

퍼지추론을 이용한 소수 문서의 대표 키워드 추출

Representative Keyword Extraction from Few Documents through Fuzzy Inference

노순억* · 김병만* · 허남철**

Sun-ok Rho, Byeong-man Kim, and Nam-chul Huh

*금오공과대학교 대학원 컴퓨터공학과

**대구미래대학 컴퓨터정보처리학과

요 약

본 논문은 사용자의 관심 내용을 포함하는 소수 문서들로부터 대표 용어들을 추출하고 가중치를 부여하는 새로운 방법을 제시한다. 대표 용어들의 추출 방법에서는 우선 예제 문서들로부터 후보 용어들을 추출하고 퍼지 추론을 적용하여 초기 대표 용어들을 선택한 후 예제 문서 내에서의 이들 용어들과 후보 용어들의 발생 빈도의 유사성을 이용하여 가중치를 재산정하고 대표 용어들을 자동 확장하였다. 제안 방법의 성능은 초기 대표 용어들을 선택하는 방법에 의해 영향을 크게 받는다. 따라서 문서집합에서 대표 용어를 추출하는 문제는 불확실성을 내포하고 있으므로 이러한 문제 해결에 효과적인 퍼지 추론을 초기 대표 용어의 선택 방법에 적용하였다. 본 논문에서 다루는 문제는 문서 집합의 중심 벡터를 계산하는 것으로 볼 수가 있다. 성능 평가를 위해 기존의 대표적인 Rocchio 알고리즘과 Widrow-Hoff 알고리즘과의 문서 분류 실험을 하였다. 실험 결과 우수한 성능을 보여줌으로서 제안 방법의 유용성을 확인 할 수 있었다.

Abstract

In this work, we propose a new method of extracting and weighting representative keywords(RKs) from a few documents that might interest a user. In order to extract RKs, we first extract candidate terms and then choose a number of terms called initial representative keywords (IRKs) from them through fuzzy inference. Then, by expanding and reweighting IRKs using term co-occurrence similarity, the final RKs are obtained. Performance of our approach is heavily influenced by effectiveness of selection method of IRKs so that we choose fuzzy inference because it is more effective in handling the uncertainty inherent in selecting representative keywords of documents. The problem addressed in this paper can be viewed as the one of calculating center of document vectors. So, to show the usefulness of our approach, we compare with two famous methods - Rocchio and Widrow-Hoff - on a number of documents collections. The results show that our approach outperforms the other approaches.

Key words : keyword extraction, fuzzy inference, user preference

1. 서 론

웹 검색 엔진 혹은 다양한 정보 검색 시스템을 이용하는 일반 사용자는 자신이 원하는 내용에 가장 적합한 정보를 찾고자 관심 대상 영역에 대한 제한된 어휘력과 전문성을 바탕으로 검색 질의어를 구성한다. 마찬가지로 정보 필터링 시스템 이용시 사용자는 적절한 정보를 추천 혹은 제공받고자 자신의 프로파일에 관심 사항을 기술한다.

검색 시스템의 경우 제공된 질의어로 검색 기능을 수행하고 그 결과에 대해서 사용자로부터 피드백을 받거나 검색된 결과를 이용해 자동으로 질의어를 수정하고 중요

도를 재산정하는 등의 부가 기능들을 수행함으로써 사용자에게 편의성을 제공하고 검색 효율을 높이고 있다[1, 5]. 정보 필터링 시스템 역시 위와 비슷한 성격의 프로파일 수정 과정들을 가진다[2, 3].

사용자에 의한 적절한 프로파일(질의어) 작성은 사용자에게 부담을 줄 수 있고 용어 불일치 문제로 인한 부적절한 필터링(검색) 결과를 가져 올 수 있다. 따라서 사용자의 프로파일 작성의 부담과 용어 불일치 문제의 수위를 낮추기 위해서 사용자로부터 관심 내용과 유사한 문서 집합을 제공 받아 이를 활용하는 것도 하나의 해결 방법이 될 수 있다. 이 경우에 발생하는 문제점은 제공된 문서 집합으로부터 사용자를 대신해서 대표 용어를 추출하고 이들에게 어느 정도의 중요도를 부여 할 것인가이다. 본 논문에서는 위의 문제 해결을 위한 새로운 방법을 제시한다.

접수일자 : 2001년 11월 1일

완료일자 : 2001년 12월 1일

본 연구는 한국과학재단 목적기초연구(2000-1-51200-008-2) 지원으로 수행되었음.

2. 관련연구

본 제안 방법은 예제 문서 집합이외의 시소러스 또는 주제 영역(domain)별 용어 사전 등과 같은 정보를 이용하여 대표 용어를 추출하거나 새로운 용어를 추가하지 않는다. 문서 집합에 대한 용어 정보(TF, DF, IDF)만을 이용하고 이들 용어 집합내에서 대표 용어를 추출한다.

