정의하기 위해 입력 데이터 내에서 용어의 사용을 분석한다. 지식공학 도구는 편집기와 브라우저를 포함하고 있으며 용어적 형태를 개념으로서 표현하는 것을 쉽게 하도록 돕는 역할을 한다(Gomez-Perez and Manzano-Macho, 2005).

TERMINAE를 이용하여 개념을 추출하는 과정은 다음과 같다. 우선, 이 도구는 용어추출 도구를 사용하여 용어의 리스트를 만들게 되며, 추출된 일련의 후보 용어들이 온톨로지스트들에게 개념으로서 제안되면, 온톨로지스트들은 관련된 용어들을 선택한다. 이후 온톨로지스트들은 이 용어들을 개념화(conceptualization)하고 의미를 정의하기 위해 각 용어가 어떻게 사용되는지를 분석한다. 그 다음 온톨로지스트가 정의된 내용을 온톨로지 언어(ontology language)로 바꾼다.

3.4 OntoBuilder⁴⁾

OntoBuilder는 이스라엘의 Technion-Israel Institute에서 개발되었다(Gal and Modica, 2004). OntoBuilder는 원래 자동 스키마 매칭 알고리즘을

평가하고, 그 알고리즘을 개선하기 위한 하나의 하위 모듈로서 개발되었다. 웹 브라우저처럼 작동하도록 고안된 OntoBuilder는 단순한 검색엔진에서부터 여러 페이지로 이루어진 웹사이트에 이르는, 다양한 웹 검색 인터페이스로부터 온톨로지를 자동적으로 추출하는 것을 지원한다. 이를 위해, OntoBuilder는 HTML 문서들의 훈련집합(training set)으로부터 학습된 휴리스틱(heuristic) 방법에 기반하여 개념과 그 관계를 추출한다.

OntoBuilder를 이용하여 온톨로지를 추출하기 위해서, 사용자는 단순히 웹 사이트의 URL을 입력하면 된다. 추출하고자 하는 웹 페이지가 OntoBuilder에 로딩이 되면, 각 페이지는 “document object model”이라고 불리는 데이터 구조로 파싱이 된다. 이후, OntoBuilder는 HTML 페이지로부터 레이블과 필드 이름을 추출하여 용어들의 사전을 만들고, 매칭 알고리즘을 이용하여 그 용어들 사이의 독특한 관계를 인식하여 관계를 추출한다.

4. 온톨로지 자동추출도구를 위한 평가지표

새롭게 개발된 시스템들은 잘 완성된 평가 방법을 사용하여 그 시스템의 유용성 및 성능에 대해 엄격하게 입증되어야 하며, 이와 같은 평가는 개발된 시스템이 실제 사용되는 비즈니스 환경 내에서의 기술적인 측면을 포함하여 면밀히 이루어져야 한다(Hevner et al., 2004). 개발된 시스템을 평가함에 있어서 적절한 평가지표들이 우선 정의되어 있어야 하며, 평가에 필요한 데이터를 모으고 분석하는 것 또한 매우 중요하다. 다시 말해 개발된 시스템의 성능을 측정하기 위한 평가지표가 적합해야만 구축된 시스템의 성능에 대해 보다 엄밀히 평가할 수 있다. Hevner 등(2004)은

3) <http://www-lipn.univ-paris13.fr/~szulman/>

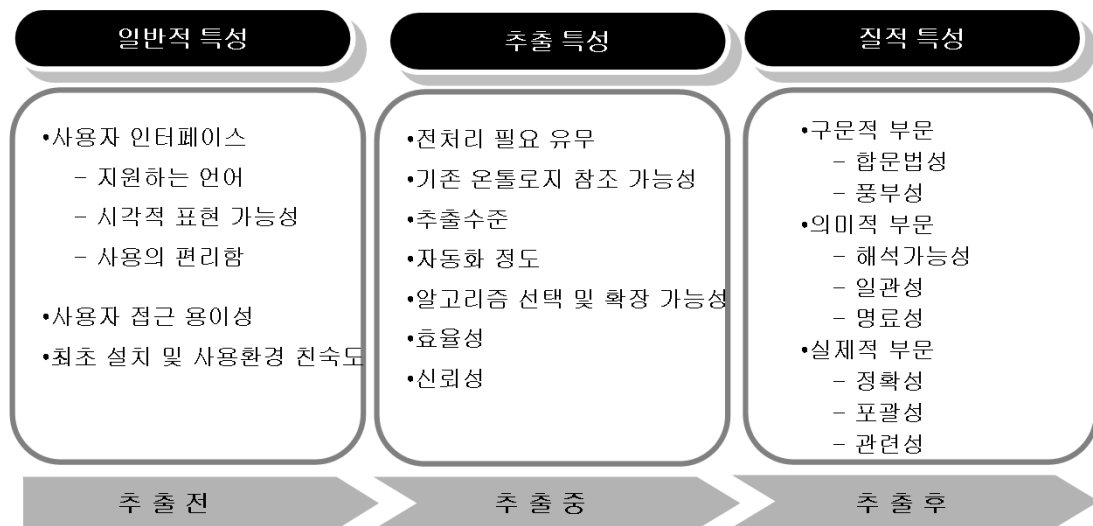
4) <http://iew3.technion.ac.il/OntoBuilder/>

정보시스템 연구 프레임워크를 제안하면서 새로운 시스템을 개발함에 있어서 시스템을 구축하는 것(building)과 함께 시스템의 정당화 및 평가(justification and evaluation)라는 두 가지 주요 활동을 정의하면서 평가지표의 중요성에 대해서 강조하고 있다.

