

서 론

최근에는 매일 신문이나 라디오 같은 미디어로부터 인터넷과 같은 전자 매체까지 다양한 경로에서 정보를 습득할 수 있게 되었다. 특히, 인터넷의 확산을 통해 이러한 여러 형태의 정보를 하나로 통합하여 사용자에게 제공함으로써 보다 쉽고 편리하게 정보를 얻고 활용하는 단계에 이르렀다. 이와 같이 인터넷이 폭 넓게 보급되어 온라인(on-line)상에서 얻을 수 있는 텍스트(text) 정보의 양이 급증함에 따라 이러한 거대한 텍스트 집합으로부터 의미 있는 지식을 찾아내는 작업은 많은 분야에서 매우 다양하게 요구되고 있다.

텍스트로부터 추출할 수 있는 유용한 정보 중에 하나가 작자가 해당 문서의 주제에 대해 표현한 감정 혹은 의견(sentiment or opinion)이다[1]. 예를 들어, 기업은 자신들의 상품에 대한 소비자들의 평판을 아는 것이 상품개발과 마케팅을 위한 유용한 정보로 사용될 수 있으며, 또한 영화배급사는 영화에 대한 관객들의 평판을 파악하여 개봉관의 수를 적절하게 조절할 수 있을 것이다. 전통적으로 이러한 평판은 비싼 비용을 지불하고 조사(survey)되어 왔으나, 근래에 들어 인터넷을 통해 상품에 대한 평가(review)를 온라인으로 손쉽게 수집할 수 있게 됨에 따라, 텍스트 문서들에서 자동으로 감정과 의견을 추출할 수 있다면, 저비용으로 그리고 자동으로 의견 조사가 가능할 것이다. 최근 외국에서는 이러한 작자의 의견이 담겨있는 문서로부터 작자의 감정을 자동으로 판별하는 연구가 활발히 진행되고 있다. 전통적인 문서 분류가 문서의 주제(topic)에 초점을 맞추었다면 감정 분류(sentiment classification)는 저자의 주제에 대한 긍정 감정과 부정 감정에 초점을 맞춘 연구 분야로서, 고객 평가의 요약, 공공 의견 조사, 고객 성향 분석 등의 응용 영역을 가지고 있다.

일반적인 문서 분류는 사람이 문서에 나타난 자질을 보고 인식하여 정해진 범주로 분류하는 과정을 수학적으로 모델링하여 기계가 동일한 과정으로 학습하여 문서를 분류하도록 하는 것[2]이다. 효과적인 문서 분류를 위해서 가장 중심이 되어야 하는 부분이 자질의 선정 방법[3]이다. 문서 감정 분류를 위한 효과적인 감정 자질의 선정을 위해 고려해야 할 사항은 감정 분류는 문서에 나타나는 단어의 형태뿐만 아니라 단어의 의미에도 기반 해야 한다는 점이다. 감정 분류는 긍정과 부

정의 감정에 초점을 두기 때문에 먼저 이를 가장 잘 표현하는 기본적인 단어인 감정 자질의 생성이 중요하다. 인간이 사용하는 말들 중엔 긍정과 감정을 나타내는 표현들이 있다. 그리고 긍정인 문서에선 긍정적인 표현이 많이 나오고, 부정인 문서에서는 부정적인 표현이 많이 나온다. 이러한 단어들을 잘 판단할 수 있다면 감정 분류를 하는데 도움이 될 것이다. 하지만, 인간이 사용하는 모든 긍정적, 부정적 표현을 다 찾아내는 일은 쉬운 일이 아니다. 그러므로, 영어권 선행 연구[4,5]를 바탕으로 대표적인 긍정, 부정을 나타내는 단어를 통해서 그 단어의 유의어를 모은다면, 각 감정을 나타내는 단어들을 모을 수 있을 것이라 판단하고 연구를 수행하였다. 자질로부터 문서의 감정 분류를 위해서 사용될 충분한 양의 감정 자질을 추출하기 위해서 사전상의 유의어 및 반의어의 의미적 정보를 활용하여 단어의 의미 확장을 시도하였으나, 한국어 사전의 유의어 및 반의어의 정보가 빈약하여 충분한 양의 감정 자질을 얻을 수가 없었다. 대안으로 영어 단어 시소러스 유의어 정보를 이용하여 단어를 확장하고 이를 한영사전을 통해 번역하여 감정 자질을 추출하였다. 확장된 감정 자질을 이용하여 기존의 문서 분류 기법을 적용하고 문서에 대한 감정을 분류하여 추출된 감정 자질의 유용성을 평가하였다.

본 논문에서는 한국어 문서의 감정을 분류하기 위한 효과적인 감정 자질 추출 방법을 제안하고, 이를 통해 문서를 표현하여 기계학습 기법 중 하나인 지지 벡터 기계를 사용하여 성능을 평가하여 추출된 감정 자질의 유용성을 평가하였다.