문서 집합에서 대표 용어를 추출하고 이들의 가중치를 부여하는 문제는 기존의 대표적인 선형 분류기인 Rocchio와 Widrow-Hoff 알고리즘들[7]이 학습 문서 집합을 대표하는 중심 벡터를 구성하는 것과 성격이 같다. 식 (1)과 식 (2)는 배치 모드(batch mode)로 동작하는 Rocchio와 온라인(on-line mode)로 동작하는 Widrow-Hoff 알고리즘을 각각 보여주고 있다.

$$w_j = \alpha w_{1,j} + \beta \frac{\sum_{i \in C} x_{i,j}}{n_c} - \gamma \frac{\sum_{i \notin C} x_{i,j}}{n - n_c} \quad (1)$$

- n : 학습 문서(training documents)의 수.
- C : 긍정적 문서(positive documents)집합.
- n_c : 긍정적 문서의 수.

$$w_{i+1,j} = w_{i,j} - 2\eta (w_i \cdot x_i - y_i) x_{i,j} \quad (2)$$

- η : 학습율(learning rate).
- w_i : i 번째 문서 처리전의 중요도 벡터.
- y_i : 클래스 라벨.

이들 알고리즘들은 본 제안 방법과 마찬가지로 용어의 가중치 산정시 발생 빈도수(TF)와 역문헌 빈도수(IDF)를 결합하는 방법을 취하고 있지만 문서내 또는 문서 집합내 용어들간의 관련성을 용어의 가중치 계산에 반영하고 있지는 않다. 따라서 TF가 높은 용어는 높은 가중치를 가지게 되는데 대표 용어로서 실제 중요하지 않는 용어임에도 문서내에 자주 발생만 되면 높은 가중치 값을 부여받을 수 있다는 단점을 지니고 있다.

본 제안 방법은 크게 두 부분 즉 초기 대표 용어의 선택 부분과 용어 가중치 재산정과 대표 용어 자동 확장 부분으로 나눌 수 있다. 첫번째 부분인 초기 대표 용어의 선택 문제를 다루기 위해서 다수의 카타고리들을 포함하고 있는 다량의 문서집합에 대하여 문서 분류기의 성능 향상과 차원 축소(dimensionality reduction)에 중점을 두고 있는 특징 선택(feature selection) 방법들[4]을 고려해 볼 수 있다. 그러나 이들 방법들은 서로 다른 카타고리들의 정보 및 부정적 문서 집합의 정보를 이용하여 특징(용어)들을 선택하고 있기 때문에 동일한 주제에 속하는 단일한 소수의 긍정적 문서집합만을 대상으로 하고 있는 본 제안 방법에 이들 방법들을 사용하기가 곤란하다. 이와 달리 위의 Rocchio 와 Widrow-Hoff 는 긍정적 문서 집합만을 이용해서도 효과적으로 용어를 추출하고 가중치를 부여 할 수 있기 때문에 초기 대표 용어 추출의 비교 방법들로 사용하였다.

두번째 부분인 용어 가중치 재산정과 대표 용어 자동 확장 부분에서는 용어간의 관련성을 용어의 중요도 계산에 반영하기 위해서 벡터 모델에 기반한 질의 용어 가중치 재산정 및 질의 용어 확장 방법[1, 5, 10] 중 하나인 [5]를 참조하였다. 이 방법에서는 피드백 문서 집합 내에서 후보용어와 사용자가 제시한 질의 용어간의 발생 빈

도수의 유사도를 계산하여 확장된 질의 용어들의 가중치 산정에 반영하고 있다. 본 제안 방법에서 초기 대표 용어 추출 부분을 제외한다면 위의 적합 피드백(relevance feedback) 문제 해결 방법과 성격이 유사하다고 할 수 있다.

3. 대표 용어 추출 및 가중치 부여

예제(학습) 문서들로부터 대표 용어 확장 및 가중치 재산정에 필요한 사용자의 관심 내용을 가장 잘 대변하는 초기 대표 용어의 선택이 무엇보다 중요하다. 이들 초기 대표 용어들과 각각의 예제 문서 내에 존재하는 후보 용어들과의 발생 빈도 유사도 계산이 서로 이루어짐으로 선택 방법과 기준이 성능에 큰 영향을 미친다. 특정 용어의 중요도 계산에 사용되는 입력 정보(예: TF, DF, IDF)들은 정량적으로 정확히 해석될 수 없는 부정확하고 불확실한 특성을 내포하고 있다. 따라서 본 논문에서는 이러한 불확실성의 문제 해결에 효과적인 퍼지 추론을 적용하여 후보 용어들의 가중치를 계산하고 이 값들에 따라 선택 우선 순위를 부여하였다(3.1참조).