본 연구에서도 평가의 중요성을 인식하고 온톨로지 자동추출도구들을 평가하기 위한 일련의 지표들을 개발했다. 이를 위해 온톨로지 구축의 자동화와 관련된 문헌연구를 수행하였으며, 온톨로지 추출도구들을 조사하고, 그것들 중 일부를 실제 사용하여 보았다. 그 결과, 온톨로지 자동추출도구들이 갖춰야 할 핵심 요소 및 기능적 특성들을 찾았으며, 이를 온톨로지 추출도구를 위한 평가 프레임워크로 제안하고자 한다. 앞서 언급했듯이, 본 연구에서는 기존에 제안된 온톨로지 병합도구나 온톨로지 편집도구의 평가 프레임워크가 단순히 표현언어, 인터페이스 같은 관련 도구의 기능적인 면의 평가에 국한된 것과는 달리, 온톨로지 추출도구

의 성능적인 측면을 평가함에 있어 온톨로지 추출도구를 사용하여 구축된 온톨로지의 질적인 평가 부분도 함께 고려한다.

본 연구가 제안하고 있는 온톨로지 추출도구를 위한 평가지표들은 다음의 <그림 1>에서 보는 바와 같이 세가지 특성으로 구분할 수 있다. 첫 번째는 ‘일반적 특성’으로, 도구의 인터페이스 및 사용의 편리함과 관련된 지표들이다. 온톨로지 추출도구가 온톨로지를 추출하는 프로세스를 진행하기 전에 이 특성이 고려되고 평가될 수 있다. 두 번째 관점은 ‘추출 특성’이다. 이것은 온톨로지 추출도구를 평가하는데 있어서 가장 중요한 지표들로 간주되며, 실제 온톨로지를 추출하는 프로세스를 진행하는 과정에서 고려될 수 있는 항목들이다. 마지막은 ‘질적 특성’이다. 온톨로지 추출이 끝난 후, 사용자는 구축된 온톨로지의 질적인 측면을 평가해야 한다. 구축된 온톨로지의 질적인 평가는 문법적(syntactic), 의미적(semantic), 그리고 실제적(pragmatic) 측면에서 평가된다.



<그림 1> 온톨로지 추출도구 평가 프레임워크

4.1 일반적 특성

위에서 언급했듯이, 이 특성은 도구의 외형적인 특성 및 사용상의 편리함과 관련된 지표들을 포함하고 있다. 일반적 특성에서, 첫 번째 평가지표는 **사용자 인터페이스**이다. 많은 연구들에서 사용자 인터페이스는 프로그램들이나 소프트웨어를 평가할 때 중요한 특성으로서 고려되었다(Duineveld et al., 1999; Lambrix and Edberg, 2003; Murshed and Singh, 2005; Angele and Sure, 2002). 사용자 인터페이스를 평가하기 위해 본 연구는 (1) 지원하는 언어, 추출된 요소의 (2) 시각적 표현 가능성, (3) 사용의 편리함이란 세가지 속성을 정의했다. 첫 번째, 지원하는 언어는 도구가 다양한 언어를 지원하는지 그리고 사용자가 인터페이스 언어를 선택할 수 있는지를 평가한다. 인터페이스가 공용어라고 할 수 있는 영어 같이 사용자에게 친숙한 언어로 표시된다면 사용자는 보다 쉽게 추출도구를 사용할 수 있을 것이다. 그러나 일부 추출 도구들은 개발자가 속한 국가의 언어로 개발된 경우도 많아서 그 언어를 이해하지 못하는 사용자의 이용을 제한하고 있다. 두 번째, 추출된 요소의 시각적 표현 가능성은 추출된 개념이나 관계의 시각적 표현 가능 여부를 평가한다. 온톨로지 추출도구를 사용해서 구축된 온톨로지를 시각적으로 표현할 수 있다면 사용자는 그 온톨로지의 구조를 좀 더 쉽게 이해할 수 있으며, 그 결과 온톨로지를 편집하는 것이 용이할 수 있다. 사용자 인터페이스를 평가하기 위한 마지막 속성으로, 사용의 편리함 측면도 고려할 필요가 있다. 온톨로지 추출도구의 인터페이스가 사용하기 복잡하게 디자인되었다면 추출을 수행하기 전에 관련된 기능을 마스터하기 위해 상당한 시간을 요할 수 있다.

일반적 특성의 두 번째 평가지표는 **사용자 접근 용이성**이다. 사용자가 큰 어려움과 노력없이 도구에 접근하여 사용할 수 있다면, 그 도구의 사용자 접근 용이성은 그렇지 않은 것들보다 높다고 할 수 있다. 특히, 무료로 도구를 다운로드할 수 있다면 그것의 사용자 접근 용이성은 매우 높다고 할 수 있다.

일반적 특성의 마지막 평가항목은 **최초 설치 및 사용환경 친숙도**이다. 이는 그 도구를 PC에 설치하기 용이한지, 설치를 위해 추가적인 플랫폼 요구사항이 존재하는지 그리고 도움말이나 매뉴얼 기능을 제공하고 있는지에 대해서 평가한다.

4.2 추출 특성

본 연구는 온톨로지 추출도구를 위한 평가 프레임워크이기 때문에, 추출 특성이 다른 특성들 보다 더 중요하게 고려되며, 그 결과 많은 평가지표들을 포함하고 있다. 추출 특성은 앞에서 언급했듯이 온톨로지 추출 프로세스 자체와 직접 관련이 있다.

