

바이오 문서에서 지지 벡터 기계를 이용한 문법관계 분석

고려대학교 컴퓨터학과
박경미[†] · 황영숙 · 임해창

Grammatical Relation Analysis using Support Vector Machine in BioText

Kyung-Mi Park,[†] Young-Sook Hwang, Hae-Chang Rim

Department of Computer Science & Engineering, Korea University, Seoul, Korea

요 약

동사와 기본구 사이의 문법관계 분석은 품사부착과 기본구 인식이 수행된 상태에서, 동사와 의존관계를 갖는 기본구를 찾고 각 구의 구문적, 의미적 역할을 나타내는 기능태그를 인식하는 작업이다. 본 논문에서는 바이오 문서에서 단백질과 단백질, 유전자와 유전자 사이의 상호작용관계를 자동으로 추출하기 위해서 제안한 문법관계 분석 방법을 적용하고 따라서 동사와 명사구, 전치사구, 종속 접속사의 관계만을 분석하며 기능태그도 정보 추출에 유용한 주어, 목적어를 나타내는 태그들로 제한하였다. 기능태그 부착과 의존관계 분석을 통합해 수행하였으며, 지도학습 방법 중 분류문제에서 좋은 성능을 보이는 지지 벡터 기계를 분류기로 사용하였고, 메모리 기반 학습을 사용하여 자질을 추출하였으며, 자료부족문제를 완화하기 위해서 저빈도 단어는 품사 타입 또는 워드넷의 최상위 클래스의 개념을 이용해서 대체하였다. 실험 결과 지지 벡터 기계를 이용한 문법관계 분석은 실제 적용시 빠른 수행시간과 적은 메모리 사용으로 상호작용관계 추출에서 효율적으로 사용될 수 있음을 보였다.

서 론

동사와 기본구 사이의 문법관계 분석은 동사와 의존관계를 갖는 기본구를 찾고 각 구의 구문적, 의미적 역할을 나타내는 기능태그를 알아내는 일이다. 동사와 기본구 사이의 문법관계 분석은 여러 가지 분야에서 응용될 수 있는데 이 논문에서는 바이오 문서에서 단백질과 단백질, 유전자와 유전자 사이의 상호작용관계를 자동으로 추출하는데 적용하였다. 현재, 바이오와 관련된 다양한 연구 결과가 발표되고 있으며, 관련 문헌수의 양적인 증가는 점차 가속화되고 있다. 이처럼 생물학분야에서는 새로운 형태의 단백질 혹은 유전자 명칭들이나, 이들간의 관계에 관한 새로운 연구관련 문헌이 끊임없이 쏟아지고 있기 때문에 일선 분야의 연구자들은 점차 원하는 정보를 얻기가 어려워지고 있

다. 따라서 바이오 관련 문헌 데이터베이스에서 유의미한 정보를 자동으로 추출해내는 정보추출 기술의 중요성은 점점 더 강조되고 있다. 즉, 텍스트로부터 단백질 간의, 유전자 간의 상호작용관계를 자동으로 추출해 데이터베이스에 저장하면, 특정 단백질이나 유전자에 대한 검색을 통해 그와 상호작용관계를 갖는 모든 단백질 및 유전자 정보를 그 래프 등으로 볼 수 있고 각각이 어떤 관계를 갖는지를 알 수 있게 된다. 이를 통해 생물학관련 연구자는 여러가지로 도움을 받을 수 있다.^{1,2)}

본 논문에서는 상호작용관계를 자동으로 추출하기 위해 동사와 주격, 목적격 등의 문법관계를 갖는 기본구를 인식한다. 따라서, 기본구 사이의 의존관계 분석뿐만 아니라 주어, 목적어를 나타내는 기능태그 부착까지 고려해야 하는데 기본구의 구문적, 의미적 역할을 표현하는 기능태그는 영어에 대한 구문분석 말뭉치인 Penn Treebank에 정의되어 있다. 총 20개의 기능태그가 정의되어 있는데 이 중 주어를 나타내는 SBJ 태그가 41%를 차지한다.³⁾ 말뭉치에서 목적어를 나타내는 OBJ 태그는 정의되어 있지 않은

[†]E-mail : kmpark@nlp.korea.ac.kr
E-mail : yshwang@nlp.korea.ac.kr
E-mail : rim@nlp.korea.ac.kr

데, 몇가지 규칙¹에 의해서 부착될 수 있다. 상호작용관계를 추출하는데 필요한 만큼의 문법관계 분석만을 수행하기 위해 21개의 기능태그 중 본 논문에서 고려하는 기능태그는 SBJ(주어), OBJ(목적어), LGS(논리적 주어)로 제한한다. 이 태그들은 구문적 역할을 표현하는 기능태그들로 문장에서 이벤트²를 확인하는데 유용하다.

