

U-WIN의 구문관계 자동구축 방법 1)

임지희¹ 김동명¹ 최호섭² 윤화목² 옥철영¹
울산대학교 컴퓨터정보통신공학과 한국어처리연구실¹
한국과학기술정보연구원 정보기술개발단 정보시스템개발팀²
arisu80@ulsan.ac.kr, spelldm@nate.com,
{hschoe, hmyoon}@kisti.re.kr², okcy@ulsan.ac.kr

Automatic Construction of Syntactic Relation in U-WIN

Jihui Im¹ Dongmyoung Kim¹ Cheolyoung Ock¹
Korean Language Processing Laboratory,
Dept. of Computer Engineering and Information Technology,
University of Ulsan

Hoseop Choe² Hwa-mook Yoon²
Information System Development Team,
Korean Institute of Science and Technology Information

요 약

일반적인 어휘망이 의미 관계에 의한 연결 구조를 중심으로 연구 개발된 것과는 달리, U-WIN은 의미관계를 비롯하여 개념 관계, 형태 관계, 구문 관계 등과 같이 의미 관계의 범위를 확장한 어휘 관계를 적용하여 구축하고 있다. 본 연구에서는 U-WIN의 어휘 관계 중의 하나인 구문관계를 자동으로 구축하는 방법을 제시하고자 한다. 먼저, 용언의 용례에서 문형 정보를 기준으로 구문관계를 형성할 수 있는 후보명사를 추출하였으며, 추출한 후보명사는 용언의 세분화된 의미별로 정확하고 다양하게 추출할 수 있었다. 그러나 U-WIN은 단어의 뜻풀이 하나하나를 개별적인 어휘로 구분하여 구축하였으므로, 어휘 간의 구문관계를 설정하기 위해서는 후보명사의 여러 의미 중에서 하나의 의미로 결정해야 한다. 그래서 본 연구에서는 용례 매칭 규칙, 구문패턴, 의미 유사도 등을 차례로 적용하여 후보명사의 의미를 분별하였으며, 또한 구문패턴의 빈도 정보를 이용하여 용례에 나타나지 않지만 구문관계를 형성할 수 있는 명사를 추출하여 구문관계를 확장하고자 하였다. 이러한 연구는 명사 중심의 어휘망이 용언과의 구문관계 구축을 통해 형태소 분석, 구문 분석, 의미 분석 등에 광범위하게 활용할 수 있는 어휘망의 기반을 다지는 작업이 될 수 있을 것이다.

1. 서 론

최근 의미적 언어자원에 대한 연구 및 관심이 증가하여, 의미부류, 어휘망, 시소러스, 온톨로지 등이 핵심 연구 대상으로 부각되고 있다. 그중에서 어휘망은 한 어휘가 다른 어휘와 가지는 다양한 관계를 망(Network) 형태로 나타내어 데이터베이스화한 것으로, 어휘망의 활용은 언어자원의 효율적인 관리, 의미적 자연언어처리 기술 향상 등의 기대효과가 있다.

본 연구팀이 2002년부터 개발 중인 한국어 사용자 어휘지능망 (User-Word Intelligent Network, 이하 U-WI

N)은 한국어의 공통적이고 개별적인 속성을 바탕으로 한국인의 보편적인 인지 체계와 개념 관계를 파악하여 이를 어휘의 의미적·개념적 네트워크로 형성한 대규모 어휘망이라 할 수 있다. 특히 일반적인 어휘망이 의미 관계에 의한 연결 구조를 중심으로 연구 개발된 것과는 달리, U-WIN은 의미관계를 비롯하여 개념 관계, 형태 관계, 구문 관계 등과 같이 의미 관계의 범위를 확장한 어휘 관계를 적용함으로써 어휘망의 확장적 형태를 모색하고자 한다. 예를 들어 <그림 1>은 동사 ‘먹다’를 중심으로 상하 관계·동의관계의 의미관계, 사동관계·피동관계·방언관계의 형태 관계, 술목관계의 구문관계 등으로 어휘 간의 다양한 관계를 표현하고 있다. 형태관계와 구문관계는 U-WIN을 구성하는 어휘 집합이 모든 품사를 대상으로 함으로써 고려한 어휘 관계로서 형태소 분석 및 구문 분석에 활

1) 이 논문은 2007년 한국과학기술정보연구원에서 시행하는 위탁연구과제로 수행되었다.