온톨로지 추출도구를 평가할 때 우선 입력 데이터와 관련해서 두 가지 평가지표들을 제안한다. 첫째가 **전처리 필요 유무**이다. 도구가 온톨로지를 추출할 수 있는 소스 데이터의 선택폭이 좁거나 특정 데이터 형태만을 입력 데이터로서 받아들이는지 평가해 볼 필요가 있다. 만약 사용하고자 하는 소스 데이터는 순수한 텍스트 파일인데, 사용하는 온톨로지 추출도구가 입력 데이터로서 사용할 수 있는 데이터 형태를 HTML으로 제한한다면 데이터 형태를 바꾸기 위한 추가작업이 필요할 수밖에 없다. 또한 일부 온톨로지 추출도구들은 언어학적으로 주석처리가 되어 있어야지만 추출 과정을 진행할 수 있는 것들도 있다. 이처럼 온톨로지

추출도구가 추출 과정을 진행하기 전에 입력 데이터들에 추가적인 전처리가 선행되어야 하는 경우는 그렇지 않은 도구들 보다 사용상에 있어 더 많은 제약을 가지게 된다. 뿐만 아니라 기존에 이미 구축된 온톨로지들의 개념을 참조하여 사용할 수 있는 지도 평가해보아야 한다(**기존 온톨로지 참조 가능성**). 현재 여러 도메인에 다양한 형태의 온톨로지가 구축되어 있으므로 기존에 존재하는 온톨로지를 활용하면 보다 쉽게 온톨로지를 구축할 수 있다.

또한, 온톨로지 추출도구를 평가할 때 각 도구의 **추출수준**을 조사할 필요가 있다. 온톨로지는 개념과 그 개념의 관계로 구성되는데, 추출수준은 도구가 이런 개념과 관계 모두를 자동(혹은 반자동)적으로 추출하는지, 아니면 개념만을 자동(혹은 반자동)적으로 추출할 수 있는지에 대해서 평가한다. 또한 추출수준과 함께 온톨로지의 개념과 관계를 자동적으로 추출할 수 있는지도 평가해야 한다(**자동화 정도**). 자동적 추출도구란 추출 프로세스의 전과정 중에 있어서 인간의 개입 없이 개념이나 관계를 추출하는 것을 의미한다. 반면, 추출 규칙 등이 사용자에게 의해 미리 정의되어야 하는 도구들은 반자동 추출도구로 분류된다. 온톨로지 자동추출도구라고 이름 붙은 많은 도구들이 사실상 자동적이라고 할 수 없는 경우가 많기 때문에 자동화 정도에 대한 평가지표도 함께 고려되어야 한다.

그 동안 개념과 관계를 추출하기 위한 다양한 알고리즘들이 제안되어 왔다. 따라서 온톨로지 추출도구가 구축하고자 하는 온톨로지의 성격에 맞춰 다양한 추출 알고리즘 중에 선택하여 사용할 수 있도록 옵션을 제공하는지의 여부도 평가해야 한다(**알고리즘 선택 및 확장 가능성**). 만약 온톨

로지 추출도구가 다양한 알고리즘을 제공하고, 사용자가 구축하고자 하는 온톨로지에 적합한 알고리즘을 적용할 수 있다면 보다 효율적으로 온톨로지를 구축할 수 있을 것이다. 더 나아가 온톨로지 추출도구에 플러그인으로 사용자가 추출 알고리즘을 쉽게 추가할 수 있다면 추출도구의 유연성을 높일 수도 있다.

추출된 산출물 역시 평가해야 한다. 이를 위해 다른 소프트웨어와 마찬가지로, 온톨로지 추출도구 역시, **효율성** 측면이 평가되어야 한다. 입력 데이터를 읽어 들여 목표 산출물을 사용자에게 보여주는 시간을 계산함으로써 추출도구가 얼마나 빨리 개념이나 관계를 추출하는지 평가할 수 있다. 어떤 온톨로지 추출도구는 추출 알고리즘의 비효율성으로 인해 추출하는데 오랜 시간이 걸리는 것들도 있으므로, 만약 비슷한 추출 결과를 보여주는 도구라면 보다 효율성이 높은 도구를 선택하는 것이 합리적이다. 마지막으로, **신뢰성**이란 같은 입력 데이터를 사용하여 추출 과정을 진행한 경우에 항상 동일한 추출 결과를 보여주는 지에 대해 평가하는 항목으로, 이 지표는 산출물의 질적인 측면을 평가하는 것이 아니라 추출도구가 얼마나 일관된 결과를 산출하는지를 조사하는 지표이다.

4.3 질적 특성

온톨로지를 추출하는 과정이 끝났다면, 이제는 구축된 온톨로지의 질적인 측면을 평가해야 한다. 많은 연구들이 구축된 온톨로지의 질적인 측면을 평가하기 위해 형식 온톨로지(formal ontology)를 참조하여 왔다(Chidamber and Kemerer, 1994; Wand and Weber, 1995). 그러나 참조모델로 선택된 형식 온톨로지 역시 유효성을 평가할 방법은

없다. 뿐만 아니라 형식 온톨로지는 실제 적용에 있어서 너무 높은 수준이고 철학적이기 때문에 도메인 수준에서 구축된 온톨로지의 질적인 측면을 평가하는데 있어서 다소 무리가 있다. 반면, 관련 연구에서도 살펴보았듯이, Burton-Jones 등(2005)은 온톨로지의 질적인 측면의 평가를 위해 기호학적 이론(semiotic theory)에 기초하여 일련의 평가 지표들을 도출하였으며 실제 적용과 관련된 평가 지표들을 포함하고 있다. 따라서 본 연구에서는 온톨로지 추출도구에 의해 구축된 온톨로지의 질을 평가함에 있어서 그들의 잘 정의된 지표들을 참조한다. 그들의 프레임워크에는 구문적, 의미적, 실제적, 사회적 부분의 지표들을 포함하고 있다. 그러나 사회적 부분은 온톨로지가 얼마나 사용되는지, 다른 온톨로지들에 의해서 얼마나 참조되는지의 등을 평가하는 것으로서 지속적인 모니터링과 조사가 필요하기 때문에 본 연구에서는 포함하지 않았다. 따라서 본 연구에서는 구축된 온톨로지의 질적인 측면을 평가하기 위해 구문적, 의미적, 실제적 부분만을 평가한다.

