

# 온톨로지 기반의 웹 페이지 분류 시스템

최 재 혁<sup>\*</sup> · 서 혜 성<sup>\*\*</sup> · 노 상 옥<sup>\*\*\*</sup> · 최 경 희<sup>\*\*\*\*</sup> · 정 기 현<sup>\*\*\*\*\*</sup>

## 요 약

본 논문은 온톨로지(ontology)에 기반한 자동화된 웹 페이지 분류 시스템을 제안한다. 웹 페이지의 분류를 위하여 첫 번째 단계에서는 각 웹 페이지가 속한 범주(category)를 대표할 수 있는 단어를 선정하며, 이를 위하여 단어빈도와 문서빈도를 곱한 값을 계산한다. 두 번째 단계에서는 첫 번째 단계에 의해 선택된 단어의 정보이득(information gain)을 계산해 분류 확률이 높은 단어를 우선적으로 선정한다. 두 단계를 통하여 선정된 단어들과 웹 페이지의 분류 정보를 가지고, 기계학습에 의하여 컴파일된 규칙(compiled rules)을 생성한다. 생성된 규칙은 임의의 웹 페이지들을 도메인 온톨로지에 의해 정의된 범주 별로 분류할 수 있도록 한다. 본 논문의 실험에서는 주어진 웹 페이지 집합에서 각 범주 별로 평균 240개의 단어로부터 78개의 단어를 결과적으로 선정하였으며, 이를 바탕으로 웹 페이지 분류 규칙을 생성하였다. 실험 결과에서 제안한 시스템의 평균 분류 정확도는 약 83.52%로 측정되었다.

## Web Page Classification System based upon Ontology

Jaehyuk Choi<sup>\*</sup> · Haesung Seo<sup>\*\*</sup> · Sanguk Noh<sup>\*\*\*</sup>  
Kyunghee Choi<sup>\*\*\*\*</sup> · Gihyun Jung<sup>\*\*\*\*\*</sup>

## ABSTRACT

In this paper, we present an automated Web page classification system based upon ontology. As a first step, to identify the representative terms given a set of classes, we compute the product of term frequency and document frequency. Secondly, the information gain of each term prioritizes it based on the possibility of classification. We compile a pair of the terms selected and a web page classification into rules using machine learning algorithms. The compiled rules classify any Web page into categories defined on a domain ontology. In the experiments, 78 terms out of 240 terms were identified as representative features given a set of Web pages. The resulting accuracy of the classification was, on the average, 83.52%.

**키워드 :** 웹 페이지 분류(Web Page Classification), 온톨로지(Ontology), 정보이득(Information Gain), 기계학습(Machine Learning)

### 1. 서 론

현대 사회는 인터넷의 발달로 새로운 방식의 사무형태를 자리잡아 가고 있다. 점차 다양하고 신속한 정보를 필요로 하는 시대적 요구에 부응하여 나타난 인터넷이라는 매체는 저렴한 비용으로 다양한 정보를 제공하며, 여타의 매체보다 신속하게 정보를 전달할 수 있다. 따라서 기업들은 인터넷의 사용이 필수 불가결하게 되었으며 이에 따른 투자를 점차 늘리고 있다. 그러나 그 정보의 다양함 가운데는 업무를 방해하는 유해한 요인들도 있어 이 때문에 오히려 업무의 효율성이 저하되는 사례들이 발견되고 있다[3]. 이런 요인들은 자원 운용 측면에

서 볼 때, 불필요한 컴퓨팅 자원의 낭비와 통신 대역의 손실을 유발한다. 때문에 많은 기업들은 업무에 적합하지 않은 웹 페이지들을 탐지하고 차단하기 위한 방법들에 대한 관심을 가지게 되었다. 본 논문에서는 일반적인 네트워크 환경을 통해 사내망과 같은 지역 네트워크 안으로 전달되는 스트림들로부터 부적합한 웹 페이지들을 판별하여 온톨로지(ontology)[10, 11]에 기반한 클래스로 분류하는 방법을 소개한다.

본 논문이 제안하는 웹 페이지 분류 시스템은 크게 세 개의 모듈로 구성되어 있다. 그 세 개의 모듈은 패킷 수집 모듈, 프로토콜 분석 모듈, 콘텐츠 분류 모듈이다. 패킷 수집 모듈은 TCP 기반의 네트워크에서 서버와 클라이언트 간의 세션으로부터 패킷을 수집한다. 프로토콜 분석 모듈은 수집된 패킷의 헤더 정보를 분석하고 패킷이 사용되는 응용 프로그램에 적합한 형태로 재구성한다. 예를 들어, 웹 서비스의 경우 웹 서버가 지원하는 압축 방식으로 웹 페이지를 제공하기도 하는데, 이와 같은 웹 페이지를 분석하기 위해서 프로토콜 분석 모듈이 웹 서버가 지원하는 압축 방식을 식별하여 해제

\* 이 논문은 2003년도 한국학술진흥재단의 지원과 국가중요연구실 사업의 지원에 의해 연구되었음(KRF-2003-041-D20465, KISTEP National Research Laboratory program).

<sup>\*</sup> 정 회 원 : 삼성전자 네트워크 사업부

<sup>\*\*</sup> 정 회 원 : 아주대학교 정보통신전문대학원 정보통신공학과

<sup>\*\*\*</sup> 정 회 원 : 가톨릭대학교 컴퓨터정보공학부 교수

<sup>\*\*\*\*</sup> 정 회 원 : 아주대학교 정보통신전문대학원 교수

<sup>\*\*\*\*\*</sup> 정 회 원 : 아주대학교 전자공학부 교수

논문접수 : 2004년 4월 29일, 심사완료 : 2004년 8월 12일

하는 역할을 한다. 콘텐츠 분류 모듈은 프로토폴 분석 모듈이 재구성한 웹 페이지의 내용을 분석하여 적합한 범주로 분류한다.

