

개최장소 추출을 위한 LGG의 구축

김경렬[○], 최동현, 김은경, 최기선
한국과학기술원, 시맨틱웹 첨단연구센터
{barnabas, cdh4696, kekeeo, kschoi}@world.kaist.ac.kr

Construction of LGG for Extracting Meeting Location

Kyoung-Ryol Kim[○], Dong-Hyun Choi, Eun-Kyung Kim, Key-Sun Choi
Semantic Web Research Center, KAIST

요 약

본 논문에서는 회의공지 이메일을 대상으로 하는 개최장소 추출시스템에 대하여 소개한다. 개최장소 추출시스템은 두 단계로 구성되는데, 첫 번째 단계는 본문에 포함된 개최장소의 추출이고, 두 번째 단계는 추출된 개최장소의 Geocoding이다. 개최장소의 추출을 위하여 문맥 패턴을 분석하여 개최장소가 포함된 문장 주변의 패턴을 반영하는 Local-Grammar Graph를 구축하며, 개최장소의 Geocoding을 위하여는 Addr2Geocode API를 사용한다. 본 논문은 일정공지메일의 개최장소를 추출하기 위한 LGG 방법론 기반의 어휘-통사적 언어 정보를 기술하는 것을 목적으로 한다.

주제어: Information Extraction, Geocode, Local-Grammar Graph

1. 서론

최근 애플의 iPhone OS, 구글의 Android와 같은 스마트폰 운영체제에서는 일정관리를 효과적으로 할 수 있도록 여러 가지 정보인식 및 추출 기술을 제공하기 시작했다. 예를 들어, iPhone에서는 이메일에 간단한 시간 표현이 발견되면 자동으로 이벤트를 생성하여 캘린더에 추가할 수 있는 링크를 생성해주며, 간단한 주소에 대하여서도 자동으로 인식하여 지도와 연결되는 링크를 생성해 주고 있다[1].

본 논문에서는 회의공지 이메일을 대상으로, 본문에 포함되어 있는 회의정보를 추출하여 캘린더에 자동으로 추가하는 회의정보 추출시스템의 모듈인 개최장소 추출시스템에 대하여 소개한다. 개최장소 추출시스템은 두 단계로 구성되는데, 첫 번째는 본문에 포함된 개최장소를 추출하는 단계이고, 두 번째 단계는 추출된 개최장소의 Geocoding이다. 본 논문은 LGG(Local-Grammar Graph) 방법론에 기초하여 어휘-통사적 언어 정보를 기술하고, 이로부터 유한상태 변환기(Finite-State Transducer)를 구축하여 본문에 포함된 개최장소를 추출한다.

회의공지를 대상으로 회의정보를 추출하는 연구는 영어권에서 활발히 진행되었는데 CMU에서 제작된 Seminar Announcement Corpus[2]가 대표적이다. 규칙을 기반으로 추출하는 연구가 꾸준히 연구되어 왔고 [3][4][5][6], 최근에는 통계 모델의 학습을 통하여 일정정보를 추출하는 연구가 주류를 이루게 되었다. Hidden Markov Models을 이용한 연구[7][8][9], Maximum Entropy Models기반 연구[10][11][12], Conditional Random Fields를 이용한 연구들이 [13][14][15] 대표적이다. 각 방법론은 데이터가 가지

는 특성에 따라 다른 성능을 얻을 수 있어, 주어진 문제에 적절한 방법론을 이용하는 것이 필요하다. 영어권의 다양한 연구결과에 비하여 한국어를 대상으로 일정정보를 추출하는 연구는 상대적으로 많이 부족하다. 공개된 일정공지 말뭉치가 존재하지 않아 연구자들이 직접 말뭉치를 수집하여 연구해왔으며, 연구의 수도 많지 않다.

본 논문에서는 개최장소의 추출문제에 적절한 방법론을 찾기 위하여, 회의공지이메일에 포함된 개최장소 주변의 어휘-통사적 언어 정보를 분석 및 기술하여 LGG를 구축한다. LGG로부터 유한상태 변환기를 얻을 수 있는데, 이를 이용하여 개최장소를 추출하는 방법을 설명한다. 추출된 개최장소 정보로부터 Geocode를 얻는 과정도 간략히 소개한다.

