

SNS와 뉴스기사의 감성분석과 기계학습을 이용한 주가예측 모형 비교 연구

김동영* · 박제원** · 최재현***

A Comparative Study between Stock Price Prediction Models Using Sentiment Analysis and Machine Learning Based on SNS and News Articles

Dongyoung Kim* · Jeawon Park** · Jaehyun Choi***

■ Abstract ■

Because people's interest of the stock market has been increased with the development of economy, a lot of studies have been going to predict fluctuation of stock prices. Latterly many studies have been made using scientific and technological method among the various forecasting method, and also data using for study are becoming diverse. So, in this paper we propose stock prices prediction models using sentiment analysis and machine learning based on news articles and SNS data to improve the accuracy of prediction of stock prices. Stock prices prediction models that we propose are generated through the four-step process that contain data collection, sentiment dictionary construction, sentiment analysis, and machine learning. The data have been collected to target newspapers related to economy in the case of news article and to target twitter in the case of SNS data. Sentiment dictionary was built using news articles among the collected data, and we utilize it to process sentiment analysis. In machine learning phase, we generate prediction models using various techniques of classification and the data that was made through sentiment analysis. After generating prediction models, we conducted 10-fold cross-validation to measure the performance of them. The experimental result showed that accuracy is over 80% in a number of ways and F1 score is closer to 0.8. The result can be seen as significantly enhanced result compared with conventional researches utilizing opinion mining or data mining techniques.

Keyword : Data Mining, Stock Price Prediction, Sentiment Analysis, SNS, Big Data,
Machine Learning

1. 서 론

경제의 발전과 더불어 사람들의 주식시장에 대한 관심이 증가함에 따라 수시로 변하는 주가를 예측하기 위한 많은 연구들이 진행되고 있다. 주가를 예측하기 위한 방법론은 크게 3가지로, 기본적인 분석(fundamental analysis)과 기술적 분석(technical analysis) 그리고 과학기술적 방법(technological methods)이 있다. 기본적인 분석은 재무 정보나 수익률과 같은 기업의 과거 성과를 평가하여 주가를 예측하는 방법으로 펀드 매니저나 주식 투자자 등이 주로 사용하는 방식이다. 기술적 분석은 과거주식가격의 동향을 파악해 미래주식가격을 예측하는 방식으로 EMA(exponential moving average) 등과 같은 통계적인 기법을 활용한다. 과학기술적 방법은 수많은 연산을 짧은 시간에 처리할 수 있는 컴퓨터 기술과 이와 관련된 이론들의 발전으로 인해 가능하게 된 방법이다. 대표적인 기법으로는 기계학습을 바탕으로 하는 인공신경망, 유전자 알고리즘 등이 있다. 여러 예측방법들 가운데 최근 과학기술적 방법을 이용한 많은 연구들이 활발하게 이루어지고 있다(Bollen et al., 2011; Kim, Kim, and Jung, 2012; Sagong, 2012; Kim, 2013; Kim, Lee, and Lee, 2013; Chun, 2013). 과학기술적 방법을 사용한 연구들에서는 기계학습, 데이터 마이닝, 오피니언 마이닝 등의 다양한 기법들을 이용해 연구를 진행 하였고, 예측에 사용하는 데이터 또한 다양해지고 있다. 이는 빅데이터의 등장으로 인한 것으로 여러 분야에서 기존의 사용하지 않던 데이터를 연구에 이용하기 시작했고, 주가 예측 분야에서도 이러한 현상이 반영된 것이다. 때문에 기존의 주가 예측에 사용하던 기업의 재무정보나 투자정보뿐만 아니라 뉴스기사(Ahn, 2010; Kim, Nam, Jo, and Kim, 2012; Kim, Kim, and Jung, 2012; Kim, 2013; Chun, 2013)나 SNS(Bollen, Mao, and Zeng, 2011; Kim, Jung, and Lee, 2014)의 데이터를 이용하는 등 다양한 데이터를 활용하여 주가를 예측하기 위한 연구를 진행하고 있다.

이에 본 논문에서는 주가 예측의 정확도를 높이기 위해 SNS와 뉴스기사 데이터를 동시에 이용한 다수의 주가예측 모형을 생성한 후 정확성을 비교하는 연구를 진행하였다. 우선 주가예측에 사용될 표본 회사를 추출한 후 일정 기간 동안 표본 회사가 언급된 SNS 데이터와 뉴스 기사를 수집한다. 수집한 데이터를 형태소 분석기를 이용해 분해한 후에 감성분석 과정을 거쳐 기사와 SNS 데이터의 긍정지수를 계산한다. 데이터들의 긍정지수와 빈도수를 예측에 사용되는 기법들을 기계학습 시키기 위한 데이터로 사용하여 예측 모형을 생성한다. 기계학습을 통해 생성된 다수의 예측 모형들과 주가와와의 관계를 통계적인 기법을 이용해 분석하고 모형들 간의 예측 정확도를 비교하여 연구를 진행한다.

2. 관련 연구

감성분석이란 자연어 처리와 텍스트 분석, 전산 언어학 등을 이용해 텍스트 내에서 주관적인 정보를 확인하고 추출하는 기법으로 ‘오피니언 마이닝’이라고도 한다. 감성분석의 기본 작업은 텍스트의 극성을 긍정, 부정, 중립 등으로 분류하는 것이다. 하지만 본 논문에서는 감성분석을 주가예측에 직접적으로 이용하지 않고 데이터 마이닝의 예측기법들을 기계학습 시키기 위한 데이터를 생성하기 위해 사용하고자 한다. 때문에 감성분석의 과정에서 형태소 분석기를 사용해 텍스트를 분해한 후 감성사전과의 비교를 통해 텍스트의 긍정지수를 도출하는 과정까지만 진행하였고 전체 텍스트의 극성을 분류하는 과정은 생략하였다.

