

# 상대적 출현 빈도를 이용한 조사/어미 사전의 구성

강 승 식  
한성대학교 정보전산학부

## A Construction of Josa/Eomi Dictionary using Relative Frequency

Kang, Scung-Shik  
School of Information and Computer Engineering, Hansung University

### 요 약

한글 문서에서는 일부 조사와 일부 어미가 자주 출현하며 그 외의 조사/어미는 출현 빈도가 낮을 것으로 추측되고 있다. 본 연구에서는 실험에 통해서 이러한 사실을 확인하고 자주 출현하는 통합형 조사와 어미의 빈도를 구하기 위하여 한국어 말뭉치에서 통합형 조사와 통합형 어미의 상대적 출현 빈도를 조사하였다.

통합형 조사의 상대적 출현 빈도를 조사한 결과 말뭉치의 분야에 따라 약간의 차이가 있으나 평균 상위 9개의 통합형 조사가 전체 조사의 70%를 차지하고 있으며 상위 20개, 32개, 69개의 통합형 조사가 각각 90%, 95%, 99%를 차지하고 있음을 확인하였다. 통합형 어말어미의 경우에는 상위 10개의 통합형 어말어미가 전체 어말어미의 70%를 차지하고 상위 33개, 54개, 117개의 통합형 어미가 각각 90%, 95%, 99%를 차지하고 있다.

본 논문에서는 조사, 어미의 상대적 출현 빈도에 따라 계층적으로 조사/어미 사전을 구성함으로써 형태소 분석 효율을 높이고 형태소 분석기가 다양한 응용 분야에 쉽게 적용할 수 있도록 하는 방법을 제안한다. 또한 통합형 조사, 어미의 상대적 출현 빈도는 미등록어 추정을 용이하게 하거나 형태론적 모호성을 해결할 때에도 유용하게 활용될 수 있음을 보인다.

### 1. 서 론

한국어 단어의 구조는 '어휘형태소+문법형태소'로 이루어지는데 어휘형태소의 수에 비해 문법형태소의 수는 매우 적다[1,2,3]. 따라서 어휘형태소에는 복합어와 신조어, 외래어 등 사전 미등록어가 많지만 문법형태소에는 미등록어가 거의 없는 편이다. 형태소 분석 과정에서 문법형태소는 형태소를 분리할 때 중요한 역할을 하는데, 문법형태소의 수가 어휘형태소에 비해 매우 적다는 사실과 미등록 문법형태소는 거의 없다는 사실을 이용하면 문법형태소는 미등록 어휘형태소를 인식하는데 핵심적인 요소로 활용할 수 있다[4,5,6].

문법형태소는 한국어 분석시에 형태소 분석뿐만 아니라 구문 분석이나 의미 분석에도 중요한 역할을 하고 있다[7]. 그러나 문법형태소에 관한 연구는 그 중요성에 비해 추이 불 때 미미한 실정이다. 현대 국어에서 사용되는 문법형태소의 수가 몇 개인지에 대해서도 검증된 바가 없고 단지 추정치만 알려져 있을 뿐 정확한 숫자가 알려져 있지 않다.

본 논문에서는 우선 현대 국어의 문어체에서 사용되는 문법형태소 중에서 자주 출현하는 통합형 조사와 통합형 어말어미를 발견하고 상대적인 출현빈도를 구한다. 통합형 조사와 통합형 어말어미의 상대적 출현 빈도를 구함으로써 자주 출현하는 문법형태소에 가중치를 부여할 수 있고 그 결과를 형태소 분석을 비롯한 한국어 분석 시스템에서 체계적으로 활용할 수 있다.

통합형 조사와 통합형 어말어미의 출현 빈도를 구하기 위하여 분야별로 말뭉치를 구성하였다. 실험에 사용된 말뭉치는 약 160만 어절로 전산학 및 문헌정보학 논문 요약집 14만 어절(KT Test Set), 신문 기사에서 추출한 어휘 66만 어절, 국민학교 교과서에서 추출한 어휘 49만 어절, PC 통신에 공개된 소설 등 문학 작품에서 추출한 어휘 34만 어절이다[8].

