

개연성 규칙과 문장추상화를 활용한 문서요약

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

Text Summarization with Abductive Rules and Sentence Abstraction

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

요 약

본 논문에서는 문장추상화와 문장간 개연적 연결상황을 활용한 문단 기준의 문서요약을 생각하였다. 구상한 문단기준 문서요약 방법론은 다음과 같은 절차로 구성되어 있다: (1) 문단의 문장들을 추상화시킨다, (2) 문장구성성분들의 문장간 개연적 연결상황을 확인한다, (3) 연결집중도가 상대적으로 높은 문장을 문단의 화제를 담고 있는 것으로 인정한다. 본 논문에서는 이 과정에서 문장추상화에 필요한 구문분석기와 온톨로지를 구체화하였고, 문장추상기로 설화문장 추상화를 하였다. 그 후 개연성 규칙을 적용하여 문단의 주제문을 선별하였다.

1. 서론

인터넷의 대중화로 온라인화된 각종 문서의 양이 급격히 증가하고 있다. 이에 비례하여 문서요약이나 검색의 필요성이 절실해지고 있다. 문서요약은 새로운 이슈가 아니며, 지금까지 다양한 방법론과 그 작동화에 대한 많은 연구가 진행되어 왔다. 그럼에도 불구하고, 최신의 정보를 정확하고 신속하게 제공하기 위한 효과적인 문서요약의 중요성은 더욱 부각되고 있다.

문서요약 방법[1, 6]은 원문표현 방식과 요약과정의 각도에서 대별할 수 있다. 원문표현 방식으로 보면, 언어학적인 지식과 단서를 활용하는 방법과 이해한 원문의 내용을 형식화한 지식으로 표현하는 방법으로 나눌 수 있다. 한편, 요약과정에서 핵심적으로 활용하는 개념으로 보면, 본문문법(Text Grammar), 수사관계(Rhetorical Relation), 어휘사슬(Lexical Chain) 등에 기초한 방법론으로 분류할 수 있다. 그리고, 통계학적인 측면에서 보면, 원문에 나타나는 용어나 문장의 빈도수(Frequency), 단어 가중치(Word Weight) 등에 기반한 방법론들이 있다.

일반적인 원문의 표면적 구조로 보았을 때 문단은 작문의 단위로서, 저자가 이야기하고자 하는 화제가 통상 한 개씩 들어간다[8]. 이 점에 착안하여 문단요약을 원문요약의 출발점으로 삼는다.

문단 요약과정은 일종의 추상화 과정이다[9]. 추상화에는 추출할 것에 대한 선별작업이 수반된다. 문단 추상화의 경우, 화제와 밀접한 관계가 있는 문장을 선택해야 하고 이렇게 발췌된 문장들에 대한 또 한번의 추상화 작업이 필요하다. 그러나 화제관련성이 높은 문장을 선정하는 작업이, 이해를 수반한다는 점에서, 그렇게 용이하지 않다. 그래서 전자와는 반대방향으로 문장 추상화에서 시작하여 문단 추상화로 접근해 가는 방식을 고려할 수 있다. 이와 같은 문단 추상화를 본 논문에서는 문장 추상화를 통해 실현하는 방법을 고안하고, 개연성 규칙과 함께 그것을 설화문단 추상화와 주제문 선별에 활용하였다.

2. 설화용 존재론: OFN

본 논문에서는 설화를 이해하기 위한 존재론으로

OfN을 사용한다. 온톨로지 OfN은 다음의 7가지 범주로 구성되어 있다: 등장인물(Character), 심상(Affect State), 사건(Event), 상태(State), 시간과 공간의 변화(Delta-(Time, Space)), 담화표지(Discourse Marker). 이렇게 설정한 OfN을 구축하기 위해서 먼저 Roget 시소러스[2]의 범주를 심상, 시간과 공간, 사건, 그리고 상태 등으로 재편성하였다. 등장인물 유형에 속하는 어휘들은 고유명사 자원[7]을 이용하여 선정하였다. 담화표지의 경우는 수사구조의 연구결과[6]를 활용하였다. 이와는 달리 시공의 변화는, 구문분석 후 문장의 구성성분간의 상호작용에 의하여 확인되는 유형인 바, 그 기본유형은 시간과 공간이다.