본 논문은 다음과 같이 구성된다. 2장에서는 LGG방법론과 이를 작성하는 도구인 UNITEX에 대하여 소개하고, 3장에서는 개최장소 추출방법에 대하여 설명하고, 4장에서는 추출된 개최장소의 Geocoding 방법을 설명한다. 5장에서는 실험 및 결과에 대하여 논의하고, 6장에서 결론 및 추후 연구과제에 관하여 논한다.

2. LGG와 UNITEX

2.1. Local-Grammar Graph

Local-Grammar Graph는 프랑스의 전산언어학자인 모리스 그로스에 의해 제안된 언어 기술 모델로써, 특정 영역별로 부분적인 언어 정보를 유한상태 오토마타(Finite-State Automata) 문법의 형태로 구현하고 이를 이용하여 자연 언어 텍스트에 대한 자동 분석 및 생성, 정보 추출 등을 수행하는 것을 목적으로 한다. 언어 지식을 형식화하는 문법을 최대한 어휘화함으로써 시스템

3가지로 분류되는데, 첫 번째는 1개의 개최장소를 포함하는 패턴, 두 번째는 N(>1)개의 개최장소를 포함하는 패턴, 세 번째는 집결장소, 예상 장소, 장소 미정과 같이 개최장소이지만 속성이 개최장소와 차이가 있는 장소를 포함하는 패턴이다. 각 분류 아래에는 해당 패턴으로 추출가능한 장소의 타입별로 구분된 목록을 가진다. 표 4는 LGG의 분류를 나타낸다. 개최장소1_1, 개최장소1_2과 같이 표현된 형태는 1개의 개최장소의 일부 장소정보가 2개로 분리되어 작성된 경우를 표현한다. 예를 들어, 표 3과 같이, '무역협회 중회의실'과 '삼성동 트레이드 타워 51층'은 같은 장소를 나타내고 있지만, 예문에서는 굳이 괄호로 구분을 하고 있다.

표 3. 개최장소 정보가 2개로 분리되어 있는 예

1. 일시 및 장소 : 2010. 5. 12(수) 14:00~16:00, 무역협회 중회의실 (삼성동 트레이드 타워 51층)

표 4. 개최장소 LGG의 분류

- 1. 개최장소 1개
 - 1.1. 개최장소(국문명)
 - 1.1.1. 개최장소
 - 1.1.2. 개최장소1_1 | 개최장소1_2
 - 1.2. 개최장소(국문명) + 주소
 - 1.2.1. 개최장소 | 주소
 - 1.2.2. 개최장소1_1 | 개최장소1_2 | 주소
 - 1.3. 개최장소(국문명) + 랜드마크
 - 1.3.1. 개최장소 | 랜드마크
 - 1.3.2. 개최장소1_1 | 개최장소1_2 | 랜드마크
 - 1.3.3. 개최장소 | 랜드마크1 | 랜드마크2
 - 1.3.4. 개최장소 | 랜드마크 | 주소
 - 1.4. 개최장소(국문명, 영문명)
 - 1.4.1. 개최장소 | 개최장소(영문)
 - 1.4.2. 개최장소1_1 | 개최장소 (영문) | 개최장소1_2
 - 1.4.3. 주소 | 개최장소1_1 | 개최장소(영문) | 개최장소1_2
- 2. 개최장소 N개 (N>1)
 - 2.1. 개최장소 2개
 - 2.2. 개최장소 3개
 - 2.3. 개최장소 4개
- 3. 기타 개최장소
 - 3.1. 집결장소
 - 3.1.1. 집결장소
 - 3.1.2. 개최장소 | 개최장소(예전이름) | 집결장소
 - 3.2. 예상장소
 - 3.2.1. 예상장소
 - 3.2.2. 예상장소1 | 예상장소2
 - 3.3. 장소미정

특별히 표 4의 분류 1.2, 1.3.4를 이용하여 개최장소

와 주소를 함께 추출할 수가 있다. 예를 들어, 표 5의 개최장소는 '울산광역시 울주군 상북면 등역리 27번지'라는 주소에 위치한 '먹고쉬었다가'라는 음식점이며, 표 4의 분류 1.2.1에 해당하는 예이다. 이처럼 개최장소의 주소를 함께 추출할 수 있는 경우는 4장에서 소개할 Addr2Geocode API를 통하여 바로 Geocode를 얻을 수 있다.