본 논문의 구성은 다음과 같다. 2장에서는 앞서 연구된 관련 연구에 대해 언급하였으며, 3장에서는 본 논문에서 제안하는 문서 감정 분류를 위한 자질 추출 방법과 자질들의 가중치 선정 방법에 대해 논의한다. 4장에서는 본 논문에서 제안하는 방법의 유용성을 평가하고 마지막 장에서는 결론 및 향후 과제에 대해서 기술한다.

관련 연구

상품에 대한 평가와 영화에 대한 관객들의 평론에서 나타나는 주관적, 감정적 표현을 여러 기계 학습 방법과 자연어 처리 기술을 통해 문서를 분류하는 연구가

진행되고 있다[3,4,6,7].

특히 문서 감정 분류 시스템은 문서 분류의 특화된 분야이기 때문에 문서 분류에서 사용되어온 여러 가지 기계 학습 기법들이 문서 감정 분류에도 적용되어 왔다. 영화 평론과 상품 평가와 같은 특정 영역에서 나타나는 감정적 표현을 나이브 베이즈(naive bayes), 최대 엔트로피(maximum entropy), 지지 벡터 기계(support vector machine) 등의 기계 학습을 통해 문서를 긍정과 부정의 범주로 분류하는 연구가 진행되어 왔다[3,4,6,7]. 감정 분류에 대한 응용 영역으로는 먼저 상품에 대한 고객들의 평가에 들어있는 감정을 분류하여 내용을 요약하는 응용 분야(customer review) [6,8]와 공공의 의견을 조사하여 요약하는 응용 분야(public opinion survey)[9,10] 그리고, 고객들의 성향을 분석(trend analysis)[11]하는 분야 등 폭넓은 응용 영역을 가지고 있다.

또한, 분류의 대상이 문서뿐만 아니라, 문장[12,13], 구(phrase)[14,15], 토론의 연결가[16], 그리고 문장의 감정 패턴 분석을 통해 문장의 여러 감정적 표현을 인식하고 분류하는 연구도 수행되었다[5,17].

그리고, 국내에서는 최근 많은 네티즌들이 사용하고 있는 메신저 프로그램내의 대화 내용의 감정을 파악하여 자동으로 그림말을 붙여주는 시스템과 최근 많은 문제가 되고 있는 인터넷의 악성 댓글을 판별하는 시스템[18]에 관한 연구도 진행되고 있다. 하지만, 지금까지 국내의 연구결과가 아직 기초적인 연구에 머무르고 있기 때문에 영어권 선행 연구 결과[4,5]를 바탕으로 한국어에 적용하였다.

영어권의 연구와 본 연구는 언어와 실험 데이터도 다르기 때문에 직접적인 비교가 어렵다. 하지만, 한국어권에서는 공개되어 있는 WordNet[19]와 같은 어휘자원이 없기 때문에 영어권 시소러스와 공개된 전자사전을 이용해서 어떻게 효과적으로 한국어 감정자질을 추출하고, 이를 바탕으로 한국어 감정분류를 수행하는지에 대해서 연구하였다.

한국어 문서 감정 분류 시스템

일반적인 문서 분류 시스템에서의 자질 선정 방법은 학습 문서에서 형태소 분석

을 통해 내용어(content word)를 추출하고 추출된 대상 자질에 대해 가중치를 부여하는 것이 일반적이다. 하지만 아래의 [그림 1]처럼 감정 분류 시스템에서는 의미적 문서 분류를 위해서 먼저 긍정과 부정을 나타내는 어휘 즉, 감정 자질(emotion feature)들을 따로 추출하여야 한다.