초기 대표 용어들은 선택 우선 순위에 따라서 각각의 예제 문서가 적어도 하나의 대표 용어를 포함할 때까지 확장되어 진다(3.2참조). 위의 결과로 초기 대표 용어들이 생성되면 다음 단계로 각각의 예제 문서 내의 후보 용어들과의 발생 빈도 유사도를 이용한 후보 용어들의 가중치 재산정이 수행된다(3.3참조). 기존의 학습 문서 집합을 대표하는 용어들의 가중치 부여 방법들(Rocchio, Widrow-Hoff)에서는 구성 용어들간의 어떠한 관련성을 계산에 반영하고 있지 않다. 이에 초기 대표 용어들이 사용자의 관심 내용을 가장 잘 나타내고 있다는 가정 아래 이들 용어들과 후보 용어들과의 발생 빈도의 유사성을 하나의 관련성으로 두고 중요도 계산에 이를 반영함으로써 분류 성능의 향상을 도모하였다[5].

대표 용어들의 자동 확장 단계에서는 기존 비교 방법들과 동일하게 모든 후보 용어들을 확장에 포함시켰으며 용어 가중치 재산정의 결과로 수정된 각각의 예제 문서들의 가중치 벡터들과 정보 검색 분야의 질의 용어 가중치 계산식[5, 8]을 사용한 초기 대표 용어들의 가중치 벡터를 전부 합산하여 최종적으로 예제 문서 집합의 대표 가중치 벡터를 구성하였다(3.4 참조). 그림 1은 위에서 설명한 과정들을 나타내고 있다

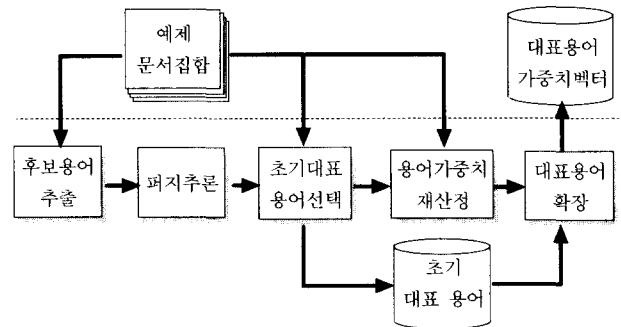


그림 1. 대표 용어 추출을 위한 시스템 구성도
Fig. 1 System architecture for Representative - Keyword Extraction

3.1 퍼지 추론을 이용한 대표 용어 중요도 계산

예제 문서들은 불용어 처리 그리고 스템밍 과정(Porter stemmer 사용)에 의해 후보 용어들의 집합으로 변형되며 이 집합으로부터 각각의 용어들의 TF(Term Frequency), DF(Document Frequency), IDF(Inverse Document Frequency) 정보가 구해진다. 이들 정보들이 퍼지 추론을 위한 퍼지 제어의 퍼지입력값으로 이용된다[6].

■ TF(Term Frequency)

각 용어의 발생 빈도수는 퍼지 계산에 사용되어지기 위해 정규화(NTF) 되어야 하며 아래의 식(3)을 사용하였다.

$$NTF_i = \frac{\frac{TF_i}{DF_i}}{\text{Max}_j \left[\frac{TF_j}{DF_j} \right]} \quad (3)$$

TF_i : 예제 문서 집합에서 i 번째 단어의 발생 빈도수
 DF_i : 예제 문서 집합에서 i 번째 단어를 포함하는 문서의 수

■ DF(Document Frequency)

각 용어의 예제 문서 집합 내에서의 문서 발생 빈도수를 나타내며 TF와 마찬가지로 아래의 식(4)을 사용하여 정규화(NDF) 하였다

$$NDF_i = \frac{\frac{DF_i}{TD}}{\text{Max}_j \left[\frac{DF_j}{TD} \right]} \quad (4)$$

TD : 예제 문서의 수
 DF_i : 예제 문서 집합에서 i 번째 단어를 포함하는 문서의 수

■ IDF(Inverse Document Frequency)

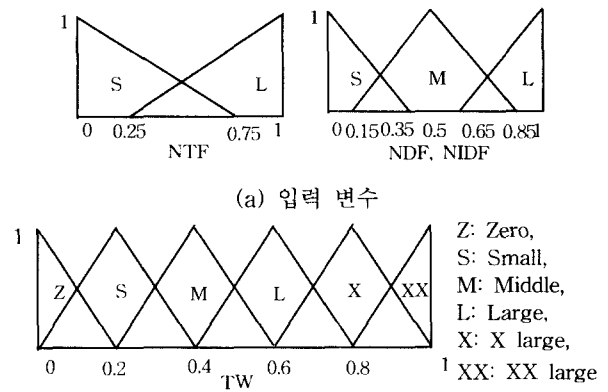
각 용어의 전체 예제 문서 집합 내에서의 역문헌 빈도수를 나타내며 아래의 식(5)을 사용하여 정규화(NIDF) 하였다.

$$NIDF_i = \frac{IDF_i}{\text{Max}_j [IDF_j]} \quad (5)$$

IDF_i : i 번째 단어의 역문헌 빈도수

그림 2는 퍼지 추론을 위하여 사용된 입출력 변수들을 나타내고 있다. 용어별로 구해진 NTF, NDF, NIDF 값들을 퍼지 추론에 적합한 형태로 퍼지화 시켜야 한다. 본 논문에서는 그림 2와 같은 삼각형 형태의 퍼지 수를 사용하였다. 그림 2(a)에서 NTF 입력 변수 값은 S(Small)과 L(Large)로 2개의 소속 함수 부분으로 나누었고 NDF와 NIDF 들은 S(Small), M(Middle), L(Large)로 하였다. 그림 2(b)에서 중요도를 나타내는 퍼지 출력 변수인 TW(Term Weight)는 6개의 소속 함수 부분으로 나누었다.

