

Automatic Extraction of Web Search Interfaces for Interface Schema Integration

Hai He, Weiyi Meng
Dept. of Computer Science
SUNY at Binghamton
Binghamton, NY 13902

{haihe, meng}@cs.binghamton.edu

Clement Yu
Dept. of Computer Science
Univ. of Illinois at Chicago
Chicago, IL 60607

yu@cs.uic.edu

Zonghuan Wu
Center for Adv. Compu. Studies
Univ. of Louisiana at Lafayette
Lafayette, LA 70504

zwu@cacs.louisiana.edu

ABSTRACT

This paper provides an overview of a technique for extracting information from the Web search interfaces of e-commerce search engines that is useful for supporting automatic search interface integration. In particular, we discuss how to group elements and labels on a search interface into attributes and how to derive certain meta-information for each attribute.

Categories & Subject Descriptors

H.3.5 [Information Storage and Retrieval]: Online Information Services – Commercial Services, Web-based Services. H.5.2 [Information Interface and Representation]: User Interfaces – User Interface Management Systems.

General Terms

Algorithms, Management, Performance, Design.

Keywords

Search interface representation, Search interface extraction, Search engine, Metasearch engine.

1. INTRODUCTION

This paper provides an overview of a technique for extracting information from the Web search interfaces of e-commerce search engines (ESEs) that is useful for constructing e-commerce metasearch engines (EMSEs). More specifically, the aim is to extract information that supports automatic search interface integration. In [1], we presented a model to capture the information on a search interface and a tool (WISE-Integrator) that can *automatically* build a unified search interface over multiple heterogeneous ESE search interfaces of the same product domain based on the model. However, in [1], the information in the model was manually obtained from each search interface. In this paper, we outline our technique for *automatically* constructing the model.

A Web search interface for e-commerce typically contains some HTML form control **elements** such as *textbox* (i.e., a single-line text input), *radio button*, *checkbox* and *selection list* (i.e., a pull-down menu) that allow a user to enter search information. Each *element* usually has a **label** – a descriptive text – associated with it. An element may have one or more **values**. For example, a selection list usually has a list of values (options) for users to select, and a radio button/checkbox usually has a single value. Logically, elements and their associated labels together form different **attributes** of the products in the underlying database of the search engine. Often an attribute may consist of one or more labels and elements. For example, the *author* attribute in Figure 1

has four elements including a textbox and three radio buttons. The label of an attribute is referred to as **attribute name**, and element(s) of the attribute are treated as **attribute domain**. If an attribute contains multiple elements, these elements may be related in some way. For example, among the four elements of *author* in Figure 1, the textbox is treated as **domain element** while the three radio buttons are treated as **constraint elements** since each of them specifies a constraint on the domain element. In addition to such explicit composition information of attributes, each attribute is also *implicitly* associated with a set of meta-information such as the domain type (e.g., finite, range) and value type (e.g., date, currency), which is essential for enhancing attribute matching [1].

The screenshot shows the Amazon.com search interface. It features several search fields: 'Author:', 'Title:', 'Subject:', 'ISBN:', and 'Publisher:'. Each field has a corresponding radio button with a label: 'First name/initials and last name', 'Start of last name', 'Exact name', 'Title word(s)', 'Start(s) of title word(s)', 'Exact start of title', 'Subject word(s)', 'Start of subject', and 'Start(s) of subject word(s)'. There is a 'Search Now' button in the top right. Below the search fields is a 'Refine your search (optional):' section with various filters: 'Used Only' (checkbox), 'Format' (dropdown), 'Reader age' (dropdown), 'Language' (dropdown), 'Publication date' (dropdown and text input), and 'Sort results by' (dropdown). The 'Publication date' dropdown is set to 'All dates' and the text input shows '(e.g. 1999)'. The 'Sort results by' dropdown is set to 'Featured Items'.

Figure 1. The book search interface of amazon.com.

2. INTERFACE EXTRACTION

In our work, interface extraction consists of two major steps: (1) *Attribute extraction*: given all search interfaces of a domain, extract labels and elements appearing in each form and then group them into logic attributes; (2) *Attribute analysis*: analyze the labels and elements of each attribute to derive meta-information of the attribute such as domain type and value type.

2.1 Attribute Extraction

We observe that labels and elements of the same attribute have a certain layout and are usually close to each other, and that in most cases they share some similar information. The layout of labels and elements can be captured as an *interface expression (IEXP)*. For a given search interface, its IEXP is a *string* consisting of three basic items 't', 'e' and '|', where 't' denotes a label/text, 'e' denotes an element, and '|' denotes a row delimiter which represents a physical row border in the search interface. IEXP provides a high-level description of the layout of different labels and elements on the interface while ignoring the details like the values of the elements and the actual implementations of laying out labels and elements. For example, the IEXP of the search