형태소는 언어의 형태론적 수준에서의 최소단위를 뜻하며, 형태소 분석이란 문장의 어절을 형태소 단위로 분리한 다음 각 형태소에 맞는 범주를 부여하는 것으로 정의할 수 있다(Shim and Yang, 2004). 본 논문에서 필요한 데이터를 추출하기 위해서는 수집한 뉴스기사와 SNS 데이터를 형태소 단위로 분석하는 과정을 수행해야한다. 본 연구에서는 서울대학교 IDS 연구실에서 개발한 꼬꼬마

형태소 분석기를 사용하여 수집한 데이터들의 형태소 분석을 진행하였다.

형태소 분석을 통해 분리된 형태소의 극성을 분류하기 위해서는 개별 언어의 내용을 대상으로 단어 또는 단어의 의미 수준에서 긍정/부정/중립의 평가를 해놓은 감성사전을 사용한다(Yu, Kim, Kim, and Jeong, 2013). Song and Lee(2011)는 도메인의 특징을 고려하여 구축한 감성사전을 이용하는 것이 감성 평가의 정확도를 향상시킴을 나타내었고, 이를 바탕으로 Yu, Kim, Kim, and Jeong(2013)은 주식 도메인에 특화된 감성사전을 구축하고 활용하는 방안을 제시하였다. 본 연구에서는 Yu, Kim, Kim, and Jeong(2013)의 방법을 바탕으로 주식시장에 맞는 감성사전을 구축하였고 이를 이용해 뉴스와 SNS 데이터의 감성분석을 실시하였다.

지금까지 주가를 예측하기 위한 많은 연구들이 진행되어 왔다. Song(2002)은 뉴스의 발생이 주가에 유의한 영향을 미침을 나타내었고, 이를 바탕으로 Kim, Kim, and Jung(2012)은 뉴스 콘텐츠를 감성분석하여 주가 지수를 예측하는 투자의사결정 모형을 제시하였다. 제시한 모형과 주가 지수간의 관계를 통계기법을 이용하여 분석한 결과 통계적으로 유의한 관계를 가지는 것으로 확인되었다. 또한 Bollen, Mao, and Zeng(2011)은 SNS의 한 종류인 Twitter 사용자들의 데이터를 감성분석하여 주가 지수와의 연관성에 대한 실증적인 연구를 진행하였다. 사용자들의 데이터를 감성분석하기 위해 사용한 2가지 툴 중 하나인 Google-Profile of Mood States는 사용자의 감성을 6개의 수준(Calm, Alert, Sure, Vital, Kind, and Happy)으로 분류하였고, 이중 Calm으로 분류된 감정상태의 변화가 주가 지수의 변동과 유사한 흐름을 나타내고 있음을 발견하였다.

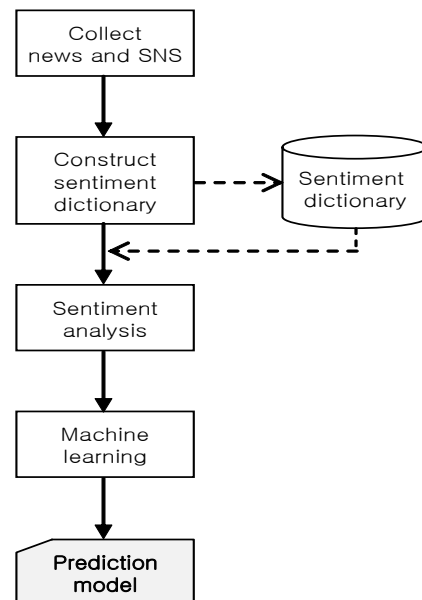
이에 본 논문에서는 뉴스기사와 SNS 사용자들의 데이터 모두를 주가예측을 위한 데이터로 사용하여 예측의 정확도를 높이하고자 한다. 또한 오피니언 마이닝을 통해 뉴스나 SNS에 담겨있는 의미를 파악하여 주가예측을 하는 방법(Kim, Kim, and

Jung, 2012; Kim, 2013)이 아닌, 오피니언 마이닝의 과정 중에 발생하는 데이터를 이용해 기계학습의 과정을 거쳐 주가를 예측하고자 한다. 또한 기존의 연구들과는 다르게 주가지수의 등락을 예측하는 것이 아닌 개별 기업 주가의 등락을 다룬다. 이를 위해 본 연구에서는 KOSPI 상장 기업 중 7개의 회사를 표본으로 추출하여 실험을 진행하였다.

3. 주가예측 모형

3.1 SNS와 뉴스기사의 감성분석과 기계학습을 통한 주가예측 모형

본 논문에서는 주가예측의 정확도를 높이하고자 기존의 연구에서 많이 사용되어오던 뉴스기사와 최근 많이 대두가 되고 있는 SNS의 데이터를 동시에 이용하였다. 또한 감성분석 기법을 적용하여 뉴스와 SNS의 데이터를 정제하였고, 이를 데이터 마이닝 분야의 여러 예측기법들을 기계학습 시키는 데에 사용하여 주가예측 모형을 생성하였다.



〈Figure 1〉 Flow of Stock Price Prediction Model Generation

제안하는 주가예측 모형은 <Figure 1>과 같으며, 데이터 수집, 감성사전 구축, 감성분석, 기계학습의 4단계를 거쳐 생성된다. 데이터 수집단계에서는 일정기간동안 해당 기업과 연관되어 발생한 뉴스기사와 SNS 데이터를 수집한다. 다음단계에서는 수집한 데이터 중 뉴스 기사를 분석하여 도메인에 맞는 감성사전을 구축한다. 감성분석 단계에서는 구축한 감성사전으로 수집한 데이터의 감성분석을 수행하여 기계학습을 위한 데이터를 생성한다. 기계학습 단계에서는 감성분석과정을 통해 생성된 데이터를 기반으로 기계학습 과정을 거쳐 최종적으로 예측 모형을 생성한다.