측정하였다.

그리고 Wu and Palmer[10]와 Hirst and St. Onge[11]은 각각 계층구조의 깊이, 관계종류를 기준으로 유사도를 측정하였다.

(2) 정보량(Information Content:IC) 기반 측정방법

정보량(Information Content)은 대용량 말뭉치 내 개념의 발생 빈도를 기반으로 MLE(Maximum Likelihood Estimate)방법으로 얻는다. 많은 정보량이 할당된 개념은 특정 주제에 매우 세부적인 개념이고, 적은 정보량이 할당된 개념은 더 일반적인 개념으로 판단할 수 있다.

$$IC(\text{concept}) = -\log(P(\text{concept})) \dots\dots\dots (2)$$

의미 주석 말뭉치가 없을 경우에 개념 발생 빈도는 단어별 빈도를 해당 단어의 동형어의어/다의어 개수로 나누어 할당하거나, 단어별 빈도를 해당 단어의 동형어의어/다의어에 그대로 할당하는 방법을 사용한다.

대표적인 정보량 기반 측정 방법에는 Resnik[9], Jiang and Conrath[7], Lin[8] 등이 있다. Resnik[9]은 정보량을 사용하여 수식(3)에 의해 유사도를 측정한다.

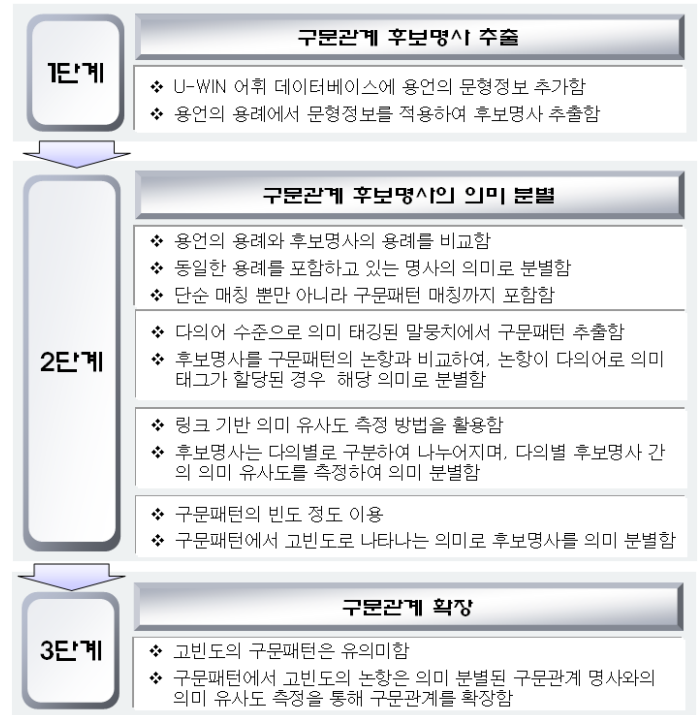
$$Sim_{res}(c_1, c_2) = IC(lcs(c_1, c_2)) \dots\dots\dots (3)$$

$lcs(c_1, c_2)$ 는 개념 c_1 과 c_2 의 공통 상위어 중에서 가장 하위에 위치한 개념³⁾(LCS : lowest common subsumer)를 의미한다. 수식(3)에 의해, 부모 노드가 같은 개념들의 유사도는 최소 공통 상위어가 같아서 항상 같은 값을 가진다. 그러나 주로 계층이 큰 덩어리 형태로 이루어진(coarse-grained) 동사 어휘망은 동일한 최소 공통 상위어를 가지는 개념들이 많으므로, Resnik[9]은 가장 좋은 coarse-grained measure로 알려져 있다.