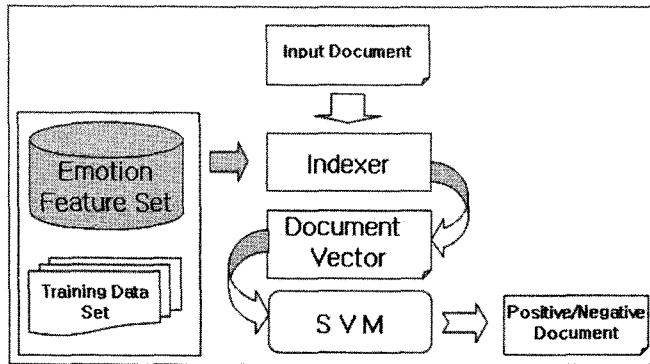


그림 1. 한국어 감정 분류 시스템 구성도

자질 추출 및 확장

감정 자질들과 일반적인 정보검색에서 사용되는 어휘들과의 가장 큰 차이점은 정보 검색에서 사용되는 어휘들의 품사는 명사, 동사가 중요하게 사용되는 반면 감정 분류에서는 형용사, 부사 등도 중요하게 사용된다는 점이다. 이러한 감정 어휘 집합을 추출하기 위해서는 여러 가지 어휘자원들이 필요한데 외국의 연구에서는 WordNet과 같은 어휘 의미망이 많이 사용되고 있다. 한국어 감정 자질들의 확장을 위하여 한국어 사전을 파싱한 후 DB(data base)를 구축하여 동의어, 반의어 정보를 획득하고자 하였으나, 부정(13개), 긍정(12개)의 감정 자질만 획득하여, 원하는 결과를 얻을 수가 없었다. 한국어 사전에서는 어휘의 동의어와 반의어의 비중이 낮다고 판단하고 영어단어 시소러스 유의어 정보¹⁾[20]를 이용하였다.

1) 폴린스 코빌드에서 출간된 유의어사전(Thesaurus)으로 최신 유의어사전을 사용하고 있으며, 가장 유용한 유의어를 설명 맨 앞에 제공하고 반의어가 유의어와 함께 제공. 총 11

한국어 감정 자질 추출을 위하여 본 논문에서는 한국어에서 긍정과 부정을 나타내는 대표 어휘를 영어권 선행 연구 결과[4,5]를 바탕으로 [표 1]과 같이 대표 어휘를 선정하고 이들 단어들의 종자 어휘로 사용하여 한국어 감정 자질을 추출하고 확장한다.

표 1. 긍정/부정 영어 단어 대표 어휘

긍정	good, correct, positive, excellent, nice, fortunate, superior
부정	bad, nasty, negative, poor, unfortunate, wrong, inferior

본 논문에서 사용하는 자질들은 생성 방식에 따라 다음의 4가지로 나눌 수 있다.

내용어(Content Word) - 자질 1

일반적인 정보 검색에서 사용되는, 형태소 분석 결과로 얻어진 명사, 동사를 사용하여 자질들을 생성하였다.

감정 내용어(Emotion Content Word) - 자질 2

한국어 감정 분류를 위해 중요하게 사용되는, 형태소 분석 결과로 얻어진 형용사, 부사를 자질1에 포함하여 자질들을 생성하였다.

감정 자질(Emotion Feature) - 자질 3

선정된 대표 어휘를 대상으로 영어 단어 유의어 시소러스 정보를 이용하여 형용사, 부사를 포함한 어휘의 의미를 [그림 2]와 같이 대표 유의어(예:harmful)를 순차적으로 확장하였다. 이러한 확장 방법을 각 감정 자질 목록이 더 이상 추가되는 단어가 없을 때까지 수행하여 감정 자질들을 생성하였다. 그 후, 사람이 직접 영한 사전을 이용하여 적절한 한국어 감정 자질만 선정하여 긍정 감정 자질(861개)과 부정 감정 자질(1,834개)들을 생성하였다.

만개 이상의 단어, 숙어가 설명되어 있음.