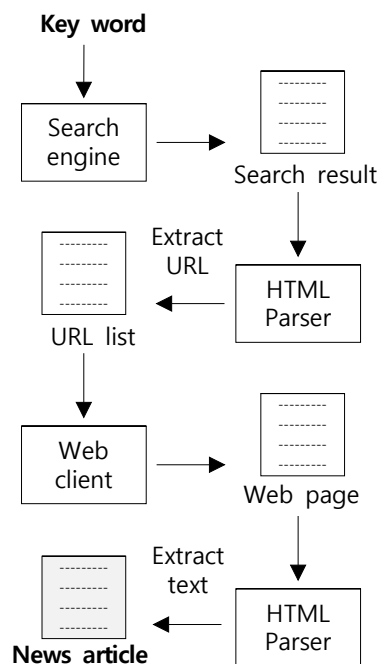
3.2 데이터 수집

뉴스기사와 SNS 데이터는 다양한 경로를 통해 대량으로 생산된다. 때문에 이를 연구에 이용하기 위해서는 자동화된 방식에 의해 필요한 정보만을 수집하여 저장해야 한다. 이를 위해 본 연구에서는 뉴스기사와 SNS 데이터를 수집하는 프로그램과 데이터베이스를 구현하였고, 이를 이용해 자동화된 방식으로 데이터를 정제하여 수집하였다.

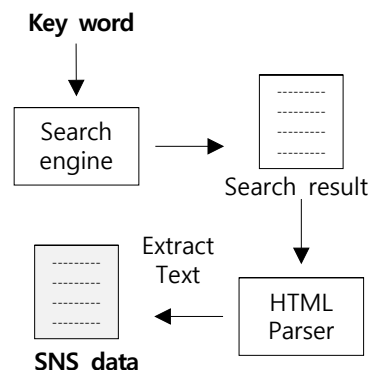
뉴스기사와 SNS 데이터의 수집과정은 <Figure 2>와 <Figure 3>에 각각 나타내었다. 수많은 뉴스기사와 SNS 데이터들 중 필요한 데이터를 수집하기 위해 검색엔진을 이용하여 해당 키워드가 들어간 데이터만을 추출한다. 뉴스기사의 경우 검색 결과에 기사 전문이 표시되지 않기 때문에 HTML Parser를 이용해 검색결과에서 해당기사의 원문이 게시되어있는 웹페이지의 URL을 추출한다. 추출한 URL을 통해 웹페이지에 접근한 후 다시 한 번 HTML Parsing 과정을 거쳐 뉴스 기사를 수집한다.

SNS의 경우 게시하는 글의 글자 수 제약이 있기 때문에 검색결과상에 글의 전문이 표시되어있다. 때문에 뉴스기사와는 다르게 검색결과에서 직접 텍스트를 추출해 데이터를 수집한다. 수집한 뉴스기사와 SNS 데이터는 데이터베이스에 저장하여 관리하며 형태소 분석과 감성분석 과정을 거쳐 주가

를 예측하기 위한 데이터로 활용한다.



<Figure 2> Process of News Article Collection



<Figure 3> Process of SNS Data Collection

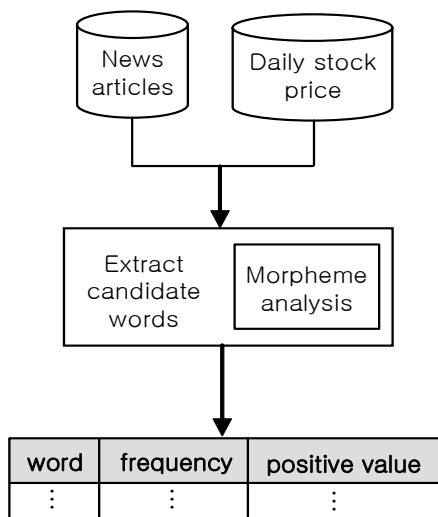
3.3 형태소 분석

텍스트로부터 작성자의 감정이나 의견을 추출하기 위해선 텍스트를 형태소단위로 분리하여 각 형태소별 극성을 파악한 후 전체 텍스트의 극성을 분류하는 방식을 사용한다. 본 논문에서는 서울대

학교 IDS 연구실에서 개발한 ‘꼬꼬마 형태소 분석기’를 사용하여 텍스트의 형태소를 분석하였다. 꼬꼬마 형태소 분석기는 한글 형태소의 품사를 ‘체언, 용언, 관형사, 부사, 감탄사, 조사, 어미, 접사, 어근, 부호, 한글 이외’의 항목으로 나누고 세부 품사를 구분하였다. 이중 체언의 경우 명사와 수사, 대명사로 구분하였고 명사는 다시 세부적으로 보통명사, 고유명사, 일반 의존 명사, 단위 의존 명사로 구분하였다. 본 연구에서는 여러 품사 중 체언에 속하며 실질적 개념을 표시하는 명사에 해당하는 형태소들을 추출하여 연구에 활용하였다.

3.4 감성사전 구축

감성분석 과정에서 극성을 분류하기 위해 사용하는 감성사전은 분석의 정확도를 높이는 데 있어 매우 큰 비중을 차지한다. 현재 한국어로 구성된 범용 감성사전이 제대로 구축되어 있지 않을 뿐만 아니라, 일반적인 단어들의 경우 주가의 변동을 제대로 반영하기에 적절하지 않다고 판단하여 주식시장에 특화된 감성사전을 구축하여 연구를 진행하였다. 감성사전의 구축과정은 <Figure 4>와



<Figure 4> Candidate Word Extraction for Building Emotional Dictionary

같으며, 이전 과정에서 수집한 뉴스기사와 일별 주가 변동 데이터를 이용한다. 우선 수집한 뉴스기사의 형태소를 분석하여 감성사전에 실릴 후보 단어들을 추출하고 단어들의 빈도수와 긍정 값을 계산한다. 빈도수는 해당 단어가 나온 기사의 수를 합산하여 계산하고, 긍정 값은 해당 단어가 들어간 기사가 게재되었을 때 익일 주가(next day stock price, NSP)가 상승한 경우의 수를 합산하여 계산한다.

