

Word2Vec을 이용한 웹 문서 클러스터링 시스템 구현*

이현석, 안성훈, 이용환, 천명재, 박혁주, 박미화, 이용규
동국대학교 컴퓨터공학과-서울
e-mail : lazuki@dongguk.edu

Implementation of a Web Document Clustering System Using Word2Vec

Hyun Seok Yi, Sung Hun Ahn, Yong Hwan Lee, Myung Jae Cheon,
Hyeok Ju Park, Mee Hwa Park, Yong Kyu Lee
Department of Computer Science and Engineering, Dongguk University-Seoul

요 약

웹 문서 추천 시스템에서는 유사한 내용의 문서임에도 불구하고 URL이 달라서 다른 문서로 인식하여 사용자에게 추천하는 데이터 희소성 문제가 있다. 여기서 기존 연구들은 이 문제에 대한 해결 방법으로 TF-IDF를 이용하였으나 비용 및 시간의 한계가 있으며 유의어 분류 문제가 있다. 본 논문에서는 Word2Vec을 이용한 웹문서 학습 시스템을 통해 문제를 해결한다. 제안 시스템은 언론사의 뉴스를 수집하고 이를 정형화된 형식으로 분석하여 가공하는 전처리 과정을 거친 후 Word2Vec 학습을 통해 문서 벡터를 생성하고 이를 K-Means 클러스터링으로 유사 문서군으로 분류한다. 이 시스템을 이용하면 데이터 희소성 문제를 해결할 뿐만 아니라 연산량이 TF-IDF에 비해 줄어들고 유의어 분류 시 유사도가 높아지는 장점이 있다.

1. 서론

추천 시스템에서 주로 사용하는 협업 필터링[1]은 현재 사용자와 유사한 패턴을 보이는 다른 사용자를 찾아 그들이 평가했던 아이템 중 현재 사용자가 좋아할만한 아이템을 찾아내는 것이다. 이러한 협업 필터링의 문제점은 사용자가 평가한 기존 아이템의 수가 적거나, 사용자에게 비해 아이템의 수가 많아 사용자들 사이의 교집합을 찾을 수 없는 것이다. 이러한 문제를 데이터 희소성[2]이라고 한다.

하지만 웹 문서 추천에서의 데이터 희소성 문제는 조금 다르다. 인터넷 뉴스 서비스의 경우, 하나의 사건에 대해 여러 언론사에서 글을 작성하기 때문에, URL(이하 ID)은 다르지만 같은 내용의 뉴스(이하 아이템)들이 다수 존재하게 된다. 이 경우 동일하거나 유사한 내용의 아이템이지만 다른 아이템으로 인식하여 사용자에게 추천해주는 문제점이 있다.

기존 시스템에서는 이러한 문제점을 해결하기 위해 문서에 속한 단어들의 TF-IDF[3] 값을 이용하여 문서 벡터를 만든다. 그 후에 각 문서들의 유사도를 계산하여 문서들을 군집화한다. 이를 통해 유사한 내용의 문서들이 같은 군집에 속하게 되고, 군집을 하나의 아이템으로 판단하면 사용자가 본 아이템들의 교집합이 커지기 때문에 데이터 희소

성 문제가 완화된다. 하지만 계산하는 데에 시간이 오래 걸리고, PCA나 SVD와 같은 차원 축소 방법이 수행되어야 한다. 또한 정확한 Stemming 작업이 이루어지지 않거나, 동일한 단어가 아닌 유의어가 사용된다면 정확한 분류가 이루어지지 않을 수 있다.

Word2Vec은 문서를 이루는 말뭉치를 얇은 신경망을 통해 학습시키는 Word Embedding 방법이다. Hierarchical Softmax기법[4]을 사용하여 연산량이 줄어들고, 유의어의 영향이 적은 장점이 있다.

본 논문에서는 Word2Vec을 이용하여 웹 문서를 학습하는 시스템을 제안한다. 이를 통해 기존의 방식보다 빠르고 유의어의 영향이 적은 방식으로 클러스터링하여 데이터 희소성을 감소시킬 수 있다.

2. 관련 연구

2.1 협업 필터링

협업 필터링은 많은 사용자들로부터 얻은 기호정보를 이용하여 사용자들에게 아이템을 추천해주는 방법이다.