### 2. 문법형태소 집합의 특성

한국어 문법형태소 중에서 국어 사전에 수록된 단위 조사의 수는 140여개이고, 어말어미의 경우에도 국어 사전

에 수록된 것은 약 400여개이다[4,9]. 그러나 단위 조사와 2개 이상의 조사가 결합된 형태를 포함한 통합형 조사의 수는 약 1,000~3,000여개로 추정되고 있으며, 통합형 어미의 경우에도 어미와 어미 혹은 어미뒤에 조사가 결합된 형태를 포함하더라도 약 1,000~3,000여개로 추정된다(선어말어미가 결합된 것은 제외한다)[10].

[가정-1] 문법형태소 집합은 유한집합이며 어휘형태소에 비해 그 수가 매우 적다.

문법형태소의 집합은 조사와 어말어미, 선어말어미, 접두사, 접미사를 포함한 허사(虛辭)들로 구성되는데, 국어사전에 모든 문법형태소가 나열되었다고 볼 수 없으며, 또한 언어는 시대에 따라 조금씩 변하므로 정확한 숫자를 파악하기가 어렵다. 그러나 문법형태소는 어휘형태소에 비해 상대적으로 변화 속도가 느리며, 국어사전에 수록된 어휘의 수를 비교해 볼 때 문법형태소의 수는 매우 적다. 또한 현재 표준어로 채택되어 있는 어휘들을 각 품사별로 분류할 때 각 품사들의 집합은 유한집합이라 할 수 있고 문법형태소의 집합도 유한집합이라 할 수 있다.

[가정-2] 통합형 조사와 통합형 어말어미의 집합은 유한 집합이다.

일반적으로 유한집합인 알파벳 T와 문법 G에 의해 생성되는 언어 L(G)는 무한집합이다[11]. 그러나 L(G)를 구성하는 string의 최대 길이를 제한하면 T\*와 T\*의 부분집합인 언어 L(G)는 모두 유한집합이 된다. 따라서 1개 이상의 조사로 이루어진 통합형 조사의 집합과 1개 이상의 어말어미와 0개 이상의 조사로 이루어진 통합형 어말어미 집합의 경우에 최대 스트링의 길이가 유한하다고 가정할 때 유한집합이 된다. 즉, 임의의 통합형 조사  $\alpha$ 와 임의의 통합형 어말어미  $\beta$ , 그리고 스트링의 최대길이 c에 대하여 스트링  $\alpha$ ,  $\beta$ 의 길이  $|\alpha|$ ,  $|\beta|$ 는  $|\alpha| \leq c$ ,  $|\beta| \leq c$ 이다.

[정리-1] 문법형태소에는 미등록어가 거의 없다.

가정-1에 의하여 문법형태소의 집합은 유한집합이며 그 수가 많지 않다. 또한 문법형태소는 언어의 변화 속도가 매우 느리므로 거의 모든 문법형태소를 문법형태소 사전에 나열할 수 있다. 가정-2에 의하여 통합형 조사와 통합형 어말어미를 문법형태소 사전에 수록하더라도 그 집합은 역시 유한집합이며 대부분의 문법형태소를 사전에 수록할 수 있다. 따라서 문법형태소에는 미등록어가 거의 없다.

1) 현대 국어에 사용된 것만 계산하고 고어는 제외한다.

[정리-2] 통합형 조사와 통합형 어미는 한글 문서에서 일부 조사와 일부 어미는 자주 사용되지만 그 외의 것은 드물게 사용된다.

한글 문서에 나타나는 통합형 조사는 수천 개이다. 그 중에서 '은/는/이/가/을/를/의/에/로' 등은 자주 사용되는 반면 대부분의 통합형 조사는 문서에서 드물게 나타난다. 통합형 어말어미의 경우에는 조사보다 집중도가 약간 떨어지지만 '아/어/는/다/ㄴ/ㄹ' 등이 자주 사용되고 대부분의 통합형 어말어미는 드물게 사용된다. 이러한 현상은 말뭉치에 대한 실험 결과에 의하여 확인된다.

[정리-3] 문법형태소는 형태소 분리 과정과 미등록 어휘 형태소의 인식 과정에서 중요한 역할을 한다.