표 5. 개최장소와 주소가 함께 공지된 예

3. 장 소 : 울산광역시 울주군 상북면 등역리 27번지
먹고쉬었다가 (052-263-1206)

4. 개최장소의 Geocoding

Geocoding이란, 토지 내 중심점의 지리적 좌표로서 토지를 구분하는 방법으로 특정지도 투영법에 의해 지표상의 위치를 X, Y 좌표로 나타내는 방법이다 [17]. 이를 표현하는 좌표계의 종류와 제정된 표준의 수는 상당히 다양하지만, WGS84, TM128등 몇 가지가 사실상 표준으로 채택되어 Google Maps, Naver 지도뿐 아니라 OpenStreetMap 등의 웹기반 지도서비스에서 사용되고 있다. 본 논문에서는 최근 가장 많이 사용되고 있는 WGS84를 사용한다.

개최장소를 Geocoding하는 방법은 크게 2가지로 나뉜다. 첫째는 주소정보가 포함되어 있는 경우로 Addr2Geocode API를 사용하며, 둘째는 주소정보가 포함되지 않은 경우로 외부 지리정보자원에 장소명으로 검색하는 방법이다. 본 논문에서는 첫 번째 방법 중에서도 번지수까지 포함하는 주소를 가지는 경우를 다룬다. 주소정보의 일부만 포함하는 경우와 두 번째 방법은 본 논문에서 다루지 않는다.

주소의 Geocoding을 위하여 Daum에서 제공하는 Addr2Geocode API를 사용하였다. API의 입력은 번지를 포함하는 주소 문자열이며, 출력으로 해당주소의 WGS84 경위도 좌표를 RSS/XML/JSON 형태로 반환한다.

5. 실험 및 결과

본 논문에서 사용한 회의공지 이메일 말뭉치는 인터넷을 통하여 '공지'라는 검색어로 1,011개의 이메일을 수집하였다. 3명의 서로 다른 어노테이터가 어노테이션 작업을 수행하였으며, 어노테이션 간 충돌이 발생하였을 때에는 또 다른 컨주게이터가 충돌 해소 작업을 진행하였다.

그 중에서 순서대로 선택된 555개의 문서에 대하여 LGG를 구축하여 실험을 진행하였다. 39개의 LGG가 작성되었으며, 그 중 7개는 개최장소와 주소를 함께 추출

할 수 있는 형태이고 나머지 32개는 개최장소만을 대상으로 추출하는 형태이다.

표 6은 작성된 LGG를 적용하여 추출된 결과이다. 'Exact'는 어노테이션된 개최장소와 시스템이 추출한 개최장소가 완전히 일치한 결과이고, 'Contain'은 시스템이 추출한 결과가 어노테이션된 개최장소를 포함하는 경우의 결과이다. Exact, Contain에 대하여 각각 93.41%, 99.41%의 높은 Recall을 보였으나 Precision에 대하여는 82.11%, 87.39%로 상대적으로 낮은 수치를 보였다. Precision이 낮은 원인으로서는 LGG가 지나치게 일반화된 경우, 조사가 개최장소 뒤에 붙어있는 경우 등을 찾을 수 있었다. 특별히, 조사가 분리된 명사를 분석하기 위하여 필요한 한국어 사전이 UNITEX에서 요구하는 방식으로 구현되어야 한다. 하지만, 라이선스의 문제로 인하여 사용이 어렵기 때문에 조사가 개최장소 뒤에 붙어 있는 경우는 후처리를 통하여 조사를 제거해주는 방법이 추가되어야 한다.

또한, 표 7은 개최장소-주소를 함께 추출 가능한 경우의 패턴을 적용하여 추출된 결과를 보여준다. 빈지수를 포함한 주소를 대상으로 하였기 때문에 그 수가 많지 않다. 19개 중 10개만이 개최장소와 주소가 함께 인식되었는데 에러를 분석하여 보면, 주소와 개최장소의 문서 내 위치가 상당히 떨어져있어 Local-Grammar만으로는 처리가 어려운 경우들이었다.