표 1은 NTF 퍼지 입력값의 소속 정도에 따라 두 부분으로 나누어 규칙들을 표현하고 있다. 작성 과정의 예를 살펴보면 NTF 입력값이 S(발생 빈도수가 낮음), NDF는 S(문서 빈도수가 낮음), 그리고 NIDF가 S(역문헌 빈도수가 낮음)일 경우 모든 특성치들이 낮은 값을



(a) 입력 변수
 (b) 출력 변수
 그림 2. 퍼지 입력 출력 변수
 Fig. 2 Fuzzy Input and Output Variables

가짐으로 중요 용어로서의 관련 정도를 Z(거의 관련 없음)로 두었다. NTF가 S, NDF가 L(문서 빈도수가 높다), 그리고 NIDF가 S 일 경우, 해당 용어가 대부분의 예제 문서들에 등장함으로 인해 관련성을 높게 평가 할 수 있지만 NTF와 NIDF 둘 모두가 낮은 값을 취함으로 관련 정도는 S(낮음)으로 설정하였다. 이와 같은 과정으로 다른 모든 규칙들의 후건부를 설정하였다.

NTF, NDF, NIDF 퍼지 입력값을 위의 결과로 생성된 18개의 추론 규칙별로 이들의 전건부의 소속 함수에 적용시킨다. 각각의 소속 정도가 구해지면 이들 중에서 최소(min)값을 취한다. 그 결과 규칙별로 하나씩의 퍼지 값이 생성되며 이 퍼지 값들을 퍼지 출력 변수 TW에 따라 6개의 그룹으로 분류하고 그룹별로 해당 그룹에 속한 퍼지 값들 중 최대(max)값을 취하여 총 6개의 퍼지 값들을 생성한다. 최종적으로 이들 6개의 퍼지 값들을 무게중심법(center of gravity)으로 비퍼지화(defuzzification)한 값이 해당 용어의 중요도 값으로 결정되어진다.

표 1. 퍼지 추론규칙
 Table 1. Fuzzy Inference Rules

NIDF \ NDF	S	M	L
S	Z	Z	S
M	Z	M	L
L	S	L	X

NTF = S

NIDF \ NDF	S	M	L
S	Z	S	M
M	S	L	X
L	S	X	XX

NTF = L

3.2 초기 대표 용어 선택

퍼지 추론을 통해 후보 용어들의 중요도 값들이 계산되어지면 이 값들에 따라 선택 우선 순위를 부여한다. 각각의 예제 문서가 적어도 1개의 초기 대표 용어를 포함해야 한다는 제약 사항을 두었다. 그림 3는 초기 대표 용어 선택 알고리즘을 보여주고 있다. 초기 대표 용어의 개수는 후보 용어와의 발생 빈도 유사도 계산의 정확성을 위해 최소로 유지하면서 위의 제약 사항에 따라 용어들을 선택하였다. 다음과 같은 예를 고려해보자.