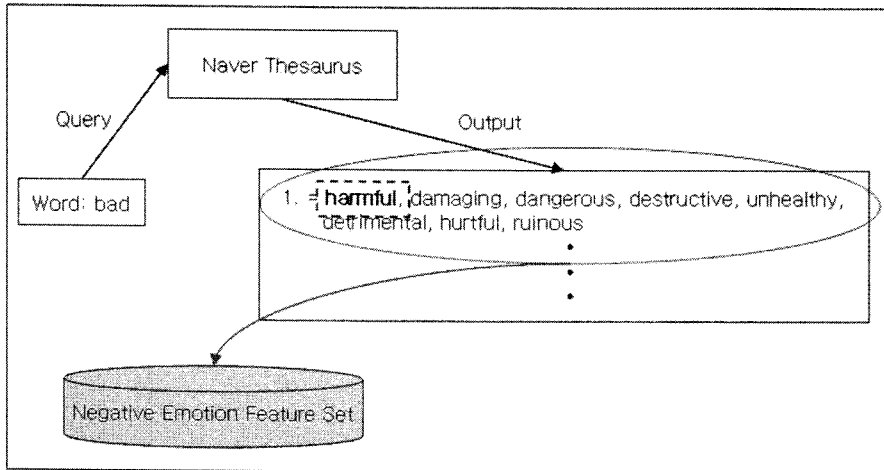


그림 2. 부정 단어(bad)의 확장 예

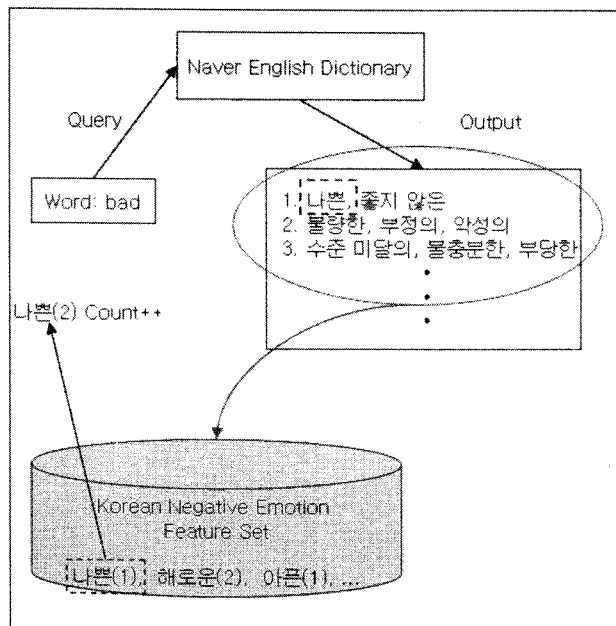


그림 3. 영한 번역시에 한국어 감정 자질 Count의 예

균형 감정 자질(Balanced Emotion Feature) - 자질4

자질3에서 부정 감정 자질(1834개)이 긍정 감정 자질(861개)에 비해 약 2.1배 많
이 생성되었기 때문에, 4.2절에 소개 될 기준 시스템(baseline system)의 실험 결과를
기반하여, [그림 3]과 같이 영한 번역 작업 시에 출현한 한국어 자질의 횟수
(Occur_cnt)가 많은 자질들을 우선하여 감정 자질들을 선정하였다. 2번 이상 출현한
자질들을 선정했을 때 부정 자질의 수가 844개로 긍정 자질(861개)의 수와 거의 균
형을 이루었기 때문에 2번 이상 출현한 자질들을 균형 감정 자질로 선택하였다.

최종적으로 생성된 감정 자질[표 2]와 자질의 출현 횟수는 문서 벡터의 생성 시
각 자질의 가중치 계산에 적용되어 사용된다.

표 2. 실험에서 사용된 자질의 구성

자질 구분	내용
내용어(자질1)	형태소 분석으로 추출한 내용어 (명사, 동사)
감정 내용어(자질2)	형태소 분석으로 추출한 내용어 (명사, 동사, 형용사, 부사)
감정 자질(자질3)	대표 어휘의 유의어 단어 집합 {긍정:861개 / 부정:1834개}
균형 감정 자질(자질4)	자질3의 부정 자질을 줄인 단어 {긍정:861개 / 부정:844개}

문서 표현 및 지지 벡터 기계(SVM)

입력 문서를 형태소 분석[21] 후, 앞 단계에서 선택된 자질을 기준으로 아래식의
TF-IDF 가중치 기법과 감정 자질의 자질 추출 시 출현 횟수(Occur-cnt)를 사용하여
가중치를 계산한다.

TF-IDF 가중치 기법은 식 (2)와 같이 문서에 어휘 t 가 나타난 어휘 빈도수(tf:term
frequency) tf_t 와 역 문서 빈도수(idf:inverse document frequency, 식(1))의 곱으로 나타
낸다.

$$idf_t = \log_2 \frac{N}{df_t} \tag{1}$$

여기서 N 은 전체 문서의 수이며, df_t 는 어휘 t 가 출현한 문서의 수이다.

$$weight_t = tf_t \cdot idf_t \quad (2)$$

각 감정 자질(t)마다 출현 횟수(Occur-cnt)를 이용하여 식 (3)을 이용하여 가중치 (W_t)를 추정한다. 각각의 긍정과 부정의 의미에 대하여 많이 출현한 단어가 의미가 더욱 강하므로 높은 가중치를 부여하며, 전체에 비례하여 평준화 한다.

$$W_t = \frac{t_{cnt}}{N_{maxcnt}} \quad (3)$$

t 는 감정 자질 단어이며, t_{cnt} 는 t 가 영한 번역시에 출현한 횟수이며, N_{maxcnt} 는 부정 감정 자질의 단어들 중 최대로 출현한 단어의 출현 횟수이다. 각 감정 자질의 가중치는 위의 식 (2)와 (3)을 이용하여 책정한다.

다음으로, 문서 분류기는 지지 벡터 기계를 사용하였다.