웹 페이지의 내용을 식별하기 위해서 본 논문에서는 단어빈도(Term Frequency : TF)와 문서빈도(Document Frequency : DF)[1, 18]를 곱한 값과 정보이득(information gain) 계산 방법[20]을 결합한 방법을 제안한다. 그리고 제안된 방법론에 의해 선택된 단어를 기반으로 하여 임의의 웹 페이지를 도메인 온톨로지 내에 있는 범주로 분류한다. 이 과정에서 실질적인 분류는 기계학습(machine learning)에 의해 생성된 분류 규칙을 사용하여 웹 페이지와 미리 정의된 범주 간의 관련성을 식별함으로써 가능하게 되는 것이다.

웹 페이지의 분류를 위하여는 웹 페이지가 속하는 범주를 우선적으로 정의해야 하며, 이를 위하여 본 논문에서는 온톨로지(ontology)를 구성한다. 온톨로지는 객체와 특성들을 계층적으로 구조화시켜 모델링하는 특징이 있으며, 개념의 확장이나 진행을 지속적으로 유지시켜 줄 수 있다. 그러나, 변화되는 온톨로지 구성 요소를 자주 갱신하지 않으면 웹 페이지 분류에 대한 정확도가 낮아지는 문제가 발생한다. 따라서 제안하는 시스템은 보다 정확한 온톨로지를 유지하기 위하여, 웹 페이지들의 범주를 계층적으로 구성하고 각 범주를 대표할 수 있는 단어들을 계속적으로 갱신한다.

제안하는 웹 페이지 분류 시스템은 실시간으로 웹 페이지를 분류하여 검사 중인 웹 페이지가 지역 네트워크에 부적합한 범주에 속할 경우 별도의 보고서를 생성할 수 있을 것으로 기대한다. 또한 웹 페이지 분류 시스템은 향후에 네트워크의 사용이 안전하게 되고 있음을 보증하는데 필수 불가결한 요소가 될 것임을 확신한다.

본 논문의 다음 장에서는 온톨로지 등을 이용한 웹 페이지 분류 방법과 관련된 연구들에 대해 논하고, 3장에서는 온톨로지를 구성하는 단어 선별 과정에 대한 자세한 설명을 한다. 4장에서는 웹 페이지 분류 시스템에 대한 시스템 구조를 설명하고, 5장에서는 본 논문에서 수행한 시뮬레이션 결과에 대해 논한다. 끝으로 결론에서 실험 결과를 정리하고 향후 연구 과제에 대해 언급한다.

## 2. 관련 연구

지금까지 온톨로지를 기반으로 하여 웹 페이지를 분류하려는 많은 연구가 있어왔다[6, 8, 12, 13]. Prabowo[13]와 그의 연구에서 온톨로지는 “특정 도메인의 개념적 인스턴스이며, 다른 인스턴스와 구별될 수 있는 개체(entity)”라고 정의했다. Prabowo는 듀이 십진 분류법(DDC)과 미 의회 도서관 분류법(LCC)을 고려하여 웹 페이지를 분류했는데 이 과정에

서 온톨로지를 구축하였다. 이 방법은 기존의 접근법에 비해 몇 가지 특징이 있다. 그 특징은 분류를 위해 백과사전(thesaurus)이나 사전을 이용하지 않고, 온톨로지를 사용했다는 점, 그리고 웹 페이지의 텍스트로부터 온톨로지를 구축했다는 점, 그리고 온톨로지와 분류체계의 연관 관계를 구성한 점 등을 들 수 있다. 이 접근법의 단점은 비록 표준 분류 방법을 따르고 있지만, 분류에 대한 사용자의 복잡한 요구가 적절히 반영되기 어렵다는 것이다. 반면에, 본 논문에서 제안하는 방법은 온톨로지의 구성요소를 지속적으로 갱신하기 때문에 사용자의 요구에 보다 유연하게 대처하는 장점을 갖는다.

분류 규칙을 만들기 위해서는 기계학습에 적당한 크기의 특성들을 입력값으로 제공해야 한다. 제공할 특성을 선별하기 위한 방법으로 Ruger와 Gauch[17]의 연구에서는 단어빈도의 최소, 최대값 만의 임계값을 설정해 놓고, 그 사이의 단어들을 선택했다. 그들의 연구는 문서빈도에 의존적이어서 특성으로 선택된 단어들이 웹 페이지들의 특정 범주에서 자주 발견되기는 하지만, 그것을 다른 범주와 그것이 속한 문서의 범주를 구분하는 요소로 사용할 수는 없다. 그러나 본 논문은 문서빈도뿐만 아니라 단어빈도를 사용하여 각 범주를 대표하면서도 중요한 의미를 갖는 단어들을 선별할 수 있었다. 주목할 사항은 기존의 정보 검색에서 단어의 가중치를 계산하기 위해서 단어빈도(TF)와 역문서 빈도(IDF)를 이용하여 계산하나 본 시스템은 이를 변형한 TF/DF로써 각 단어의 가중치를 계산한다. TF/IDF가 각 문서마다 고유한 단어들에 가중치를 부여하는 반면 TF/DF는 각 문서 집합에 고유한 단어들에 가중치를 부여한다.

웹 페이지를 분류하기 위해서 대부분의 접근 방식들은 기계학습에 의한 접근 방식을 취하고 있다[16, 19, 22, 23]. Sebastiani[19]는 텍스트 기반의 분류법에서의 기계학습의 역할을 정리했다. Sebastiani에 의하면 기계학습 알고리즘을 통해 높은 성능을 얻기 위해서는 입력값으로서 선택되는 특성들의 크기를 줄이는 것이 관건이다. 본 논문의 관점에서는 Sebastiani의 연구가 도메인 온톨로지를 구축하지 않기 때문에 각 범주를 대표하기 위한 특성들이 중요하면서 가치 있는 것으로 선정되었는지 알 수 없다. 한편으로 Riboni[16]의 연구에서는 특성추출을 위해 정보이득, 단어 수(word frequency) [16] 그리고 문서빈도를 이용하였다. Riboni는 정보이득, 단어 수, 문서빈도 각각에 의해 추출된 특성의 수를 변경하면서 기계학습 알고리즘을 적용하였다. 그 결과 정보이득에 의한 특성추출 방법이 분류 성능에 가장 좋은 영향을 미쳤다. 그러나 Riboni의 방법은 특성추출을 위해 사용한 각각의 방법을 개별적으로 사용하였기 때문에 결합된 방법이 분류 성능에 미칠 영향에 대해서는 연구가 미흡했다. 따라서 본 논문

에서는 특성들의 크기를 조절하여 사용할 수 있게 한 점에 대해서는 방법을 같이 하되, 분류의 정확도를 높이기 위하여 문서빈도, 단어빈도, 그리고 엔트로피 계산을 결합한 메커니즘을 제안한다.

