

AC 머신을 이용한 유해 사이트 차단에 관한 연구

정현수* · 정규철* · 김후남* · 박기홍*

A Study about interception on Hurtfulness Site using Aho-Corasik machine

Hyun-Soo Jung* · Kyu-Cheol Jung* · Hoo-Nam Kim* · Kihong Park*

요 약

지식정보사회로의 변천은 우리의 삶을 보다 편리하고 윤택하게 해주고 있지만 미처 예상하지 못했던 부작용 또한 적지 않게 발생하고 있어 이에 대한 해결방안을 모색하는 것이 시급한 실정이다. 아직 인성이 발달되지 못한 청소년들이 정보통신망을 통해서 유통되고 있는 수많은 음란물과 폭력물과 같은 유해정보에 무방비상태로 노출되고 있는 것도 정보화사회의 대표적인 역기능 중의 하나라고 할 수 있다. 이 문제점 해결을 위해 본 논문에서는 사이트 상에서 제공되어지는 내용을 AC 머신을 이용하여 유해 단어를 추출하고 유해 정보 데이터베이스를 이용해서 유해 단어에 가중치를 부여했다. 그 결과로 유해 정보를 포함한 사이트는 90%의 차단률을 보여 효율적인 시스템으로 판명되었다.

ABSTRACT

Change is doing our life more conveniently and abundantly by knowledge information society, but side effect and that is happening considerable and gropes solution in reply that did not expect in advance is urgent real condition. It can be called one of representative dysfunction of information-oriented society that human nature is revealed in open state to great many objectionable material and poisonous information such as violence kind that teenagers who do not grow are gotten abroad through Information network system yet. So, to solve these fallacy, word-weighting process, where several harmful words which can be obtained in internet site are discriminant and weighted, is utilized by using AC machine. At the result, the isolation rate of harmful site rose up to 90%, which means this process is greatly efficient.

키워드

유해, 유해사이트, 유해어 추출기, 차단 시스템

1. 서 론

정보통신기술의 발달로 우리는 인터넷이라는 가상공간(cyber space)을 통하여 지구 반대편에 있는 사람과 실시간으로 쌍방향 의사전달을 할 수가 있고, 굳이 외출하지 않고도 집안에 앉아서 전자상거래를 통해 각종 재화와 서비스를 공급받을 수 있는 편리한 시대에 살고 있다. 아마도 21세기는 정보통신기술의 발달과 더불어 그 어느 때보다도 변화의 속도를 더해 갈 것이 분명하며, 지식과 정보가 결합되어 우리의 삶을 보다 편리하고 윤택하게 해 줄 지식정보사회로 이행되어 나갈 것으로 전망된다.

그러나 산이 높으면 골이 깊듯이 지식정보사회의 도래가 항상 밝은 면만을 우리 인류에게 가져다 주는 것은 아니다. 지식·정보격차의 발생, 사이버 테러의 성행, 음란물의 범람 등 기존 산업사회에서는 예상하지 못했던 많은 문제점을 던져주고 있는 것이 현실이다. 이러한 문제점들 가운데 인터넷을

통해 제공되는 음란, 엽기, 자살, 도박, 폭력, 마약 등에 대한 유해한 정보를 차단하기 위한 도구 개발에 대한 관심이 고조되고 있으나 전 세계의 어느 나라에서나 인터넷을 사용하기 때문에 어느 한 특정 국가에서 법·제도로서 규제하기 힘들다. 그렇기 때문에 기술적으로 유해한 정보를 차단할 수 있는 소프트웨어를 개발 및 보급하고 있다. 이러한 소프트웨어들의 대표적인 예가 등급표시제와 차단 목록 기반의 유해 정보 차단 소프트웨어이다. 등급표시제는 내용물을 제공하는 사업자로부터 하여금 사이트에 등급을 표시해서 성인용과 청소년용으로 나누는 것이며 차단 목록 기반의 유해 정보 차단 소프트웨어는 차단 목록 데이터베이스를 이용해서 데이터베이스에 있는 사이트를 접속하면 차단하는 소프트웨어이다. 등급표시제는 등급을 표시하는 기관이나 사람에 의한 오류와 사전 검열이라는 문제점이 있으며 차단 목록 기반의 유해 정보 차단 소프트웨어는 차단 목록의 계속적인 업데이트가

이루어져야 하는 단점이 있다.