지지 벡터 기계는 두 개의 범주를 구분하는 문제를 해결하기 위해 1995년에 Vapnik에 의해 소개된 학습 기법으로 두 개의 클래스의 구성 데이터들을 가장 잘 분리해 낼 수 있는 초평면(optimal hyperplane)을 찾는 모델이다[22]. 지지 벡터 기계에서의 초평면은 식 (4)와 같이 나타낼 수 있다.

$$\vec{w} \cdot \vec{x} - b = 0 \quad (4)$$

여기서 \vec{x} 는 분류하고자 하는 문서의 벡터이며 \vec{w} 와 b 는 학습 데이터로부터 학습되어 나온 결과이다. 학습 문서 집합을 $D = \{(y_i, \vec{x}_i)\}$ 과 같이 나타냈을 때, 각각의 학습 문서 벡터(\vec{x}_i)가 임의의 범주에 속한 문서이면 y_i 의 값에 +1을 할당하

고, 범주에 속하지 않은 문서에는 -1을 할당한다. 결국 지지 벡터 기계는 식 (5)와 (6)을 만족시키는 \vec{w} 와 b 를 찾는 문제이다.

$$\vec{w} \cdot \vec{x}_i - b \geq +1 \text{ for } y_i = +1 \quad (5)$$

$$\vec{w} \cdot \vec{x}_i - b \leq -1 \text{ for } y_i = -1 \quad (6)$$

위의 수식들에 따르면 두 개의 클래스를 구분하는 초평면은 무수히 많이 존재하는데, 이들 초평면들 중에서 최적의 초평면은 두 클래스를 구분하는 거리(margin)가 최대가 되는 초평면을 정의할 수 있다. [그림 4]는 벡터를 2차원으로 표현한 한 예로서, 각 x축, y축은 자질들을 나타낸다. 실선은 두 개의 클래스를 구분하는 초평면이고, 점선은 이들 초평면들 중에서 최적의 초평면으로 두 클래스를 구분하는 거리(margin)가 최대가 되는 초평면을 나타낸다. 두 클래스를 나누는 초평면 중에서 초평면들 사이의 거리(d)가 최대인 초평면을 보여주고 있다.

지지 벡터 기계는 직선으로 나눌 수 있는 문제(linearly separable problem)에 사용되는 알고리즘이지만, 다차원의 부드러운 곡선을 이용하여 초평면을 설정하거나, 실제 데이터 벡터를 새로운 자질을 포함한 새로운 벡터 공간에 매핑하는 방법을 통해서 직선으로 나눌 수 없는 문제도 해결 할 수 있다. 지지 벡터 기계 모델을 문서 범주화에 적용되어 좋은 성능을 보여 왔다[2].

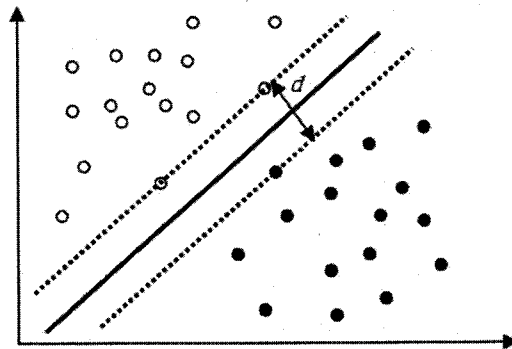


그림 4. 초평면의 거리(Margin)

실험 및 결과

실험 데이터 구성과 실험 환경

실험 데이터

실험에 사용된 데이터는 [표 3]과 같이 총 2,479개의 감정 문서이며, 현재 공개되어 있는 감정 분류를 위한 실험 데이터가 없는 관계로 3명의 학생이 3개의 분야를 나누어 수집하였다. 사이트 상²⁾에서 나타난 {찬성/반대}, {추천/비추천}의 정보를 바탕으로 직접 문서를 읽고 2명 이상이 판단한 감정을 문서의 감정으로 하여 데이터를 구축하였다. 최종적으로 구축된 말뭉치는 신문 기사 729개, 제품 리뷰 395개, 영화 리뷰 1,355개의 문서이다.

표 3. 실험에 사용한 테스트 말뭉치

분야	긍정	부정	총합
신문 기사	417	312	729
제품 리뷰	205	190	395
영화 리뷰	703	652	1,355
총합	1,325	1,154	2,479

성능 평가 방법

본 논문에서는 다양한 자질 단어와 가중치 책정 방법을 사용하여 10-fold cross validation 방법으로 실험을 하였으며, 인터넷 사이트상에서 수집된 문서 집합의 평가 방법으로는 정보 검색 분야에서 일반적으로 사용되는 정확율(precision)과 재현율(recall)을 사용하였다.

정확율은 다음 식 (7)과 같이 표현된다.

2) <http://cineast.kr/>, <http://news.naver.com/>, <http://www.gmarket.co.kr/>

$$\text{정확율} = \frac{\text{시스템에 의해 판단된 적합 문서수}}{\text{시스템이 적합하다고 판단한 문서수}} \quad (7)$$

재현율은 다음 식 (8)과 같이 표현된다.

$$\text{재현율} = \frac{\text{시스템에 의해 판단된 적합 문서수}}{\text{적합 문서수}} \quad (8)$$

정확율과 재현율을 하나의 값으로 표현해주기 위해서 다음 식 (9)와 같이 $F_1 - measure$ 를 사용하였다.

$$F_1(r, p) = \frac{2 \cdot r \cdot p}{r + p} \quad (9)$$

식 (9)에서 r 은 재현율에 해당하고 p 는 정확율에 해당한다. 아래 [표 4, 5, 6, 7, 8]은 $F_1 - measure$ 값을 표시하였다.