첫째, **구문적(syntactic) 부문**은 온톨로지가 쓰여진 방식에 따라서 온톨로지의 질적인 측면을 평가한다. 이 부문은 합문법성과 풍부성이란 두 개의 속성을 가진다. 합문법성은 구축된 온톨로지의 구문이 작성된 언어 규칙에 어느 정도까지 맞게 생성되었는지를 평가한다. 즉, 온톨로지의 구문(syntax) 정확성을 평가한다. 풍부성은 구축된 온톨로지가 개념들과 공리들을 모두 포함하는지 아니면 개념들만을 가지고 있는지를 평가한다. 단순히 개념들만을 가지고 있는 온톨로지보다는 공리들까지 함께 가지고 있는 온톨로지가 더 풍부하다고 할 수 있다.

둘째, **의미적(semantic) 부문**은 구축된 온톨로

지에서 표현된 용어의 의미를 조사하며, 해석가능성, 일관성, 명료성이란 세가지 속성을 가진다. 해석가능성은 온톨로지 내에서 표현된 용어가 현실에서도 실제 사용되는 올바른 의미를 가지고 있는지를 체크한다. 예를 들어 신체와 관련된 온톨로지에 보조개 대신 볼우물로 표현되어 있는 경우, 현실에서는 볼우물이란 말이 거의 사용되지 않기 때문에 해석가능성이 떨어진다고 할 수 있다. 일관성은 그 용어가 온톨로지 내에서 일관된 의미를 가지고 있는지를 평가한다. 예를 들어, 한 온톨로지 내에서 A가 B의 부분집합이고, 또한 B가 A의 속성이 된다고 정의되어 있다면, 이는 일관된 의미를 갖지 않으며, 어떤 의미적 가치도 없다. 명료성은 용어의 맥락이 분명한지를 평가한다. 예를 들어, 배라는 클래스가 신체라는 속성을 가진다면, 시스템은 그 배가 과일인 한 종류인 먹는 배가 아닌 사람의 배를 의미하는 것임을 분명히 알 수 있다.

마지막으로, **실제적(pragmatic) 부문**은 사용자나 그들의 에이전트들에게 있어 온톨로지의 유용성을 평가한다. 이 부문은 정확성, 포괄성, 관련성이란 세가지 속성을 가진다. 정확성은 구축된 온톨로지에 표현되고 있는 개념이나 관계가 실제 사실인지에 대해 평가한다. 예를 들어, 온톨로지에 기혼자가 성인으로 표현이 되어있다면, 이는 현실 세계에서는 기혼자중 미성년자가 존재하기 때문에 이는 사실과 다르므로 정확성에 위배된다. 포괄성은 구축된 온톨로지가 커버하는 범위를 평가한다. 온톨로지가 커버하는 범위가 크다면 도메인을 보다 포괄적으로 표현하고 있다고 판단할 수 있으며, 그 결과 사용자에게 보다 유용한 정보를 제공할 수 있을 것으로 기대할 수 있다. 관련성은 온톨로지가 사용자의 특정 필요사항을 만족하는지를 평가한다.

5. 온톨로지 자동추출도구의 실험

본 연구에서는 제 3장에서 소개된 OntoLT, Text-To-Onto, TERMINAE, OntoBuilder라는 네 개의 온톨로지 추출도구들을 제 4장에서 제안한 평가 프레임워크를 사용하여 평가했다. 온톨로지 구축 경험이 있는 네 명의 온톨로지 전문가들이 의학 논문의 데이터 집합을 사용하여 실제 온톨로지를 구축해 보고 각 도구들을 평가했다.

5.1 일반적 특성

사용자 인터페이스의 첫 번째 속성인 지원하는 언어 측면의 평가결과, TERMINAE를 제외하고는 모두 같았다. OntoLT, Text-To-Onto, OntoBuilder는 모두 영어로 인터페이스가 지원되는 반면, TERMINAE는 영어와 프랑스어 중에서 선택할 수 있도록 되어있다. 두 번째 속성인 추출된 요소의 시각적 표현 가능성은 네 도구 모두 제공하고 있다. 또한 사용자 인터페이스의 마지막 속성인 사용의 편리함이란 평가지표는 Text-To-Onto와 OntoBuilder를 제외하고는 그다지 사용하기 쉽게 설계되지 않은 것으로 평가되었다. 특히 TERMINAE의 경우, 인터페이스가 사용하기 복잡해서 그 도구들의 기능을 마스터하는데 많은 시간을 요한다고 평가되었다. 또한 메뉴를 클릭하면 각 기능들이 팝업창의 형태로 나타나므로 TERMINAE를 사용하여 작업하는 경우 화면이 클릭 몇 번으로도 충분히 복잡해지게 되어 사용하기 편하게 설계되지 않았다고 평가되었다.

또한, OntoLT, Text-To-Onto, OntoBuilder, TERMINAE 모두 웹에서 다운로드하여 사용할

수 있었으므로 사용자 접근 용이성이 모두 높았다. 마지막으로 최초 설치 및 사용환경 친숙도 측면에서는 네 개의 도구 모두 자바에 기반하여 개발되었기 때문에 Java 프로그램이 돌아가기 위한 시스템의 기초환경으로서 Java 가상머신이 미리 설치되어 있어야 했다. 또한 다른 도구들은 모두 영어로 도움말이나 설치 매뉴얼이 제공되는 반면, TERMINAE의 경우 그것들이 모두 프랑스어로 제공되기 때문에 비프랑스어 사용자의 경우는 설치나 기능을 마스터하는데 있어 매우 복잡하다고 평가되었다.

5.2 추출 특성

추출특성의 첫 번째 지표로서 전처리의 필요유무에 대해서 평가해 본 결과, OntoLT와 TERMINAE는 모두 전처리가 필요한 도구들로 평가되었다. OntoLT를 사용하여 온톨로지를 추출하기 위해서는 데이터집합이 모두 언어학적으로 주석처리가 되어있어야 하며, TERMINAE 역시 YATEA⁵⁾와 같은 용어추출기로 전처리가 되어 있어야 한다. 또한 기존 온톨로지 참조 가능성 부분에 있어서는 Text-To-Onto만이 기존 온톨로지를 참조하여 추출과정을 진행할 수 있는 도구로 평가되었다.