Jiang and Conrath[7]과 Lin[8]는 각각 Resnik 기반의 명사들의 유사도 측정방법과 문서간의 유사도 측정방법 중에 하나인 Dice Coefficient를 이용한 방법이다.

3. 구문관계 자동구축 방법

명사와 용언 간의 구문관계를 자동으로 구축하는 개괄적인 과정은 <그림 2>와 같다.



<그림 2> U-WIN의 명사-용언 구문관계 자동구축 과정

U-WIN 어휘 사전 데이터베이스에서 문형정보는 구문정보를 제공하며, 용례는 선택 제약이나 결합 관계 등의 정보를 추출할 수 있는 말뭉치 역할을 담당하고 있다. 그러므로 용언의 용례에서 문형정보에 해당하는 논항을 추출하여, 해당 용언과 구문관계를 형성할 수 있는 후보명사로 이용함으로써, 용언의 세분화된 의미별로 다양하고 명확한 후보명사를 추출할 수 있다. 그러나 U-WIN은 다의어를 개별 어휘로 판단하고 있어, U-WIN의 구문관계 자동구축을 위해서는 후보명사를 다의어 수준으로 의미 분별을 해야 한다.

그래서 본 연구에서는 용례 규칙, 의미 주석 말뭉치에서 추출한 구문패턴, 의미 유사도 등을 차례로 적용하여 후보명사의 의미를 분별하였으며, 특히 구문패턴은 다의어 수준으로 의미 태그가 부착된 말뭉치에서 추출함으로써 의미 태그가 부착된 논항은 후보 명사의 의미 분별 및 구문관계 확장에 활용하였다.

3.1 구문관계 후보명사 추출

U-WIN 어휘 사전 데이터베이스는 [표준국어대사전]을 기반으로 <표 1>과 같이 어휘 내적 정보를 설계하여 구축하였으며, 이러한 어휘 내적 정보 구조는 명사를 중심으로 설계되었다.

3) 이후 개념 c_1 과 c_2 의 공통 상위어 중에서 가장 하위에 위치한 개념을 '최소 공통 상위어'라고 한다.

<표 1> U-WIN 어휘 사전 데이터베이스의 어휘 내적 정보 구조

'어휘 내적 정보' 항목	설명
식별자(Identifier : ID)	각 어휘가 가지는 식별자
의미 코드(Sense Code)	각 어휘가 가지는 의미 관리 코드
어휘소(Lexme)	어휘의 형태
의미 표지(Sense Tag)	동형이의어/다의어 정보
일차 분석 어휘소 (Analyzed lexeme)	일차 형태소 분석 결과 및 띄어쓰기 정보
뜻풀이(Definition)	어휘의 뜻풀이 정보
품사(Parts-Of-Speech)	어휘의 품사 정보
한자(Chinese Character)	어휘의 한자 정보
원어(Original Word)	외래어의 원어 정보
대역(Translation)	어휘의 대역 정보
전문분야(Domain)	전문용어 분야 정보
용례(Example)	어휘의 용례 정보
출처(Source)	어휘의 출처 정보

본 연구에서는 후보명사를 추출하는 기준으로 각 용언의 문형 정보를 활용하였다. 문형정보는 문장을 구성하는 데 있어서, 각 용언이 요구하는 필수적 성분들, 즉 논항의 순서 있는 목록으로서, <그림 3>과 같이 U-WIN 어휘 사전 데이터베이스에 문형정보 항목을 추가하고, 용언의 의미별 문형정보를 추출하여 문형정보 항목에 할당하였다.