국어 단어는 어휘형태소와 문법형태소로 이루어져 있고 형태소 분석은 각 어휘형태소와 문법형태소들을 분리-인식하는 것이다. 그런데 어휘형태소는 그 수가 많고 모든 어휘를 사전에 수록하는 것이 불가능하므로 일반적으로 미등록어가 다수 발생한다. 이에 비해 문법형태소의 수는 많지 않으며, 거의 모든 어휘를 사전에 수록할 수 있으므로 미등록어가 거의 없다. 또한 국어 단어의 구조는 어휘 형태소 하나에 문법형태소 0개 이상으로 구성되는데 문법 형태소가 2개 이상으로 구성되는 유형이 많다[12]. 따라서 형태소를 분리하고 인식하는 과정에서 문법형태소는 중요한 역할을 하게 된다.

미등록 어휘형태소의 인식은 단어의 구조를 파악하여 단어를 구성하고 있는 가능한 모든 형태소들을 인식한 후에 인식되지 않은 어휘형태소 부분을 미등록어로 추정하는 과정이다. 이 때 단어의 구조에 따라 미등록 어휘형태소를 제외한 다른 형태소(문법형태소)를 인식하기 위하여 문법형태소는 매우 중요한 역할을 한다. 이 과정에서 미등록 어휘형태소와 이웃한 문법형태소를 잘못 인식하면 미등록어 인식 오류가 발생하게 되므로 문법형태소의 인식 결과는 미등록 어휘형태소의 인식률을 좌우하는 핵심적인 요소이다.

### 3. 문법형태소의 출현빈도

#### 3.1 말뭉치의 구성 및 실험 환경

문법형태소 중에서 통합형 조사와 통합형 어말어미의 출현빈도를 구하고 전체 조사/어미에 대하여 각각의 상대적 출현빈도를 구하기 위하여 한글 문어체 말뭉치에 대하여 실험하였다. 말뭉치의 성격에 따라 문법형태소의 출현 빈도가 달라질 수 있으므로 말뭉치를 두 가지 유형으로 구성하였다. 첫번째 유형은 전문 분야에 대한 기술적인 문

시로 문장 유형에 제약이 있는 문서이다. 이러한 유형의 문서 집합으로 전산학과 문헌정보학 분야의 논문 요약집(KT Test Set), 그리고 신문 기사에 대한 말뭉치를 구성하였다. 두번째 유형은 다양한 유형의 단어와 문장 형태가 사용되는 문서이다. 이 유형의 문서 집합으로는 국민학교 교과서와 소설 등 문학 작품에 대한 말뭉치이다. 각 말뭉치의 크기는 표1과 같다.

표1. 조사/어미 추출 실험에 사용된 말뭉치

논문요약	신문기사	교과서	문학작품	전체
14만 어절	66만 어절	49만 어절	34만 어절	163만 어절

표1의 각 말뭉치에서 통합형 조사와 통합형 어말어미를 추출하고 출현 빈도를 구하는 과정은 다음과 같다.

- ① 한국어 형태소 분석기를 이용하여 말뭉치의 각 단어에 대하여 형태소 분석을 한다.
- ② 형태소 분석된 단어 중에서 추정된 결과 및 분석실패 어절을 제외하고, 분석 성공 및 복합명사 추정된 단어에 대하여 통합형 조사와 통합형 어미를 추출한다.
- ③ 추출된 통합형 조사를 출현 빈도에 따라 sorting하여 각 조사의 출현 횟수(단순 빈도) 및 모든 조사의 출현 횟수 합계에 대한 각 조사의 상대적 출현 빈도(상대 빈도)를 계산한다. 통합형 어말어미에 대해서도 같은 방법으로 단순 빈도와 상대 빈도를 계산한다.

단계①에서 한국어 형태소 분석기는 PC 통신망에 실험 및 연구용으로 공개된 'HAM version 1.51'을 사용하였다. 본 실험에서는 형태소 분석 결과만을 이용하여 통합형 조사와 통합형 어말어미를 추출하였기 때문에 분석 실패한 단어, 형태론적 모호성이 있는 단어 등 문장의 구조 및 의미에 적합하지 않는 분석 결과가 존재할 수 있다.