- i) 예제 문서 집합 : $DS=\{d1, d2, d3, d4, d5, d6\}$
 각각의 문서들은 다음과 같은 용어들을 포함하고 있다.
 $d1=\{a, b, f, \dots\}$, $d2=\{a, c, d, \dots\}$, $d3=\{d, e, f, \dots\}$,
 $d4=\{d, f, \dots\}$, $d5=\{b, c, e, \dots\}$, $d6=\{e, f, \dots\}$.
- ii) 후보 용어 집합 :
 $TS=\{(a, 0.9), (b, 0.8), (c, 0.7), (d, 0.6), (e, 0.5), (f, 0.4)\}$
 편의상 퍼지 추론의 결과값인 중요도를 함께 기입하였다.
 위의 예에 그림 3의 알고리즘을 적용시키면 우선 그림 3의 2~4줄의 수행문들에서 임시변수의 적절한 초기화가 이루어진다. 다음으로 DS에서 1개의 임의의 문서요소 d를 선택하고(6줄) 후보 용어들의 중요도 값에 따라 내림차순으로 용어를 선택해 가면서 해당 용어가 문서 내에 존재하는지 확인한다. 만약 이들 중에서 11줄의 수행문을 만족하는 즉 문서(d)에 나타나는 용어(t)가 발견되면 이 용어는 해당 문서(d)를 대변하는 초기 대표 용어로 선택됨과 동시에 문서(d)에 대한 처리가 끝나게 된다. 이처럼 모든 문서로부터 초기 대표 용어들이 추출될 때까지 위의 과정을 반복해서 수행하게 된다(5~14 줄). 만약 위의 DS 집합에 나열된 문서 순으로 반복 수행된다고 가정하면 1번 수행후(d1 처리후) 출력값인 초기 대표 용어 집합(ITS)은 {a}가 되며 2번 수행후(d2 처리후)까지 변함이 없다가 3번 수행후(d3 처리후)에는 용어(d)가 추가되어 {a, d}가 된다. 4번 수행후 ITS는 변함이 없으며 5번 수행후 용어(b)가 추가되며 문서 d6를 처리한 뒤인 6번 수행후에는 마지막으로 용어(e)가 추가되어진다. 따라서, 문서의 개수만큼의 반복 수행이 완료(15줄)되면 ITS는 {a, b, d, e}로 확장되어진다. 이들 용어들과 중요도 값을 보아 알 수 있듯이 알고리즘(그림 3)이 앞에서 언급한 제약 사항을 잘 따르면서 용어들을 선택하였음을 확인 할 수 있다.

```

Input: DS ( Example Documents Set )
      TS ( Candidate Terms Set )

1] Procedure get_ITS(DS, TS)
2] ITS: Initial Representative Terms Set,
   initialized to empty.
3] TS': Temporary Terms Set, initialized to TS.
4] d, t: Document and Term element respectively.
5] Repeat
6]   Select a document element as d from DS.
7]   Repeat
8]     Select the highest element as t in TS'
       according to the weight.
9]     If t appears in d and not member in ITS
       Then Add t to ITS.
10]    Remove t from TS.
11]  Until t appears in d.
12]  Remove d from DS.
13]  Assign TS to TS'.
14] Until DS is empty.
15] Return ITS.
    
```

그림 3. 초기 대표 용어 선택 알고리즘

Fig. 3 The algorithm for selection of Initial-Representative Terms

3.3 용어 가중치 재산정과 대표용어 자동확정

초기 대표 용어들의 선택이 완료되면 이들 용어들과 각 문서내에서의 후보 용어와의 발생 빈도 유사도를 아

래의 식 (6)를 이용하여 계산하게 된다. 문서 집합내의 문서들은 사용자의 관심 내용 즉 특정 주제를 공통으로 표현하고 있다. 달리 말하자면 각각의 문서가 특정 주제에 대해 일반적(general)거나 세부적(specific)이거나 혹은 부분적인(partial) 내용을 표현한다고 볼 수 있다. 따라서 초기 대표 용어를 그 주제의 핵심 용어들로 간주할 경우 이들과 관련이 있는 용어들(전문용어, 유의어 등)은 당연히 함께 등장할 기회를 많이 가질 것이며 서로의 발생 빈도수(TF) 또한 비슷할 것으로 판단된다. 이러한 유사성을 후보 용어의 가중치 계산에 반영함으로써 중요도를 보다 정확하게 부여할 수 있다. 식 (6)에서는 초기 대표 용어(핵심 용어)들과 후보용어의 발생 빈도수의 차이들을 누적시킨 뒤 그 값이 작을 수록 핵심 용어와의 관련 정도를 높게 산정하고 있으며 반대의 경우 낮게 산정하고 있다.

$$Rd_{il}(K, t_i) = 1 - \log_p \left(\sqrt{\frac{\sum_{j=1}^n (kf_{jl} - tf_{il})^2}{n}} \right) \quad (6)$$

- $Rd_{il}(K, t_i)$: 문서 l에서 후보 용어 t_i 와 초기 대표 용어 간의 관련 정도
- kf_{jl} : 문서 l에서 초기 대표 용어 j의 발생 빈도수
- tf_{il} : 문서 l에서 후보 용어 i의 발생 빈도수
- K : 전체 초기 대표 용어들
- n : 초기 대표 용어들의 개수
- p : 조정 상수

후보 용어와 초기 대표 용어들간의 관련 정도가 산정되면 다음 단계로 이를 반영하면서 문서 집합에 대한 후보 용어의 가중치를 식(7)를 이용하여 계산하게 된다[5].

$$wt_i = \sum_{l=1}^m (wt_{il} \times Rd_{il}) \quad (7)$$

$$wt_{il} = TF_{il} \times IDF_i$$

- wt_i : 후보 용어의 문서 집합에서의 가중치
- wt_{il} : 문서 l에서 후보 용어 i의 가중치
- Rd_{il} : 문서 l에서 후보 용어 i와 초기 대표 용어들간의 관련 정도
- TF_{il} : 문서 l에서 후보 용어 i의 발생 빈도수
- IDF_i : 후보 용어 i의 역문헌 빈도수
- n : 문서 집합내의 문서 개수

식 (7)는 문서 별로 초기 대표 용어와 문서내에서의 후보 용어와의 관련 정도를 계산하고 그 정도를 일반적인 가중치 계산 방법인 TF와 IDF를 결합한 문서내의 용어 가중치 산정에 반영하면서 이 결과값들을 모두 합산하는 방식을 취하고 있다.