의미 표지	뜻풀이	문형정보
먹다_001000	키나 코가 막혀서 채 기능을 하지 못하게 되다.	(-을)
먹다_002001	음식 (따위를) 입문 통하여 배 속에 들어보 내다.	-을
먹다_002002	담배나 아편 (따위를) 피우다.	-을
먹다_002003	연기나 가스 (따위를) 들이마시다.	-을
먹다_002004	어떤 마음이나 감정을 품다.	-을
먹다_002005	일정한 나이에 이르거나 나이를 더하다.	-을
먹다_002006	욕, 민간 따위를 듣거나 달하다.	-을
먹다_002007	(속되게) 뇌물을 받아 가지다.	-을
먹다_002008	수익이나 이윤을 차지하여 가지다.	-을
먹다_002009	물이나 습기 따위를 빨아들이다.	-을
먹다_002010	어떤 등급을 차지하거나 흡수를 띠다.	-을
먹다_002011	구기 경기에서, 점수를 잃다.	-을
먹다_002012	(속되게) 여자의 정조를 유린하다.	-을
먹다_002013	매 따위를 맞다.	-을
먹다_002014	남의 재물을 다루거나 맡은 사람이 그 재물을 부당하게 자기의 것으로 만들다.	-을
먹다_002015	날이 있는 두구가 소해를 겪거나 자르거나 갈거나 하는 작동을 하다.	-에
먹다_002016	바르르 물결이 배어들거나 고무 껌지다.	-에
먹다_002017	발레, 군 따위가 꺾어 넘어지거나 꺾지다.	-에
먹다_002018	종이나 빨지 따위가 뭉개나 쓰이다.	-에

<그림 3> 문형정보 추출 및 할당(예-'먹다')

그리고 표준국어대사전 편찬 지침[12]를 살펴보면 “용례는 뜻풀이를 이해하거나 실제로 사용하는 데에 도움을 줄 수 있고 가급적 선택 제약이나 결합 관계 등을 보여주는 전형적인 것이어야 하며, 특히 명사 표제어라면 함께 쓰이는 서술어의 대표 예가, 동사 표제어라면 함께 쓰이는 명사의 대표 예가 충분히 제시될 필요가 있다”고 기술되어 있으므로, 용례는 논항 정보를 추출할 수 있는 대표적인 말뭉치라고 할 수 있다.

또한 용언 표제어의 용례는 제시된 문형정보를 충분히 반영하고 있으므로, 용례에서 문형정보를 기준으로 구문

관계를 형성할 수 있는 후보명사를 추출하였다.

