

## 온톨로지를 이용한 정보 추출

김인수, 이복주

단국대학교 전자 컴퓨터 공학과

insu@ai.dankook.ac.kr, blee@dankook.ac.kr

### Information Extraction Using the Ontology

Insu Kim and Bogju Lee

Department of Computer Engineering, Dankook University, Korea

#### 요 약

정보 추출은 텍스트로 되어 있는 비 정형화된 데이터로부터 정형화된 정보를 추출하는 분야이다. 기존의 정보 추출이 구문 중심의 방법인데 비해 본 논문에서는 시맨틱 웹과 온톨로지를 이용한 의미 기반의 정보 추출을 시도한다. 또한 본 논문에서는 기존의 정보 추출 모델을 분류해 보고 반자동 정보 추출이라는 새로운 모델을 제시한다. 이 모델에 기반하여 개인 정보를 자동으로 정형화 시켜주는 정보 추출 도구를 개발하고 이를 소개한다.

#### 1. 서론

회원 가입이나 전자상거래시 배송지 정보와 같은 개인 신상 정보를 입력해야 하는 페이지 또는 이력서나 다른 일반 문서를 작성할 때 매번 똑같은 개인 신상 정보를 입력해야 하는 번거로움이 있다. 이러한 개인 신상 정보를 자동으로 채워지게 하여 간편하게 만들면 편하게 될 것이다. 여기에 정보 추출(Information Extraction)을 이용할 수 있다. 정보추출은 텍스트로 되어 있는 비 정형화된 데이터(unstructured data)로부터 정형화된 데이터(structured data)를 추출하는 분야이다. 기존의 정보 추출이 구문 중심의 방법인데 비해 본 논문에서는 시맨틱 웹과 온톨로지를 이용한 의미 기반의 정보 추출을 시도한다.

기존의 정보 추출 모델은 두 가지로 분류할 수 있다. 즉, 수동 정보 추출과 자동 정보 추출이 그것이다. 수동 정보 추출은 텍스트로부터 사람이 필드와 필드 값을 인지하여 추출하는 방식이고 자동 정보 추출은 그것을 전부 컴퓨터 프로그램에 맡기는 것이다. 지금까지의 연구 성과로도 완전한 자동 정보 추출은 달성하지 못한다. 본 논문에서는 이 두 가지 정보 추출 사이에 존재하는 반자동 정보 추출이라는 새로운 모델을 제시한다. 즉 프로그램으로 하여금 일단 정보 추출을 하게 한 다음 사용자가 확인하고 수정하게 하는 방식이다. 현실적으로 이 방식이 완전 자동 정보 추출보다 훨씬 유용하다.

본 논문에서는 이 모델에 기반한 정보 추출 도구를 개발하고 이를 실험하였다. 정보 추출 도구의 개발 도메인은 웹 사이트에서 흔히 사용하는 사용자의 인적 정보에 관한 것이다. 상업적 웹 사이트 또는 일반 웹 사이트에 회원 가입할 때 사용자는 흔히 신상 정보를 웹에 입력하게 된다.

또는 이력서나 다른 일반 문서를 작성할 때도 테이블을 신상 정보로 채우게 된다. 이럴 때 텍스트 기반으로 되어 있는 만들어진 원래 신상 정보 파일이 있고 이로부터 자동적으로 정보 추출이 되어 해당 목표 페이지나 테이블에 "로드" 된다면 매우 편리할 것이다. 로드 된 후 사용자는 혹시 잘못 옮겨진 필드를 수정하면 될 것이다.

#### 2. 기존 연구

현재 인터넷에는 수많은 웹사이트가 존재하고 이들 중 대부분은 사용자 확인을 위해 회원가입을 요구한다. 그래서 사용자는 수많은 웹사이트를 일일이 가입하고 각 사이트의 계정 정보를 관리해야 하는 번거로움이 생긴다. 이러한 문제를 해결하기 위해 마이크로소프트에서는 닷넷 패스पोर्ट를 제안하였다.

닷넷 패스पोर्ट는 <http://www.passport.net> 에서 등록하여 사용할 수 있다. 닷넷 패스पोर्ट에 등록된 하나의 계정으로 MSN메신저와 닷넷 패스पोर्ट를 채택한 사이트의 로그인할 수 있다. 윈도우에서는 사용자 별로 닷넷 패스पोर्ट를 저장하여 사용 가능하다. 또 MSN메신저에 로그인한 상태라면 메신저의 계정 정보로 웹사이트의 로그인이 이루어진다.

MSN메신저는 기본적으로 닷넷 패스워드를 계정으로 사용한다. 웹사이트의 경우, 로그인 방식을 닷넷 패스पोर्ट로 채택한 사이트에 한하여 간단히 로그인 버튼을 클릭하는 것만으로 자동으로 로그인이 된다. 그러나 닷넷 패스पोर्ट를 채택한 웹사이트는 MSN홈페이지(<http://www.msn.co.kr>) 이외에는 찾아보기 힘들다. 데브피아(<http://www.devpia.com>)는 닷넷 패스पोर्ट 로그인을 서비스 하였지만 현재 중단한 상태이다.

이러한 방식은 회원관리부분을 닷넷 패스포트에 맡김으로써 편리하지만, 모든 웹사이트가 이 방식을 채택하기란 사실상 힘들다. 수많은 웹사이트에 독립적으로 사용자의 정보를 관리하여 로그인 폼이나 회원가입 폼을 채워주는 프로그램이 좀 더 현실적인 방법이 될 것이다.