빈도수(frequency)와 긍정 값(positive)을 수식으로 표현하면 다음과 같다.

$$word(i, j) = \begin{cases} 1 & \{ \text{기사 } j \text{에 단어 } i \text{가 포함된 경우} \\ 0 & \{ \text{그 외의 경우} \} \end{cases}$$

$$frequency(i) = \sum_{j=1}^n word(i, j)$$

$$NSP(j) = \begin{cases} 1 & \{ \text{기사 } j \text{가 게재된 후 익일} \\ & \text{주가가 상승한 경우} \\ 0 & \{ \text{그 외의 경우} \} \end{cases}$$

$$positive(i) = \sum_{j=1}^n \{ word(i, j) \times NSP(j) \}$$

추출한 후보단어들 중 상대적으로 빈도수가 너무 작은 단어들은 주가 변동을 제대로 반영할 수 없다고 판단하여 평균 빈도수 이하의 단어들은 제거하였다. 마지막으로 추출한 단어들의 긍정지수를 계산하여 감성사전을 완성한다. 긍정지수는 긍정 값을 빈도수로 나누어 나타내며, 식으로 표현하면 다음과 같다.

$$P(i) = \frac{\sum_{j=1}^n \{ word(i, j) \times NSP(j) \}}{\sum_{j=1}^n word(i, j)}$$

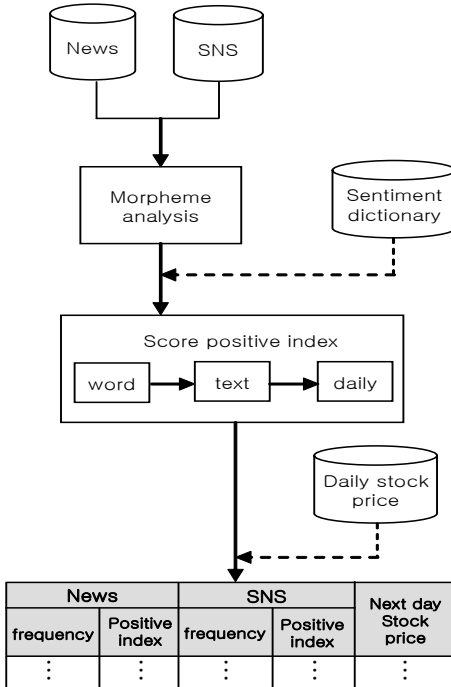
본 논문에서 구축한 감성사전은 단어와 긍정지수의 두 가지 속성으로 이루어져 있다. 단어는 명사만을 추출하여 구성하였고, 긍정지수는 0에서 1 사이의 값으로 1에 가까울수록 긍정의 의미를 나타낸다.

3.5 감성분석

구축한 감성사전을 이용하여 뉴스기사와 SNS 데이터의 감성분석을 <Figure 5>와 같은 과정을 거쳐 진행한다. 우선 수집한 데이터의 형태소를 분석하여 명사를 추출한 후 추출한 명사와 감성사전의 단어들을 비교해 해당 텍스트의 긍정지수를 계산한다. 텍스트의 긍정지수(positive index of text, PT)는 해당 텍스트에서 추출한 명사들의 긍정지수를 합해 그 개수로 나눈 산술평균값으로 나타내며 수식으로 표현하면 다음과 같다.

$$match(i, j) = \begin{cases} 1 & \left\{ \begin{array}{l} \text{텍스트 } i \text{에 포함된 명사 } j \text{가} \\ \text{감성사전에 존재 할 경우} \end{array} \right. \\ 0 & (\text{그 외의 경우}) \end{cases}$$

$$PT(i) = \frac{\sum_{j=1}^n \{match(i, j) \times P(j)\}}{\sum_{j=1}^n match(i, j)}$$



<Figure 5> Process of Sentiment Analysis

텍스트의 긍정지수를 계산한 후 일별 긍정지수를 계산한다. 일별 긍정지수(daily positive index, DP)는 해당일자에 게재된 텍스트들의 긍정지수를 합해 그 개수로 나눈 산술평균값으로 나타낸다. 일별 긍정지수를 수식으로 표현하면 다음과 같다.

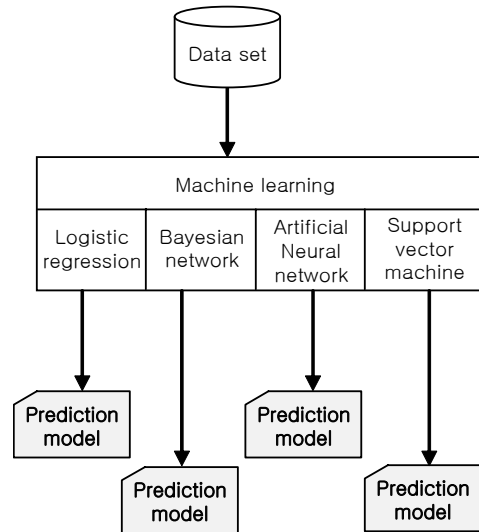
$$DP(i) = \frac{\sum_{j=1}^n PT(j)}{n}$$

$n = \text{number of text in } i$

최종적으로 일별 뉴스기사의 빈도수와 긍정지수, SNS의 빈도수와 긍정지수를 도출한 후 익일 추가변동 항목을 추가해 기계학습을 위한 데이터로 사용한다. 익일 추가변동 항목은 주식가격의 오르고 내림의 정도를 표시하지 않고 오르고 내림의 여부를 표현하기 위해 up, 0, down 세 가지로 나타낸다.

3.6 기계학습

기계학습을 위한 데이터는 속성 집합과 클래스 레이블로 이루어져 있으며, 기계학습은 속성 집합



<Figure 6> Flow of Machine Learning

본 연구에서는 <Figure 6>에서 나타낸 바와 같이 기계학습이론을 바탕으로 하는 분류기법들을 사용하여 예측 모형을 생성하였고, 여러 분류기법들 중 로지스틱 회귀분석, 베이지안 네트워크, 인공신경망, 서포트 벡터 머신 방법을 이용하였다.