본 논문에서 문법관계 분석은 생물학자에 의해서 이벤트를 표현하는 동사³로 정의된 것에 대해서만 분석이 이루어진다. 또한 이벤트성 동사와 모든 기본구가 아닌 명사구(NP), 전치사구(PP), 종속접속사(SBAR)와의 문법관계만을 분석한다. 이 기본구들은 동사와 주격, 목적격 등의 문법관계를 갖는 기본구들이다.

이 논문에서 문법관계 분석은 테스트할 동사와 기본구가 주어졌을 때 이것을 8개의 클래스 중 하나로 분류하는 작업이다. 여기서 8개의 클래스는 NOFUNC,⁴ NP,⁵ NP-SBJ, NP-OBJ, PP, PP-LGS, SBAR, SBAR-OBJ 등이다. 실제 말뭉치에는 여러 기본구와 기능태그의 결합으로 573개의 다양한 클래스가 존재하지만 본 논문에서는 바이오 문서에서 이벤트를 추출하는데 유용한 문법관계 분석만을 목적으로 하기 때문에 분석하는 기본구와 기능태그를 제한한다.

앞으로 2장에서는 바이오 문서에서 문법관계 분석에 대해 설명하고, 3장에서는 문법관계 분석과 관련된 기존 연구를 살펴본다. 4장에서는 지지 벡터 기계를 이용해 문법관계 분석을 수행하는 방법을 살펴보고 5장에서는 문법관계 분석과 관련된 실험 결과를 보인다. 그리고 6장에서 결론을 맺는다.

바이오 문서에서 문법관계 분석

바이오 문서에서 단백질과 단백질, 유전자와 유전자 사이의 상호작용관계는 이벤트성 동사와 주격, 목적격 등의 문법관계를 갖는 기본구를 인식함으로써 분석된다. Fig. 1은 바이오 문서⁶에서 문법관계 분석의 예를 보이고 있다. 품사 부착과 개체명 인식, 기본구 인식을 수행한 결과가 입력으로

	품사 부착	개체명 인식	기본구 인식	
19	,		0	0
20	and	CC	0	0
21	also	RB	0	I-ADUP
22	indicate	VB	0	I-UP
23	two	CD	0	I-NP
24	mechanisms	NNS	0	I-NP
25	by	IN	0	I-PP
26	which	WDT	0	I-NP
27	distinct	JJ	0	B-NP
28	cytokines	NNS	B-protein	I-NP
29	can	MD	0	I-UP
30	activate	VB	0	I-UP
31	the	DT	0	I-NP
32	same	JJ	0	I-NP
33	Stat	NN	B-protein	I-NP
34	protein	NN	I-protein	I-NP
35	.	.	0	0

Fig. 1. 바이오 문서에서 문법관계 분석의 예.

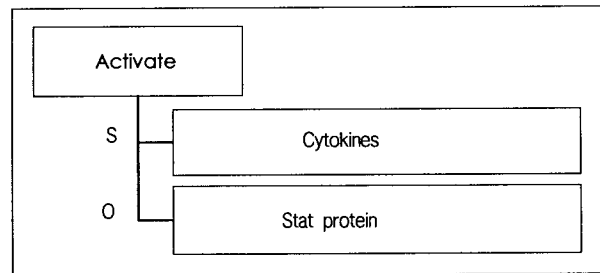


Fig. 2. 단백질간의 상호작용관계를 나타내는 이벤트.