최종적으로 문서 집합의 대표 용어들과 이들의 가중치 벡터를 구하기 위해서 다음과 같은 단계를 거친다.

- i) 초기 대표 용어 가중치 벡터를 구성한다. 용어의 가중치 계산은 식 (8)을 사용한다.
- ii) 식 (7)를 이용하여 구성된 후보 용어 벡터들을 i)단계로 생성된 벡터와 합산하여 초기 용어로부터 확장된 전체 대표(중심) 벡터를 구한다.

$$w_i = \left(0.5 + \frac{0.5 \text{ freq}_i}{\text{Max}_j [\text{freq}_j]} \right) \times \log \left(\frac{N}{n_i} \right) \quad (8)$$

freq_i : 문서 집합에서의 초기 대표 용어 i의 발생 빈도수
 n_i : 초기 대표 용어 i가 나타나는 문서 개수

3.4 대표 벡터의 구성 예

i) 초기 대표 용어들과 후보 용어간의 관련 정도 산정
 사용자가 제시한 문서 d1에서 초기 대표 용어(K) k1, k2, k3 가 4, 3, 1 회 발생하고, 후보 용어 t1이 2 회 발생할 경우, 식 (6)을 사용한 초기 대표 용어들과 후보 용어의 관련정도는 아래와 같이 계산되어진다.

$$Rd_{11}(K, t_1) = 1 - \log_{10} \left(\sqrt{\frac{(2)^2 + (1)^2 + (-1)^2}{3}} \right) = 1 - 0.15 = 0.85$$

ii) 후보 용어의 가중치 산정

후보 용어 t1 이 1.0 의 IDF 값을 가지고 문서 d1, d2, d3에 3, 2, 1 회 발생하고, 식 (6)을 사용한 각 문서에서의 관련정도 값들이 0.4, 0.5, 0.7이라고 가정하면 식 (7)을 사용한 후보 용어 t1의 가중치는 아래와 같이 계산되어진다.

$$wt_1 = ((3.0 \times 0.4) + (2.0 \times 0.5) + (1.0 \times 0.7)) = 1.82$$

iii) 대표 용어 확장

초기 대표 용어 집합(K)와 식(8)를 사용하여 구성된 가중치 벡터(WK), 그리고 후보 용어 집합(T)와 식(7)를 사용하여 구성된 가중치 벡터(WT)가 각각 아래와 같을 경우,

$$K = \{ t_1, t_3, t_4 \}, WK = \{ 3.0, 2.0, 1.0 \}, \\ T = \{ t_1, t_2, t_3, t_4, t_5 \}, WT = \{ 5.0, 4.0, 3.0, 2.0, 1.0 \}$$

최종적으로 구성되는 대표 용어 집합(P)과 가중치 벡터(WP)는 다음과 같다. P={t1, t2, t3, t4, t5}, WP={8.0, 4.0, 5.0, 3.0, 1.0}

3.5 실험 및 결과

본 논문의 제안 방법의 유용성을 평가하기 위해서 기존의 대표적인 선형 분류기인 Rocchio, Widrow-Hoff 알고리즘들과의 문서 분류 성능을 비교해 보았다. 예제 문서 집합으로부터 대표 용어들을 추출하고 이들에게 가중치를 부여하는 문제는 위의 비교 대상 알고리즘들이 학습 문서 집합의 중심(center) 벡터를 구성하는 것과 성격이 같다.

실험 문서 집합으로는 Reuters-21578을 선택하였다. 본 논문에서는 Reuters-21578의 TOPICS 범주들을 선택하였으며 ApteMod 버전을 사용했고 라벨이 없는 문서들은 제외시켰다. 실험 대상으로 소수 예제 문서 집합들을 준비하고자 테스트 문서 집합과 학습 문서 집합에 적어도 하나의 문서를 각각 포함하고 있는 범주(category)들을 선택(총 90개)한 후 이 중에서 학습 문서 개수가 10개~30개인 범주 21개를 마지막으로 선별했다.

테스트 문서 집합은 3019개의 문서들을 포함하고 있다. 용어의 역문헌 빈도수(IDF)값을 구하기 위해 90개의 범주들에 속하는 7770개의 학습 문서 집합으로부터 문서 빈도수 정보를 이용하였다. 사용자는 자신의 관심 사항

에 부합하는 긍정적 문서 집합(positive documents)만을 제공한다는 가정하에 알고리즘 수행시 부정적 문서(negative documents)들의 정보 이용은 모두 제외시켰다.