단계②에서 통합형 조사와 통합형 어미를 추출할 때 형태론적 모호성에 의하여 조사와 어미로 2가지 이상의 분석 결과가 존재하는 단어에 대해 각 분석 결과를 조사, 어미의 출현 횟수로 중복하여 계산한다(한 단어에 대한 통합형 조사 혹은 통합형 어말어미가 2개 이상인 경우에는 1개로 계산한다). 실험에 사용된 HAM version 1.51은 체언과 용언의 형태소 분석 결과로 품사 정보를 단순화시켜 명사, 대명사, 의존명사, 수사 등 체언은 'noun', 용언은 모두 'verb'로 출력하므로 어휘형태소의 품사 모호성이 거의 없다.

단계③에서 통합형 조사에 대한 상대적 출현 빈도의 계산 방법은 다음과 같다. 말뭉치에 출현한 통합형 조사의 수가  $n$ 개이고 각 통합형 조사를 단순 빈도(term frequency)가 높은 순으로  $j_1, j_2, \dots, j_n$ 이라 할 때 조사  $j_i$

에 대한 단순 빈도는  $j_i$ 의 출현 횟수  $n(j_i)$ 이다. 조사  $j_i$ 의 상대 빈도(relative frequency)를  $p(j_i)$ 라 할 때  $p(j_i)$ 는 다음과 같이 계산한다.

$$p(j_i) = \frac{n(j_i)}{\sum_k n(j_k)} \times 100$$

통합형 어말어미를 단순 빈도가 높은 순으로  $o_1, o_2, \dots, o_m$ 이라 할 때 임의의 통합형 어말어미  $o_j$ 에 대한 단순 빈도는  $o_j$ 의 출현 횟수  $n(o_j)$ 이고 상대 빈도는  $p(o_j)$ 이다.

### 3.2 실험 결과

4가지 분야의 말뭉치에 대한 통합형 조사와 통합형 어미의 실험 결과는 부록 1, 부록 2와 같다. 부록에서 각 말뭉치에 대한 기술 내용은 순서대로 조사(또는 어미), 단순 빈도, 상대 빈도, 그리고 상대 빈도에 대한 누적 빈도이다.

표2는 통합형 조사에 대한 실험 결과(부록 1)에서 누적 빈도가 70%, 90%, 95%, 99%인 고빈도 통합형 조사의 수이다.

표2. 출현 빈도가 높은 상위 n개의 통합형 조사

말뭉치 %	논문요약	신문기사	국민학교 교과서	문학작품	평균
70%	8	9	9	9	9
90%	16	20	20	22	20
95%	25	31	32	39	32
99%	49	65	68	93	69

표2에서 평균 9개의 통합형 조사가 전체 출현 빈도의 70%를 차지하고 있으며 평균 69개의 조사가 99%를 차지하고 있음을 알 수 있다. 즉, 고빈도 통합형 조사 9개로 형태소 분석을 한다고 할 때 실제 문서에서 약 70%의 조사를 분석해 낼 수 있으며, 상위 69개의 조사만으로 전체 조사의 99%를 분석할 수 있음을 알 수 있다.

표3. 출현 빈도가 높은 상위 n개의 통합형 어미

말뭉치 %	논문요약	신문기사	국민학교 교과서	문학작품	평균
70%	6	8	14	14	10
90%	16	28	40	47	33
95%	27	43	66	80	54
99%	60	99	139	170	117

표3은 통합형 어말어미에 대한 실험 결과(부록2)에서 누적 빈도가 70%, 90%, 95%, 99%인 고빈도 통합형 어말어미의 수이다. 표3에서 평균 11개의 통합형 어말어미가 전체 출현 빈도의 70%를 차지하고 있으며 평균 117개의 어말어미가 99%를 차지하고 있음을 알 수 있다. 즉, 고빈도 통합형 어말어미 11개로 형태소 분석을 한다고 할 때 실제 문서에서 약 70%의 어말어미를 분석해 낼 수 있으며, 69개의 어말어미만으로 99%의 어말어미를 분석할 수 있음을 알 수 있다.