다음으로 추출수준 측면에서는, Text-To-Onto와 OntoLT, OntoBuilder는 데이터 집합으로부터 개념과 그것의 관계까지 추출할 수 있는 반면, TERMINAE는 개념만을 추출할 수 있는 것으로 평가되었다. 또한 TERMINAE는 용어추출기를

5) YATEA(<http://search.cpan.org/~thamon/Lingua-YaTeA-0.1/>)는 무료로 사용 가능한 용어추출기로서, 소스데이터로부터 용어들의 후보집합을 찾고, 그 중에 적합한 용어들을 자동적으로 추출하는 도구이다.

포함하고 있지 않으므로 용어추출기 없이는 개념을 추출할 수 없다. 그러나 다른 용어추출기와 상호운영이 가능하며, 본 연구에서는 용어추출기로 YATEA를 사용했다. YATEA는 무료로 이용 가능한 용어추출기로 데이터 집합에서 개념의 후보자 집합을 뽑아낸다. 그 결과 TERMINAE는 온톨로지스트들이 YATEA에 의해 추출된 용어들을 가지고 반자동적으로 온톨로지를 구축할 수 있도록 돕는다.

각 도구들의 자동화 정도 역시 다르다. Text-To-Onto는 개념과 그 관계들을 반자동 또는 자동으로 추출할 수 있다. 즉, Text-To-Onto를 사용하여 온톨로지를 구축하는 동안 사용자는 자동적 추출과 반자동적 추출 사이에서 선택권을 가지게 된다. 만약 사용자가 반자동 방법을 선택하면 추출된 결과가 맞는지 온톨로지 편집기능에 보내기 전에 판단해 볼 수 있다. 그러나, OntoLT를 사용한 경우 관계는 반자동의 방법으로만 추출된다. 따라서 사용자는 추출된 결과가 유용한지, 의미를 가졌는지 반드시 확인해야만 한다.

OntoLT와 Text-To-Onto는 다른 도구들에 비해 알고리즘을 선택하는데 있어서 더 유연성을 가지고 있다. OntoLT는 플러그인으로 많은 매핑 규칙을 포함하고 있으며, 사용자가 매핑 규칙을 추가로 정의할 수도 있다. Text-To-Onto는 역시 다양한 알고리즘들 중에서 선택하여 사용할 수 있으며, 새로운 추출 알고리즘을 추가하여 확장할 수도 있다.

마지막으로 효율성과 신뢰성 측면에 대한 평가 결과는 다음과 같다. 추출도구가 얼마나 빨리 개념이나 관계를 추출하는지를 평가하는 효율성의 경우, 2초 미만의 시간에 추출 결과를 보여주는 경우 '높음'으로 판단했으며, 2~5초의 시간이 걸

리는 경우는 '보통,' 그리고 5초 이상의 시간이 필요한 경우는 '낮음'으로 판단했다. 네 도구들에 대한 평가결과, OntoLT와 OntoBuilder는 입력 데이터를 읽어 들여 목표 산출물을 사용자에게 보여주는 시간이 2초 미만으로 매우 빨리 추출 결과를 보여줬다. 반면, Text-To-Onto는 다량의 데이터에 대해서 추출작업을 진행하는 경우, 3~5초로 다소 더디게 온톨로지를 추출했으며, TERMINAE는 그 속도가 Text-To-Onto 보다 훨씬 느렸다. 각 도구들의 신뢰성 평가를 위해, 같은 입력데이터에 대해 일관된 결과를 산출한 경우 '높음'으로, 그렇지 않은 경우는 '낮음'으로 판단했다. 평가 결과, 네 개의 도구 모두 같은 입력데이터를 가지고 반복적으로 실험을 수행한 결과, 각 도구들이 모두 일관된 결과를 산출하였기 때문에 '높음'으로 평가되었다.

5.3 질적 특성

각 추출도구를 사용하여 구축된 온톨로지를 질적인 측면에서 평가해본 결과, 우선 구문적 부분의 합문법성에 대해서는 실험에 사용된 모든 추출 도구가 문법에 맞는 구문을 가진 온톨로지를 생성하는 것으로 나타났다. 또한, 풍부성 측면에서 Onto-Builder는 클래스, 도메인, 속성, 공리 같은 제한된 구문적 어휘를 사용하기 때문에 나머지 다른 도구들이 비해 덜 풍부한 표현력을 가졌다. 반면 나머지 도구들은 클래스, 하위클래스, 속성, 하위속성, 공리, 관계차수 등 다양한 유형의 구문적 어휘를 통해 온톨로지를 표현하고 있다.