표 6. 39개 Path의 LGG로 555개의 문서에 대한 개최장소 추출 결과

	Exact	Contain
Relevant	683	
Retrieved	777	
Ret. & Rel.	638	679
Recall	93.41%	99.41%
Precision	82.11%	87.39%
F-measure	87.40%	93.01%

표 7. 개최장소-주소를 함께 추출 가능한 7개 Path의 LGG로 555개의 문서에 대한 추출 결과

	Exact
Relevant	19
Retrieved	10
Ret. & Rel.	10
Recall	47.62%
Precision	100.00%
F-measure	64.52%

6. 결론

본 논문에서는 회의공지 이메일을 대상으로 개최장소를 추출하는 문제를 해결하기 위하여 LGG를 구축하는 방법을 소개하였다.

555개의 이메일에 포함된 개최장소를 추출하기 위하여 39개의 패턴만이 사용되어 Exact, Contain 일치에 대하여 각각 F-measure 87.40%, 93.01%의 높은 추출 성능을 보였다. 세밀한 일반화 및 후처리 작업 등을 통하여 Precision이 더 상승할 수 있을 것으로 기대한다. 또한, 이러한 어휘-통사론적 분석 결과는 추후 다른 방법론을 적용하기 위한 근거자료로써 활용도 가능할 것으로 예상된다.

추후 연구할 과제는, 충분히 분석된 LGG 결과를 바탕으로 개최장소 온톨로지를 구축하여 정보 타입 간의 온톨로지 추론을 통하여 개최장소 추출 및 주소 매핑을 수행하는 연구이다. 이에 앞서, 수집된 이메일의 50% 정도를 대상으로 실험하였기 때문에 나머지 데이터에 대한 LGG 구축과 일반화 작업이 필요하다. 뿐만 아니라, 본 연구는 한국어 개최장소만을 대상으로 실험을 수행하였지만, 특정언어에 의존적인 방법론이 아니기 때문에 추후 영어 등의 다른 언어에 대하여도 같은 방법으로 실험하여 언어에 독립적인 방법론으로 일반화 해볼 수 있다.

감사의 글

본 논문은 지식경제부 산업원천기술개발사업(10035348, 모바일 플랫폼 기반 계획 및 학습 인지 모델 프레임워크 기술 개발)의 지원으로 수행되었음

참고문헌

- [1] Mail Data Detection, Apple Inc., 2008, <http://itunes.apple.com/podcast/mac-quick-tips/id257243321>
- [2] Seminar Announcements Dataset, CMU, <http://www.cs.cmu.edu/~dayne/SeminarAnnouncements/>
- [3] Muggleton, S., Buntine, W., "Machine invention of first-order predicates by inverting resolution", Proceedings of the 5th International Conference on Machine Learning, 1988
- [4] Riloff, E., "Automatically constructing a dictionary for information extraction tasks", Proceedings of the 11th National Conference on Artificial Intelligence (AAAI), 811-816, 1993
- [5] Kim, J., Moldovan, D., "Acquisition of linguistic patterns for knowledge-based information extraction", IEEE Transaction Knowledge Data Engineering, 1995
- [6] Chai, J., Biermann, A., Guinn, C., Two dimensional generalization in information extraction, Proceedings of the 16th AAAI National

Conference on Artificial Intelligence (AAAI), 1999

[7] Seymore, K., Mccallum, A., Rosenfeld, R., "Learning hidden Markov model structure for information extraction", Proceedings of the 16th AAAI National Conference on Artificial Intelligence (AAAI), 1999

[8] Freitag, D., Kushmerick, N., "Boosted wrapper induction", Proceedings of the ECAI Workshop on Machine Learning for Information Extraction, 2000

[9] Ray, S., Craven, M., "Representing sentence structure in hidden Markov models for information extraction", Proceedings of the 17th International Joint Conference on Artificial Intelligence (IJCAI01), 2001

