# Chapter 7.12 Search Engine-Based Web Information Extraction

**Gijs Geleijnse** Philips Research, The Netherlands

Jan Korst Philips Research, The Netherlands

### ABSTRACT

In this chapter we discuss approaches to find, extract, and structure information from natural language texts on the Web. Such structured information can be expressed and shared using the standard Semantic Web languages and hence be machine interpreted. In this chapter we focus on two tasks in Web information extraction. The first part focuses on mining facts from the Web, while in the second part, we present an approach to collect community-based meta-data. A search engine is used to retrieve potentially relevant texts. From these texts, instances and relations are extracted. The proposed approaches are illustrated using various case-studies, showing that we can reliably extract information from the Web using simple techniques.

# INTRODUCTION

Suppose we are interested in *'the countries where Burger King can be found', 'the Dutch cities with a university of technology'* or perhaps *'the genre of the music of Miles Davis'*. For such diverse factual information needs, the World Wide Web in general and a search engine in particular can provide a solution. Experienced users of search engines are able to construct queries that are likely to access documents containing the desired information. However, current search engines retrieve Web pages, not the information itself<sup>1</sup>. We have to search within the search results in order to acquire the information. Moreover, we make implicit use of our knowledge (e.g. of the language and the domain), to interpret the Web pages.

Web Corpus	Newspaper Corpus
<b>Redundancy</b> . Because of the size of the Web, we can expect information to be duplicated, or formulated in various ways. If we are interested in a fact, we have to be able to identify just one of the formulations to extract it.	<b>No or fewer redundancy.</b> Especially for smaller corpora, we cannot expect that information is redundantly present.
<b>Temporal and unreliable</b> . The content of the Web is created over several years by numerous contributors. The data is thus unreliable and may be out-dated.	<b>Constant and reliable</b> . In corpus-based IE, it is assumed that the information in the corpus is correct and up-to-date.
<b>Multilingual and heterogeneous.</b> The Web is not restricted to a single language and the texts are produced by numerous authors for diverse audiences.	<b>Often monolingual and homogeneous.</b> If the author or nature (e.g. articles from the Wall Street Journal) of the corpus is known beforehand, it is easier to develop heuristics or to train named entity recognizers.
<b>No representative annotated corpora.</b> As no representative annotated texts are available, the Web as a corpus is currently less suited for supervised machine learning approaches.	Annotated test corpora available. In order to train supervised learning based named entity recognizers (NERs), test corpora are available where instances of a limited number of classes are marked within the text.
<b>Dynamic.</b> The contents of the Web changes continuously, results of experiments may thus also change over time.	<b>Static.</b> Experimental results are independent of time and place as the corpora are static.
<b>Facts and opinions.</b> As a multitude of users contribute to the Web, its contents are also suited for opinion mining.	<b>Facts only.</b> Information Extraction tasks on Newspaper corpora mainly focus on the identification of facts.

Table 1. Comparison between the Web as a corpus and 'traditional' corpora

Apart from factual information, the Web is the de-facto source to gather community-based data as people with numerous backgrounds, interests and ideas contribute to the content of the Web. Hence the Web is a valuable source to extract opinions, characterizations and perceived relatedness between items.

In this chapter, the focus is on gathering and structuring information from the 'traditional'Web. This structured information can be represented (and shared) using the standard Semantic Web (SW) languages. Hence, this chapter focuses on the automatic creation of content for the SW. For simplicity, we abstract from the SW standards RDF(S)/OWL.

# The Web-as-a-Corpus vs. Traditional Text Corpora

Information extraction (IE) is the task of identifying instances (or *named entities*) and relations between those instances in a collection of texts, called a text corpus.

In the nineties, the Message Understanding Conferences (MUC) focused on the recognition of named entities (such as names of persons and organizations) in a collection of texts (Chinchor, 1998). Initially, this work was mostly based on rules on the syntax and context of such named entities. For example, two capitalized words preceded by mr. will denote the name of a male person. As the creation of such rules is a laborious task, approaches became popular where named entities were recognized using machine learning (Mitchell, 1997), for example in (Zhou & Su, 2002; Brothwick, 1999; Finkel, Grenager, & Manning, 2005). However, such approaches typically make use of annotated training sets where instances (e.g. 'Microsoft') are labeled with their class ('Organization').

Traditional information extraction tasks focus on the identification of named entities in large text corpora such as collections newspaper articles or biomedical texts. In this chapter however, we focus on the Web as a corpus. In Table 1 the most important differences between the two can be found. 32 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage: www.igi-global.com/chapter/search-engine-based-web-information/37729

# **Related Content**

#### A Location-Aware Access Control Model for Mobile Workflow Systems

Michael Decker (2009). International Journal of Information Technology and Web Engineering (pp. 50-66). www.irma-international.org/article/location-aware-access-control-model/4030

#### **Explaining Semantic Web Applications**

Deborah L. McGuinness, Vasco Furtado, Paulo Pinheiro da Silva, Li Ding, Alyssa Glassand Cynthia Chang (2010). *Web Technologies: Concepts, Methodologies, Tools, and Applications (pp. 2304-2327).* www.irma-international.org/chapter/explaining-semantic-web-applications/37739

#### The Role of Physical Affordances in Multifunctional Mobile Device Design

Sorin Adam Matei, Anthony Faiola, David J. Wheatleyand Tim Altom (2010). *International Journal of Information Technology and Web Engineering (pp. 40-57).* www.irma-international.org/article/role-physical-affordances-multifunctional-mobile/49199

#### Applying Ontology Similarity Functions to Improve Software Agent Communication

Jairo Francisco de Souza, Sean W.M. Siqueiraand Rubens N. Melo (2012). *Models for Capitalizing on Web Engineering Advancements: Trends and Discoveries (pp. 43-57).* www.irma-international.org/chapter/applying-ontology-similarity-functions-improve/61899

#### From User's Goal to Semantic Web Services Discovery: Approach Based on Traceability

Houda el Bouhissi, Mimoun Malkiand Mohamed Amine Sidi Ali Cherif (2014). International Journal of Information Technology and Web Engineering (pp. 15-39).

www.irma-international.org/article/from-users-goal-to-semantic-web-services-discovery/123182