### 3. 웹 페이지 분류를 위한 온톨로지

본 논문에서는 개념에 대한 계층적 구조를 갖는 온톨로지를 이용하여 웹 페이지를 모델링하고 분류하였다. 제한한 웹 페이지 분류 시스템의 온톨로지는 전체 개념의 계층구조에서 각각의 분류 범주를 정의하고 각 분류 범주를 대표하는 단어들을 포함한다. 따라서, 시스템의 분류 범주 각각이 주어진 웹 페이지를 어떤 기준(또는 단어)으로 분류하는지를 제시하며, 온톨로지의 분류 범주에 포함된 단어를 이용하여 기계 학습 규칙을 생성한다. 그러므로 주어진 도메인의 온톨로지를 정의하기 위하여 특정 범주(category)에 연관된 단어의 집합을 어떻게 선택하는가가 중요한 요인이 된다. 3장에서는 이러한 단어 집합의 선택 방법에 대하여 구체적으로 설명한다.

#### 3.1 특성 추출

도메인 온톨로지를 구축하는 첫 번째 단계는 특정 범주를 대표할 수 있는 적절한 단어의 집합을 식별하는 일이다. 직관적으로 볼 때, 어떤 단어가 특정 범주에 속하는 한 웹 페이지 내에서 자주 발견되고 또한 그 특정 범주에 속하는 대부분의 웹 페이지들에서도 발견이 된다면 그 단어는 특정 범주를 대표한다고 할 수 있을 것이다. 예를 들어, “은행과 금융”이라는 범주를 포함하는 도메인 온톨로지가 있고, ‘account’라는 단어가 앞서 설명한 대로 “은행과 금융” 범주 내의 특정 웹 페이지에 집중적으로 나타나고 또한 “은행과 금융” 범주 전체에서 골고루 발견된다면 ‘account’라는 단어는 “은행과 금융” 범주에 속하는 웹 페이지들을 대표하는 단어가 될 수 있다.

한 페이지에서 발견되는 단어의 빈도수와 특정 단어가 전체 문서에서 발견되는 빈도수는 정보검색 분야에서 사용되는 단어빈도(TF)와 문서빈도(DF)[1, 18]라는 용어로 대체하여 사용할 수 있다. 앞서 직관적인 개념에 의존했던 설명을 단어빈도와 문서빈도를 이용한 수식으로 설명하고자 한다면 다음과 같이 할 수 있을 것이다. 임의의 웹 페이지에 속하는 단어의 집합을  $I = \{1, 2, \dots, m\}$ 라 하고, 특정 범주에 속하는 웹 페이지들의 집합을  $J = \{1, 2, \dots, n\}$ 라 하자. 단어빈도와 문서빈도는 각 범주를 대표하는 단어에 대해 비례관계에 있으므로 이를 일반화하여 표현하면 아래 식과 같다.

$$W_{i,j} = \frac{TF_{i,j}}{\max_{k \in I} TF_{k,j}} \times \frac{DF_i}{n} \quad (1)$$

- $W_{i,j}$ 는 웹 페이지  $j \in J$ 에 속하는 단어  $i \in I$ 의 가중치이다.
- $TF_{i,j}$ 는 웹 페이지  $j \in J$ 에 속하는 단어  $i \in I$ 의 빈도이다.
- $DF_i$ 는 단어  $i$ 가 존재하는 웹 페이지의 개수이다.

위의 식 (1)은 특정 단어가 하나의 웹 페이지에서 다른 단어들 보다 많이 나타나고 또한 그 단어가 범주를 구성하는 대부분의 웹 페이지에서 나타난다면 그 가중치는 다른 단어들에 비해 크다는 것을 의미한다. 다른 단어들에 비해 더 큰 가중치를 가지는 단어는 그 단어가 속하는 범주에 대한 대표성이 더 크다는 것을 의미한다. 단어의 가중치를 계산함으로써 단어와 범주 간의 관련성을 측정해 볼 수 있다. 따라서 높은 가중치의 단어일수록 계층적 온톨로지 상에 특정 범주를 표현하는데 사용될 가능성이 높게 된다.

#### 3.2 엔트로피의 계산

우리는 특정 범주에 연관된 단어들을 가지고 임의의 웹 페이지를 온톨로지 기반으로 나뉜 범주들 중 하나의 범주로 분류해야 한다. 그러나 3.1절에서 계산된 단어들의 가중치는 그 단어가 해당 범주를 대표한다는 것을 나타내는 데는 유용하게 사용될 수 있지만, 범주를 분류하는데 있어 어느 정도 유용한지는 알 수 없다. 또한 유사 범주들간에는 같은 단어가 많이 나타나기 때문에 각 범주를 잘 설명하는 단어가 곧 해당 범주를 다른 범주와 잘 구분하는 단어가 될 가능성 또한 작아진다. 따라서 본 논문은 웹 페이지의 분류를 위해서 [20]에서 제안된 방법을 통해 단어의 엔트로피를 계산하는 방식을 도입한다. 단어의 엔트로피는 정확한 분류를 하는 데 필요한 정보의 기대 값을 제공한다. 낮은 엔트로피 값은 웹 페이지를 분류하는데 필요한 정보의 양이 적다는 것을 의미한다. 그러므로 보다 효율적인 분류를 위해서는 낮은 값의 엔트로피를 갖는 속성들을 사용해야 한다. 엔트로피의 계산은 다음과 같은 식으로 표현할 수 있다.

$$E^S = - \sum_{i=1}^n p_i \log_2(p_i) \quad (2)$$

- $E^S$ 는 확률  $p_1, \dots, p_n$ 의 클래스들의 집합  $S$ 의 엔트로피이다.