본 논문에서는 사이트 상에서 제공되어지는 내용을 가지고 명사를 추출하여 AC(Aho-Corasik)머신과 유해 정보 데이터베이스를 이용해 사이트의 유해성을 판단하는 시스템을 구축했다.

본 논문에서는 유해 정보를 차단하기 위한 전 단계작업으로 유해 정보의 목적과 유해정보의 유형에 대해 1장에서 다루지며 유해 정보 차단 시스템의 설계와 구현에 대해 2장에서 다루어진다. 마지막으로 결론에서는 유해 정보 차단 시스템의 테스트 결과와 향후 추가되어야 할 사항들에 대해서 제시한다.

II.. 본 론

1. 유해 정보의 차단 목적과 유해 정보의 유형

1.1 유해 정보의 차단 목적

인터넷을 통하여 음란, 엽기, 자살, 도박, 폭력, 마약 등의 유해 정보들이 아무런 여과 없이 청소년들에서부터 어른들에 이르기까지 무분별하게 유해 정보의 제공이 이루어지고 있다. 이러한 유해 정보로 인하여 청소년들의 폭력성향에 대한 행동을 강화시킬 수 있고, 유해정보에 반복하여 접촉하게 될 경우, 성과 관련된 가치와 규범이 모호하게 되고, 성 규범을 해체시켜 탈 억제적 행동을 초래할 가능성이 있기 때문에 더 이상은 유해 정보에 청소년들이 노출되지 않도록 막고자 하는 것이 목적이다.

1.2 유해정보의 유형

유해 정보로서 많이 나타나는 유형은 다음과 같다. 포르노 동영상 같은 야동(야한동영상), 누드 또는 포르노 사진과 같은 야사(야한사진), 섹스소셜, 강간소설 등등의 야설(야한소설)과 성폭력이나 살인을 조장, 살인청탁을 받는 홈페이지나 카페나 커뮤니티, 인터넷 자기, 웹 자키가 방송을 이끌어가는 성인용 인터넷 방송, 마약류에 대한 정보로 코카인이나 헤로인, 아편, 대마초 등에 대한 사용을 조장하는 사이트나 카페 등으로 좀 더 세부적으로 유해 정보의 유형을 나눈다면 다음과 같이 나눌수 있다.

1.2.1 P2P(peer to peer)파일 공유서비스

인터넷이용자는 P2P파일공유서비스에 접속하여 디지털콘텐츠를 자유롭게 이용할 수 있고, 특정한 경우를 제외하고 거의 무료 디지털콘텐츠 등의 정보를 이용할 수 있을뿐만 아니라 자신의 PC에 그대로 옮겨 놓을 수 있다. 이를 이용하여 야동이나 야사 야설 등의 유해 정보를 무료로 자신의 컴퓨터에 내려 받을수 있다.

1.2.2 www

대부분의 음란사이트들은 일단 접속하면 사용자 가 18세를 넘었는지 확인한다. 그런데 이것은 형식

에 불과한 절차여서 사용자가 아무리 어려도 '예(yes)'를 클릭하기만 하면 홍보용 무료사진을 볼 수 있으며 결제용 신용카드번호나 핸드폰 요금결제와 비밀번호의 입력을 요구받게 된다. 물론 일부 사이트들은 사용자의 주민등록번호를 확인하는 절차를 거친 후에 홍보용 무료사진을 보여주기도 하지만 이는 그나마 양심적인 경우에 속한다.