표2와 부록1로부터 논문 요약집에 사용되는 통합형 조사의 수는 107개, 누적 빈도 99%에 속하는 통합형 조사의 수 49개로 다른 말뭉치에 비해 적게 나타난다. 이에 비해 소설 등 문학 작품에서 사용되는 통합형 조사의 수는 245개, 누적 빈도 99%에 속하는 통합형 조사의 수 93개로 다른 말뭉치에 비해 많이 나타난다. 즉, 논문과 같이 특정 유형의 단어와 특정 유형의 문장이 많이 사용되는 문서와 문학 작품처럼 다양한 유형이 많은 문서는 통합형 조사와 어말어미의 수 및 출현 빈도에 차이가 있음을 확인할 수 있다. 이러한 현상은 통합형 어말어미의 경우에도 비슷하다(표3, 부록2 참조).

#### 4. 조사/어미 사전의 구성

형태소 분석을 위한 조사 사전과 어미 사전을 구성하는 방법으로는 단위 조사와 어미를 사전의 항목으로 구성하고 조사간의 결합, 어미간의 결합은 접속정보표를 참조하여 검사하는 방법과 통합형 조사, 통합형 어말어미를 사전의 항목으로 구성하는 방법이 있다[13,14,15,16]. 첫번째 방법은 사전의 크기가 작은 대신 접속정보표를 사용해야 하기 때문에 새로운 결합 유형을 등록하려면 접속정보표를 갱신해야 하므로 추가 및 수정이 용이하지 않다.

두번째 방법은 통합형을 모두 사전에 수록하므로 사전의 크기가 커지는 단점이 있지만 수정 및 삭제가 용이하다. 일반적으로 사전의 크기가 커지면 사전 탐색 부담이 커지게 된다. 그러나 본 논문의 실험에 사용된 160만 어절에 출현한 통합형 조사와 통합형 어말어미는 각각 250여개(부록1), 440여개(부록2)로 단위 형태소를 수록한 경우 140여개, 400여개와 비교할 때 통합형을 모두 수록하더라도 사전의 크기가 커지는 부담은 거의 없다. 따라서 조사/어미 사전은 통합형을 사전에 수록하는 방식을 취하는 것이 바람직하다.

통합형 조사/어미 사전을 구성할 때는 통합형 조사와 통합 어말어미의 출현 빈도에 근거하여 고빈도어와 저빈도어로 나누어 계층적으로 구성하면 사전 탐색의 효율을 높일 수 있다. 이 때 고빈도어와 저빈도어를 구별하는 기준점을 발견하기 위하여 부록1, 부록2의 누적 빈도를 도표로 나타내면 각각 그림1, 그림2와 같다.

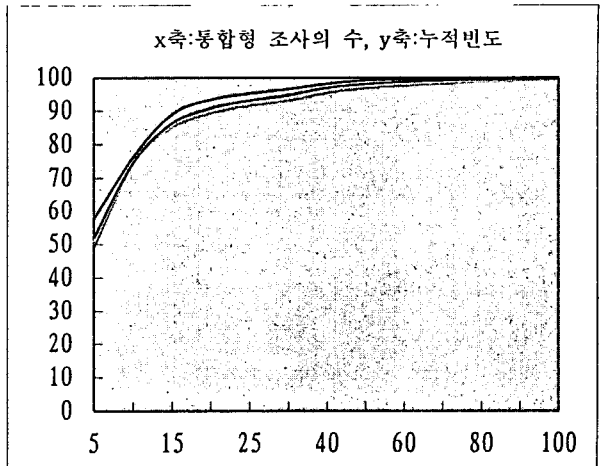


그림 1. 통합형 조사의 누적 빈도

통합형 조사 사전은 그림1에서 처리 증가율이 둔화되는 점과 처리율을 고려하여 초고빈도어 15개(조사그룹1), 고빈도어 35개(조사그룹2), 그 외의 조사(조사그룹3)로 세 개의 그룹으로 나누어 계층 구조로 사전을 구성하면, 조사그룹1이 약 86%, 조사그룹2가 12%, 조사그룹3이 나머지 2%를 처리하게 된다.