실험은 KOSPI 상장사 중 7개의 기업을 대상으로 진행하였으며, 2013년 1월 2일부터 2013년 12월 30일까지의 뉴스기사와 SNS 데이터, 주가변동 데이터를 수집하여 실험에 활용하였다.

실험에 필요한 데이터는 자체적으로 제작한 프로그램과 데이터베이스를 사용하여 수집하였다. 뉴스 기사의 경우 경제관련 기사를 중점적으로 다루는 11개의 언론사에서 경제 섹션에 게재한 기사들을 대상으로 검색하였고, 이 중 해당기업과 관련된 기사만을 추출하여 총 132,123건의 기사를 수집하였다.

〈Figure 7〉 Portion of Collected News Articles and
Twits

〈Table 1〉 Number of Collected Data for Each Company

| Company | Number of News Data | Number of SNS Data |
|---------|---------------------|--------------------|
| A | 41,343 | 96,146 |
| B | 28,410 | 20,714 |
| C | 10,879 | 6,561 |
| D | 15,214 | 26,851 |
| E | 16,923 | 34,076 |
| F | 7,136 | 4,141 |
| G | 12,218 | 39,771 |

기업별로 수집한 뉴스기사와 주가변동 데이터를 이용하여 감성분석에 필요한 감성사전을 각각 구축한다. 우선 뉴스 기사를 형태소분석하여 명사를 추출한 후 그 빈도수를 계산하여 후보단어들을 도출한다. 이 중 뉴스에 포함 된 해당기업명, 언론사명 등 불필요한 단어와 평균 빈도수 이하의 단어들을 제거한 후, 추출한 단어들의 빈도수와 긍정값으로 긍정지수를 계산하여 감성사전의 구축을

| word | positive_index | word | positive_index |
|------|----------------|------|----------------|
| 절벽 | 0.147321 | 패소 | 0.195122 |
| 민사 | 0.203252 | 증발 | 0.222222 |
| 시동 | 0.307692 | 촉구 | 0.308333 |
| 투자 | 0.442915 | 전자 | 0.446831 |
| 증권 | 0.449651 | 상승 | 0.451437 |
| 국산화 | 0.514793 | 테크 | 0.514811 |
| 기어 | 0.649275 | 회장단 | 0.626667 |
| 시리마 | 0.746032 | 베를린 | 0.710462 |

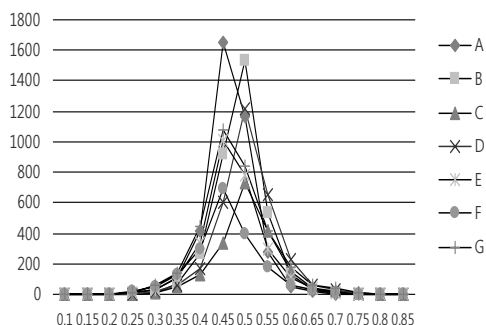
〈Figure 8〉 Portion of Sentiment Dictionary

마무리한다. 감성사전의 일부를 <Figure 8>에 나타냈으며, <Table 2>에 기업별 감성사전의 단어 수를 나타내었다.

<Table 2> Number of Words on Sentiment Dictionary for Each Company

| Company | Number of Candidate Words | Number of Words on Sentiment Dictionary |
|---------|---------------------------|---|
| A | 31,789 | 3,740 |
| B | 29,108 | 3,513 |
| C | 16,541 | 1,884 |
| D | 24,246 | 3,041 |
| E | 22,490 | 2,730 |
| F | 15,429 | 1,851 |
| G | 22,966 | 3,277 |

감성사전의 긍정지수는 0에서 1사이의 값으로 표현되며 1에 가까울수록 긍정의 의미를 나타낸다. 실험을 통해 구축한 각 기업 별 감성사전의 긍정지수 분포를 <Figure 9>에 표현하였다. 긍정지수의 분포는 기업별로 근소한 차이를 보이고 있지만 0.4에서 0.5사이 값에 집중되어 있는 거의 동일한 형태를 이루고 있다.



<Figure 9> Positive Index Distribution of Sentiment Dictionary for Each Company

4.3 감성분석

구축한 감성사전을 이용하여 각 기업별로 수집한 뉴스기사와 트윗들의 감성분석을 진행한다. 감성분

석과정을 거쳐 생성된 데이터 셋은 일별 뉴스기사의 빈도수와 긍정지수, 일별 트윗의 빈도수와 긍정지수, 익일 주가변동 항목으로 구성되며, <Table 3>에 데이터 일부를 나타내었다.

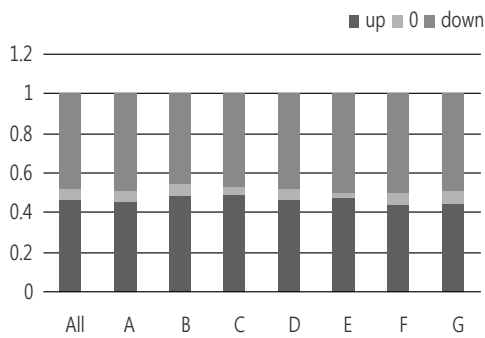
<Table 3> Portion of Data Set

| Attribute set | | | | Class label |
|---------------|----------------|-----------|----------------|----------------------|
| News | | SNS | | Next day stock price |
| frequency | positive index | frequency | positive index | |
| 117 | 0.442 | 369 | 0.438 | up |
| 183 | 0.440 | 871 | 0.424 | down |
| 212 | 0.440 | 542 | 0.430 | down |
| 265 | 0.440 | 651 | 0.435 | 0 |
| 260 | 0.445 | 492 | 0.435 | up |

수집한 뉴스기사와 트윗들은 매일 발생하지만 주식시장은 매일 개장하는 것이 아니기 때문에 익일 주가변동 데이터가 존재하지 않는 날이 발생하게 된다. 때문에 익일 주가변동 데이터가 존재하지 않을 경우 익일 주가변동 데이터가 존재할 때까지 데이터들을 합산하여 데이터 셋을 생성하였다. 예를 들어 금요일에 게재된 뉴스기사의 경우 다음날인 토요일의 주가에 영향을 미친다고 가정했지만 토요일엔 주식시장이 개장되지 않는다. 때문에 금요일에 게재된 뉴스기사는 토요일과 일요일의 기사들과 같이 월요일 주가에 영향을 미친다고 가정하고 이들 데이터를 합산하여 하나의 데이터로 구성하였다.