문형정보와 용례를 이용한 후보명사 추출

의미 표지	뜻풀이	문형정보	용례	후보명사
먹다_001000	키나 코가 막혀서 채 기능을 하지 못하게 되다.	(-을)	갑은 주위를 배다 / 주위 (때때로는 구기) 먹어 / 귀를 먹었는지 아무도 몰라도 그냥...	코, 귀
먹다_002001	음식 (따위를) 입문 통하여 배 속에 들어보 내다.	-을	밥을 먹다 / 술을 먹다 / 약을 먹다 / 물을 먹다 / 음식을 배 불러 먹다 / (아이) 모이를 먹다 / 물이 약해진 누나는 반항을 못 지냈나 먹어도 잘 풀을거였다.	밥, 술, 약, 물, 음식, 모이, 보약
먹다_002002	담배나 아편 (따위를) 피우다.	-을	담배를 먹다 / 아편을 먹다.	담배, 아편
먹다_002003	연기나 가스 (따위를) 들이마시다.	-을	연탄가스를 먹다 / 태내를 먹다.	연탄가스, 태내
먹다_002004	어떤 마음이나 감정을 품다.	-을	안심을 먹고 두서를 하다 / 세상살이란 마음 먹기에 달려 있다 / 한번 먹은 마음이 변하지 않도록 하자. / 나는 마음을 흔하게 먹고 크기를 요원하였다.	안심, 마음
먹다_002005	일정한 나이에 이르거나 나이를 더하다.	-을	내 살 먹은 아이 / 나이를 먹다 / 내년이면 십삼을 먹는다.	나이
먹다_002006	욕, 민간 따위를 듣거나 달하다.	-을	하루 종일 욕만 피게 먹었다 / 그래도 그는 속이는 소리를 하다가 가끔 편지를 듣는 것이었다. <오기말, 할>	욕, 편지
먹다_002007	(속되게) 뇌물을 받아 가지다.	-을	뇌물을 먹다 / 뇌물을 먹고 양심을 눈감아 주다.	뇌물
먹다_002008	수익이나 이윤을 차지하여 가지다.	-을	남은 이익은 모두 내가 먹어야 / 시세가 마지 않은 것 같아 새 일을 찾았단 것인데 전 영을 먹기는 고사하고	이익
먹다_002009	물이나 습기 따위를 빨아들이다.	-을	기름 먹은 종이 / 김이 습기를 먹어 눅눅해졌다. / 숨이 물을 마려 먹었다.	기름, 습기, 물
먹다_002010	어떤 등급을 차지하거나 흡수를 띠다.	-을	시름을 먹다 / 우승을 먹다 / 100점을 먹다 / 체육 대회에서 유린기가 일등을 먹었다.	우승
먹다_002011	구기 경기에서, 점수를 잃다.	-을	상대편에게 먼저 한 골을 먹었다.	골
먹다_002012	(속되게) 여자의 정조를 유린하다.	-을	그는 벌써 여러 여자를 먹었다.	여자
먹다_002013	매 따위를 맞다.	-을	상대의 세 주먹을 한 방 먹고 나가떨어졌다.	주먹
먹다_002014	남의 재물을 다루거나 맡은 사람이 그 재물을 부당하게 자기의 것으로 만들다.	-을	관인 직원이 회사의 공금을 먹었다.	공금
먹다_002015	날이 있는 두구가 소해를 겪거나 자르거나 갈거나 하는 작동을 하다.	-에	이 고기에는 날이 잘 먹이 않는다. / 대패가 잘 먹는다.	고기
먹다_002016	바르르 물결이 배어들거나 고무 껌지다.	-에	뽕껍에 물이 잘 먹어야 다음틀기가 좋다. / 원술에 화강이 잘 먹지 않고 뭉는다.	뽕껍, 원술
먹다_002017	발레, 군 따위가 꺾어 넘어지거나 꺾지다.	-에	시뮬에 발레가 많이 먹었다. / 옷에 흙이 먹어 못 가게 되었다.	시뮬, 옷
먹다_002018	종이나 빨지 따위가 뭉개나 쓰이다.	-에	공사에 종이 색지가 많이 먹어 걸린다. / 낮은 집 수련에는 지선 지는 것보다 비호기 더 먹을 수 있다.	공사, 수련

<그림 4> 문형정보와 용례를 이용한 후보명사 추출(예-'먹다')

예를 들어 <그림 4>와 같이 '먹다02'의 첫 번째 다의어(예)는 문형정보("...을")를 기준으로 용례에서 {밥, 술, 약, 물, 음식, 모이, 보약}의 술목관계를 형성할 수 있는 후보명사를 추출하였다.

3.2 구문관계 후보명사의 의미 분별

(1) 용례 매칭 규칙 적용

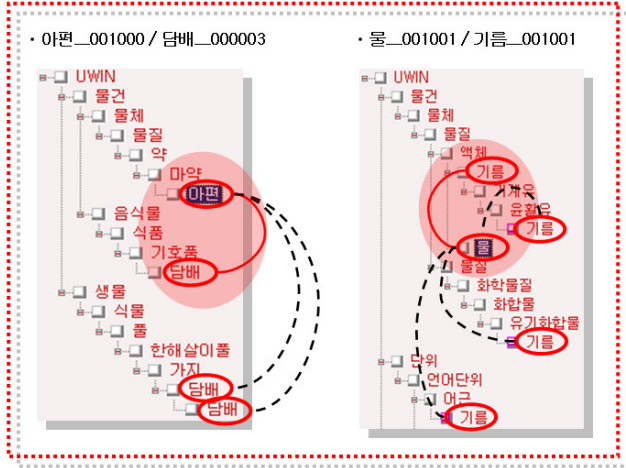
명사/용언 표제어의 용례에는 표제어와 함께 쓰이는 대표적인 서술어/명사를 포함하고 있으므로, 이렇게 함께 쓰이는 대표적인 명사와 용언은 각각 같은 용례를 가지고 있는 경우가 많다.