비교 대상 알고리즘에 사용된 벡터들의 가중치는 용어의 TF × IDF로 계산하였다. 실험시 사용된 조정 상수(parameter)들의 설정값들은 Rocchio의 경우 α=0, β=1, γ=0로 두었고 Widrow-Hoff의 경우 η=0.25, y_i=1로 두었으며[9] 본 제안 방법에서 사용되는 식 (6)의 p는 10으로 두었다. 유사도 계산식은 cosine 식을 이용했으며 문서 분류 성능의 척도로서 standard recall, precision를 사용했다[8].

표 2는 11 standard recall levels 에 해당하는 보간(interpolation)절차를 통해 얻어진 정확율들의 평균값들과 비교 알고리즘들에 대한 제안 방법의 성능 향상율들을 각 범주별로 보여주고 있다. 본 제안방법(R.K.E.F : Representative Keyword Extraction through Fuzzy Inference)의 성능이 Rocchio 보다 평균 17%, Widrow-Hoff 보다는 평균 14% 향상되었음을 알 수 있다.

본 제안 방법에서는 초기 대표 용어들의 추출 방법으로 퍼지 추론을 이용하였다. 퍼지 추론의 유용성을 확인하고자 비교 알고리즘들을 초기 대표 용어 추출에 이용하여 각각의 성능을 비교해 보았다. 표 2의 R.K.E.R 와

표 2. 21개 범주들에 대한 성능 및 초기 대표 용어 - 선택 방법 별 성능 : R.K.E.F(Representative - Keyword Extraction through Fuzzy Inference), R.K.E.R(Rocchio), R.K.E.W(Widrow-Hoff).

Table 2. Performance on 21 categories and -Performance of each method for Initial Representative Keywords selection

	R.K.E.F				
	averaged precision	over Rocchio	over W.H	over R.K.E.R	over R.K.E.W
lumber	0.550	+51.9%	+25.8%	+34.4%	+49.8%
dmk	0.084	+52.7%	+61.5%	+ 0.0%	+ 0.0%
sunseed	0.451	+18.6%	+18.3%	+302.6%	+302.6%
lei	0.363	+ 0.0%	+ 0.0%	+ 0.0%	+ 0.0%
soy-meal	0.772	+40.6%	+43.7%	+101.0%	+101.0%
fuel	0.518	+40.3%	+58.8%	+36.3%	+51.0%
heat	0.626	+ 9.6%	+10.0%	+ 2.6%	+ 2.6%
soy-oil	0.323	+20.5%	+ 0.6%	+62.3%	+ 76.5%
lead	0.614	+13.9%	+10.6%	+76.4%	+16.2%
strategic-metal	0.120	+15.3%	+11.1%	+21.2%	+20.2%
hog	0.485	- 5.6%	- 9.1%	- 1.2%	- 11.8%
orange	0.975	+ 5.1%	+ 4.3%	+ 0.0%	+ 0.0%
housing	0.352	- 9.0%	- 5.6%	- 10.2%	- 1.6%
tin	0.986	+ 2.0%	+ 1.7%	+ 0.0%	+ 0.0%
rapeseed	0.575	+24.1%	+28.3%	- 0.8%	- 0.8%
wpi	0.728	- 8.1%	- 9.4%	- 4.9%	- 4.9%
pet-chem	0.308	+ 5.4%	-16.9%	+ 4.0%	-11.7%
silver	0.770	+61.4%	+50.0%	+ 2.2%	+ 2.2%
zinc	0.921	+16.2%	+ 8.2%	+ 0.0%	+ 0.0%
retail	0.194	+ 2.6%	+ 1.5%	+ 3.7%	+102.0%
sorghum	0.591	+16.3%	+18.2%	- 4.0%	+94.4%
Average	0.530	+17.8%	+14.8%	+29.8%	+37.5%

표 3. lumber 범주에서 추출된 적용방법 별 초기 대표 용어들
Table 3. Initial Representative Keywords of the "lumber"
category extracted by each method

Methods	Initial Representative Keywords
Fuzzy	lumber, softwood, wood, agre
Rocchio	lumber, softwood, timber, guarante, forest
W.H	lumber, timber, guarante, champion, wood

R.K.E.W는 초기 대표 용어의 추출 방법으로 각각 Rocchio와 Widrow-Hoff 를 사용한 대표 용어 추출 방법들을 나타내고 있다. 표 3은 각 방법 별로 lumber 문서 집합에서 추출되는 초기 대표 용어들의 예를 보여주고 있다. 실험결과, 퍼지 추론을 이용한 대표 용어 추출 방법(R.K.E.F)에 비해 두 비교 알고리즘을 각각 응용한 대표 용어 추출 방법들(R.K.E.R, R.K.E.W)은 뚜렷한 성능향상을 보여주지 않았다. 즉 적절한 핵심용어를 선별하지 못했음을 알 수 있었다. 표 2의 오른쪽 결과를 통해 비교 알고리즘들을 사용한 대표 용어 추출 방법들 보다 퍼지 추론을 사용한 대표 용어 추출 방법이 보다 나은 성능을 보여주고 있음을 알 수 있다.