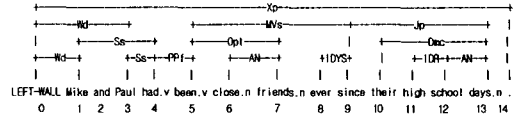
OfN은 문장추상화 과정에서 추출할 문장구성성분에 대한 선택기준을 제공한다. OfN과 함께 건설한 구문분석기도 활용하는데 그 과정은 다음과 같다: (1) 주어진 문장을 구문분석하여, 구성성분에 대한 구문상 중요도를 파악한다. (2) 중요 구성성분에 대한 OfN 유형을 확인한다. (3) 확인된 OfN 유형을 토대로, 구문상 중요도를 평가한다. (4) OfN 유형으로 확인된 것만을 추상화된 문장의 구성성분으로 채택한다.

3. 구문분석기: LGPI+

구문분석기로는 LGPI(Link Grammar Parser Interface)[5]를 확장시킨 LGPI+를 이용하였다. LGPI+는 Link Grammar Parser[4]에 대한 SWI-Prolog API를 제공한다. 입력문장에 대한 LGP 구문분석 결과는, 표식고리의 집합으로 문장의 통사구조가 표현된다. 표식고리는 한 쌍의 단어를 연결하며 그것들의 문법적인 기능을 표시한다.

다음 문장을 생각해 보자: *Mike and Paul had been close friends ever since their high school days.* 이에 대한 구문분석 결과는 그림 1과 같다: (1) *Xp*는 문장의 마침표와 좌벽을 연결한다. (2) *MVs*는 동사를 접속사와 연결시킨다. (3) *Jp*는 전치사와 그것의 복수형 목적어를 연결한다. (4) *Opt*는 *be* 동사를 복수명사에 연결하고 *there* 구문 후처리에 참가한다. (5) *Dmc*는 정관사와 복수명사를 연결한다. (6) *Wd*는 평서문에서 주부를 좌벽에 연결한다. (7) *Ss*는 단수명사를 단수동사형과 연결한다. (8) *PPf*는 *have* 동사를 과거분사형과 연결하고 *it*나 *there* 구문 후처리에 참가한다. (9) *A*는 한정형용사를 명사와 연결한다. (10) *ID[X][Y]*는 속어단어를

일렬로 연결한다. 여기에서 *X*, *Y*는 임의의 영문자이다. 마지막으로 (11) *AN*은 수식명사를 명사와 연결한다.



```
[link({m}, connection(3-4, s-[s], paul(_G1443), had(v))),
link({m}, connection(0-3, w-[d], 'left-wall'(_G1413), paul(_G1415))),
link({m}, connection(11-12, idr-[], high(_G1383), school(_G1385)),
link({m}, connection(12-13, an-[], school(_G1356), days(n))),
link({m}, connection(10-13, d-[m, c], their(_G1329), days(n))),
link({m}, connection(9-13, j-[p], since(_G1296), days(n))),
link({m}, connection(8-9, idys-[], ever(_G1266), since(_G1268))),
link({m}, connection(6-7, an-[], close(n), friends(n))),
link({m}, connection(5-7, o-[p, t], been(v), friends(n))),
link({m}, connection(5-9, mv-[s], been(v), since(_G1181))),
link({m}, connection(4-5, pp-[f], had(v), been(v))),
link({m}, connection(1-4, s-[s], mike(_G1119), had(v))),
link({m}, connection(0-1, w-[d], 'left-wall'(_G1089), mike(_G1091))),
link([], connection(0-14, rw-[], 'left-wall'(_G1059),
'right-wall'(_G1061)))
```