예를 들어 A와 B라는 두 개의 클래스로 구성되는 집합이 있을 때, 집합 내의 모든 개체 수에 대한 A 클래스에 속하는 개체 수의 비율을  $p_1$ 이라 하자. 마찬가지로 B 클래스에 대한  $p_2$ 가 있다. 그러면 엔트로피는  $-(p_1 \log_2 p_1 + p_2 \log_2 p_2)$ 에 의해 계산된다.

위에서 계산된 엔트로피를 가지고 속성의 분류 확률을 나타내는 정보이득 값을 계산할 수 있다. 정보이득에 대한 계산식은 아래와 같다.

$$Gain(a) = E^S - \sum_{j=1}^m \left( \frac{|S_{a_j}|}{|S|} \times E^{S_{a_j}} \right) \quad (3)$$

- $Gain(a)$ 는 속성  $a$ 에 대한 정보이득 값이다.
- $S_{a_j}$ 는  $j \in \{1, \dots, m\}$ 의 값을 갖는 속성  $a$ 의 집합이다.

위 식에 사용된 속성이라는 것은 특정한 범주에 속하는 단어를 의미한다. 계산된 속성들의 정보이득 값을 통해 단어가 갖는 분류 확률의 순위를 매길 수 있으며, 정보이득 값에 근거해서 단어들을 선택할 수 있게 된다.

### 3.3. 단어 구성

웹 페이지를 분류하기 위해 속성들을 하나의 튜플로 조합시킬 필요가 있다.  $A$ 를 속성들의 집합이라고 한다면, 속성의 선택은 아래와 같은 식을 적용하여 수행한다.

$$H_\mu = \{a \mid a \in A, Gain(a) \geq \mu\} \quad (4)$$

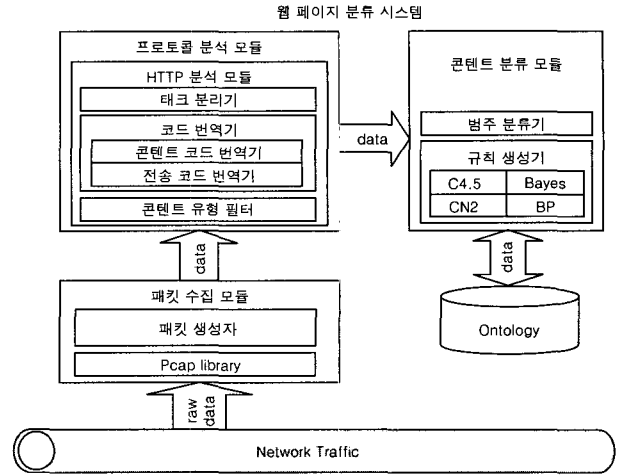
임계값  $\mu$ 는 각각의 범주에 있는 단어들 중 상대적으로 중요하지 않은 단어들을 걸러주는 역할을 한다. 위의 식을 수행한 결과로 만들어진 집합  $H_\mu$ 는  $\mu$ 가 높을수록 해당 범주와 관련이 깊은 속성들의 집합이 된다.

요컨대 자동화된 웹 페이지 분류를 위해서는 첫 번째로 TF/DF를 계산하여 해당 범주를 대표할 수 있는 단어들을 선정한다. 두 번째로 첫 번째 과정을 통해 선정된 단어들의 정보이득을 계산한 값을 이용하여 속성들의 우선순위를 정해 튜플을 구성한다. 세 번째로 구성된 튜플을 가지고 여러 기계학습 알고리즘[5,9]에 대입해 분류 규칙의 집합을 만든다. 결과적으로 생성된 규칙을 통해 임의의 웹 페이지를 도메인 온톨로지가 정의하는 범주로 분류할 수 있도록 한다.

## 4. 분류 시스템의 구조

이번 장에서는 웹 페이지 분류 시스템에 대한 시스템 구조를 설명한다. (그림 1)은 웹 페이지 분류 시스템의 구조를 나타낸다. 웹 페이지 분류 시스템은 크게 세 개의 모듈 즉, 패킷 수집 모듈, 프로토콜 분석 모듈, 콘텐츠 분류 모듈로 구성되어 있다. 시스템의 전반적인 흐름은 네트워크 트래픽 중 일부 트래픽을 패킷 수집 모듈로 수집한 후에 프로토콜 분석 모듈이 분석 가능한 형태로 바꾸고 그것을 콘텐츠 분류 모듈이 온톨로

지 데이터베이스를 기반으로 미리 정의된 범주 별로 분리하는 것이다.



(그림 1) 웹 페이지 분류 시스템의 구조

패킷 수집 모듈은 libpcap[7]을 이용하여 네트워크 트래픽 중에 TCP 80 포트를 통해 흐르는 스트림을 수집하여 패킷 생성자를 통해 구조체로 표현한 후 프로토콜 분석 모듈로 넘긴다.

<표 1> HTTP 프로토콜의 Content-Encoding 유형

인코딩 방식	설 명
gzip	GNU zip 형식의 압축 포맷
compress	UNIX 파일 압축 프로그램의 압축포맷. LZW 라고도 함.
deflate	zlib과 deflate라는 압축 메커니즘이 결합된 형태
identity	기본 인코딩. 변환이 필요 없음.

프로토콜 분석 모듈은 현재 HTTP 프로토콜에 대한 분석만이 가능하므로 패킷 수집 모듈에서 수집된 데이터는 HTTP 분석 모듈로 전달된다. HTTP 분석 모듈에서는 HTTP 프로토콜을 분석하여 최종적으로 분석 가능한 형태를 만든다. 최종적으로 분석 가능한 형태란 텍스트로만 이뤄진 데이터로 이를 얻기 위해서는 다음과 같은 과정을 거친다. 먼저 콘텐츠 유형 필터에 의해 자료의 유형을 파악해야 한다. HTTP 프로토콜을 통해 전달되는 자료의 유형은 IANA[15]에 의해 등록되어 있다. 이 유형 중 우리가 취해야 할 것은 웹의 콘텐츠의 대부분을 차지하는 text/plain 유형과 text/html 유형이다. 이렇게 선별된 유형의 자료를 판독 가능한 형태로 바꾸는 작업을 해야 한다. 일부 HTTP 서버들은 <표 1>과 같은 코딩 방식 중 하나를 선택하여 인코딩하여 보내기도 한다.