2013년 한 해 동안 주식시장은 247일 개장되었고, 이 중 1월 2일을 제외한 246일에 대한 감성분석을 진행하였다. 때문에 감성분석 과정을 거쳐 생성된 대부분의 데이터 셋은 246개의 데이터로 구성되어 있지만, G기업의 경우 경제관련 뉴스기사나 트위터 데이터가 발생하지 않은 날들이 포함되어 있어 이를 제외한 238개의 데이터로 구성하였다. 또한 일일 감성지수를 계산하는 과정에서 텍스트의 감성지수가 0인 데이터들은 제외하고 일일 감성지

수를 계산하여 데이터를 생성하였다. 생성한 전체 데이터 중 익일주가가 상승한 경우는 798개, 하락한 경우는 828개, 전날과 동일한 경우는 88개로 나타났고, 전체 데이터와 기업별 데이터의 분포를 <Figure 10>에 표현하였다. 주가가 전날과 동일한 경우의 데이터 빈도수가 비교적 적지만 주가가 상승한 경우와 하락한 경우의 데이터 수는 한쪽에 치우치지 않게 분포되어 있는 것으로 판단된다.



<Figure 10> Distribution of Data for Each Class Label

4.4 기계학습

본 논문에서는 기계학습 과정을 거쳐 예측 모형의 성능을 측정하기 위해 뉴질랜드 와이카토 대학에서 개발한 오픈소스 소프트웨어인 'WEKA'를 사용하였다. 2013년 1월 2일부터 9월 1일까지의 데이터는 예측 모형의 성능을 검증하기 위한 데이터(training set, test set)로 사용하였고, 9월 2일부터 12월 30일까지의 데이터는 예측 모형간의 성능을 비교하기 위한 데이터(validation set)로 사용하였다. 예측 모형의 성능을 검증하기 위한 데이터를 학습 데이터 셋과 평가 데이터 셋으로 나누어 기계학습을 수행하기에는 충분하지 않다고 판단하여 교차 검증으로 예측 모형의 성능을 측정하였다(Ian, Eibe, and Mark, 2011). 분류기법은 로지스틱 회귀분석, 베이지안 네트워크, 인공신경망, 서포트 벡터 머신 총 4가지 방법을 이용하였고, 성능 측정은 10중 교차 검증(10-fold cross-validation) 방식을 사용하였다.

분류기법 중 서포트 벡터 머신의 경우 SMO(Sequential Minimal Optimization) 알고리즘을 사용하였고 Polynomial kernel을 사용하였다. SMO 알고리즘의 complexity 파라미터는 1부터 10까지의 수 중 최적의 성능을 내는 파라미터를 시행착오법을 통해 결정하였다. 인공신경망의 경우 4개의 입력노드와 3개의 출력노드를 사용하였으며 은닉층의 수는 서포트 벡터 머신의 경우와 같이 시행착오법을 통해 결정하였다. 각 기업별 complexity 파라미터와 은닉층의 수는 <Table 4>에 나타내었다.

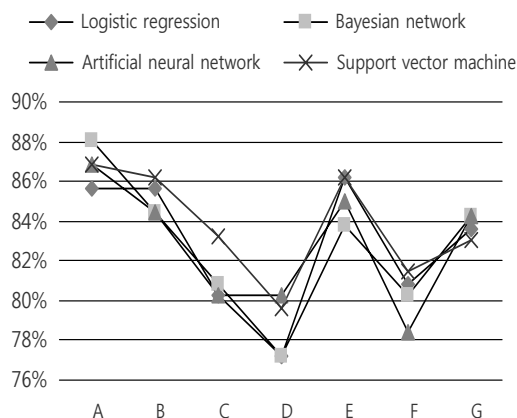
<Table 4> Parameters for Each Company

| Company | Complexity parameter | Hidden layer |
|---------|----------------------|--------------|
| A | 3 | 2 |
| B | 6 | 1 |
| C | 2 | 1 |
| D | 5 | 5 |
| E | 9 | 1 |
| F | 3 | 9 |
| G | 1 | 1 |

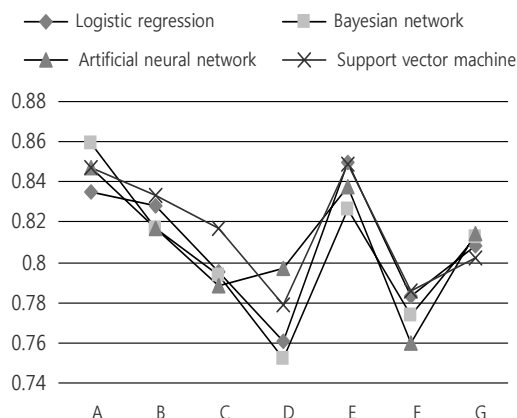
4.5 결과 분석

학습 데이터 셋과 평가 데이터 셋을 이용한 10중 교차 검증 실험 결과 7개의 기업 모두 다수의 방법에서 80%가 넘는 정확도를 보였고 0.8에 가까운 F1 score를 나타내었다. 또한 예측 모형들의 정확도와 F1 score의 산술평균값을 비교한 결과 서포트 벡터 머신을 이용하였을 때 정확도 83.8%와 F1 score 0.8161로 가장 좋은 성능을 보였다. 기업별 정확도와 F1 score는 <Figure 11>와 <Figure 12>에 나타내었고 수치화된 결과를 <Table 5>에 나타내었다.