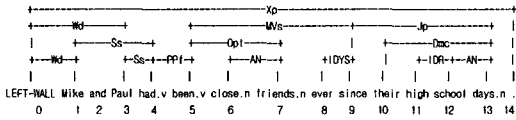
그림 1. 예문에 대한 LGP와 LGPI+의 출력결과

4. 문장추상기: SABOT

문장추상화 과정에서는 구문분석된 결과에서 주어, 동사, 목적어, 그리고 동사 수식어 등과 같은 어구들을 고려하되, 최상구(Top-Level Phrase)의 주요어(Head Word)가 우리의 주된 관심대상이 된다. 전치사구의 경우, 전치사와 그 목적어만 고려한다. 속어(Multiword Unit)는 한 단어로 취급한다. 동사가 의미상 변화를 내포하고 있거나 검토중인 단어가 심상에 연관되어 있을 때에는, 검토심도를 한 단계 깊게 한다. 이 경우, 문제단어의 목적어구나 수식어구의 주요어를 검토대상으로 삼는다. 만일 어떤 단어의 OfN 범주값이 2개 이상이면 최소값을 선택한다. 본동사가 의미상 변화를 포함하고 있고 이 동사구의 구성성분 유형이 시간(또는 공간)이라면, 구성성분의 OfN 범주는 *delta-time*(또는 *delta-space*)이 된다. 이러한 방식으로 우리는 한 문장 안에서 추상화에 참가할 요점어(Pivot Word)를 선택할 수 있다. 요점어란 OfN 범주가 확인된 문장추상화의 후보단어이다.

이러한 문장추상화 과정을 앞서 본 예문에 적용한 결과가 그림 2에 나타나 있다. Prolog로 구현한 문장추상기 *SABOT*가 요점어 *mike, paul, been,*

friends, 그리고 ever since 등을 선별해내었다. 문장 표식 $sent(X,Y)$ 에서 X 는 문단 내에서 위치를 그리고 Y 는 절의 위치를 각각 나타낸다. 술어 $affect_state$ 와 cue_phrase 는 각각 OfN 의 심상과 담화표지에 대응한다.



```
[sent(1, 1/2):[affect_state([friend, social, sympathetic]):friends/
(mike<->paul),
state([identity, absolute, relation]):been/ (mike<->paul)],
sent(1, 2/2):[cue_phrase([temporal, durative]):[ever, since]/ (mike<->paul)]]
```

그림 2. 추상화된 예문

문장추상기 **SABOT**을 문단을 구성하고 있는 각 문장에 적용하여 추상화된 문단을 얻는다. Mike와 Paul에 대한 이야기[1]의 경우를 보자. 그림 3에 추상화시킨 한 문단이 있다. **SABOT**이 문단을 처리한 결과가 그림 4에 보인다.

Mike and Paul had been close friends ever since their high school days. But now Mike wanted Paul out of town for a few days so that he could build a patio in Paul's backyard as a surprise birthday present. He suggested to Paul that he get away for a weekend, but Paul said he wasn't interested. On another occasion Mike casually spoke about the joys of fishing or camping trips. But Paul told him he enjoyed puttering around the house much more. Paul was getting very settled in his old age.

그림 3. 추상화시킬 문단

```
[sent(1, 1/2):[affect_state([friend, social, sympathetic]):friends/
(mike<->paul),
state([identity, absolute, relation]):been/ (mike<->paul)],
sent(1, 2/2):[cue_phrase([temporal, durative]):[ever, since]/
(mike<->paul)]],
[sent(2, 1/2):[affect_state([requirement, conceptional,
prospective]):wanted/ (paul<mike),
delta(space){out, of, tom}/ (paul<mike),
delta(time):days/ (paul<mike),
cue_phrase([temporal, repetitive]):[but, now]/
(paul<mike)],
sent(2, 2/2):[affect_state([wonder, contemplative]):surprise/
(paul<mike),
affect_state([giving, intersocial]):present/ (paul<mike),
delta(space):patio/ (paul<mike),
delta(space):backyard/ (paul<mike),
event([production, power, causation]):build/ (paul<mike),
cue_phrase([causal, specific, purpose]):[so, that]/
(paul<mike)]]
```

그림 4. 추상화된 문단

5. 개연성 규칙: Abductive Rules

표층적 원문 이해의 유용한 언어학적 도구로 개연성 규칙을 고안하였다. 이는 문장간 구성성분들의

개연적인 결속성을 나타내며 문장내 구성성분들이 가지는 OfN 정보로서 표현된다. 개연성 규칙의 일반적인 모습은 다음과 같다.

$$\text{Ante} \leq \text{Post} \{ \{ = + \}, = - \}, = * \} \text{Cons}$$

여기에서 (1) Ante, Post, Cons는 $pred(args)$ 의 형태를 가진다, (2) $pred(args)$ 는 OfN에 명시된 개념이다, (3) $= + \}$ and $= - \}$ 는 Post와 Ante에 제한사항이 있음을 나타내고, (4) $= * \}$ 는 담화표지가 있음을 나타낸다. 그림 5는 Mike와 Paul의 이야기[1]를 처리할 때 사용한 개연성 규칙의 예이다.