1.2.3 대화방(chat room)

인터넷에서 무료로 이용할 수 있는 서비스 중 하나가 대화방이다. 대화방도 이기적인 사람들에게 의하여 악용되고 있으며 자신의 성적 욕망을 채우려는 사람들이 대화방 이곳저곳을 기웃거리며 회생양을 찾아 헤매고 있다. 그런데 문제는 이들을 제재할 만한 방안이 거의 없다는 것이다. 예전엔 텍스트 위주의 대화방이 서비스 되었었지만 요즘은 동영상 카메라를 이용한 화상 채팅 서비스를 지원하고 있어서 컴퓨터에 설치된 동영상 카메라 앞에서 성에 관련된 행위를 하면 다른 사람들은 그 상황을 실시간으로 볼 수 있다.

1.2.4 성인용 전자게시판

자장 많은 사람들이 이용하는 전자게시판 중 한 군데는 무려 2만 5천 건에 달하는 사진 파일을 게시하고 있는데 대부분이 폭력성이 짙은 사진이나 성기가 노출된 사진이다. 이 게시판의 운영자는 로버트 토마스라는 사람으로 현재 구속 수감 중이다. 다행이 이러한 전자게시판들은 대부분 상업적인 목적으로 운영되기 때문에 경제적인 능력이 없는 아이들이 접속하기는 쉽지 않다.

2. 유해 정보 차단 시스템의 설계와 구현

2.1 유해 정보 차단 시스템의 설계

2.1.1 차단 시스템의 구성도

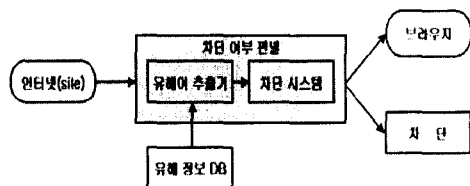


그림 1. 유해정보 차단 시스템 구성도

Fig 1. Construction of interception system on harmful information

시스템의 기본 구성은 사이트를 접속하면 인터넷 브라우저 상에서 바로 사이트를 나타내지 않고 접속하려는 사이트내의 내용의 단어를 추출한후 유해어 추출기와 유해 정보 데이터베이스를 이용해 사이트의 차단 여부를 판단하도록 하는 것이다. 이때 AC머신을 이용해서 태그를 제거한 소스로부터 명사 목록을 추출해낸다[10][11].

2.1.2 가중치 부여

AC 머신을 이용해 추출한 명사 목록을 유해 정보 데이터베이스에 있는 단어들과의 매치로 확인하고 추출된 유해 단어들의 가중치를 부여해 준다. 유해 단어에 가중치를 부여한 기준은 다음과 같다 [6].

표 1. 가중치 부여 기준표
Table 1. Standard table of grant with a weight

가중치	가중치 부여 기준
1	일상 용어(예:학교, 친구)
2	성에 관한 연상 용어(예:물건, 음수, 구멍)
3	신체 관련 연상 용어(예:음모, 질구, 가슴)
4	성에 관한 직접적인 용어(예: 강간, 애무)

2.2 구현 과정

다음의 구현 시스템은 브라우저의 주소 입력창을 통해서 입력된 주소와 해당주소의 소스를 파일로 저장하고 사이트의 소스를 분석할 수 있도록 태그 등을 제거해 유해어 추출기의 데이터로 활용한다. 유해어 추출기와 유해 정보 데이터베이스를 이용해 접속한 사이트의 유해성 가중치를 계산해 유해성 여부를 판단하여 차단하도록 한다.

실행시에 생성되는 파일들은 주소를 저장한 파일과 소스를 저장한 파일, 소스로부터 태그를 제거한 파일이 생성된다.

AC 머신으로 구성된 유해어 사전을 이용해 태그를 제거한 소스로부터 유해어를 추출해서 각 단어가 가진 가중치를 계산해낸다. 유해 정보의 판단 근거의 단어들은 경험치에 의해서 추출했으며 가중치 또한 경험치에 의해서 1~4 까지의 값을 부여했고 태그를 제거한 소스를 분석할 때 가중치의 평균을 구해서 경험치에 의한 평균값 1.5 이상을 넘을 경우 그 사이트는 유해 정보 사이트로 판명해서 차단하고 1.5 미만은 청소년들이 접속해도 무해한 사이트로 규정 하였다.