그리고 <표 1>과는 성격이 다른 인코딩 방식이 사용되기

도 하는데 이는 문서의 안전한 전송을 위해 사용하는 인코딩 방식으로 일반적으로 “chunked”라고 불린다. 모든 인코딩이 해제되었다면 본 시스템에서는 <body> 태그 내의 내용만을 취급하므로 태그 분리에 의해 태그로 사용되는 문자들을 모두 제거하고 본 내용만 추출한다. 위 과정을 통해 얻어진 콘텐츠는 콘텐츠 분류 모듈로 넘겨 해당 웹 페이지가 어떤 범주에 속하는지 분류한다.

콘텐츠 분류 모듈은 분류 규칙을 만드는 부분과 콘텐츠의 범주를 분류하는 부분으로 구성되어 있다. 분류 규칙을 만드는 부분은 온톨로지 데이터베이스로부터 각 범주의 정보를 얻어와 3장에서 기술한 바와 같이 네 가지 기계학습 알고리즘을 통해 생성하도록 되어 있고, 콘텐츠의 범주를 분류하는 부분은 생성된 규칙을 임의의 웹 페이지에 적용하여 온톨로지에 의해 정의된 범주 별로 분류하도록 되어 있다. 생성된 규칙 중 하나의 예로 CN2[2]가 생성한 규칙은 (그림 2)와 같다. (그림 2)에서 사용된 단어 ‘java’, ‘string’, ‘galaxi’ 등은 3장에서 소개한 특성 추출 방법과 정보이득 계산에 의해 선출된 단어들로 이들을 기계학습의 입력값으로 사용한 것이다. CN2는 ‘IF ... THEN ... ELSE’ 구조로 분류 규칙을 표현한다. 즉, 임의의 웹 페이지가 포함하는 단어들 중에 ‘java’, ‘string’이 없고 ‘galaxi’가 있다면 그것은 ‘SC’ 즉, ‘과학(Science)’ 범주에 속한다는 것을 뜻한다. 그 밖의 ‘SP’, ‘BF’, ‘PL’은 각각 ‘스포츠’, ‘은행과 금융’, ‘프로그래밍 언어’를 지칭한다.

```

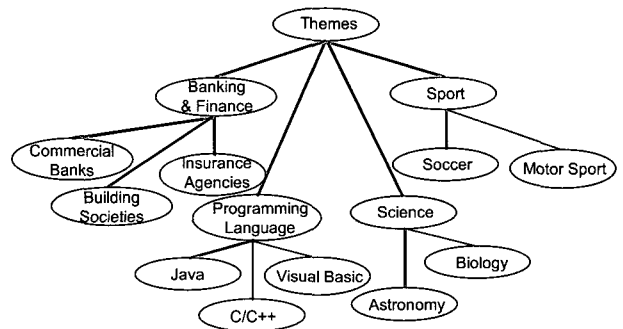
CN2 Decision List :
IF "java" = "false"
AND "string" = "false"
AND "galaxi" = "true"
THEN class = "SC" [0 0 402 0]
ELSE
IF "account" = "false"
AND "astronom" = "false"
AND "hole" = "false"
AND "footbal" = "true"
THEN class = "SP" [0 0 0 583]
ELSE
IF "bank" = "true"
AND "loan" = "true"
AND "distanc" = "false"
THEN class = "BF" [405 0 0 0]
ELSE
IF "financi" = "false"
AND "java" = "true"
AND "code" = "true"
AND "sky" = "false"
AND "manchest" = "false"
THEN class = "PL" [0 552 0 0]
.....

```

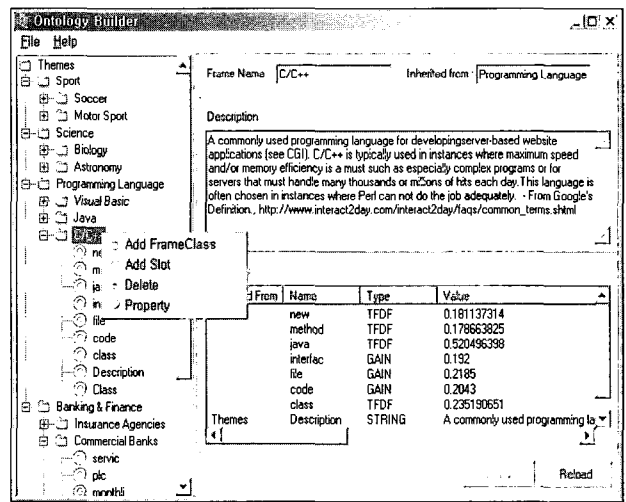
(그림 2) CN2에 의해 생성된 분류 규칙의 예

### 5. 실험 결과

우리는 자동화된 웹 페이지 분류 시스템의 분류 정확도에 대한 성능을 검증하기 위하여 Sinka[21]의 연구에서 제공하는 데이터 집합을 사용하였다. Sinka[21]가 사용한 데이터 집합은 웹 페이지 클러스터링 알고리즘을 평가하기 위해 제안한 평가 문서 집합으로 본 시스템의 웹 페이지 분류의 평가에도 적합하다고 판단된다. 실험을 위하여 1.6GHz의 펜티엄 IV PC를 사용하였다.