기준 시스템(Baseline System) 구성

본 실험에서는 Weka에서 제공된 SVM toolkit[23]을 사용하였으며, 다음의 자질과 분류 기법에 따라 다음과 같은 기준 시스템을 사용하였다.

자질1과 자질2의 TF-IDF 가중치 - SVM

일반적인 정보 검색 분야에서 사용하는 내용어 기반의 자질1과 이 자질들에 한국어 감정 분류를 위해 중요하게 작용하는 품사인 형용사, 부사를 포함한 자질2에 TF-IDF 가중치 기법을 사용하여 SVM 분류기로 분류하였다.

자질3과 자질4의 Occur-cnt 가중치를 이용한 분류기

문서에 출현한 감정 자질 단어의 Occur-cnt 가중치 식 (3)를 이용하여 식 (10)을 사용하여 분류하였다. 즉, 단순히 문서에 감정 자질의 가중치의 합이 긍정, 부정 중에 어느 쪽이 더 큰가를 사용해서 문서의 감정을 분류하는 가장 간단한 형태의

감정분류기이다. 앞으로 이 분류기를 Occur-sum 분류기로 표시한다.

$$\begin{aligned} & \text{if } \sum_{t \in d} \{W_t | t \in P\} - \sum_{t \in d} \{W_t | t \in N\} > 0, \\ & \text{then } d_{result} = \text{positive.} \\ & \text{else } d_{result} = \text{negative.} \end{aligned} \quad (10)$$

d는 입력 문서이며, t는 감정 자질 단어, P와 N은 긍정/부정 감정 자질들의 집합이다.

실험 결과

기준 시스템을 사용한 자질의 성능 및 특성 비교

[표 4]에서 알 수 있듯이, 형용사, 부사의 품사를 포함한 자질2의 성능이 자질1의 성능보다 6.27% 더 향상된 성능을 보였다. 이를 통해 한국어 문서 감정 분류를 위해서는 형용사, 부사의 품사 중요한 역할을 한다는 것을 알 수 있었다.

표 4. 기준 시스템(SVM) 성능

분류기	가중치 기법	자질	긍정	부정	평균
SVM	TF-IDF	자질1	72.80	67.60	70.20
		자질2	77.18	75.75	76.47

다음으로, [표 5]에서 식 (10)을 이용한 분류기(Occur-sum)의 성능을 관찰하면, 자질4를 사용한 결과에 비해 자질3을 사용한 결과에서 부정의 성능이 긍정의 성능보다 많은 차이로 높다. 그 이유는 자질3의 감정 자질에서 긍정 자질보다 부정 자질의 수가 월등히 많았기 때문으로 분석할 수 있다. 이 결과로 긍정, 부정 감정 자질의 수의 균형이 필요하다는 것을 확인할 수 있다.

표 5. 기준 시스템(Occur) 성능

분류기	가중치 기법	자질	긍정	부정	평균
Occur-sum	Occur-Cnt	자질3	57.32	64.83	61.08
		자질4	60.03	61.83	60.93

향후 본 논문의 기준 시스템은 일반적인 정보 검색에서 사용되는 내용어 기반인 자질1을 사용하여 SVM 분류기로 분류한 성능을 기준으로 한다.

기준 시스템과 추출된 감정 자질(자질3)을 사용한 시스템의 비교 성능 평가

실험은 각 자질들과 가중치 기법들을 혼합하여 수행함으로써 그 성능을 비교하였다. 아래 표에서 모든 실험 결과에서 사용한 문서 분류기는 SVM이며, 자질과 가중치 기법에 따라서 성능을 비교하였다. [표 6, 7]에서 가중치 기법은 SVM을 사용하기 위한 문서 표현에서 각 자질의 가중치 방법을 나타내고 있으며, Occur-cnt는 식(3)에서 계산된 감정 자질의 가중치를 이용해서 문서를 표현하는 방식이며, TF-IDF는 자질의 가중치를 문서 분류와 정보 검색에서 일반적으로 쓰이는 TF-IDF 기법을 사용해서 계산하는 방식이며, tfidf • cnt는 TF-IDF 가중치와 Occur-cnt 가중치를 곱한 가중치를 사용한 방식을 표현한다.