(1)	% 마음이 통하면 주고 싶어진다 affection(sympathetic) <= affection(offer)
(2)	% 풀이 죽으면 소극적이 된다 event(inactivity) <= event(descent).
(3)	% 장소를 바꾸고 싶을 때 여행을 한다 delta(space) <= event(journey) => affection(prospective).
(4)	% 싫은 것을 권하면 흥미를 끌지 못한다 affection(advice) <= affection(cause(pleasure)) => cue_phrase(adversative)

그림 5. 개연성 규칙

6. 실험 결과

표 1은 Mike와 Paul의 이야기뿐만 아니라 실험 원문으로 Dear Abby(<http://www.DearAbby.com>)에서 선택한 23편의 상담문 문단과 각 문단별 문장 및 단어의 구성 내용이다.

예피 소드	문단 개수	문단의 문장개수					계	문단의 단어개수					계
		1	2	3	4	5		1	2	3	4	5	
00-00	4	13	6	13	13		45	111	45	123	11		392
01-01	1	15					15	98					98
01-03	1	17					17	114					114
01-04	3	12	10	8			30	96	93	78			267
01-05	5	6	13	5	6	6	36	46	110	43	5	3	291
01-06	3	6	8	7			21	45	54	68			167
01-07	3	14	17	10			41	86	114	62			262
01-08	5	14	7	9	4	8	42	94	55	68	4	5	313
01-09	3	16	8	10			34	128	64	79			271
01-10	3	6	15	16			37	54	96	111			261
02-01	1	15					15	148					148
02-02	2	8	5				13	94	45				139
02-06	3	12	9	8			29	116	97	87			300
02-07	4	11	5	6	8		30	125	57	76	8		344
02-09	2	10	10				20	121	94				215
02-10	2	12	10				22	153	60				213
03-01	3	10	11	4			25	84	95	14			193
03-03	2	10	13				23	96	148				244
03-09	1	14					14	166					166
04-03	1	12					12	108					108
04-10	2	12	5				17	119	60				179
05-07	2	11	5				16	133	60				193
05-08	2	10	4				14	101	31				132
합계	58	266	161	96	31	14	568	2,436	1,378	809	296	91	5,010

표 1. 실험 대상 설화의 문단별 문장구성표

예과 소드	문단 개수	문단에 대한 Sentence Recall					계	문단에 대한 Topic Hit Ratio					계
		1	2	3	4	5		1	2	3	4	5	
00-00	4	3/3	2/2	5/6	5/7		3.54	3/3	2/2	6/6	7/7		4.00
01-01	1	2/5					0.40	5/5					1.00
01-03	1	3/6					0.50	4/6					0.67
01-04	3	2/4	2/3	2/3			1.84	4/4	3/3	3/3			3.00
01-05	5	0.50	0.67	0.67			1.84	1.00	1.00	1.00			3.00
01-06	3	1/2	2/3	2/2	1/1		3.67	2/2	3/3	2/2	2/2	1/1	5.00
01-06	3	1/2	1/2	3/3			2.00	2/2	1/2	3/3			2.50
01-07	3	0.50	0.50	1.00			1.85	1.00	0.60	1.00			2.60
01-08	5	4/4	1/2	1/2	1/2	3/3	3.50	4/4	2/2	1/2	2/2	3/3	4.50
01-09	3	2/4	3/3	1/3			1.83	3/4	3/3	3/3			2.75
01-10	3	3/3	4/5	1/2			2.30	3/3	5/5	2/2			3.00
02-01	1	8/9					0.89	8/9					0.89
02-02	2	2/4	0/2				0.50	3/4	1/2				1.25
02-06	3	2/5	2/3	3/4			1.82	3/5	2/3	3/4			2.02
02-07	4	4/4	1/2	2/3	2/4		2.66	4/4	1/2	2/3	2/4		2.66
02-09	2	4/5	4/5				1.60	5/5	4/5				1.80
02-10	2	1/5	0/5				0.20	2/5	2/5				0.80
03-01	3	2/3	2/4	0/1			1.17	2/3	4/4	1/1			2.67
03-03	2	0.67	0.50	0.00			1.17	0.67	1.00	1.00			1.51
03-09	1	3/5	4/7				0.67	4/6					0.67
04-03	1	2/5					0.40	4/5					0.80
04-10	2	2/5	3/3				1.40	3/5	3/3				1.60
05-07	2	3/5	1/2				1.10	3/5	1/2				1.10
05-06	2	0.60	0.50				1.66	0.60	0.50				1.66
합계	58						36.67						48.45