사용자 정보를 관리해 주는 프로그램 중 하나로 '알프레드'가 있다. '알프레드'는 각 사이트마다 로그인시 미리 저장된 아이디와 패스워드를 이용해 자동으로 로그인하는 기능과 회원 가입시 사용자 정보를 채워주는 기능을 한다. '알프레드'를 처음 시작하면 개인정보를 입력하게 되고, 이 정보를 회원 가입시 이용하게 된다. 웹브라우저에 회원 가입 폼이 나타나게 되면 '알프레드'는 이를 감지하여 저장된 사용자 정보를 입력할 것인지 묻는 팝업창을 띄우게 된다.

사용자 정보를 입력한 후에도 그림 1 과 같이 '휴대폰 번호' 필드는 채워지지 않은 것을 확인 할 수 있다. 이는 각 필드의 명칭이 정확히 매치가 되지 않아서, 요구하는 정보가 있음에도 불구하고 필요한 곳에 입력되지 않은 것이다.

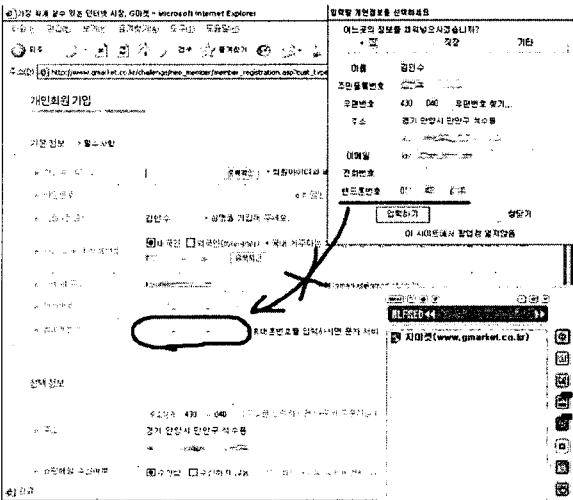


그림 1. '알프레드' 입력 예

이를 해결하기 위해서는 각 사이트가 요구하는 정보의 필드명과 '알프레드'에 저장되어 있는 정보의 필드명이 정확히 매치되도록 추가적인 정보가 필요하다. 각 사이트마다 회원가입페이지의 폼이 다르게 작성되어 있기 때문에, 추가적인 정보의 양은 사이트의 수만큼 필요하게 된다. 또 새로운 사이트에 대해서도 적용될 수 있도록 지속적인 업데이트가 요구되어 해결방법으로는 적합하지 않을 수 있다.

이러한 상황에서 온톨로지를 사용하여 요구되는 필드명이 어떤 필드와 일치하는지 추론하여 매치 시킨다면 추가적인 정보 없이 매치 가능할 것이다.

### 3. 시스템 소개 및 구조

본 논문에서 소개할 시스템은 미리 사용자의 신상 정보

를 가지고 있는 파일이 있고, 정보 입력 페이지의 field name에 따라 해당되는 정보를 파일로부터 가져와서 채워주게 된다. 이때 서로 다른 field name을 가지고 있지만 같은 정보를 요구하는 field의 경우, 서로 다른 field name 이 같은 field 라는 것을 알려주어 해당 정보를 채워야 한다. 여기서 Ontology를 사용하여 field가 같다는 것을 명시해준다. 예를 들어 '성명' 과 '이름' 은 같은 의미를 갖는 field 이므로 같이 취급해야 한다.

Ontology 를 사용하여 채워진 정보는 정확하지 않은, 관계없는 정보가 채워질 수 있기 때문에 사용자가 채워진 정보를 확인하고 수정한 후 Submit 하는 반자동 시스템으로 한다. 본 논문에서는 이 시스템을 Personal Information Manager라고 명칭 하기로 한다.

다음 그림 2 는 Personal Information Manager 의 전체적인 구조를 나타낸 그림이다.

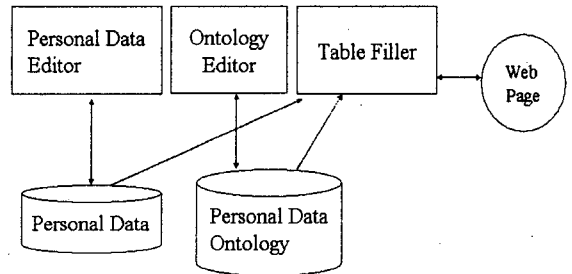


그림 2. Personal Information Manager 전체구조

각 구성요소를 살펴보면 크게 Personal Data Editor, Ontology Editor 그리고 Table Filler로 구성되어 있다.

Personal Data Editor 는 사용자의 신상정보를 입력 및 수정하여 Personal Data를 작성하는 역할을 한다. Personal Data 는 사용자의 컴퓨터 내에 저장되어있다가 필요에 따라 사용되게 된다.

Ontology Editor 는 Personal Data Ontology를 관리하기 위한 목적으로 사용된다. Personal Data Ontology 에는 사용자의 신상정보에 대한 각 필드의 유사어의 정보와 각 필드간의 관계 정보가 저장되어 있다. 그림 3 은 Ontology 내의 Field Name 간의 관계를 트리 구조로 보여주고 있다.

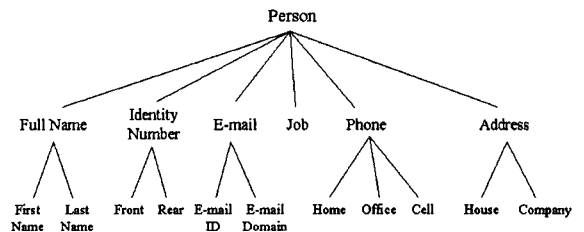


그림 3. Personal Data Ontology

