

# N3WS : 키워드 및 요약문장 추출을 이용한 인터랙티브 신문기사 탐색

조희정\*, 손지연\*, 윤별이\*, 조아현\*, 김 명\*, 박은정\*\*

\*이화여자대학교 컴퓨터공학과

\*\*네이버 파파고

heybb816@ewhain.net, gyouns@ewhain.net, stardl@ewhain.net, ccho0513@ewhain.net, mkim@ewha.ac.kr,  
lucy.park@navercorp.com

## N3WS : Interactive Newspaper Article Navigation Using Keyword and Summary Extraction

Hee-Jeong Cho\*, Ji-Youn Son\*, Byeol-Yi Yoon\*, A-Hyun Cho\*, Myung Kim\*, Eun-Jeong Park\*\*

\*Dept. of Computer Science & Engineering, Ewha Womans University

\*\*Dept. of Papago, Naver Corp.

### 요 약

최근 인터넷 기사 중에는 부정확한 제목이나 자극적인 단어를 사용하는 경우가 많아 구독자에게 불편함을 준다. 본 논문에서는 이러한 기사들의 헤드라인을 삭제하고, 기사의 내용을 3 문장으로 요약해 주어, 구독자가 원하는 기사를 효율적으로 파악할 수 있게 하는 시스템을 제안한다. 제안하는 본 시스템은 파이썬 언어의 KoNLPy 패키지를 사용하여 기사의 단어들을 형태소 단위로 분석하며, 추출된 키워드를 토대로 워드 클라우드를 생성한다. 사용자가 클라우드의 특정 단어를 선택하면, 해당 신문기사들의 본문을 분석하여 각 신문 기사만의 핵심적인 문장을 3 문장으로 출력해 준다.

### 1. 서론

언론 매체들의 수익은 구독자 수에 비례한다. 구독률을 높이기 위해, 일부 매체들은 본문의 내용과 상관 관계가 없는 자극적이고 선정적인 기사 제목으로 구독자를 유인하곤 한다. 예를 들어, 2017년 1월~3월 사이에 신문 윤리 위원회의 제재를 받은 기사 중에는 낚시성 제목을 갖는 기사가 4분의 1 이상 된다는 연합뉴스의 기사 [1]가 있었다. 낚시성을 띠지 않는 헤드라인이라고 해도, 제목만을 보았을 때 해당 기사의 내용을 쉽게 유추하기 어려운 기사들도 다수 존재한다. 독자들을 속이거나 정확한 사실을 전달하지 않는 기사 제목으로 인해 구독자들은 많은 불편함을 느끼게 된다. 이는 또한 해당 언론매체의 신뢰성 하락으로 이어질 수 있다.

현대 사회는 다양한 기술의 급속한 발전에 따라 '더 빠르고, 더 쉬운 것' 을 요구한다. 종이 신문을 읽는 것보다 스마트폰으로 인터넷 뉴스를 보는

현대인이 더 많다. 또한 긴 글의 본문을 읽기보다는 길이가 짧은 헤드라인만을 읽는 경향이 크다. 이러한 현상이 더욱 심화되고 있는 현 사태에서 신문 기사의 제목은 큰 중요성을 띠고 있으며, 현재 제시된 문제들은 반드시 보완되어야 할 점들이다.

해당 문제점들을 해결하려면, 언론에서는 독자들이 원하는 기사를 제대로 읽을 수 있는 환경을 제공해야 하며 빠른 시간 내에 기사 본문의 내용을 정확하게 파악할 수 있게 하는 방법을 제시해 주어야 한다. 본 논문은 이러한 문제점을 해결하는 방안으로, 혼란을 일으킬 수 있는 헤드라인을 삭제하고 각 기사의 본문을 3 줄로 요약하여 원하는 기사를 효율적으로 찾아 읽을 수 있는 시스템을 제안한다.

본 논문은 다음과 같이 구성된다: 2 절에서 자연어 처리를 사용한 기존의 문서 요약 사이트들을 살펴보고, 3 절과 4 절에서는 본 논문에서 제안하는 N3WS 시스템에 관해 기술한다. 끝으로 5 절에서 본 논문의 결론을 맺는다.

## 2. 기존 연구 사례 분석

### 2.1 국외 사례

텍스트 요약(text summarization)에 대한 다양한 선행 연구가 수행되어 있고 [2] [3] [4], 하나의 긴 영문 글을 원하는 길이로 요약해주는 웹 서비스 [5] [6]도 구현되어 있다. ‘Smmry’ [5]의 경우 영문으로 된 글을 텍스트 또는 PDF 파일 형태로 입력 받아 사용자가 지정하는 문장의 수로 해당 글을 요약해 준다. 이와 유사한 프로그램으로 ‘Text compactor’ 라는 프로그램 [6]도 있다.