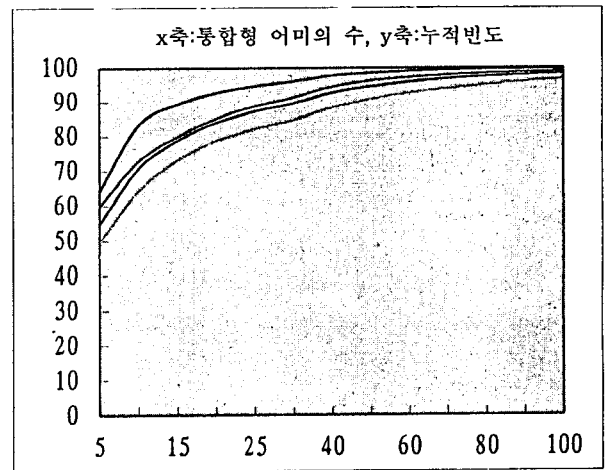


그림 2. 통합형 어말어미의 누적 빈도

통합형 어말어미 사전은 그림2에서 처리 증가율이 둔화되는 점과 처리율을 고려하여 초고빈도어 10개(어미그룹1), 고빈도어 50개(어미그룹2), 그 외의 어미(어미그룹3)로 세 개의 그룹으로 나누어 계층 구조로 사전을 구성하면, 어미그룹1이 약 71%, 어미그룹2가 25%, 어미그룹3이 나머지 4%를 처리하게 된다.

## 5. 결론

국어의 조사와 어미는 서로 다양하게 결합된 통합형이 존재한다. 형태소 분석을 위해서는 다양한 통합형 조사와 통합형 어말어미를 문법형태소 사전으로 구성해야 하는데 이 때 사전의 효율적인 구축을 위해 분야별로 160만 어절의 말뭉치에서 출현한 조사와 어말어미의 상대적 출현 빈도를 조사하였다.

상대적 출현 빈도에 대한 누적 빈도를 조사한 결과 한글 문서에서는 소수의 조사와 어말어미만이 자주 출현할 뿐이고 대부분의 조사와 어말어미는 드물게 출현함을 확인하였다. 누적 빈도의 증가율이 둔화되는 점과 처리 범위를 고려하여 조사/어미 사전을 구축할 때 각각 초고빈도 사전(Ultra-High Frequency Dictionary: UHFD), 고빈도 사전(High Frequency Dictionary: HFD), 기타로 나누어 계층적으로 구성하였다.

통합형 조사 사전을 초고빈도어 15개, 고빈도어 35개, 기타 조사로 나누어 계층 구조로 구성하면, 초고빈도 조사가 약 86%, 고빈도 어미가 12%, 기타 어미가 나머지 2%를 처리하게 된다. 통합형 어미 사전은 초고빈도어 10개, 고빈도어 50개, 기타 어미로 구성하면, 초고빈도 어미가 약 71%, 고빈도 어미가 25%, 기타 어미가 나머지 4%를 처리하게 된다.

조사와 어미의 상대적 출현 빈도는 형태소 분석사에 미등록 어휘형태소를 추정하거나 한국어 tagging 시스템에서 형태론적 모호성을 해결할 때도 유용한 정보로서 활용될 수 있을 것이다.