의미적 부문에 있어 각 도구를 사용하여 구축된 온톨로지를 평가한 결과 모두 해석가능성, 일관성, 명료성 측면에서 좋게 평가되었다. 그러나 의미적

<표 1> 실험 결과 종합

평가지표	OntoLT		Text-To-Onto		TERMINAE		OntoBuilder	
	지원하는 언어	영어	영어	영어	영어, 프랑수어	영어	영어	영어
일반적 특성	사용자 인터페이스	시각적 표현 기능성	표현 기능함	표현 기능함	표현 기능함	표현 기능함	표현 기능함	표현 기능함
	사용자 접근 용이성	사용의 편리함	Prolog6 사용이 익숙한 사용자의 경우 쉽게 사용할 수 있음	사용하기 쉬운 인터페이스	사용하기 쉬운 인터페이스	사용하기 사용하기 복잡함	각 기능과 구성이 사용하기 복잡함	각 기능과 구성이 사용하기 복잡함
추출 특성	최초 설치 및 사용환경 친숙도	사용자 접근 용이성	용이함	용이함	용이함	용이함	용이함	용이함
	전처리 필요 여부	진행 가능성	Java 가상머신이 설치되어 있어야 하며 영어로 간단히 작성된 설치 매뉴얼만이 제공되어 다소 숙련도를 요함	Java 가상머신이 설치되어 있어야 하며 프랑수어로 간단히 작성된 설치 매뉴얼 제공되어 다소 숙련도를 요함 특히 비프랑수어 사용자는 사용하기 어려움	Java 가상머신이 설치되어 있어야 하며 프랑수어로 간단히 작성된 설치 매뉴얼 제공되어 다소 숙련도를 요함 특히 비프랑수어 사용자는 사용하기 어려움	Java 가상머신이 설치되어 있어야 하며 프랑수어로 간단히 작성된 설치 매뉴얼만이 제공되어 다소 숙련도를 요함	Java 가상머신이 설치되어 있어야 하며 영어로 간단히 작성된 설치 매뉴얼만이 제공되어 다소 숙련도를 요함	
	기존 온톨로지 참조 가능성	추출 수준	참조할 수 없음	개념과 관계	개념과 관계	참조할 수 없음	참조할 수 없음	참조할 수 없음
	자동화 정도	자동화 정도	자동적으로 개념 추출 자동적으로 타사노미 추출 반자동적으로 관계 추출	자동적 또는 반자동적으로 개념 추출 자동적으로 타사노미 추출 자동적 또는 반자동적으로 관계 추출	자동적 또는 반자동적으로 개념 추출 자동적으로 타사노미 추출 자동적 또는 반자동적으로 관계 추출	자동적으로 추출 자동적으로 관계 추출	자동적으로 추출 자동적으로 관계 추출	자동적으로 추출 자동적으로 관계 추출
	알고리즘 선택 및 확장 가능성	효율성	높음	보통	높음	높음	높음	높음
온톨로지 특성	구문적 부문	합문법성	문법에 맞음	문법에 맞음	문법에 맞음	문법에 맞음	문법에 맞음	문법에 맞음
	의미적 부문	명료성	명료함	명료함	명료함	명료함	명료함	명료함
		정확성	정확함	정확함	정확함	정확함	정확함	정확함
	실제적 부문	포괄성	큰 온톨로지를 구축하는 것이 가능	큰 온톨로지를 구축하는 것이 가능	큰 온톨로지를 구축하는 것이 가능	큰 온톨로지를 구축하는 것이 가능	큰 온톨로지를 구축하는 것이 가능	큰 온톨로지를 구축하는 것이 가능
		관련성	관련성	관련됨	관련됨	관련됨	관련됨	관련됨
	알고리즘 선택 및 확장 가능성	확장성	높음	보통	높음	높음	높음	높음
	신뢰성	신뢰성	높음	높음	높음	높음	높음	높음
	구문적 부문	합문법성	문법에 맞음	문법에 맞음	문법에 맞음	문법에 맞음	문법에 맞음	문법에 맞음
	의미적 부문	명료성	명료함	명료함	명료함	명료함	명료함	명료함
		정확성	정확함	정확함	정확함	정확함	정확함	정확함
실제적 부문	포괄성	큰 온톨로지를 구축하는 것이 가능	큰 온톨로지를 구축하는 것이 가능	큰 온톨로지를 구축하는 것이 가능	큰 온톨로지를 구축하는 것이 가능	큰 온톨로지를 구축하는 것이 가능	큰 온톨로지를 구축하는 것이 가능	
	관련성	관련성	관련됨	관련됨	관련됨	관련됨	관련됨	

부문의 평가에 있어서 TERMINAE에 의해 구축된 온톨로지는 평가하지 않았다. TERMINAE는 반자동의 도구로, 전문가나 온톨로지스트가 구축하는 전 과정 동안에 추출된 개념이나 관계를 평가하기 위해 지속적으로 개입하는 것이 필요하다. 따라서 TERMINAE를 사용하여 구축된 온톨로지의 의미적 부문은 도구 자체보다는 인간 개입에 의해 더 많은 영향을 받기 때문에, 평가 결과가 객관적일 수 없다.

마지막으로 실제적 부문의 평가결과 구축된 온톨로지들이 모두 정확하고, 관련성이 높은 것으로 평가되었다. 또한 포괄성 측면에서도 OntoBuilder를 제외하고는 큰 온톨로지를 구축하는 것이 가능했다. 그러나 OntoBuilder는 HTML의 레이블과 필드 이름으로부터 개념을 추출하기 때문에 OntoBuilder에 의해 구축된 온톨로지의 크기는 제한적이라고 평가되었다. 즉 OntoBuilder의 실제적 부문의 포괄성은 한계가 있었다.

실험 결과를 종합한 <표 1>에서 보여지듯이, 실험을 위해 선택된 도구들 중에 Text-To-Onto가 상대적으로 가장 좋은 온톨로지 추출도구라고 판단하는 것이 가능하다. Text-To-Onto는 개념과 그 관계를 자동적으로 추출할 수 있을 뿐만 아니라 다양한 유형의 데이터 형태를 입력 데이터로 사용할 수도 있다. 또한 사용자로 하여금 다양한 알고리즘들 중에서 적절한 알고리즘을 선택하여 사용할 수 있도록 하여, 사용상의 유연성도 높다.

실험결과, OntoLT와 Text-To-Onto는 대용량 데이터를 입력 데이터로 사용하는 경우, 그리고 시간적 여유가 없고 관련 도메인 지식도 없는 경우 사용하는 것이 좋다고 판단되며, TERMINAE는 온톨로지를 구축하는데 충분한 시간을 가지고 있으면서 구축된 온톨로지의 도메인이 사용자와 친

숙한 도메인인 경우 사용하는 것이 더 좋다고 결론 내릴 수 있다.