들어왔을 때 먼저, 문장에서 동사⁷를 찾는다. 단어번호 22의 "indicate"와 단어번호 30의 "activate" 중에서 이벤트성 동사는 "activate"이기 때문에 이것과 문법관계를 갖는 기본구를 찾는다. Fig. 1은 이벤트성 동사인 "activate"와 문법관계를 갖는 기본구를 나타내고, 각각 주어와 목적어 관계임을 보이고 있다. 이벤트성 동사와 주어, 목적어 관계를 갖고 단백질 이름을 포함하고 있기 때문에 Fig. 2와 같은 상호작용관계를 자동으로 추출할 수 있다.

기존연구

[Blaheta, 2000]은 완전 구문분석(full parsing)이 수행된 상태에서 기능태그를 할당하는 문제를 다루고 있다. 방법은 Penn Treebank로부터 학습한 확률을 이용해 각 기능태그로 부착될 확률을 모두 구한 후 가장 확률이 높은 기능태그로 할당하는 것이다.³⁾ 사용하는 자질은 구문분석의 결과인 상위 부모 노드의 표지(label), 중심어(head) 등이다. 이 논문의 경우 기능태그 부착으로는 좋은 성능을 보였지만 실제 적용시 완전 구문분석에 의해 복잡도가 높아지고, 구문분석의 오류가 다음 단계인 기능태그 부착에

7 동사구의 중심어를 의미함.

1 규칙의 예 : 구구조인 Penn Treebank 말뭉치에서 동사구, 전치사구의 하위 노드에 위치하면서 기능태그를 갖지 않는 명사구는 목적어로 간주함.
 2 단백질 "cytokines"가 단백질 "Stat protein"을 activate(활성화시킨다)와 같은 상호작용관계 정보가 이벤트에 해당함.
 3 이벤트성 동사의 예 : activate, inhibit, induce,
 4 동사와 기본구 사이에 문법관계가 없음을 나타냄.
 5 동사와 명사구 사이에 의존관계는 있지만 주어나 목적어는 아님을 나타냄.
 6 Fig. 1의 예문은 개체명이 부착되어 있는 GENIA 말뭉치에 있는 문장의 일부임.

전파될 가능성이 높은 단점이 있다.

[Buchholz, 2002]는 동사구의 중심어와 다른 기본구의 중심어 사이의 문법관계를 분석하기 위해 메모리 기반 학습 (MBL : memory-based learning)을 이용한다.^{4,5)} 의존구조로 변형된 Penn Treebank로부터 학습자료로 총 361,342⁸개의 인스턴스를 추출하는데, 인스턴스는 동사와 기본구 주변의 문맥정보로부터 추출한 21개의 자질정보와 문법관계를 표현하는 1개의 클래스로 구성된다. 테스트 과정은 메모리 기반 학습을 이용해 학습집합으로부터 테스트 인스턴스와 가장 자질값의 차이가 적은 9개의 인스턴스를 추출해 w_j 값⁹이 가장 큰 클래스로 할당한다. 이 논문에서 문법관계 분석은 하나의 클래스로 분류하는 문제로 간주되기 때문에 의존관계 분석과 기능태그 부착을 동시에 수행한다. 바이오 문서에서 이벤트를 추출하는데 필요한 만큼의 문법관계(주격, 목적격 등)만을 추출하기 위해서 최소한의 구문분석을 수행하는데 응용할 수 있다.

지지 벡터 기계를 이용한 문법관계 분석

1. 사용하는 자질 정보

본 논문에서 자질 집합은 메모리 기반 학습을 사용하여 추출한 것으로 Table 1과 같다. 어휘자질의 경우 자질값의 종류가 2만개를 넘는 경우까지 있기 때문에 자료 부족 문제가 심각하다. 실제로 학습집합에서 빈도수 5이하로 발생한 단어의 수는 8번째 자질(기본구의 중심어)의 경우 21,111개 중 13,962개이다. 메모리 기반 학습을 사용한 기존 연구에서는 빈도수 5이하의 저빈도 자질값을 모두 하나의 값으로 대체하는데, 어휘 자질의 경우 저빈도 단어의 수가 상당히 많기 때문에 이 값들 사이의 구별을 위해 품사 타입이나 워드넷 최상위 클래스의 개념으로 대체할 필요가 있다.

Table 1에서 마지막 열은 각 자질의 가중치인 Gain Ratio를 의미한다. 대체적으로 어휘 자질들의 경우 자질값의 종류가 많아 낮은 값을 보였다.