표 2. 실험 결과

표 2는 문장추상화와 개연성 규칙을 적용한 결과이다. 문단에 대한 Sentence Recall은 사람이 선택한 중요 문장을 문장추상화와 개연성 규칙을 적용하여 회상해 낼 수 있는 비율이다. 문단에 대한 Topic Hit Ratio는 회상해 낸 문장이 문단의 주제와 직접적으로 관련이 있는지를 나타내는 비율이다. 실험 결과, 본 논문에서 적용한 설화의 전체 문단에 대한 Sentence Recall은 63%, 각 문단에 대한 Topic Hit Ratio는 평균적으로 84%로 나타났다.

7. 결론

본 논문에서는 문서요약의 한 방법으로 문장추상화에서 시작하여 문단추상화로 접근해 나가는 방식을 취하였다. 그 과정은 다음과 같다: (1) 문단 안의

문장들을 추상화시키고, (2) 이들의 논리적 연결상황을 확인한다, (3) 연결집중도가 상대적으로 높은 것을 그 문단의 화제를 담고 있는 문장으로 인정한다, (4) 추상화된 문단들을 대상으로 (2, 3)의 과정을 적용하여 원문의 주제전개 상황을 파악한다. 이러한 과정에서 존재론 *OJN*과 실용적인 구문분석기 *LGPI+*, 그리고 문장추상기 *SABOT*, 원문 이해의 언어학적 도구로 개연성 규칙 등을 활용하였다.

본 논문에서 제시한 문장추상화와 개연성 규칙을 문단의 주제문을 선별하는데 활용해 본 결과, 주요 문장 회상도가 약 66%, 그리고 선별된 문장의 주제 관련성이 약 86% 정도임을 확인할 수 있었다. 그 결과 문장추상화와 개연성 규칙을 문서요약의 효과적인 도구로 활용할 수 있음을 알았다.

참고문헌

- [1] Bae, J.-H. J. and Lee, J.-H. "Topic Sentence Selection with Mid-Depth Understanding." Proc. of ICCPOL, pp. 199-204, 2001.
- [2] Roget's Thesaurus. <http://promo.net/cgi-promo/pg/t9.cgi?entry=22&full=yes&ftpstite=ftp://ibiblio.org/pub/docs/books/gutenberg/>.
- [3] 양재균, 배재학. "온톨로지 정보를 이용한 범주 재편성: Roget 시소러스의 경우." 한국정보처리학회, 제 9권, 제 1호, pp.515-518, 2002.
- [4] Sleator, D. and Temperley, D. "Parsing English with a Link Grammar." Third International Workshop on Parsing Technologies, August 1993. <http://www.link.cs.cmu.edu/link/>.
- [5] SWI-Prolog. <http://www.swi-prolog.org/>.
- [6] 배재학. "언어학적인 방법론을 취하는 자동 문서요약에 대한 연구." 공학 연구논문집, 제 29권 2호, pp.351-363, 울산대학교, 1998.
- [7] Proper Names Wordlist. <http://clr.nmsu.edu/cgi-bin/Tools/CLR/clrcat#14>.
- [8] C. Fellbaum (ed.), WordNet: An Electronic Lexical Database (MIT Press, 1998)
- [9] Bae, J.-H. J. and Lee, J.-H. Another Investigation of Automatic Text Summarization: A Reader-Oriented Approach. In Proceedings of ANZIS '94 (Australian and New Zealand Conference on Intelligent Information Systems), pp. 472-476, 1994.