

구문분석기의 어휘확장

김민찬, 김곤, 배재학
울산대학교 컴퓨터·정보통신공학부
e-mail:{tomatuli, gonkim, jhjbae}@ulsan.ac.kr

Lexical Expansion of Sentence Parsers

Min-Chan Kim, Gon Kim, Jae-Hak J. Bae
School of Computer Engineering and Information Technology,
University of Ulsan

요 약

본 논문에서는 구문분석기의 어휘확장을 통해 구문분석의 성공률을 높이고자 하였다. 구문분석은 문장 내 구성성분들이 가지는 통사적인 관련성을 파악하는 작업이다. 구문분석 실패의 가장 빈번한 원인 중의 하나는 미등록 어휘의 출현이다. 결여된 어휘문제를 해결하는 것은 구문분석의 성공률을 높이고, 원문 이해 시스템을 보다 더 견고하게 하는데 관건으로 작용한다. 이를 위하여, 본 논문에서는 구문분석기 LGPI+의 어휘 사전에 존재하지 않는 단어들을 또 다른 어휘자원인 WordNet을 이용하여 해결하고자 하였다. 구체적으로는, (1) 미등록 어휘를 WordNet에서 찾고, (2) 그 유의어 정보를 파악하여, (3) LGPI+ 어휘사전에 추가한다. 실험을 통하여 구문분석의 실패를 해결하고, 정확도와 성공률을 높일 수 있음을 확인하였다.

1. 서론

자연어처리 분야에서는 어휘자체에 대한 의미분석에서부터 구문분석, 의미해석으로 이어지는 단어, 문장, 글의 의미분석에 대한 연구가 계속되고 있다. 문장의 통사구조분석과 의미분석은 완전히 독립적인 단계에서 이루어 질 수 없다. 즉, 문장의 구성성분들이 가지는 문법적인 기능과 내포하는 의미는 서로 의존적인 관계에 놓여 있다. 따라서, 구문분석의 정확도와 성공률 향상은 의미분석의 중의성을 해소하고, 원문 이해 시스템을 보다 더 견고하게 할 수 있다.

구문분석은 문장의 구성성분들이 가지는 통사적인 관련성을 파악하는 작업이다[1]. 이러한 구문분석의 오류를 야기하는 주된 원인으로서는 알려지지 않은 어휘가 문장에서 나타나는 경우이다. 미등록 어휘문제

를 해결하기 위한 연구로는, (1) LDOCE(Longman Dictionary of Contemporary English)를 활용하여, 계층적 구문 유형의 규칙들로부터 어휘의 의미정보를 파악하는 방법[2], (2) 특화된 전문용어의 선택적인 어휘 의미들 중에서 선별하는 방법[3] 등이 있다.

기존의 연구방법들은 미리 구축된 어휘사전에서 규칙 및 의미정보 분류를 활용한 방법이다. 본 논문에서는 이와는 달리, 문장에서 나타난 어휘를 찾을 수 없을 때, 또 다른 어휘자원을 활용하여 해결하는 방법을 모색해 보았다.

2. 어휘확장 알고리즘

구문분석기는 분석에 필요한 어휘 사전을 가지고 있다. 포함된 어휘의 수는 제한적이며, 정확한 구문 분석을 위해서는 지속적인 어휘확장이 필요하다. 본

분석 시 사전에 등록되어 있지 않은 단어의 목록을 나타낸다. ③은 ②에서 나타난 어휘들에 대하여, WordNet 유의어를 검색하고, 찾은 유의어들을 LGPI+ 어휘사전에서 확인한 후, 사전에 등록하는 기능이다.

4. 결론

LGPI+의 어휘사전은 각 어휘별 품사 및 특성에 따라 분류된 파일로 구성되어 있다. 따라서, 어휘확장 시 정확한 파일과 위치에 기록되어야 구문분석의 성공률과 정확도가 높아진다.

그림 2의 homemade가 사전에 없을 경우에는 28개의 linkages를 보이고, 그림 7의 사전에 추가한 경우에는 19개의 linkages를 보여준다. 이를 통하여, 어휘정보를 추가하였을 경우 linkages의 수가 감소함을 알 수 있다.

본 논문에서는 구문분석 시 미등록된 어휘를 확장하는 방법으로 WordNet의 유의어 정보를 이용하였다. 실험을 통해 구문분석기의 성공률과 정확도를 높일 수 있음을 알 수 있었다. 정확한 구문분석은 언어번역, 정보검색, 문법분석, 음성처리, 원문처리 및 원문이해 등 자연언어처리 분야의 전반적인 연구 분야뿐만 아니라 전문용어 및 다양한 전문가 시스템이나 학제 간 연구의 중요한 기초자료가 될 수 있을 것으로 기대된다.

<Acknowledgements>

본 연구는 한국과학재단 목적기초연구 R05-2004-000-12362-0 지원으로 수행되었음. 또한 본 연구는 정보통신부 및 정보통신연구진흥원의 대학 IT연구센터 육성·지원사업의 연구결과로 수행되었음.

[참고문헌]

- [1] 김곤, 배재학. "문서요약을 위한 문장추상화." 한국정보처리학회 춘계 학술대회 논문집, 제 9권, 제 1호, pp.531-534, 2002.
- [2] Hiyani Alshawi. "Processing Dictionary Definitions with Phrasal Pattern Hierarchies." Computational Linguistics, Volume 13, Numbers 3-4, 1987.
- [3] Walker, D. and Amsler, R. "The use of Machine Readable Dictionaries in Sublanguage Analysis." In Analyzing Language in Restricted Domains, Lawrence Erlbaum Associates, 1986.
- [4] SWI-Prolog. <http://www.swi-prolog.org/>.
- [5] Sleator, D. and Temperley, D. "Parsing English with a Link Grammar." Third International Workshop on Parsing Technologies, August 1993.
- [6] WordNet. <http://wordnet.princeton.edu/>.
- [7] 김곤, 김민찬, 배재학. 이종혁, "ISAAC: 문장 분석용 통합시스템 및 사용자 인터페이스", 한국정보처리학회 논문지B, 제11-B권, 제 1호, pp. 107-116, 2004.