예를 들어 <그림 5>와 같이, 용례 “약을 먹다”는 ‘약_007001’과 ‘먹다_002001’의 용례에 각각 나타난다. 그러므로 동사 ‘먹다_002001’의 용례와 의미 분별되지 않은 후보명사 ‘약’의 용례를 차례로 비교함으로써, 동일한 용례 “약을 먹다”를 포함하고 있는 ‘약_007001’의 의미로 후보명사의 의미를 결정할 수 있다.

또한 명사와 용언이 동일한 용례를 포함하는지 비교하는 과정은 용례 전체를 비교하는 단순 매칭뿐만 아니라, 용례에서 추출한 구문패턴의 매칭도 함께 고려해야 한다. 용언 표제어의 용례는 용언의 활용형이 포함되거나 또는 “숟가락으로 밥을 떠 먹다”의 문장의 ‘숟가락으로’ 등의 수의적 성분을 포함하는 경우도 있기 때문이다.

4) ‘먹다02’의 첫 번째 다의어는 ‘[표제어]_[동형이의어번호][다의어번호]’의 표기 형식을 따라 ‘먹다_002001’로 표기한다.

__001001'}이 되었다. 그리고 '기름'은 '물__0001001'과 'U-WIN>물건>물체>물질>액체>'의 최소 공통 상위어를 가지며, 가장 큰 값의 의미 유사도를 가지는 '기름__001001'로 의미를 결정할 수 있다.



<그림 8> 링크 기반 의미 유사도 측정에 의한 후보명사의 의미 분별

(4) 구문패턴의 빈도 정보 이용

의미 유사도에 의한 후보명사의 의미 분별은 다음과 같은 경우에는 적용할 수 없다. 후보명사가 U-WIN의 계층구조로 표현되어 있지 않거나, 다의어별 용례에서 추출한 후보명사가 '뇌물', '우승', '이익', '골', '여자', '주먹' 등과 같이 하나만 존재하는 경우이다. 또한 '기름(물건)', '습기(작용)'과 같이 두 후보명사의 의미 유사도가 작은 경우에도 후보명사의 의미를 결정하지 못한다.

이와 같은 경우 앞서 추출한 구문패턴에서의 고빈도 의미를 선택함으로써, 의미가 결정되지 않은 나머지 후보명사의 의미를 분별하였다. 예를 들어, '떡다__002006'의 후보명사 '육'과 '편찬' 중에서 '육'은 1단계 의미 분별 과정에서 '육__002001'의 의미로 결정되었지만, '편찬'은 '떡다'의 구문패턴에 나타나지 않고, U-WIN의 계층구조로 표현되어 있지 않아 의미 유사도를 구할 수 없다. 이러한 경우 전체 구문패턴에서 고빈도로 나타난 '편찬__001000'으로 의미를 결정하도록 하였다.

3.3 구문관계 확장

구문패턴에서 고빈도로 나타난 명사를 해당 용언과 구문관계를 설정함으로써, 용례에서 추출할 수 없는 구문관계를 추가적으로 구축할 수 있다. 즉, 말뭉치에서 출현한 특정 어휘가 용례에 나타나지 않지만 해당 용언과 빈번하게 사용되는 경우가 있다. 예를 들어 '떡다'의 구문패턴 명사 목록에서 '돈__001001', '겹__005000' 등의 고빈도 명사는 용언 '떡다'와 구문관계를 설정할 수 있을 것이다.