### 2.2 국내 사례

국내 사례로는 국민대학교 한글 공학 정보 검색 연구소에서 개발한 한국어 처리 시스템인 ‘Korean Language Processing System’ [7]이 있다. 해당 프로그램은 형태소 분석, 명사 추출, 자동 문서 분류 등 한국어를 이용한 형태소 분석, 합성 등 한국어의 맥락을 이해하고 내용을 분류해 주는 기능이 있다. 본 연구 프로젝트와 가장 큰 유사성을 갖는 프로그램인 ‘세줄 요약기’ [8]는 장문의 글을 3 문장으로 요약해주는 사이트이다. 그림 1은 ‘세줄 요약기’가 글을 요약해 주는 예제이다.

#### 3줄 요약 v2

참고로 좀 노려요  
지원 언어: 한국어, 영어

#### 요약할 문서

강아지나 고양이 등의 반려동물들 기르는 가구가 많아지면서 수제간식, 영양제 등 반려동물 건강에 대한 관심도 높아지고 있다. 반려동물도 사람처럼 다양한 질병에 걸리는데 대부분 질환 초기에는 특별한 이상이 나타나지 않아 치료에 어려움을 겪는 경우가 많다.

특히 반려동물의 기침은 가볍게 지나서는 안되는 요소 중 하나이다. 기침이 장기간 지속된다면 ‘만성기관지염’에 걸린 상태일 수 있기 때문이다. 특별한 질환이 없는데도 두 달 이상 계속 기침을 한다면 병원을 찾아 검진을 받는 것이 좋다.

검진 결과 특별한 이상이 없더라도 흡연이나 먼지 등 기침을 자극하는 요인을 최소화하기 위해 주변환경관리에 세심히 신경써야 한다. 또 평소엔 기관지 건강에 도움을 주는 수제간식이나 영양제를 급여하는 것도

요약

#### 결과

1. 또 평소엔 기관지 건강에 도움을 주는 수제간식이나 영양제를 급여하는 것도 도움이 된다.
2. 반려동물 전용 영양제 전문 브랜드 마이펫닥터는 평소 강아지나 고양이의 호흡기 건강 관리에 도움을 주는 보양식 ‘노하브레스’를 선보이고 있다.
3. ‘노하브레스’는 호흡기 세포보호를 돕는 코엔자임Q10, 호흡기내 삼출물(가래)을 분해해주는 상황버섯과 브로멜라인, 활성 산소를 제거하여 호흡기 건강 유지에 도움을 주는 비타민E 등이 함유되어 있다.

(그림 1) 세줄 요약기의 요약 사례.

한글은 영어와는 달리 명사의 복수 표현, 동사의 과거, 현재, 미래 표현 등이 복잡하기 때문에 분석하기가 쉽지 않다. 따라서 영문 요약 시스템들과는 달리 한글 요약 프로그램들은 찾아보기 힘들다. 본 연구에서 제안하는 시스템은 국내의 유사

사례들과 달리 사용자의 사용 시간에 따라 자동으로 여러 개의 글을 크롤링하여 수집하고, 분석한 후 사용자가 원하는 글들을 요약해서 보여준다는 면에서 차별점을 갖는다.

## 3. N3WS : 인터랙티브 신문기사 탐색 시스템

본 프로젝트에서 제안하는 N3WS 시스템은 그림 2와 같은 구조를 갖는다. 이 시스템은 사용자의 접속 시간을 기준으로 신문 기사를 크롤링하여 워드클라우드로 보여주는 부분과 사용자가 선택하는 키워드에 관한 기사들의 세줄 요약물을 보여주는 부분으로 구성된다. 이제 각 부분에 대해 구체적으로 설명하고자 한다.



(그림 2) N3WS 시스템 구성도.

### 3.1 기사 크롤링 및 워드 클라우드 생성

사용자가 N3WS 웹사이트에 접속하면, 그 접속 시간을 기준으로 최근 기사들을 크롤링한다. 초기 데이터베이스로 약 2,000 개의 기사를 크롤링 하며, 10 분 간격으로 새로 업로드 된 기사를 검색하여 자동으로 저장한다.

크롤링 된 기사들에서 키워드를 추출한 후, 해당 기사의 주소와 본문을 함께 뉴스 정보를 저장하는 데이터베이스(이하 NEWS DB)에 저장한다. 이 때, 전체 기사를 대상으로 빈도수가 높은 키워드를 추출하였을 때 보다, 각 기사당 높은 빈도수를 가진 3 개의 키워드를 추출하고 그 빈도수를 다시 계산하는 것이 정확도가 높았다. 따라서, NEWS DB에는 키워드 1, 키워드 2, 키워드 3을 저장하며 해당 단어들을 워드 클라우드에 사용하도록 설정한다. “이대”, “이화여대”와 같이, 동일한 대상을 지칭하는 단어들은 덕셔너리를 사용한 함수를 이용하여 미리 통일시키며, 동일한 함수를 통해 “일보”, “것” 등 키워드로 인식하지 않을 단어들을 사전에 제거한다. 그림 3은 초기 NEWS