(그림 3) Sinka [21]가 제공한 네 가지 범주 : 은행과 금융, 프로그래밍 언어, 과학, 스포츠



(그림 4) 온톨로지 생성기를 이용하여 프로그래밍 언어 범주를 정의한 예제 화면

Sinka의 데이터 집합을 상위레벨에서의 범주로 나타내면 (그림 3)과 같다. 데이터 집합은 은행과 금융, 프로그래밍 언어, 과학, 스포츠로 총 네 개의 범주로 구성되어 있다. 주어진 데이터 집합에서 각 범주를 계층구조로 하는 온톨로지를 정의하기 위하여, (그림 4)와 같은 온톨로지 생성기(ontology builder)를 설계 및 구현하였다. 온톨로지 생성기를 이용하여 각 분류 범주에 속한 단어들을 정의하고 유지하









위 결과를 정리해 보면 무조건 많은 단어를 선택하는 것이 좋은 분류성능을 가져오지 않음을 알 수 있다. 이와 같은 현상을 <표 10>과 (그림 8)에 정리하였다. <표 10>은  $\mu$  값을 0.01부터 0.20까지 0.01씩 증가시킬 때 증가되는 분류 정확도를 추가로 선택되는 단어 개수로 나눈 값을 나타낸다. 굵은 글씨체로 강조된 수치는 기울기가 1.00이상인 것들을 표시한다.

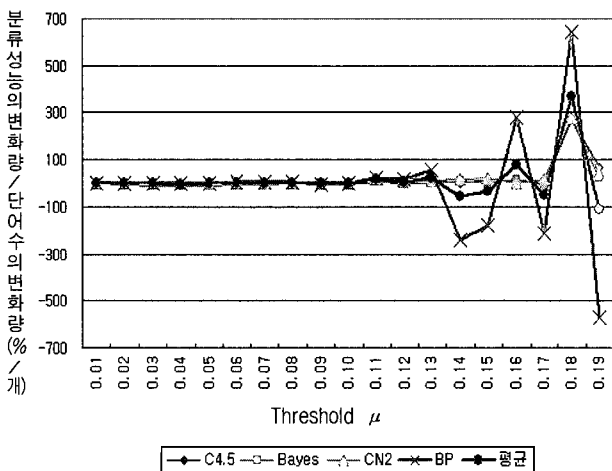
선택된 단어의 수는 비용과 반비례 관계에 있기 때문에 단어의 수는 적으면 적을수록 좋다. 그러나 단어의 수를 줄이기 위해서는 임계값  $\mu$ 를 높여야 하는데 이것은 분류 성능의 감소를 가져온다. 따라서 이를 절충하기  $\mu$ 값에 따른 분류 성능의 변화 량이 가장 큰 구간 [0.11~0.12]에서  $\mu$  값 0.11을

선택하여 실험하는 것이 분류 성능과 비용 측면에서 유리하다고 판단된다. 따라서 우리는 기계학습 알고리즘의 입력값으로 쓰일 단어를 선택하기 위해  $\mu$  값 0.11을 선택하여 실험을 하였다.

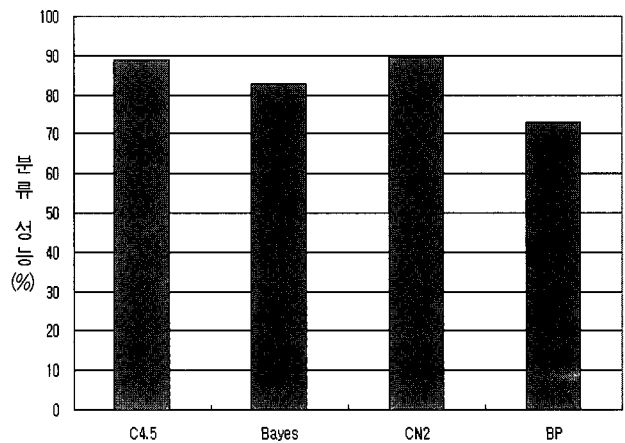
순위가 매겨진 단어들에 대해 식 (4)를 적용해 임계값  $\mu$ 를 0.11로 정하고, 임계값보다 큰 정보이득 값을 갖는 단어들을 선택한다. 그 결과 범주 당 13개, 총 78개의 단어가 선택되었다. 마지막으로 선택된 단어들로 조합을 만들어 각각의 조합이 어떤 범주에 해당하는지 파악해 기계학습에 사용될 수 있게 만든다. 기계학습의 입력값으로 사용될 이 튜플의 크기는 엔트로피 계산에 의해 그 크기를 줄일 수 있는 것이다.

<표 10> 통합 범주 분류 상에서  $\mu$ 값의 구간별 분류 성능의 변화량

$\mu$ 구간	단어수	단어 수의 변화량	분류 성능의 변화량(분류 성능의 변화량/단어 수의 변화량 %/개)				
			C45	Bayes	CN2	BP	평균
0.01 ~ 0.02	1,172	267	0.12	-1.62	0.06	0.39	-0.26
0.02 ~ 0.03	905	224	-0.02	-6.16	0.11	-0.29	-1.59
0.03 ~ 0.04	681	179	0.08	-10.75	-0.08	0.70	-2.51
0.04 ~ 0.05	502	110	0.40	-13.00	-0.13	<b>1.30</b>	-2.86
0.05 ~ 0.06	392	106	0.05	-9.10	0.70	<b>1.02</b>	-1.83
0.06 ~ 0.07	286	57	-0.68	-3.28	-0.26	<b>3.74</b>	-0.12
0.07 ~ 0.08	229	56	0.96	-3.18	0.10	<b>8.16</b>	1.51
0.08 ~ 0.09	173	41	<b>1.67</b>	-1.45	<b>1.85</b>	3.71	<b>1.45</b>
0.09 ~ 0.10	132	28	0.13	<b>1.46</b>	-0.07	-2.93	-0.35
0.10 ~ 0.11	104	26	<b>3.70</b>	-4.91	<b>1.28</b>	0.73	0.20
<u>0.11 ~ 0.12</u>	<u>78</u>	<u>19</u>	<u><b>9.75</b></u>	<u><b>7.21</b></u>	<u><b>15.40</b></u>	<u><b>24.26</b></u>	<u><b>14.16</b></u>
0.12 ~ 0.13	59	17	<b>2.18</b>	<b>2.07</b>	<b>3.49</b>	<b>15.71</b>	<b>5.86</b>
0.13 ~ 0.14	42	12	<b>14.20</b>	-1.85	<b>11.73</b>	<b>57.75</b>	<b>20.46</b>
0.14 ~ 0.15	30	1	-1.85	<b>9.26</b>	<b>14.82</b>	-241.00	-54.69
0.15 ~ 0.16	29	7	<b>18.78</b>	<b>2.38</b>	<b>23.55</b>	-176.57	-32.96
0.16 ~ 0.17	22	4	<b>10.18</b>	<b>18.06</b>	0.00	<b>278.75</b>	<b>76.75</b>
0.17 ~ 0.18	18	2	<b>1.86</b>	-6.48	<b>18.52</b>	-215.50	-50.40
0.18 ~ 0.19	16	4	<b>276.90</b>	<b>270.42</b>	<b>282.00</b>	<b>644.00</b>	<b>368.33</b>
0.19 ~ 0.20	12	2	64.81	23.15	44.44	-570.50	-109.52