3.4 구문관계 자동 구축 예

우선 구문관계를 구축하기 위해서는 구문관계를 형성할 수 있는 후보명사를 용언의 문형정보를 기준으로 용례에서 추출하였다. 동사 '떡다'의 14개의 의미('떡다__002001'~'떡다__002014')를 대상으로 문형정보 '...을'을 기준으로 구문관계를 형성할 수 있는 26개의 후보명사를 추출하였다.

그런 다음 <그림 9>와 같이, 4단계의 후보명사 의미 분별을 통해 구문관계를 설정하고 고빈도 어휘를 이용하여 구문관계를 확장하였다. 그 결과 26개 후보명사 중 '골__004005'와 '공금__001001'을 제외한 24개의 명사가 정확하게 의미가 결정되었고, 구문패턴에서 추출한 고빈도 어휘 중에서 '아침밥', '겹__005000', '돈__001001'이 추가적으로 구문관계를 설정하였다.



<그림 9> 후보명사의 의미 분별을 통한 구문관계 자동 구축 및 확장

4. 결론 및 향후 연구방향

본 연구는 U-WIN의 어휘 관계 중의 하나인 구문관계를 자동 구축하는 방법을 제시하였다. 용언의 의미별 용례와 문형정보를 이용하여 다양한 후보명사를 추출할 수 있었으며, 용례 규칙, 구문패턴, 의미 유사도 등을 이용한 명사 의미 분별을 수행하여 명사와 용언 간의 구문관계를 자동으로 구축할 수 있었다.

이러한 방법을 통해 뜻풀이 하나하나를 개별어휘로 구분하여 구축한 U-WIN 내부에 명사와 용언 간의 구문관계를 설정할 수 있을 것이며, 이러한 정보는 WSD, 격틀사전 구축, 정보검색, 클러스터링, 구문분석, 의미분석 등의 다양한 자연언어처리 분야에서의 활용을 기대해 볼 수 있을 것이다.

참고문헌

- [1]홍재성 외, “21세기 세종 계획 <전자사전 개발> 연구보고서”, 국립국어원, 2005.
- [2]전문용어언어공학센터[KORTERM], 『다국어 어휘망』 총3권, KAIST Press, 2005.
- [3]최호섭, “한국어 명사 개념망 구축-경제용어를 중심으로”, ETRI 지식정보검색연구팀 경제개념망 구축 결과보고서, 2001.
- [4]최호섭, “대규모 사용자 어휘지능망 구축과 활용”, 울산대 박사학위논문, 2007.
- [5]R. Rada, H. Mili, E. Bicknell, M. Blettner, Development and application of a metric on semantic nets, IEEE Transactions on Systems, Man and Cybernetics 19 (1) 17-30, 1989.
- [6]C. Leacock, M. Chodorow, Combining local context and WordNet similarity for word sense identification, in: C. Fellbaum (Ed.), WordNet: An electronic lexical database, MIT Press, pp. 265-283, 1998.
- [7]J. Jiang, D. Conrath, Semantic similarity based on corpus statistics and lexical taxonomy, in: Proceedings on International Conference on Research in Computational Linguistics, Taiwan, pp. 19-33, 1997.
- [8]D. Lin, Using syntactic dependency as a local context to resolve word sense ambiguity, in: Proceedings of the 35th Annual Meeting of the Association for Computational Linguistics, Madrid, pp. 64-71, 1997.
- [9]P. Resnik, Using information content to evaluate semantic similarity in a taxonomy, in: Proceedings of the 14th International Joint Conference

onArtificial Intelligence, Montreal, pp. 448-453, 1995.

- [10]Wu, Z., Palmer, Verb semantics and lexical selection, 32nd Annual Meeting of the Association for Computational Linguistics, New Mexico State University, LasCruces, New Mexico, 1994.
- [11]Hirst, G. and D. St.Onge. Lexical Chains as Representations of Context for the Detection and Correction of Malapropisms. WordNet. C. Fellbaum. Cambridge, MA, The Mit Press, 1995.
- [12]국립국어연구원, “<표준국어대사전> 편찬 지침 I-II”, 국립국어연구원, 2000.