

Deactivation of Unwelcomed Deep Web Extraction Services through Random Injection

Varun Bhagwan, Tyrone Grandison

IBM Almaden Research Center, 650 Harry Road, San Jose, California 95120 USA
{vbhagwan, tyroneg}@us.ibm.com

Abstract

Websites serve content both through Web Services as well as through user-viewable webpages. While the consumers of web-services are typically 'machines', webpages are meant for human users. It is highly desirable (for reasons of security, revenue, ownership, availability etc.) for service providers that content that will undergo further processing be fetched in a prescribed fashion, preferably through a supplied Web Services. In fact, monetization of partnerships within a services ecosystem normally means that website data translate into valuable revenue. Unfortunately, it is quite commonplace for arbitrary developers to extract or leverage information from websites without asking for permission and or negotiating a revenue sharing agreement. This may translate to significant lost income for content providers. Even in cases where website owners are happy to share the data, they may want users to adopt dedicated Web Service APIs (and associated API-servers) rather than putting a load on their revenue-generating websites. In this paper, we introduce a mechanism that disables automated web scraping agents, thus forcing clients to conform to the provided Web Services.

1. Introduction

Many companies take the prudent approach when it comes to the development of systems that use online data from external parties, i.e. they undergo a formal documented legal process where consent is granted and a revenue package is agreed upon. However, these same companies tend to be frustrated when it comes to their own data being leveraged by other parties, who have less to lose and that do not take this prudent approach.

These companies normally employ web-scraping [1] to harvest their information. Formally defined, web-scraping [1] is the act of going through the content of a website for the purpose of extracting information from it. It is typically implemented by means of authoring an automated agent that makes appropriate HTTP requests to the website with the desired content, and 'scrapes' the said content from the result of the HTTP request (related issues have been dealt in detail in [2]). The scraping (or extraction or harvesting) is used to collect content such as user-data, image-links, user-comments, email addresses or any other data of potential value from the source website.

In the most malicious of cases, it can involve copying entire websites to direct traffic away from the source website and onto the (typically spam and or ad infested) malicious website. In this paper, we introduce a mechanism, Random Injection-based Deactivation (RID), to forcibly disallow automated web-scraping agents from harvesting or collecting data from a website.

2. Background

Figure 1 shows the HTML code that web scrapers would need to navigate in order to obtain image links for artists.



Figure 1. Image Harvesting.

Another example of the information that screen scrapers may be interested in is depicted in Figure 2, which shows the HTML hierarchy that the scrapers would navigate when extracting user data and comments.

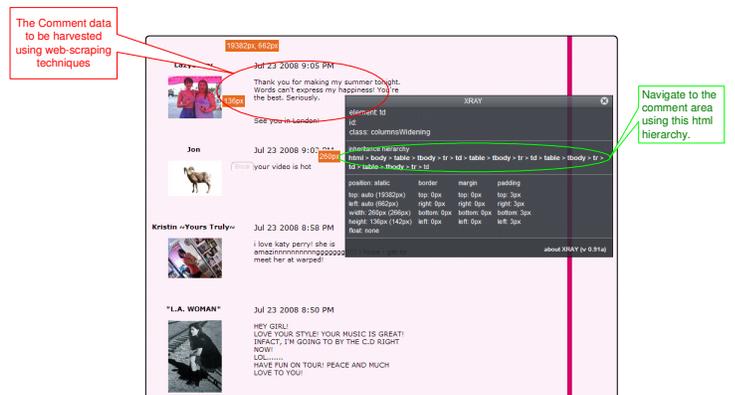


Figure 2. User Comment Harvesting

As stated previously, even if web-scraping is being done completely legitimately and for acceptable reasons, source websites may wish to divert traffic away from their main servers, and or to encourage such 'scrapers' to switch to using the provided Web Service APIs instead of scraping the (HTML) source code, whether for technical or