검증 데이터 셋을 이용한 실험 결과 5개의 기업(B, C, D, F, G)에서는 다수의 방법에서 80%가 넘는 정확도와 0.8에 가까운 F1 score를 보였고, 2개의 기업(A, E)에서는 75% 이상의 정확도와 0.75에 가까운 F1 score를 나타내었다. 또한 예측 모형들의 정확도와 F1 score의 산술평균값을 비교한 결과



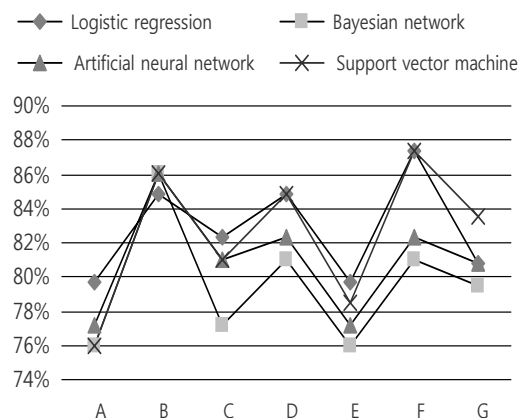
〈Figure 11〉 Accuracy of 10-Fold Cross-Validation



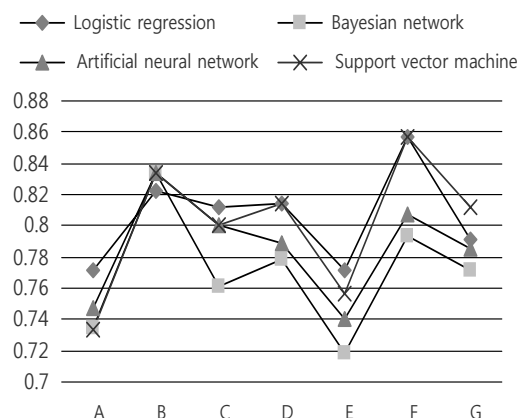
〈Figure 12〉 F1 Score of 10-Fold Cross-Validation

로지스틱 회귀분석(정확도-82.7%, F1 score-0.805)과 서포트 벡터 머신(정확도-82.4%, F1 score-0.801)을 이용하였을 때 비슷한 수준으로 좋은 성능을 보였다. 검증 데이터 셋을 이용한 실험의 기업별 정확도와 F1 score는 〈Figure 13〉와 〈Figure 14〉에 나타내었고 수치화된 결과를 〈Table 6〉에 나타내었다.

본 논문의 실험결과는 주가예측을 위해 뉴스 콘텐츠와 오피니언 마이닝을 이용하거나(Kim, 2013; Chun, 2013), 데이터 마이닝 기법들을 이용한 기존 연구들(Ahn, 2010; Park and Shin, 2011; Sagong, 2012)의 결과와 비교했을 때 상당히 진전된 결과로 볼 수 있다.



〈Figure 13〉 Accuracy of Using Validation Set



〈Figure 14〉 F1 Score of Using Validation Set

5. 결 론

본 논문에서는 개별 기업 주가의 등락을 예측하기 위해 뉴스기사와 SNS 데이터를 바탕으로 감성분석과 기계학습기법을 사용한 예측 모형을 제시하였다. 데이터의 수집은 자체적으로 개발한 프로그램으로 수행하였으며, 수집한 데이터 중 뉴스 기사를 이용해 주식도메인에 맞는 감성사전을 구축하였다. 구축한 감성사전을 사용한 감성분석 과정을 거쳐 기계학습을 위한 데이터를 생성하였고, 이를 이용해 다수의 기법들을 기계학습 시켜 각각의 예측 모형을 생성하였다. 생성한 예측 모형의 성능을 측

〈Table 5〉 Results of 10-Fold Cross Validation

| Classifier | Measure | A | B | C | D | E | F | G |
|---------------------------|-------------|---------|---------|---------|---------|---------|---------|---------|
| Logistic regression | Accuracy(%) | 85.6287 | 85.6287 | 80.2395 | 77.2455 | 86.2275 | 80.8383 | 83.6364 |
| | F1 score | 0.835 | 0.828 | 0.795 | 0.761 | 0.85 | 0.783 | 0.808 |
| Bayesian network | Accuracy(%) | 88.024 | 84.4311 | 80.8383 | 77.2455 | 83.8323 | 80.2395 | 84.2424 |
| | F1 score | 0.859 | 0.817 | 0.794 | 0.752 | 0.826 | 0.774 | 0.813 |
| Artificial neural network | Accuracy(%) | 86.8263 | 84.4311 | 80.2395 | 80.2395 | 85.0299 | 78.4431 | 84.2424 |
| | F1 score | 0.847 | 0.817 | 0.788 | 0.797 | 0.838 | 0.76 | 0.814 |
| Support vector machine | Accuracy(%) | 86.8263 | 86.2275 | 83.2335 | 79.6407 | 86.2275 | 81.4371 | 83.0303 |
| | F1 score | 0.847 | 0.833 | 0.817 | 0.779 | 0.849 | 0.786 | 0.802 |

〈Table 6〉 Results of Using Validation Set

| Classifier | Measure | A | B | C | D | E | F | G |
|---------------------------|-------------|---------|---------|---------|---------|---------|---------|---------|
| Logistic regression | Accuracy(%) | 79.7468 | 84.8101 | 82.2785 | 84.8101 | 79.7468 | 87.3418 | 80.8219 |
| | F1 score | 0.771 | 0.822 | 0.812 | 0.814 | 0.771 | 0.857 | 0.791 |
| Bayesian network | Accuracy(%) | 75.9494 | 86.0759 | 77.2152 | 81.0127 | 75.9494 | 81.0127 | 79.4521 |
| | F1 score | 0.735 | 0.834 | 0.761 | 0.779 | 0.719 | 0.794 | 0.771 |
| Artificial neural network | Accuracy(%) | 77.2152 | 86.0759 | 81.0127 | 82.2785 | 77.2152 | 82.2785 | 80.8219 |
| | F1 score | 0.747 | 0.834 | 0.8 | 0.789 | 0.74 | 0.807 | 0.785 |
| Support vector machine | Accuracy(%) | 75.9494 | 86.0759 | 81.0127 | 84.8101 | 78.481 | 87.3418 | 83.5616 |
| | F1 score | 0.734 | 0.834 | 0.8 | 0.814 | 0.756 | 0.857 | 0.812 |

정하기 위한 실험을 진행하였으며, 기존의 연구들에 비해 개선된 결과를 도출 하였다.