[10] Chieu, H. L., Ng, H. T., "A maximum entropy approach to information extraction from semistructured and free text", Proceedings of the 18th National Conference on Artificial Intelligence (AAAI), 2002

[11] Kambhatla, N., "Combining lexical, syntactic, and semantic features with maximum entropy models for extracting relations", Proceedings of the 42nd Annual Meeting of the Association for Computational Linguistics (ACL), 2004

[12] Turmo, J., Ageno, A., Catala, N., "Adaptive information extraction", Journal of ACM Computing Surveys (CSUR), Vol.38-2, 2006

[13] Lafferty, J., Mccallum, A., Pereira, F., "Conditional random fields: probabilistic models for segmenting and labeling sequence data", Proceedings of the 18th International Conference on Machine Learning (ICML), 2001

[14] Cox, C., Nicolson, J., Finkel, J., Manning, C., Langley, P., "Template sampling for leveraging domain knowledge in information extraction", First PASCAL Challenges Workshop, 2005

[1] Lee, C. K., Hwang, Y. G., Oh, H. J., Lim, S. J., Heo, J., Lee, C. H., Kim, H. J., Wang, J. H., Jang, M. G., "Fine-Grained Named Entity Recognition Using Conditional Random Fields for Question Answering," Proceedings of the 30th annual international ACM SIGIR conference on Research and development in information retrieval, 2007

[15] 남지순, "웹문서 의미 지식 추출을 위한 LGG의 구축", 프랑수어문교육 제25집, 105-112, 2007

[16] 강영욱, 이동연, "위치정보제고를 위한 주소표시제도 개선방안", 국토계획 제31권 제6호, 5-314, 1996

[17] Google Maps, <http://maps.google.com/>

[18] Naver 지도, <http://map.naver.com/>

[19] OpenStreetMap, <http://www.openstreetmap.org/>

[20] Daum 주소좌표변환 API, <http://dna.daum.net/apis/maps/reference>

부록 1. 분류별 정보타입

분류	정보 타입	설명
장소	locDst	장소 : 개최예정장소
	locAre	장소 : 지역 (장소+dirWrd 를 포함하는 지역)
	locAdr	장소 : 주소
	locGtr	장소 : 집결장소
	locGnr	장소 : 일반장소
	locLmk	장소 : 랜드마크
	locMtg	장소 : 개최장소
	locOld	장소 : 장소의 이전이름
	locSwy	장소 : 지하철역
	locSlr	장소 : 여러 장소 중 하나로 개최장소 선택예정
	locUnk	장소 : 장소 미정
	locWay	장소 : 거리 (way)
	locCnf	장소 확정 지시어
	locDstCnf	예정 장소 지시어
	locEqv	동일한 장소 표현 접속사 및 기호
	locTyp	장소의 타입 (업종)
swySta	지하철역명	
dirWrd	방향성을 나타내는 단어	
disWrd	거리(distance)를 나타내는 단어	
시간	timDur	시간 : 기간
	timEnd	시간 : 종료시간
	timGnr	시간 : 일반시간 (시작,종료 시간이 아닌 시간표현)
	timStr	시간 : 시작시간
주제	titMtg	모임제목
	tpcMtg	모임주제
행위주제	perLst	사람 목록
	perNam	사람이름
	perTit	사람 직책
	orgNam	기관이름