(그림 8) 통합 범주 분류 상에서  $\mu$ 값의 구간별 분류 성능의 변화량



(그림 9) 기계학습 알고리즘의 분류 성능

자동화된 웹 페이지 분류 시스템을 구현하기 위해서 위에서 생성한 튜플들을 기계학습의 학습 데이터로 사용한다. 본 논문에서 사용한 기계학습 알고리즘은 C4.5[14], 베이저언 분류법(naïve Bayesian classifier)[4], CN2[2], 역전파(back-propagation) 신경망[24] 알고리즘으로 총 네 개이다. 이들 알고리즘에 의해 생성된 규칙들은 실시간으로 임의의 웹 페이지들을 분류하는데 사용된다.

최종 분류정확도를 측정하기 위하여, 위에서 생성한 규칙을 3,600개의 웹 페이지에 적용해 (그림 3)에서 제시한 네 개의 범주 중 하나로 분류하는 실험을 하였다. (그림 9)는 각각의 기계학습 알고리즘에 의해 생성된 분류 규칙이 갖는 분류 정확도 성능 측정값을 그래프로 나타낸 것이다. CN2의 성능은 89.05%로 네 개의 알고리즘 중에 가장 좋은 성능을 보여줬고, 가장 나쁜 성능을 보인 베이저언 분류법도 73.11%의 분류 성능을 보였다. 따라서 3,600개의 웹 페이지에 대해 실시한 분류 실험은 평균적으로 83.52%의 분류 성능을 가졌다고 할 수 있다.

## 6. 결론 및 향후과제

본 논문에서는 온톨로지에 기반한 자동화된 웹 페이지 분류 시스템을 제안하였다. 이러한 시스템을 개발하기 위하여, 온톨로지의 각 분류 범주에 속하는 단어들이 가져야 할 두 가지 속성을 (1) 각 분류 범주를 대표할 수 있으며, (2) 각 분류 범주를 다른 분류 범주와 뚜렷이 구분할 수 있는 특성으로 구체화 하였다. 따라서, 첫 번째 속성을 위하여 단어들의 가중치를 계산하여 각 분류 범주에 속하는 문서 집합을 대표하는 단어들을 선택하였으며, 두 번째 속성을 구현하기 위하여 선택된 단어들의 정보이득 값을 계산하여 우선순위를 결정하였다. 그리고 추출된 특성 즉, 이와 같이 선정된 단어들을 가지고 기계학습 알고리즘을 적용하여 분류 규칙을 생성하였다.

실험에서는 본 논문이 제안한 웹 페이지 분류 시스템의 분류 정확도에 대한 성능을 측정하였다. 그 결과 제안된 방법이 기계학습의 입력값으로 쓰이는 단어의 수를 효과적으로 줄일 수 있음을 알 수 있었다. 분류된 웹 페이지들의 집합으로부터 대표성을 갖는 단어를 추출하는 실험에서, 평균적으로 250개의 단어로부터 분류 범주를 대표하는 78개의 단어를 선정할 수 있었다. 이때의 분류 정확도에 대한 성능은 약 83.52%를 보였다.

본 연구의 결과를 다양한 도메인의 웹 페이지 분류를 위하여 적용할 수 있다고 판단한다. 특성을 추출할 때 HTML

문서의 구조적 특징이나 하이퍼링크와 같은 웹 페이지의 특징을 추가하여 이 특징들에 의한 성능 향상을 실험해 볼 수 있으며, 여러 가지 다른 도메인의 웹 페이지 분류를 수행할 수 있을 것으로 기대한다. 또한 분류 시스템에 새로운 웹 페이지들이 지속적으로 입력될 때, 이에 의한 온톨로지와 분류규칙의 갱신 및 성능 변화를 실험해 볼 수 있을 것이다.

본 연구를 토대로 향후에는 다양한 도메인에 대한 웹 페이지 분류를 해볼 계획이며, 특히 자체적으로 수집한 30,000개의 웹 페이지들에 대하여 웹 페이지의 유해성 여부를 분류하여 시스템의 유용성을 분석할 계획이다. 또한 새로운 웹 페이지들을 이용하여 분류 시스템의 규칙을 갱신할 때, 본 논문에서 제안한 단어 선택 방법의 유용성을 검증해 볼 것이다. 즉, 제안된 방법을 사용하여 갱신된 분류 규칙이 새로운 웹 페이지의 특성을 어느 정도 반영하는지 분석하며, 갱신된 규칙을 기존의 웹 페이지 집합에 적용하여 성능을 평가할 계획이다.

## 참 고 문 헌

- [1] R. Baeza-Yates and B. Ribeiro-Neto, *Modern Information Retrieval*, ACM Press/Addison-Wesley, New York, 1999.
- [2] P. Clark and T. Niblett, "The CN2 Induction algorithm," *Machine Learning Journal*, Vol.3, No.4, pp.261-283, 1989.
- [3] C. Ding, C. Chi, J. Deng and C. Dong, "Centralized content-based Web filtering and blocking : how far can it go?," Proc. of 1999 IEEE International Conference on Systems, Man and Cybernetics, Vol.2, pp.115-119, October, 1999.
- [4] R. Hanson, J. Stutz and P. Cheeseman, *Bayesian Classification Theory*, Technical Report FIA-90-12-7-01, NASA Ames research Center, AI Branch, 1991.
- [5] L. Holder, ML v2.0, *Machine Learning Program Evaluator*, available on-line, <http://ranger.uta.edu/~holder/courses/cse6363/ml2.0.tar.gz>.
- [6] C. Jenkins, M. Jackson, P. Burden and J. Wallis, "Automatic RDF metadata generation for resource discovery," Proc. of 8th International WWW Conference, Toronto, pp.11-14, May, 1999.
- [7] Lawrence Berkeley National Labs Network Research Group, *libpcap*, available on-line, <http://ftp.ee.lbl.gov>.
- [8] Y. Ng, J. Tang and M. Goodrich, "A binary-categorization

approach for classifying multiple-record Web documents using application ontologies and a probabilistic model," Proc. of 7th International Conference on Database Systems for Advanced Applications, pp.58-65, April, 2001.