2. 학습 및 테스트 과정

지지 벡터 기계의 학습 과정은 자질 집합을 결정한 후 학습 데이터를 이용해 NOFUNC을 제외한 각 클래스별로

8 Penn Treebank의 section 10부터 19까지에서 추출한 결과임. 이 중 동사와 문법관계를 갖는 인스턴스는 115,650개임. 학습 인스턴스 중에서 동사와 의존관계가 없는 것이 전체의 68%임. 이로 인해, 의존관계가 있음에도 불구하고 의존관계가 없다고 판단한 경우가 12.6%였음.

9 $w_j = e^{-30d_j}$ 여기서 d_j 는 테스트 인스턴스와 학습 인스턴스의 거리를 의미함.

Table 1. 자질 집합

	자 질	자질값의 예	Gain Ratio	
1	동사구-1 ¹¹	중심어	6,775개의단어	0.0374
2	동사구	품사 ¹²	MD, TO, VB, ...	0.0185
3		동사	5,785개의동사	0.0467
4	기본구-2	중심어	16,850개의단어	0.0537
5	기본구-1	중심어	15,898개의단어	0.0784
6		기본구 타입	ADJP, ADVP, ...	0.1115
7	기본구	전치사 ¹³	158개의전치사	0.1723
8		중심어	21,111개의단어	0.0763
9	기본구	중심어의 품사	36개의품사	0.1354
10		기본구 타입	ADJP, ADVP, ...	0.2952
11	기본구+1 ¹⁴	기본구 타입	ADJP, ADVP, ...	0.0834
12	기본구	방향	+, -	0.2847
13		거리	1, 2, 3,	0.1204
14	기본구	동사구의 수	0, 1	0.1250
15		품사의 수	0, 1, 2,	0.1101
16	기본구	품사 CC의 수	0, 1, 2,	0.0705
17		SBAR의 수	0, 1, 2,	0.0705
18	사이의 정보	명사구의 수	0, 1, 2,	0.1253
19		뒤큰따옴표의 수	0, 1, 2,	0.0931
20	사이의 정보	콜론의 수	0, 1, 2,	0.0711
21		앞큰따옴표의 수	0, 1, 2,	0.0448

7개의 분류 모델을 만드는 것이다.⁶⁾ 벡터의 차원은 자질값의 수에 따라 결정된다.¹⁰⁾

지지 벡터 기계의 테스트 과정은 7개의 이진 분류 모델을 이용해 테스트 벡터의 클래스를 결정하는 것이다. 7개의 이진 분류기의 prediction값 중 가장 값이 큰 클래스로 분류되는데, 모든 prediction값이 0보다 작은 경우 NOFUNC 클래스를 할당받는다.

실험 및 평가

1. 실험 환경

실험 데이터는 의존구조로 변형된 Penn Treebank의 section 10부터 19까지를 사용했는데, 동사와 기본구의 관계에 대한 259,205개의 벡터를 추출했다. 동사와 기본구 사

10 Penn Treebank section 11에서 19까지를 학습 데이터로 이용하고 빈도수 5이하의 저빈도 단어는 품사 타입으로 대체한 경우 21,547 차원임.

11 동사구의 왼쪽 첫 번째 기본구를 의미함.

12 동사구의 의미적 중심어가 아닌 구문적 중심어(조동사 등)의 품사를 의미함. 8개의 자질값을 가짐.

13 전치사구 다음에 명사구가 나타난 경우 이들은 하나의 구로 취급함. 여기서 명사구의 중심어가 전체의 중심어가 됨. "전치사" 자질은 전치사구의 중심어를 의미함.

14 기본구의 오른쪽 첫 번째 기본구를 의미함. 유일하게 사용되는 오른쪽 문맥 정보임.

Table 2. 기계학습 방법에 따른 실험결과

학습방법	MB ¹⁵	m : s ¹⁶	정확률	재현율	F _β
MBL	65	44 : 20	90.68%	89.30%	89.99%
SVM	5	7 : 12	89.39%	86.67%	88.01%

이에 나타난 동사구의 수가 특정값 이하인 경우만 고려했다. 기본구가 동사구의 왼쪽에 나타난 경우 동사구의 수는 1이하로 제한하고 오른쪽에 나타난 경우는 0인 경우까지만 고려한다. 이렇게 제한할 경우 전체 동사와 기본구 사이의 문법관계 중 96.5%를 커버하는 것으로 조사됐다.⁴⁾ 또한, 동사와 명사구, 전치사구, 종속접속사 사이의 문법관계만을 분석하고 SBJ, OBJ, LGS 등 3개의 기능태그만을 고려해 8개의 클래스 중 하나로 분류한다. 실험은 10-fold cross validation을 수행한다.