레이블	adrLbl	주소 레이블	
	agdLbl	안건 레이블	
	ancLbl	진행자 레이블	
	appLbl	신청 레이블	
	brgLbl	준비물 레이블	
	chgLbl	변경정보 레이블	
	dirLbl	오는길 레이블	
	feeLbl	참가비 레이블	
	gtrLbl	모임 집결 레이블	
	hstLbl	주최 레이블	
	inqLbl	문의 레이블	
	lcaLbl	지역 (지역 : 의 형태) 레이블	
	lctLbl	개최장소 및 시간 레이블	
	locLbl	개최장소 레이블	
	nopLbl	인원 레이블	
	mtgLbl	모임 레이블	
	phnLbl	전화번호 레이블	
	prsLbl	발표자 레이블	
	refLbl	참고자료 레이블	
	sclLbl	규모 레이블	
	symEmp	강조 심볼	
	spnLbl	후원 레이블	
	tgtLbl	대상 레이블	
	timLbl	시간 레이블	
	tlcLbl	시간 및 개최장소 레이블	
	tprLbl	교통편 레이블	
	tpcLbl	주제 레이블	
	webLbl	홈페이지 레이블	
	이음 문자열	decPrd	서술성 종결사
		locLmk2locLmk	랜드마크 사이의 문자열
locGnr2locAre		일반장소와 지역사이의 문자열	
locGnr2titMtg		일반장소와 제목 사이의 문자열	
locGnrAft		일반장소 뒷쪽에 나오는 문자열	
locMtg2locGnr		개최장소와 일반장소 사이의 문자열	
locMtg2locMtg		개최장소 사이의 문자열	
locMtg2perNam		개최장소와 사람이름사이의 문자열	
locMtg2phnLbl		개최장소와 전화 레이블 사이의 문자열	
locMtg2timStr		개최장소와 시작시간 사이의 문자열	
locMtg2titMtg		개최장소와 모임제목 사이의 문자열	
locMtg2tpcMtg		개최장소와 모임주제 사이의 문자열	
locMtgAft		개최장소 뒷쪽에 나오는 문자열	
locMtgBef		개최장소 앞쪽에 나오는 문자열	

	locAdr2locMtg	주소와 개최장소 사이의 문자열
	locAdr2timStr	주소와 시작시간 사이의 문자열
	locWay2locMtg	거리(way) 와 개최장소 사이의 문자열
	orgNam2orgNam	기관이름과 기관이름 사이의 문자열
	orgNam2locGnr	기관이름과 일반장소 사이의 문자열
	orgNam2timGnr	기관이름과 일반시간 사이의 문자열
	orgNam2tpcMtg	기관이름과 모임주제 사이의 문자열
	disWrd2locWay	거리를 나타내는 단어와 거리 사이의 문자열
	locSwy2locAre	지하철역과 랜드마크 장소 사이의 문자열
	locSwy2locSwy	지하철역 사이의 문자열
	locWay2dirWrd	거리(way)와 방향을 나타내는 단어 사이의 문자열
	perNam2locAre	사람이름과 지역 사이의 문자열
	perNam2locMtg	사람이름과 개최장소 사이의 문자열
	perNam2titMtg	사람이름과 모임제목 사이의 문자열
	perNam2tpcMtg	사람이름과 모임주제 사이의 문자열
	perNamAft	사람이름 뒷쪽에 나오는 문자열
	perNamBef	사람이름 앞쪽에 나오는 문자열
	phnNbr2adrLbl	전화번호와 주소레이블 사이의 문자열
	timStr2locLbl	시작시간과 개최장소 레이블 사이의 문자열
	timStr2locMtg	시작시간과 개최장소 사이의 문자열
	timStr2timEnd	시작시간과 종료시간 사이의 문자열
	timStr2titMtg	시작시간과 제목 사이의 문자열
	timStrAft	시작시간 뒷쪽에 나오는 문자열
	titMtg2locAdr	모임제목과 주소 사이의 문자열
	titMtg2locMtg	모임제목과 개최장소 사이의 문자열
	titMtg2perNbr	모임제목과 사람이름 사이의 문자열
	titMtg2timStr	모임제목과 시작시간 사이의 문자열
	titMtgBef	모임제목 앞쪽에 나오는 문자열
	tpcMtg2locMtg	모임주제와 개최장소 사이의 문자열
	tpcMtg2titMtg	모임주제와 모임제목 사이의 문자열
tpcMtg2perLst	모임주제와 사람목록 사이의 문자열	
tpcMtgAft	모임주제 뒷쪽에 나오는 문자열	
timEnd2locMtg	종료시간과 개최장소 사이의 문자열	
locAdr2locAdr	주소 사이의 문자열	
locAre2perNam	장소(지역)과 사람이름 사이의 문자열	
기타 정보	webAdr	홈페이지 주소
	phnNbr	전화번호
	phnNbrAft	전화번호 뒷쪽에 나오는 문자열
	addInf	추가정보
-	LF	<Enter>, line feed