협업 필터링은 사용자 기반 방식과 아이템 기반 방식으로 구분할 수 있다. 사용자 기반 방식은 비슷한 패턴을 보이는 사용자들끼리 군집화하여 현재 사용자가 속한 군집의 사용자들이 높은 점수를 준 아이템을 추천해주는 방식이다. 이와 달리, 아이템 기반 방식은 비슷한 아이템으로 군집화하여 군집을 생성하고 그 군집 안에서 현재 사용자의

* 본 연구는 미래창조과학부 및 정보통신기술진흥센터의 서울어코드 활성화지원사업의 연구결과로 수행되었음(IITP-2016-R0613-16-1147)

* 본 연구는 미래창조과학부 및 정보통신기술진흥센터의 SW중심대학지원 사업의 연구결과로 수행되었음(R7116-16-1014)

패턴을 이용하여 다른 아이템을 추천해주는 방식이다[5].

본 논문에서는 주제에 적합한 사용자 기반의 협업 필터링을 사용한다.

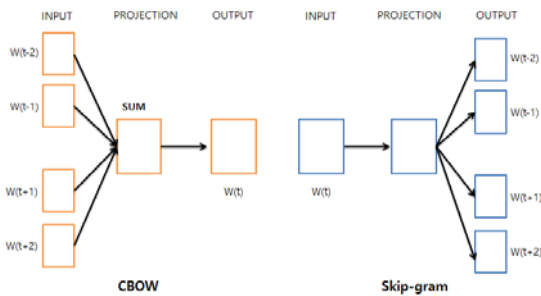
2.2 Word2Vec

Word2Vec은 Neural Network 기반의 Continuous Word Embedding 학습 모형이다. 단어가 가지는 의미 자체를 다차원 공간에서 벡터로 표현하여 벡터 연산을 통해 단어 사이의 관계를 파악하거나 추론을 내릴 수 있다.

Word2Vec의 학습 모델은 2가지가 있다. CBOW 모델[4]과 Skip-gram 모델[6]이다.

CBOW 모델은 여러 개의 단어들로부터 그 단어들과 가까운 특정 단어를 예측하는 모델이다. Skip-gram 모델은 특정 단어로부터 그 단어와 가까운 여러 단어들을 예측하는 모델이다.

(그림 1)은 CBOW 모델과 Skip-gram 모델의 구조를 나타낸다.



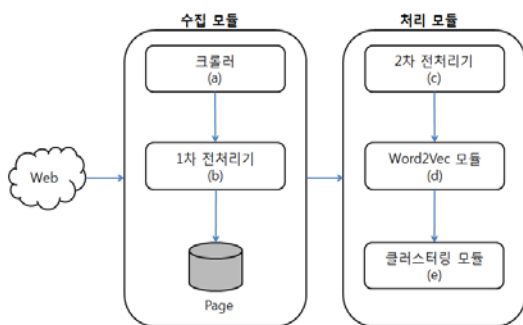
(그림 1) CBOW 모델과 Skip-gram 모델

특정문서에 대해 유사한 문서를 찾기 위해 본 논문에서는 Skip-gram 모델을 사용하였다.

3. 웹 문서 수집 및 Word2Vec 클러스터링 시스템

3.1 시스템 구성도

본 논문에서 제안하는 시스템은 (그림 2)와 같이 수집 모듈과 처리 모듈로 구성된다.



(그림 2) 시스템 구성도

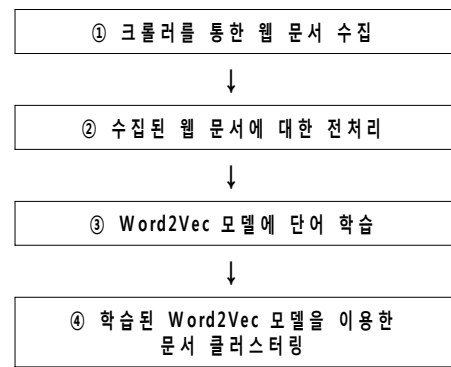
수집 모듈은 언론사 홈페이지의 뉴스 데이터를 수집하기

위한 크롤러(Crawler)와 수집된 뉴스 데이터 내의 학습이 불가능한 단어 제거, 광고 문구 제거 작업을 수행하는 1차 전처리기(Preprocessor)로 구성되어 있다. 1차 전처리를 수행한 데이터는 학습을 위해 데이터베이스에 저장된다.