[9] S. Noh, C. Lee, K. Choi and G. Jung, "Detecting Distributed Denial of Service(DDoS) Attacks Through Inductive Learning," Lecture Notes in Computer Science 2690, pp.286-295, Springer, 2003.

[10] S. Noh, H. Seo, J. Choi, K. Choi and G. Jung, "Classifying Web Pages Using Adaptive Ontology," Proc. of the IEEE International Conference on Systems, Man and Cybernetics, pp.2144-2149, Washington, D.C., October, 2003.

[11] N. F. Noy and D. L. Mcguinness, "Ontology development 101 : A guide to creating your first ontology," *Knowledge Systems Laboratory(KSL), Department of Computer Science, Stanford : Technical report*, KSL-01-05, 2001.

[12] S. Parent, B. Mobasher and S. Lytinen, "An adaptive agent for web exploration based on concept hierarchies," Proc. of 9th International Conference on Human Computer Interaction, New Orleans, August, 2001.

[13] R. Prabowo, M. Jackson, P. Burden and H. Knoell, "Ontology-Based Automatic Classification for the WEB Pages : Design, Implementation an Evaluation," Proc. of 3rd International Conference, Singapore, pp.182-191, 2002.

[14] J. R. Quinlan, *C4.5 : Programs for Machine Learning*, Morgan Kaufmann, 1993.

[15] J. Reynolds and J. Postel, "Assigned Numbers," STD 2, RFC 1700, October, 1994.

[16] D. Riboni, "Feature Selection for Web Page Classification," EURASIA-ICT 2002 Proc. of the Workshops, Shiraz, Iran, October 2002.

[17] S. M. Ruger and S. E. Gauch, *Feature Reduction for Document Clustering and Classification*, Technical report, Computing Department, Imperial College, London, 2000.

[18] G. Salton, and C. Buckley, "Term weighting approaches in automatic text retrieval," *Information Processing and Management*, Vol.24, No. 5, pp. 513-523, 1988.

[19] F. Sebastiani, "Machine learning in automated text categorization," *ACM Computing Surveys*, Vol.34, No.1, pp.1-47, 2002.

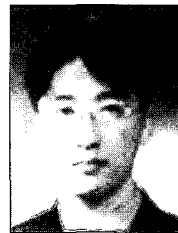
[20] C. E. Shannon, "A mathematical theory of communication," *Bell System Technical Journal*, Vol.27, pp.379-423 and 623-656, July/October, 1948.

[21] M. P. Sinka and D. W. Corne, "A large benchmark dataset for web document clustering," *Soft Computing Systems : Design, Management and Applications, Frontiers in Artificial Intelligence and Applications*, Vol.87, pp.881-890, 2002.

[22] N. Soonthornphisaj, P. Chartbanchachai, T. Pratheeptham, and B. Kijisirikul, "Web page categorization using hierarchical headings structure," Proc. of 24th International Conference on Information Technology Interfaces, Vol.1, pp.37-42, 2002.

[23] A. Sun, E. Lim and W. Ng, "Web classification using support vector machine," WIDM'02, Virginia, November, 2002.

[24] D. R. Tvetter, *Backprop Package*, available on-line, <http://www.dontveter.com/nsoft/bp042796.zip>, 1996.



최재혁

e-mail : jaehyuk04.choi@samsung.com

2002년 아주대학교 정보통신대학정보 및 컴퓨터공학부(학사)

2004년 아주대학교 정보통신전문대학원 정보통신공학과(석사)

2004년~현재 삼성전자 네트워크 사업부

관심분야 : 인터넷 보안, 인공지능, 실시간 시스템 등



서혜성

e-mail : retry@ajou.ac.kr

2003년 아주대학교 정보통신 대학 정보 및 컴퓨터공학부(학사)

2003년~현재 아주대학교 정보통신전문 대학원 정보통신공학과

관심분야 : 네트워크 보안, 분산 시스템, 인공지능 등



노상욱

e-mail : sunoh@catholic.ac.kr

1987년 서강대학교 생명과학(이학사)

1989년 서강대학교 컴퓨터공학(공학석사)

1999년 텍사스 주립대(Arlington)

컴퓨터공학(공학박사)

1989년~1995년 국방과학연구소 연구원

2000년~2002년 미조리 주립대(Rolla) 컴퓨터학과 조교수

2002년~현재 가톨릭대학교 컴퓨터정보공학부 조교수

관심분야 : Knowledge Management, Intelligent Agent, Multi-Agent System, Machine Learning, Distributed Real-Time System 등



### 최 경 희

e-mail : khchoi@madang.ajou.ac.kr

1976년 서울대학교 사범대학 수학교육과  
(학사)

1979년 프랑스 그랑데폴 Enseiht  
정보공학과(공학석사)

1982년 프랑스 Paul Sabatier 정보공학과  
(공학박사)

1982년~현재 아주대학교 정보통신전문대학원 교수

관심분야 : 운영체제, 분산시스템, 실시간 및 멀티미디어 시스템 등



### 정 기 현

e-mail : khchung@madang.ajou.ac.kr

1984년 서강대학교 공과대학 전자공학과  
(학사)

1988년 미국 Illinois 주립대 EECS(공학석사)

1990년 미국 Perdue 전기전자공학부  
(공학박사)

1991년~1992년 현대전자 반도체 연구소

1993년~현재 아주대학교 전자공학부 교수

관심분야 : 컴퓨터구조, VLSI설계, 멀티미디어 및 실시간 시스템 등