## 참고문헌

- [1] 김영근, 국어' 형태론 연구, 서울대학교 출판부, 1989.
- [2] 김진수, 국어 접속조사와 어미 연구, 탐출판사, 1987.
- [3] 서태룡, 국어 활용어미의 형태와 의미, 탐출판사, 1988.
- [4] 강승식, 음절 정보와 복수어 단위 정보를 이용한 한국어 형태소 분석, 서울대학교 컴퓨터공학과 박사학위논문, 1993년 2월.
- [5] 김덕봉, 최기선, 강재우, "한국어 형태소 처리와 사전 - 접속정보를 이용한 한글 철자 및 띄어쓰기 검사기 -", 어학연구, 26권, 1호, pp.87-113, 1990.
- [6] 이은철, 이종혁, "계층적 기호 접속정보를 이용한 한국어 형태소 분석기의 구현", 제4회 한글 및 한국어 정보처리 학술발표 논문집, pp.95-104, 1992.
- [7] 최기선, "한국어 해석을 위한 격어동 패턴의 고찰", 인지과학, 민음사, pp.364-388, 1989.
- [8] 김성혁 외 5인, "자동색인기 성능시험을 위한 Test Set 개발", 정보관리학회지, 11권 1호, pp.81-100, 1994.
- [9] 금성출판사, 뉴에이스 국어 사전, 금성출판사, 1989.
- [10] 부산대학교, 조사의 유형, Technical Report 90-1, 부산대학교 전자계산학과 인공지능 연구실, 1990.
- [11] P. Denning, J. Dennis, J. Qualitz, *Machines, Languages, and Computation*, Prentice-Hall, 1978.
- [12] 김성용, 최기선, 김길창, "Tabular Parsing 방법과 접속 정보를 이용한 한국어 형태소 분석기", 한국정보과학회 인공지능연구회 춘계 인공지능 학술발표회 논문집, pp.133-147, 1987.
- [13] 최계혁, 이상조, "양방향 최장일치법에 의한 한국어 형태소 분석기에서의 사전 검색 횟수 감소 방안", 정보과학회논문지, 20권, 10호, pp.1497-1507, 1993.
- [14] H.C. Kwon, Y.S. Chac, and G.O. Jeong, "A Dictionary-based Morphological Analysis," Proceedings of Natural Language Processing Pacific Rim Symposium, pp.87-91, 1991.
- [15] A. Copestake, "An Approach to Building the Hierarchical Element of a Lexical Knowledge Base from a Machine-Readable Dictionary," Proceedings of a Workshop on Inheritance in Natural Language Processing, 1990.
- [16] Boguraev, Btanimir and Mary Neff, *Dictionary Structure and Lexical Relations: A Study in Applied Computational Lexicography*, IBM Technical Report, IBM T.J. Watson Research Center, Yorktown Heights, New York, 1990.