6. 결론

온톨로지 추출도구는 온톨로지 편집도구와 병합도구와 함께 시맨틱 웹의 발전에 있어 중요한 역할을 한다. 특히 온톨로지 추출도구를 사용하면 전문가의 수작업을 최소화할 수 있고, 여러 전문가들의 작업 결과가 일관성을 가지게 되는 장점이 있다. 본 논문에서는 도구의 기능성과 구축된 온톨로지의 질적인 평가를 포함하고 있는 온톨로지 자동추출 도구들을 평가하기 위한 포괄적인 프레임워크를 제안하고 실제 실험을 통해 본 프레임워크의 활용 가능성을 평가해보았다. 실험결과 온톨로지 추출 도구들은 추출 과정의 자동화, 알고리즘의 확장 등 개선되어야 할 부분이 여전히 많다고 판단되었다.

본 연구에서 제안한 프레임워크를 사용하여 온톨로지 추출도구들을 평가하는 실험에 있어서 몇 가지 한계가 존재한다. 우선 본 연구에서는 여러 온톨로지 추출도구가 개발되어 있음에도 불구하고 대부분의 도구들은 접근하여 사용할 수 있도록 공개되지 않았기 때문에 단지 이용 가능한 네 개의 도구들만을 사용하여 실험을 했기 때문에 현재 개발된 온톨로지 추출도구들의 특성을 완전히 파악하고 이해하는데 한계가 있었다. 그러나 이는 여러 저널에서 효과적인 추출 성능을 가진 것으로 보고되고 있는 다른 도구들이 추후 이용 가능하면 추가적인 실험을 통해 보완할 수 있는 부분이다. 또한 평가된 모든 도구들이 완전히 자동적이라면, 각 도구를 사용하여 구축된 온톨로지의 질적인 부분 역시 객관적으로 평가하는 것이 가능했을 것이다. 그러나 여전히 몇몇 도구들이 반자동적이기

때문에 구축된 온톨로지가 사람의 개입에 의해 영향을 받았으므로 도구를 사용하여 구축된 온톨로지의 질을 객관적으로 평가하는데 한계가 있다.

현재 여러 대학과 연구소를 중심으로 온톨로지 추출 도구를 개발하기 위해 노력해왔다. 하지만 개발된 온톨로지 병합도구나 편집도구와 비교하면 개발된 온톨로지 추출도구들의 수나 그 활용도는 아직 미미하다. 따라서 온톨로지 추출도구들을 개발하기 위한 더 많은 노력이 필요하다. 특히 국내의 경우, 인터페이스까지 갖춘 완벽한 시스템의 형태로 개발된 온톨로지 추출도구로 2008년 현재 CoreOnto⁶⁾라는 추출도구가 개발되었기는 하나, 온톨로지의 보다 적극적인 개발과 응용을 위해서는 온톨로지 개발 프로세스를 효율화해 주는 온톨로지 추출도구의 개발에 보다 많은 관심과 노력이 필요하다. 온톨로지 추출 도구를 개발함에 있어서 온톨로지 추출도구가 갖춰야 할 기능적 특성에 대해 쉽게 파악할 수 있다면 개발자들은 온톨로지 추출도구의 개발을 보다 쉽게 착수할 수 있을 것이다. 이때 본 연구에서 제안하고 있는 평가지표들을 온톨로지 추출도구를 개발할 때 일종의 유용한 벤치마크나 가이드라인으로 활용할 수 있을 것이다.

참고문헌

- 구미숙, 황정희, 류근호, 홍장의, “데이터마이닝 기법을 이용한 XML 문서의 온톨로지 반자동 생성”, *정보처리학회논문지D*, 13권, 3호(2006), 299~308.
- 송도규, “대용량 OWL 온톨로지 자동구축을 위한 세종전자사전 활용 방법론 연구”, *언어와 정보*, 9권, 1호(2005), 19~34.
- 윤현주, 김영민, 이상준, 변영철, “규칙기반 온톨로지 자동생성 및 검색”, *한국정보과학회 학술 발표논문집*, 31권, 2호(2004), 655~657.
- 정도현, 김태수, “시소러스를 기반으로한 온톨로지 시스템 구현에 관한 연구”, *정보관리학회지*, 20권, 3호(2003), 155~176.
- Angele, J. and Y. Sure, “Evaluation of ontology-based tool”, In *Proceedings of the 13th International Conference on Knowledge Engineering and Knowledge Management*, 2002.
- Biebow, B. and S. Szulman, “TERMINAE : a linguistic-based tool for the building of a domain ontology”, In *Proceedings of the 11th European Workshop on Knowledge Acquisition, Modeling and Management*, 1999, 49~66.
- Bisson, G., C. Nedellec and D. Canamero, “Designing clustering methods for ontology building : The Mo’K Workbench”, In *proceedings of the Workshop on Ontology Learning, 14th European Conference on Artificial Intelligence(ECAI)*, 2000.
- Brank, J., D. Mladenic and M. Grobelnik, “Gold standard based ontology evaluation using instance assignment”, In *Proceedings of the 4th International Workshop on Evaluation of Ontologies for the Web(EON) at the 15th International World Wide Web Conference*, 2006.
- Buitelaar, P., D. Olejnik and M. Sintek, “A Protégé Plug-In for Ontology Extraction from Text Based on Linguistic Analysis”, In *Proceedings of the 1st European Semantic*