2. 실험 결과

지지 벡터 기계를 이용한 문법관계 분석은 Table 2와 같은 성능을 보였다. 비교를 위해 메모리 기반 학습을 이용했을 때의 결과도 제시했는데 두 방법 모두 같은 자질 집합을 사용했고 저빈도 자질값은 하나의 값으로 대체하는 방법을 사용했다. 표에서 메모리 기반 학습의 경우 지지 벡터 기계를 이용한 경우보다 메모리 사용량과 수행 시간에 있어서 효율성이 많이 떨어져 실제 바이오 문서에서 상호작용관계를 자동으로 추출하는 시스템에 적당하지 않음을 알 수 있다. 그런데, 실험 결과 지지 벡터 기계를 이용한 분석은 88%의 F-measure를 보여 메모리 기반 학습을 이용했을 때보다 약 2%정도 성능이 낮았다. 이러한 결과가 나온 원인 중 하나는 메모리 기반 학습에 의존적인 자질 집합을 지지 벡터 기계에 적용해서 나타난 결과이기 때문에 향후에 지지 벡터 기계에 적합한 자질들을 추가함으로써 성능을 높일 수 있다.

문법관계를 학습할 바이오 문서가 없기 때문에 분야가 다른, Wall Street Journal 기사로 이뤄진 Penn Treebank를 이용해 학습했다. 따라서 실제 바이오 문서에 적용시 사용되는 어휘가 다르기 때문에 자료 부족 문제가 심각하게 발생한다. 그래서 저빈도 단어들을 하나의 값으로 대체하는 기존의 방법보다는 저빈도 단어들의 구별을 위해 품사 타입이나 워드넷 최상위 클래스의 개념을 도입할 필요가 있다. 워드넷을 이용할 경우 최상위 클래스의 수는 명사의 경우 25개, 동사의 경우 15개이므로 저빈도 단어들의 구별이 가능하다. 학습 자료에서 중복되지 않는 17,871개의

Table 3. 저빈도 자질값의 표현방법에 따른 실험결과

학습방법	표현방법	정확률	재현율	F _β
MBL	품사	90.85%	89.34%	90.09%
	워드넷	90.86%	89.36%	90.11%
SVM	품사	89.47%	86.73%	88.08%
	워드넷	89.49%	86.66%	88.05%

Table 4. 독립적인 각 단계의 성능

단 계	클래스	정확률	재현율	F _β
의존관계 분석	NP	92.71%	92.08%	92.39%
기능태그 부착	SBJ	99.55%	99.58%	99.57%
	OBJ	97.11%	98.54%	97.82%

저빈도 단어 중 6,035개가 워드넷에 존재해서 해당 최상위 클래스로 대체했다. 그런데, 2개 이상의 최상위 클래스를 갖는 경우는 의미 중의성을 해소하지 않고 가장 처음 나타난 클래스로 대체했다. 또, 워드넷에 없는 저빈도 단어의 경우 품사로 대체했고 이 결과가 Table 3에 제시되어 있다.

기존의 방법보다 성능이 나아짐을 확인할 수 있는데, 실제 바이오 문서에 적용했을 때는 어휘 자질의 자료 부족 문제가 심각하기 때문에 품사나 워드넷을 이용해 이들을 구별하는 것이 더 유용할 것으로 보인다.

문법관계 분석은 의존관계 분석과 기능태그 부착의 2가지 단계로 나눌 수 있다. 본 논문에서는 2가지 단계를 통합해 수행하지만 단계를 분리할 경우 각 단계별로 유용한 자질을 선별해 사용할 수 있어서 더 좋은 성능을 보일 수 있다. 명사구에 대해서 각 단계를 독립적으로 수행한 결과가 Table 4에 제시되어 있다. 상대적으로 기능태그 부착 단계의 성능은 높고, 의존관계 분석 단계의 성능은 낮음을 알 수 있는데, 의존관계 분석 단계의 성능을 높일 수 있는 자질의 추가가 요구된다.