판단결과 = 가중치의 합 / 명사추출 개수의 수

3.3 평가

제안된 유해 정보 차단 시스템의 평가를 위해서 100개의 사이트를 방문해본 결과 유해 정보 사이트를 차단하는데 90% 정도의 차단률을 보였다. 그러나 유해 정보를 제공하지 않는 사이트에 대해서도 5% 정도가 차단되는 문제점이 발견되었다.

표 2. 유해 정보 차단 결과
Table 2. Result of interception on harmful information

	결과	차단률	
		적절	부적절
유해 O	73%	90%	10%
유해 X	27%	95%	5%

이로써 본 논문에서 구현한 유해 정보 차단 시스템은 유해한 정보를 갖고 사이트를 차단하는데 효율적인 시스템이다.

3.3.1 시스템 환경

- 1) System : 펜티엄 IV 2.4 GHz
- 2) Memory : 256 Mbyte
- 3) OS : Windows XP Professional
- 4) Program language : Visual C++ 6.0

III. 결 론

본 논문에서는 유해 정보를 가진 사이트를 차단하는데 접속을 시도하는 사이트의 소스를 AC머신과 유해 정보 데이터베이스를 이용해 차단하려는 시도를 했다. 아직은 미비한 점이 많지만 따로 소프트웨어를 설치하지 않아도 인터넷 브라우저를 통해서 자동으로 유해 정보 사이트를 차단할 수 있도록 했다. 기존의 차단 소프트웨어에서 볼 수 없는 내용 분석 기반의 AC머신을 이용해 유해 단어를 추출하고 유해 정보 데이터베이스를 가중치를 부여하여 효적으로 유해 정보 사이트를 차단하였다. 향후 본 논문에 추가로 인터넷 사용자들이 느끼지 못하는 아주 짧은 시간에 차단 여부를 판명할 수 있는 차단 시스템을 구현하도록 하겠다.

참고문헌

- [1] 류광재, "침입탐지에 의한 실시간 인터넷 해킹방지 연구", 명지대학교 컴퓨터공학부, 석사학위 논문, 1996
- [2] 성장열, "네트워크 모니터링과 MAPI를 적용한 웹메시징 시스템 설계에 관한 연구", 대전대학교대학원, 석사학위 논문, 2000
- [3] 김현중, "HTML을 지원하는 라이브러리를 이용한 웹 문서 생성 시스템의 설계 및 구현", 이화여자대학교 대학원 전자계산학과, 석사학위 논문, 1996
- [4] 이민구, "인터넷 유해정보 유입방지 방안에 관한 연구", 아주대학교 산업대학원 컴퓨터공학과, 석사학위 논문, 1998
- [5] 정희, "유해 정보 차단을 위한 검색 시스템의 설계와 구현", 창원대학교 대학원 전자계산학과, 석사학위 논문, 1999
- [6] 정명숙, "청소년을 위한 유해 웹 영상 차단 시스템의 구현", 경상대학교 교육대학원 전산교육, 2000
- [7] 정경수, "인터넷을 이용한 정보시스템의 전략적 활용", 경북대학교 경영대학원, 1907
- [8] Makoto Okada, Kazuaki Ando, Kazuhro Morita, Jun-ichi Aoe, "An Efficient Determination of Keywords for Compound Words", Proceedings of 18th ICCPOL, Vol

- 1, pp317-320, March 1999.
- [9] Kazuaki Ando, Toshiharu Kinoshita, Masami Shishibori, Jun-ichi Aoe, "An improvement of the Aho-Corasick machine", International Journal of Information Sciences, Vol 3, pp139-151, 1998.
- [10] A. V. Aho, M. J. Corasick, "Efficient string matching: an aid to bibliographic search", Comm. ACM, Vol.18, No.6, pp.333-340, 1975.
- [11] 이진관, "테이블을 이용한 AC 기반의 키워드 검색 기법", 군산대학교 컴퓨터정보과학과, 석사학위 논문, 2002.