표 6. 추출된 감정 자질(자질3)을 사용한 실험 결과

가중치기법	자질	긍정	부정	평균	비교
Baseline	자질1	72.80	67.60	70.20	-
Occur-cnt	자질3	79.29	74.96	77.12	+6.92
TF-IDF	자질3	84.15	82.00	83.07	+12.87
tfidf • cnt	자질3	80.25	76.40	78.32	+8.12

자질3의 감정 자질을 사용한 결과에서는 TF-IDF 가중치 기법을 사용했을 때 가장 높은 성능을 얻었으며, 그 성능은 단지 명사, 동사의 내용어만 사용한 자질1의 기준 시스템에 비해 12.87%의 성능 향상을 얻었다. 이 결과로 제안된 기법에

의한 감정 자질 추출이 한국어 문서 감정 분류에 더 유용하다는 것을 확인 할 수 있었다.

균형 감정 자질(자질4)과 불균형 감정 자질(자질3)의 비교 성능 평가

아래 [표 7]에서 보는 바와 같이 균형 감정 자질(자질4)로 실험한 결과, 자질3의 결과보다 1.23%의 더 나은 성능을 보였다. 이 결과로 균형 감정 자질이 필요하다는 것을 알 수 있다. Occur-cnt 가중치 기법을 기존의 TF-IDF 기법과 결합(tfidf · cnt)한 실험 결과는, 기존의 TF-IDF와 비교해서 성능향상의 결과를 얻을 수 없었다. 가장 성능이 좋은 가중치 기법은 자질4를 사용한 경우에도 TF-IDF 가중치 기법이였다.

표 7. 균형 감정 자질(자질4)을 사용한 실험 결과

가중치기법	자질	긍정	부정	평균	비교
TF-IDF	자질3	84.15	82.00	83.07	-
Occur-cnt	자질4	83.06	80.74	81.90	-1.17
TF-IDF	자질4	85.20	83.39	84.30	+1.23
tfidf · cnt	자질4	84.34	82.40	83.37	+0.3

자질간의 최종 성능 비교

결론적으로, 어휘 자원을 사용하여 단어의 의미에 기반한 자질의 선정이 일반 정보 검색에서 사용되는 내용어 기반의 자질 선정에 비해 더 유용하며, 형용사, 부사의 품사가 한국어 문서 감정 분류에 중요한 자질이 되는 것을 알 수 있다. 마지막으로, 균형 자질의 선정 역시 유용하다는 결론을 얻을 수 있다. 최종적으로 제안된 기법에 의해 추출된 감정 자질을 사용한 감정 분류 시스템은 기존 시스템보다 14.1%의 성능향상을 얻을 수 있었다. 또한, 형용사, 부사의 품사를 포함한 자질2를 사용하는 기존 시스템에 비해서는 7.83%의 성능 향상을 보였다.

표 8. 최종 성능 비교

구분	자질	긍정	부정	평균	비교
Baseline	자질1	72.80	67.60	70.20	-
Proposed Method	자질3	84.15	82.00	83.07	+12.87
<i>Proposed Method</i>	<i>자질4</i>	<i>85.20</i>	<i>83.39</i>	<i>84.30</i>	<i>+14.1</i>

결론 및 향후 과제

본 논문에서는 한국어 감정 분류 시스템을 위한 효과적인 자질 추출 방법을 제안하고, 그 유용성을 평가하였다. 한국어 문서 감정 분류를 위해서는 일반적인 정보 검색에서 사용하는 명사, 동사의 품사를 가진 내용어 뿐만 아니라, 형용사, 부사의 품사 역시 중요하며, 단지 형태소 분석을 통한 단어의 형태보다는 그 의미에 기반한 감정 자질의 생성 또한 중요하다. 본 논문에서 제안한 방법으로 추출한 감정 자질은 한국어 문서 감정 분류에 적용했을 때 약 84%의 높은 성능을 보였다. 이 성능은 기존의 문서 분류에서 사용한 자질을 사용한 경우보다 14.1%의 성능 향상을 얻은 것이다.

향후 연구로는 감정의 특징상 문서 전체가 아닌 특정 문장에서 강하게 표현되기 때문에 그 감정이 강하게 표현되는 문장에 관해 특별한 가중치 기법으로 가중치를 선정하는 방법에 관한 연구를 수행할 것이다. 또한, 감정 표현의 이중 부정에 관한 패턴을 파악하여 파악된 패턴을 적용할 수 있는 방법에 관한 연구도 수행할 것이다.

참고문헌

- [1] M. Rimon, "Sentiment Classification: Linguistic and Non-Linguistic Issues," Hebrew University.