4. 결론

본 논문에서는 사용자가 제시한 소수의 문서 집합으로부터 관심 내용을 대표하는 중요 용어들을 추출하고 그들의 가중치를 부여하는 문제에 관하여 퍼지 추론과 용어 발생 빈도수의 유사성을 이용한 가중치 재산정 방법을 적용했다. 방법의 유용성을 보이고자 학습 문서 집합의 중심 벡터를 구하는 대표적인 선형 분류기의 알고리즘들과 성능을 비교했다. 소수의 긍정적 학습 문서 집합들에 대해서 실험한 결과 비교적 우수한 성능향상을 보여줌으로써 본 제안 방법의 유용성을 확인할 수 있었다.

사용자의 관심 내용을 가장 잘 대변하는 초기 대표 용어(핵심 용어)들의 추출 방법은 용어 발생 빈도수의 유사성에 따른 가중치의 재산정에 직접적인 영향을 주고 있다. 사용자의 관심 내용과 거리가 먼 핵심 용어를 추출했을 경우 그 용어를 중심으로 빈도수의 유사성을 보이면서 관련 전문 용어 혹은 유의어 등의 성격을 함께 가지는 후보용어의 가중치가 높게 산정됨으로 그 결과로 생성되는 중심 벡터는 문서 집합의 대표성을 상실하게 된다. 따라서 향상된 실험 결과를 통해 핵심 용어들이 올바르게 추출 되었고 사용된 퍼지 추론 방법이 효과적이었음을 확인할 수 있었다. 본 제안 방법에 이용된 용어들간의 발생 빈도수의 유사성 이외에 또다른 여러 가지 관련 정보를 이용할 수 가 있고 이와 관련된 문제 해결을 위해 퍼지 추론 방식이 중요한 기준 혹은 기반 구축 방법으로써 효과적으로 사용 될 수 있을 것이다.

참 고 문 헌

- [1] Mitra, M., Singhal, A., and Buckley, C., "Improving Automatic Query Expansion", *In Proceedings of the 21st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 206-214, 1998.
- [2] Pazzani, M., Billsus, D. "Learning and revising user profiles: the identification of interesting Web sites" *Machine Learning*, 1997.
- [3] Seo, Y., Zhang, B., "Personalized Web Document Filtering Using Reinforcement Learning", *Applied Artificial Intelligence*, 2001.
- [4] Yang, Y., Pedersen, J. "A comparative study on feature selection in text categorization", *Proceedings of the 14th International Conference on Machine Learning*, pp. 412~420, 1997.
- [5] Byeong Man Kim, Ju Youn Kim, JongWan Kim, "Query Term Expansion and Reweighting using Term Co-Occurrence Similarity and Fuzzy Inference" *IFSA/NAFIPS*, 2001.
- [6] O, Cordó, F. Herrera, and A. Peregrin, "A Practical Study on the Implementation of Fuzzy Logic Controllers", *Intelligent Control*, 1998.
- [7] David D. Lewis, Robert E. Schapire and James P. Callan and Ron Papka, "Training algorithms for linear text classifier", *Proceedings of SIGIR-96, 19th ACM International Conference on Research and Development in Information Retrieval*, 1996.
- [8] R. Baeza-Yates, B. Ribeiro-Neto, *Modern Information Retrieval*, ACM Press, NY, USA, 1999.
- [9] K. Lam and C. Ho, "Using a generalized instance set for automatic text categorization", *In 21th Ann. Int. ACM SIGIR Conference on Research and Development in Information Retrieval*, pp81~89, 1998.
- [10] Jinix Xu and W. Bruce Croft, "Query Expansion Using Local and Global Document Analysis", *Proceeding of ACM SIGIR International Conference on Research and Development in Information Retrieval*, pp4-11, 1996.

저 자 소 개



노순억 (Sunok Rho)

1999년 : 금오공과대학교 컴퓨터공학과 학사

2000년~현재 : 금오공과대학교 컴퓨터공학과 석사과정

관심분야 : 정보검색, 인공지능

E-mail : sorho@se.kumoh.ac.kr



김병만 (Byeong Man Kim)

1987년 : 서울대학교 컴퓨터공학과 학사
1989년 : 한국과학기술원 전산학과 공학석사
1992년 : 한국과학기술원 전산학과 컴퓨터
공학박사
1992년~현재 : 금오공과대학교 부교수
1998년~1999년 : 미국 Univ. of California,
Irvine Post Doc.

관심분야 : 인공지능, 정보검색, 소프트웨어 검증
E-mail : bmkim@se.kumoh.ac.kr



허남철 (Nam Chul Huh)

1986년 : 계명대학교 전자계산학과
1988년 : 한국과학기술원 전산학과 공학석사
1988년~1991년 : 산업과학기술연구소 연구원
1991년~현재 : 대구미래대학 컴퓨터정보처
리학과 부교수

관심분야 : 인공지능, 정보검색 및 필터링, 병렬처리
E-mail : nchuh@mail.tfc.ac.kr