처리 모듈은 Word2Vec 학습을 위해 문장 단위로 문서를 재구성하고 스템밍(Stemming)을 수행하는 2차 전처리기, Word2Vec의 Skip-gram 모델로 단어를 학습시키고 학습된 단어 벡터를 이용하여 문서 벡터를 생성하는 Word2Vec 모듈, 문서 벡터를 이용하여 클러스터링을 수행하는 클러스터링 모듈로 구성되어 있다.

3.2 동작흐름

(그림 3)은 본 논문에서 제안하는 시스템에 대한 순서도이다.



(그림 3) 전체 시스템 순서도

① (그림 2)의 (a)에 해당하는 크롤러는 지정된 언론사의 뉴스를 수집한다. 수집하는 정보는 문서 작성일, 문서 내용, 문서 제목, 문서 URL이다.

② 수집된 정보 중 데이터베이스에 저장할 수 없는 단어(ex>이모티콘)나 광고 문장 등의 불필요한 정보를 제거하는 작업을 (그림 2)의 (b)에 해당하는 1차 전처리기에서 수행한다. 1차 전처리 과정이 완료되면 문서 단위로 데이터베이스에 저장한다. (그림 2)의 (c)에 해당하는 2차 전처리기는 데이터베이스에 저장된 문서들을 로드하여 각 문서를 문장별로 분리하고, Stemming 작업을 거쳐 Word2Vec 학습의 입력 데이터로 사용할 수 있도록 한다.

③ 전처리 과정을 거친 데이터를 (그림 2)의 (d)에 해당하는 Word2Vec 학습 모듈의 입력으로 사용하여 Skip-gram 모델을 학습시킨다. 학습된 단어 벡터들을 다시 문서 기준으로 그룹핑 한 후 그룹 내 단어 벡터의 합을 평균하여 문서 벡터를 생성한다.

④ (그림 2)의 (e)에 해당하는 클러스터링 모듈에서는 3번에서 생성한 문서 벡터에서 임의로 K개의 문서벡터를

5. 결론

기존 방법은 실시간으로 추가되는 뉴스 추천에서의 데이터 희소성 문제를 완화하는데 계산 및 시간 비용이 높아지는 한계를 가지고 있었다. 본 논문에서 제안하는 시스템에서는 수집된 문서를 Word2Vec 알고리즘에 의한 학습을 통해 각각의 문서벡터를 생성한다.

이를 이용하여 유사 문서를 분류한 결과, 85% 이상의 정확도를 갖는 클러스터를 생성할 수 있었다. 이 분류 기법을 웹 문서 추천에서 활용하면 유사한 내용으로 다른 URL을 가지고 있는 페이지를 다른 내용으로 인식하여 추천하는 문제를 해결할 수 있다.

뉴스 데이터에서의 스테밍 과정에서 단어가 올바르게 추출되지 않아 문서의 군집화가 제대로 이루어지지 않는 경우가 발생하고 있어 향후 정확한 단어 추출을 위한 스테밍 연구를 진행할 것이다.

참고문헌

- [1] Gediminas Adomavicius, Alexander Tuzhilin. "Toward the Next Generation of Recommender Systems : A Survey of the State-of-the-Art and Possible Extensions", Foundations and Trends in Human-Computer Interaction, Vol 4, No. 2, pp. 91, 2010.
- [2] K. Goldberg, T. Roeder, D. Gupta and C. Perkins, "Eigentaste: A Constant Time Collaborative Filtering Algorithm", Information Retrieval, Vol 4, No. 2, pp.133-151, 2001.
- [3] Su. Xiaoyuan, and Taghi M. Khoshgoftaar, "A survey of collaborative filtering techniques", Advances in artificial intelligence 2009, Vol 2009, pp.5-8, 2009.
- [4] T. Mikolov, K. Chen, G. Corrado, J. Dean, "Efficient estimation of word representations in vector space", arXiv preprint arXiv, Vol 1301, No. 3781, 2013.
- [5] Singhal, Amit, "Modern information retrieval: A brief overview", IEEE Data Eng. Bull. Vol 24, No. 4, pp.35-43, 2004.
- [6] T. Mikolov, J. Dean, "Distributed representations of words and phrases and their compositionality", Advances in neural information processing systems, 2013.