본 연구를 통해 최근 빅데이터의 대두로 주목받고 있는 SNS 데이터를 추가예측 분야에서 활용할 수 있는 한 가지 방안을 제시하였다. 또한 기계학습이론을 바탕으로 하는 예측기법에 필요한 데이터를 정제하는 과정에서 감성분석 기법을 이용할 수 있음을 나타내었다. 반면, 본 연구에서 제안한 예측 모형은 뉴스기사와 SNS 상에서 충분히 언급되지 않는 기업들에 대해서는 그 정확도가 떨어질 것으로 예상된다. 또한 감성사전을 구축하기 위해 사용한 데이터와 기계학습에 이용하는 데이터가 같아 종속변수로 설정한 데이터가 독립변수에 영향을 주었을 가능성을 가지고 있다. 향후 연구에서는 뉴스기사와 SNS 데이터의 양과 예측정확도와의 상관관계에 대한 연구와 감성사전을 위한 데이터와 기계학습을 위한 데이터를 구분하여 연구

를 진행해 예측 모형의 한계점을 보완할 수 있는 대책을 강구할 필요가 있을 것이다.

References

- Ahn, S., "Stock Prediction Using News Text Mining and Time Series Analysis", M.S. thesis, *The Graduate School of Engineering*, Yonsei Univ., Seoul, Korea, 2010.
- Bollen, J., H. Mao, and X. Zeng, "Twitter mood predicts the stock market", *Journal of Computational Science*, Vol.2, No.1, 2011, 1-8.
- Ian, H., F. Eibe, and A. Mark, *Data Mining*, Morgan Kaufmann, Burlington, 2011.
- Kim, K., G. Lee, and S., Lee, "A Comparative Analysis of Artificial Intelligence System and Ohlson model for IPO firm's Stock Price

- Evaluation", *The Journal of Digital Policy and Management*, Vol.11, No.5, 2013, 145-158.
- Kim, S., D. Nam, H. Jo, and S. Kim, "A Study on the Relation of Web News and Stock Price", *Korea Society of IT Services*, Vol. 11, No.3, 2012, 191-203.
- Kim, T., W. Jung, and S. Lee, "The Analysis on the Relationship between Firms' Exposures to SNS and Stock Prices in Korea", *Asia Pacific Journal of Information Systems*, Vol. 24, No.2, 2014, 233-253.
- Kim, Y., "News Big Data Opinion Mining Model for Predicting KOSPI Movement", *Ph.D. thesis, Graduate School of Business IT*, Kookmin Univ., Seoul, Korea, 2013.
- Kim, Y., N. Kim, and S. Jung, "Stock-Index Invest Model Using News Big Data Opinion Mining", *Journal of Intelligence and Information Systems*, Vol.18, No.2, 2012, 143-156.
- Park, K. and H. Shin, "Stock Price Prediction Based on Time Series Network", *Journal of the Korean Operations Research and Management Science Society*, Vol.28, No.1, 2011, 53-60.
- Sagong, J., "A Study on Predicting Stock Price Based on Data Mining Techniques", M.S. thesis, *Dept. Data Science*, Inje Univ., Gimhae, Korea, 2012.
- Shim, K. and J. Yang, "High Speed Korean Morphological Analysis based on Adjacency Condition Check", *Korean Institute of Information Scientists and Engineers*, Vol.31, No.1, 2004, 89-99.
- Song, C., "News and Financial Prices", *International Economic Journal*, Vol.8, No.3, 2002, 1-34.
- Song, J. and S. Lee, "Automatic Construction of Positive/Negative Feature-Predicate Dictionary for Polarity Classification of Product Reviews", *Korean Institute of Information Scientists and Engineers*, Vol.38, No.3, 2011, 157-168.
- Yu, E., Y. Kim, N. Kim, and S., Jeong, "Predicting the Direction of the Stock Index by Using a Domain-Specific Sentiment Dictionary", *Journal of Intelligence and Information Systems*, Vol.19, No.1, 2013, 95-110.
- Chun, S., "뉴스 콘텐츠의 오피니언 마이닝을 통한 매체별 주가상승 예측정확도 비교 연구", M. S. thesis, Graduate School of Business IT, Kookmin Univ., Seoul, Korea, 2013.

◆ About the Authors ◆



Dongyoung Kim (dy.kim@ssu.ac.kr)

Dongyoung Kim is currently enrolled in a Masters degree in Graduate School of Software, Soongsil University. He received his bachelor's degree in Computer Science and Engineering and Industrial Engineering from Soongsil University in 2013. His research interests are in areas of Data Mining, Big Data, Software Engineering, Cloud Computing, and Mobile.



Jeawon Park (jwpark@ssu.ac.kr)

Jeawon Park received the Ph.D. degree in Computer Science from Soongsil University in Korea, 2011. He is a professor at Graduate School of Software, Soongsil University. His research interests are in areas of Software Testing, Software Process, Web Services, and Project Management



Jaehyun Choi (jaehyun@ssu.ac.kr)

Jaehyun Choi received the Ph.D. degree in Computer Science from Soongsil University in Korea, 2011. He is a professor at Graduate School of Software, Soongsil University. His research interests are in areas of Data Processing, Service Engineering, Software Engineering, and Text Mining