부록 1. 상위 70개 통합 조사의 상대적 출현빈도

KT Data Set	신문 기사	국민학교 교과서	문학 작품
을	8944 16.97 16.97	25941 12.98 12.98	16498 13.65 13.65
의	8708 16.52 33.48	22161 11.09 24.06	16023 13.26 26.92
를	5350 10.15 43.63	20200 10.10 34.17	10296 8.52 35.44
에	4950 9.39 53.02	15021 7.51 41.68	9074 7.51 42.95
으로	2378 4.51 57.53	14937 7.47 49.15	8297 6.87 49.81
은	2366 4.49 62.02	11973 5.99 55.14	8189 6.78 56.59
이	2202 4.18 66.20	11615 5.81 60.95	6787 5.62 62.21
는	1916 3.63 69.83	11426 5.72 66.67	6109 5.06 67.27
에서	1672 3.17 73.00	8348 4.18 70.84	4239 3.51 70.77
과	1630 3.09 76.10	6547 3.28 74.12	4107 3.40 74.17
에서는	1598 3.03 79.13	5639 2.82 76.94	3368 2.79 76.96
로	1486 2.82 81.95	5244 2.62 79.56	3250 2.69 79.65
하고	1340 2.54 84.49	4984 2.49 82.06	2868 2.37 82.02
가	1332 2.53 87.01	3795 1.90 83.95	1822 1.51 83.53
와	1014 1.92 88.94	3169 1.59 85.54	1627 1.35 84.88
고	610 1.16 90.09	2905 1.45 86.99	1403 1.16 86.04
이다	444 0.84 90.94	2576 1.29 88.28	1137 0.94 86.98
어	444 0.84 91.78	1624 0.81 89.09	1030 0.85 87.83
으로써	426 0.81 92.59	1506 0.75 89.85	949 0.79 88.62
도	404 0.77 93.35	1317 0.66 90.51	762 0.63 89.25
다	230 0.44 93.79	1254 0.63 91.13	751 0.62 89.87
에서의	226 0.43 94.22	1043 0.52 91.66	615 0.51 90.38
든	192 0.36 94.58	1011 0.51 92.16	454 0.38 90.76
과를	190 0.36 94.94	957 0.48 92.64	439 0.36 91.12
보다	140 0.27 95.21	901 0.45 93.09	427 0.35 91.47
라	136 0.26 95.47	775 0.39 93.48	418 0.35 91.82
나	134 0.25 95.72	746 0.37 93.85	411 0.34 92.16
는	132 0.25 95.97	690 0.35 94.20	377 0.31 92.47
으로서	120 0.23 96.20	641 0.32 94.52	340 0.28 92.75
만	114 0.22 96.41	582 0.29 94.81	336 0.28 93.03
이나	102 0.19 96.61	535 0.27 95.08	333 0.28 93.31
며	102 0.19 96.80	532 0.27 95.34	325 0.27 93.57
로부터	100 0.19 96.99	487 0.24 95.59	310 0.26 93.83
로서	98 0.19 97.18	404 0.20 95.79	309 0.26 94.09
에게	92 0.17 97.35	380 0.19 95.98	281 0.23 94.32
까지	86 0.16 97.52	374 0.19 96.17	278 0.23 94.55
과는	84 0.16 97.67	344 0.17 96.34	268 0.22 94.77
개	80 0.15 97.83	315 0.16 96.49	246 0.20 94.97
만을	76 0.14 97.97	313 0.16 96.65	232 0.19 95.17
으로는	68 0.13 98.10	296 0.15 96.80	228 0.19 95.36
대로	66 0.13 98.22	290 0.15 96.94	215 0.18 95.53
이고	60 0.11 98.34	283 0.14 97.09	210 0.17 95.71
에도	60 0.11 98.45	262 0.13 97.22	179 0.15 95.86
과의	60 0.11 98.57	262 0.13 97.35	177 0.15 96.00
로의	56 0.11 98.67	254 0.13 97.48	161 0.13 96.13
으로부터	54 0.10 98.77	242 0.12 97.60	138 0.11 96.25
이며	50 0.09 98.87	227 0.11 97.71	118 0.1 96.35
에서도	46 0.09 98.96	227 0.11 97.82	115 0.1 96.44
로써	42 0.08 99.04	227 0.11 97.94	112 0.09 96.53
로는	40 0.08 99.11	206 0.10 98.04	111 0.09 96.63
까지의	38 0.07 99.18	186 0.09 98.13	109 0.09 96.72
요	38 0.07 99.26	164 0.08 98.22	108 0.09 96.81
이라는	28 0.05 99.31	164 0.08 98.30	107 0.09 96.89
으로의	26 0.05 99.36	163 0.08 98.38	102 0.08 96.98
이라고	24 0.05 99.40	147 0.07 98.45	102 0.08 97.06
만이	18 0.03 99.44	139 0.07 98.52	102 0.08 97.15
와의	18 0.03 99.47	124 0.06 98.58	101 0.08 97.23
마다	18 0.03 99.51	122 0.06 98.64	96 0.08 97.31
와는	18 0.03 99.54	115 0.06 98.70	96 0.08 97.39
이라	16 0.03 99.57	110 0.06 98.76	88 0.07 97.46
에만	16 0.03 99.60	107 0.05 98.81	82 0.07 97.53
라고	16 0.03 99.63	103 0.05 98.86	82 0.07 97.60
보다는	16 0.03 99.66	103 0.05 98.91	78 0.06 97.66
처럼	12 0.02 99.69	100 0.05 98.96	77 0.06 97.73
부터	12 0.02 99.71	98 0.05 99.01	76 0.06 97.79
로	10 0.02 99.73	84 0.04 99.06	73 0.06 97.85
으로서의	10 0.02 99.75	76 0.04 99.09	72 0.06 97.91
로서의	10 0.02 99.76	72 0.04 99.13	72 0.06 97.97
만으로	8 0.02 99.78	70 0.04 99.16	69 0.06 98.03
토록	8 0.02 99.80	69 0.03 99.20	67 0.06 98.08
총 107가지 52,718회 출현	총 197가지 199,906회 출현	총 219가지 195,279회 출현	총 245가지 120,825회 출현