학습 데이터의 크기에 따른 문법관계 분석의 성능 변화를 통해 실제 바이오 문서에 적용시 학습 데이터의 크기를 결정할 수 있다. 그 결과가 Fig. 3에 제시되어 있는데 클래스 NP-SBJ의 경우 “1~6¹⁷⁾ 이후부터 성능 향상이 미미함을 알 수 있다. 한편, 다른 경우들은 약간씩 성능이 향상되고 있다.

결론 및 향후연구

이 논문에서는 지지 벡터 기계를 이용한 문법관계 분석에 대해 살펴보았다. 이것은 기본구 인식까지 수행한 바이

¹⁵ 메모리 사용량을 나타냄.
¹⁶ 분과 초로 나타난 수행시간임.

¹⁷ Penn Treebank의 1개의 section을 테스트 집합으로 6개의 section을 학습 집합으로 이용함을 의미함.

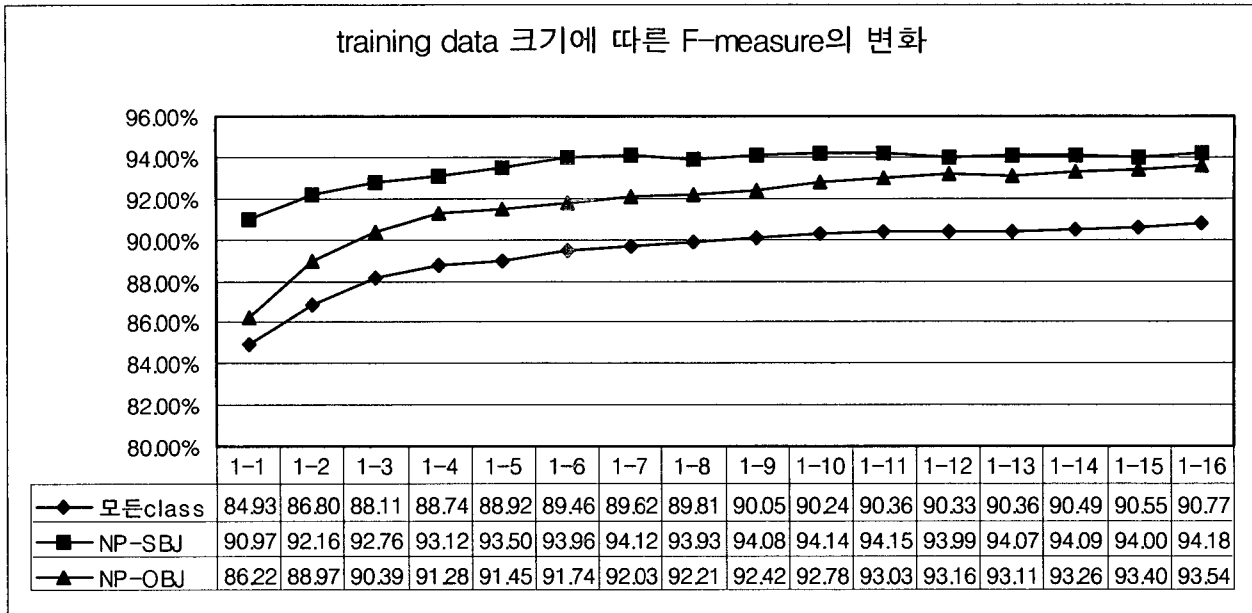


Fig. 3. 학습 데이터의 크기에 따른 F-measure의 변화.