- [2] T. Joachims, "Text Categorization with Support Vector Machines: Learning with Many relevant Features," In *Proceedings of the ECML*, pp.137-142, 1998.
- [3] J. Yi, T. Nasukawa, R. Bunescu, and W. Niblack, "Sentimental Analyzer : Extracting Sentiments about a Given Topic using Natural Language Processing Techniques," In *Proceedings of International Conference on Data Mining*, pp.427-434, 2003.
- [4] P.D. Turney and M.L. Littman, "Measuring Praise and Criticism: Inference of Semantic Orientation from Association," In *Proceedings of the ACM Transactions on Information Systems*, pp.315-346, 2003.
- [5] A. Esuli and F. Sebastiani, "Determining the Semantic Orientation of Terms through Gloss Classification," In *Proceedings of the CIKM*, pp.617-624, 2005.
- [6] B. Pang, L. Lee and S. Vaithyanathan, "Thumbs up? Sentiment Classification Using Machine Learning Techniques," In *Proceedings of the EMNLP*, pp.79-86, 2002.
- [7] N. Hiroshima, S. Yamada, O. Furuse and R. Kataoka, "Searching for Sentences Expressing Opinions by Using Declaratively Subjective Clues," In *Proceedings of the Workshop on Sentiment and Subjectivity in Text*, pp.39-46, 2006.
- [8] K. Dave, S. Lawrence, D. M. Pennock, "Mining the Peanut Gallery: Opinion Extraction and Semantic Classification of Product Reviews," In *Proceedings of the 12th WWW*, pp.519-528, 2003.
- [9] L.W. Ku, L.Y. Lee, T.H. Wu, and H.H. Chen, "Major Topic Detection and Its Application to Opinion Summarization," In *Proceedings of the ACM SIGIR*, pp.627-628, 2005.
- [10] S.M. Kim and E. Hovy, "Determining the Sentiment of Opinions," In *Proceedings of the COLING conference*, pp.1367-1373, 2004.
- [11] M. Hu and B. Liu, "Mining and Summarizing Customer Reviews," In *Proceedings of the KDD*, pp.168-177, 2004.
- [12] B. Pang and L. Lee, "A Sentimental Education: Sentiment Analysis Using Subjectivity Summarization Based on Minimum Cuts," In *Proceedings of the ACL*, pp.271-278, 2004.
- [13] Y. Mao and G. Lebanon, "Isotonic Conditional Random Fields and Local Sentiment

- Flow,” In *Proceedings of the NIPS*, 2007.
- [14] P. Turney, “Thumbs up or thumbs down? Sentiment orientation applied to unsupervised classification of reviews,” In *Proceedings of the ACL*, pp.417-424, 2002.
- [15] Y. Choi, C. Cardie, E. Riloff, and S. Patwardhan, “Identifying sources of opinions with conditional random fields and extraction patterns,” In *Proceedings of the HLT/EMNLP*, pp.355-362, 2005.
- [16] M. Thomas, B. Pang, and L. Lee, “Get out the vote: Determining support or opposition from congressional floor-debate transcripts,” In *Proceedings of the EMNLP*, pp.327-335, 2006.
- [17] E. Riloff and J. Wiebe, “Learning extraction patterns for subjective expressions,” In *Proceedings of the EMNLP*, pp.105-112, 2003.
- [18] 김묘실, 강승식, “SVM을 이용한 악성 댓글 판별 시스템의 설계 및 구현,” 한글 및 한국어 정보처리, pp.285-289, 2006.
- [19] G. A. Miller, “Nouns in WordNet: A Lexical Inheritance System,” *International Journal of Lexicography*, pp.245-264, 1990.
- [20] http://eedic.naver.com/list_thesaurus.naver 네이버 영어단어 유의어 시소러스
- [21] 강승식, *한국어 형태소 분석 및 정보 검색*, 홍릉과학출판사, 2002.
- [22] V. Vapnik, *The Nature of Statistical Learning Theory*, Springer, 1995.
- [23] E. Frank, M. Hall, and L. Trigg, *Weka 3: Data Mining Software in Java*, The University of Waikato, 2006.

1 차원고접수 : 2007. 9. 10

2 차원고접수 : 2008. 3. 10

최종게재승인 : 2008. 12. 1

(Abstract)

A Korean Emotion Features Extraction Method and Their Availability Evaluation for Sentiment Classification

Jaewon Hwang

Youngjoong Ko

Dept. of Computer Engineering, Dong-A University

In this paper, we propose an effective emotion feature extraction method for Korean and evaluate their availability in sentiment classification. Korean emotion features are expanded from several representative emotion words and they play an important role in building in an effective sentiment classification system. Firstly, synonym information of English word thesaurus is used to extract effective emotion features and then the extracted English emotion features are translated into Korean. To evaluate the extracted Korean emotion features, we represent each document using the extracted features and classify it using SVM (Support Vector Machine). In experimental results, the sentiment classification system using the extracted Korean emotion features obtained more improved performance (14.1%) than the system using content-words based features which have generally used in common text classification systems.

Keywords : Sentiment Classification, Korean Emotion Feature, Feature Expansion, SVM(Support Vector Machine)