6) <http://cscola.kaist.ac.kr/wiki/>

- Web Symposium(ESWS)*, 2004.
- Burton-Jones, A., V. Storey, V. Sugumaran and P. Ahluwalia, "A semiotic metrics suite for assessing the quality of ontologies", *Data & Knowledge Engineering*, Vol.55(2005), 84~102.
- Chaelandar, G. and B. Grau, "SVETLAN⁷-a system to classify words in context", In *Proceedings of the Workshop on Ontology Learning, 14th European Conference on Artificial Intelligence(ECAI)*, 2000.
- Chidamber, S. and C. Kemerer, "A metrics suite for object-oriented design", *IEEE transactions on software engineering*, Vol.20, No.6 (1994), 476~493.
- Duineveld, A., R. Stoter, M. Weiden, B. Kenepa and V. Benjamins, "Wondertools? A comprehensive study of ontological engineering tool", In *Proceedings of the 12th Workshop on Knowledge Acquisition, Modeling and Management*, 1999.
- Faure, D. and C. Nedellec, "Knowledge Acquisition of Predicate Argument Structures from Technical Texts Using Machine Learning : The System ASIUM", In *Proceedings of the 11th European Workshop(EKAK), LNAI 1621*, Springer-Verlag, 1999, 329~334.
- Faure, D and C. Nedellec, "A corpus-based conceptual clustering method for verb frames and ontology acquisition", In *LREC workshop on adapting lexical and corpus resources to sublanguages and applications*, 1998.
- Gal, A. and G. Modica, "OntoBuilder : fully automatic extraction and consolidation of ontologies from Web sources", In *Proceedings of the 20th International Conference on Data Engineering(ICDE)*, 2004.
- Gomez-Perez, A. and D. Manzano-Macho, "An overview of methods and tools for ontology learning from texts", *The Knowledge Engineering Review*, Vol.19, No.3(2005), 187~212.
- Hevner, A., S. March, J. Park and S. Ram, "Design Science in information systems research", *MIS Quarterly*, Vol.28, No.1(2004), 75~106.
- Lambrix, P. and A. Edberg, "Evaluation of ontology merging tools in bioinformatics", In *Proceedings of the Pacific Symposium on Biocomputing*, Vol.8(2003), 589~600.
- Maedche, A. and R. Volz, "The Text-To-Ontology Extraction and Maintenance Environment", In *Proceedings of the ICDM Workshop on integrating data mining and knowledge management*, 2001.
- Maedche, A. and S. Staab, "Discovering conceptual relations from text", In *Proceedings of ECAI-2000*, 2002.
- Murshed, A. and R. Singh, "Evaluation and ranking of ontology construction tools", Technical Report DIT-05-013, University of Trento, 2005.
- Salton, G., *Automatic Text Processing*, Addison-Wesley, 1988.
- Srikant, R. and R. Agrawal, "Mining generalized association rules", *Future Generation Computer Systems*, Vol.13(1997), 161~180.
- Velardi, P., R. Navigli and M. Missikoff, "Integrated approach for Web ontology learning

- and engineering”, *IEEE Computer*, November, 2002.
- Wand, Y. and R. Weber, “On the deep structure of information systems”, *Journal of Information Systems*, Vol.5(1995), 203~223.
- Wu, S. and W. Hsu, “SOAT : a semi-automatic domain ontology acquisition tool from Chinese corpus”, In *Proceedings of the 19th International Conference on Computational Linguistics*, 2002.

Abstract

Measurement Criteria for Ontology Extraction Tools

Jinsoo Park^{*} · Wonchin Cho^{**} · Sangkyu Rho^{*}

The Web is evolving toward the Semantic Web. Ontologies are considered as a crucial component of the Semantic Web since it is the backbone of knowledge representation for this Web. However, most of these ontologies are still built manually. Manual building of an ontology is time-consuming activity which requires many resources. Consequently, the need for automatic ontology extraction tools has been increased for the last decade, and many tools have been developed for this purpose. Yet, there is no comprehensive framework for evaluating such tools. In this paper, we proposed a set of criteria for evaluating ontology extraction tools and carried out an experiment on four popular ontology extraction tools (i.e., OntoLT, Text-To-Onto, TERMINAE, and OntoBuilder) using our proposed evaluation framework. The proposed framework can be applied as a useful benchmark when developers want to develop ontology extraction tools.

Key Words : Ontology, Automatic Ontology Construction, Ontology Extraction Tool, Semantic Web, Evaluation Framework

* Graduate School of Business, Seoul National University

** College of Business Administration, Seoul National University

저자 소개



박진수

현재 서울대학교 경영전문대학원/경영대학 부교수로 재직 중이다. The University of Arizona에서 경영정보시스템을 전공하여 경영학 박사를 취득했으며, University of Minnesota의 Carlson School of Management에서 조교수, 고려대학교 경영대학에서 조교수를 역임했다. 주요 관심분야는 온톨로지, 정보시스템 통합, 지식 공유, 시맨틱 웹, 시맨틱 모델링, 웹 정보시스템 등이 있다. MIS Quarterly, IEEE Transactions on Knowledge and Data Engineering (TKDE), IEEE Computer, ACM Transactions on Information Systems (TOIS), Information Systems Frontiers, Communications of the AIS, Journal of Global Information Technology Management (JGITM), International Journal of Electronic Business, Information Systems Review, 한국경영정보학연구, 한국전자거래학회지 등 국내외 우수 학술지에 다수의 논문을 게재하였으며, 대표 저서로는 ‘인터넷 진화의 열쇠 온톨로지 : 웹 2.0에서 3.0.’가 있다.



조원진

현재 서울대학교 경영대학에서 경영정보학 전공으로 박사과정에 재학 중이다. 아주대학교 경영학부를 졸업하고 동 대학원 경영학과에서 경영정보를 전공했으며, 데이터마이닝으로 석사학위를 취득하였다. 주요 관심분야는 시맨틱 웹과 온톨로지, 온톨로지추출도구, 데이터마이닝 등이다.



노상규

서울대학교 경영학과를 졸업하고 미국 미네소타 대학에서 MBA 및 경영학 박사학위를 취득하였으며 현재 서울대학교 경영전문대학원/경영대학 교수로 재직 중이다. 주요 연구분야로는 인터넷 비즈니스, 온톨로지, 정보시스템 모델링, 데이터마이닝 등이 있으며 IEEE Transactions on Knowledge and Data Engineering, Strategic Management Journal, Long Range Planning, Annals of Operations Research, Journal of Database Management 등 우수 학술지에 다수의 논문을 게재하였다. 저서로는 ‘인터넷 진화의 열쇠 온톨로지 : 웹 2.0에서 3.0,’ ‘한국온라인게임산업의 발전과정과 향후과제’ 등이 있다.