오 문서에 대해 이벤트성 동사와 문법관계를 갖는 기본구를 찾는데 활용할 수 있다. 그 결과 단백질과 단백질, 유전자와 유전자 사이의 상호작용관계를 자동으로 추출하는데 중요한 역할을 할 수 있다. 이러한 문법관계 분석에 대해 메모리 기반 학습을 이용한 기존연구가 있는데 이 방법은 메모리 사용량과 수행 시간의 비효율성으로 인해 실제로 적용하기 어렵다. 그래서, 메모리 기반 학습을 이용한 기존 연구의 자질 집합을 참조해 지지 벡터 기계를 이용한 분석을 수행했다. 그리고 저빈도 단어의 경우 기존 연구의 하나의 값으로 대체하는 방법 대신에 품사나 워드넷을 이용했다. 앞으로 지지 벡터 기계에 적합한 자질들을 추가해 성능을 높일 필요가 있다. 특히, 바이오 문서를 학습 데이터로 이용할 수 없기 때문에 다른 분야의 문서를 갖고 학습했을 때 발생하는 자료 부족 문제를 완화시킬 수 있는 방법에 대한 연구가 필요하다.

REFERENCES

- 1) 임해창, 황영숙, 박경미(2003) : 바이오 텍스트 마이닝 시스템 개발. 정보과학회지 21(6) : 60-68
- 2) Pustejovsky J, Castano J, Sauri R, Rumshinsky A, Zhang J, Luo W (2002) : "Medstract : Creating large-scale information servers for biomedical libraries," Proc. of the Association for Computational Linguistics (Workshop on Natural Language Processing in the Biomedical Domain), pp85-92
- 3) Blaheta D, Charniak E(2000) : "Assigning function tags to parsed text," Proc. of the 1st Conference of NAACL, pp165-174
- 4) Buchholz S(2002) : "Memory-based Grammatical Relation Finding", PhD. thesis, Tilburg University
- 5) Timbl, <http://ilk.uvt.nl/>
- 6) SVM-Light, <http://svmlight.joachims.org/>

부 록 1

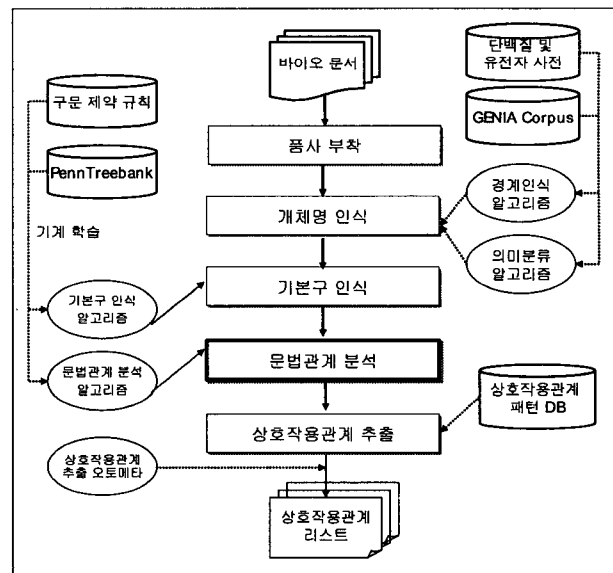


Fig. 4. 상호작용관계 자동추출을 위한 바이오 텍스트 분석 시스템 구성도.

본 논문의 지지 벡터 기계를 이용한 문법관계 분석 방법은 실제 상호작용관계 자동추출을 위한 텍스트 분석 시스템의 자연어처리 단계 중 하나로 사용됐다. Fig. 4는 전체 바이오 텍스트 분석 시스템의 구성도이다. 문법관계 분석을 위해 Penn Treebank의 16개의 section을 학습데이터로 사용해 분류모델을 생성했는데 빈도수 50번 이하의 저빈도 단어들은 품사로 대체했다. 문법관계 분석기의 성능

Table 5. 문법관계 분석기의 성능

클래스	①	②	③	정확률	재현율	F_{β}
모든 클래스	37,578	38,879	33,677	89.62	86.62	88.09
NP	2,531	3,038	2,128	84.08	70.05	76.42
NP-OBJ	8,920	8,656	8,025	89.97	92.71	91.32
NP-SBJ	15,488	15,810	14,503	93.64	91.73	92.68
PP	8,367	9,102	7,039	84.13	77.33	80.59
PP-LGS	545	543	495	90.83	91.16	90.99
SBAR	1,006	976	809	80.42	82.89	81.63
SBAR-OBJ	721	754	678	94.04	89.92	91.93

을 측정하기 위해 section 20~24를 이용해 실험한 결과 정답의 수, ③은 정확하게 예측한 수이다.
가 아래의 Table 5와 같다. 여기서 ①은 예측한 수, ②는