**Computers & Security**

# Stable web spam detection using features based on lexical items

CrossMark

Marcin Luckner [a,*], Michał Gad [b], Paweł Sobkowiak [c]

[a] Faculty of Mathematics and Information Science, Warsaw University of Technology, Koszykowa 75, 00-662 Warszawa, Poland
[b] EO Networks, ul. Głuszycka 5, Warszawa, Poland
[c] Sensi Soft, ul. Głuszycka 5, Warszawa, Poland

## ARTICLE INFO

## ABSTRACT

Web spam is a method of manipulating search engines results by improving ranks of spam pages. It takes various forms and lacks a consistent definition. Web spam detectors use machine learning techniques to detect spam. However, the detectors are mostly verified on data sets coming from the same year as the learning sets. In this paper we compared Support Vector Machine classifiers trained and tested on WEBSPAM—UK data sets from different years. To obtain stable results we proposed new lexical-based features. The HTML document — transformed into a text without HTML tags, a set of visible symbols, and a list of links including the ones from tags — gave information about weird combinations of letters; consonant clusters; statistics on syllables, words, and sentences; and the Gunning Fog Index. Using data collected in 2006 as a learning set, we obtained very stable accuracy among years. This choice of the training set reduced the sensitivity in 2007, but that can be improved by managing the acceptance threshold. Finally, we proved that the balance between the sensitivity and the specificity measured by the Area Under the Curve (AUC) is improved by our selection of features.

© 2014 Elsevier Ltd. All rights reserved.

## 1. Introduction

The detection of spam emails was — and still is — a serious problem for the Internet community (Carpinter and Hunt, 2006). The new spam filtering approaches have raised the hope of reducing this problem (Filasiak et al., 2014), but the detection of new types of spam — SMS spam (Xu et al., 2012), MMS spam (Yoon et al., 2010), video spam, and web spam (Potdar et al., 2010) — is still one of the most challenging issues. The most common type in the list is web spam that exploits vulnerabilities and gaps in the web 2.0 to inject links to spam content into dynamic and sharable content such as blogs, comments, reviews, or wikipages.

Web spam takes various forms and lacks a consistent definition. Therefore, web spam detectors use machine learning techniques to create a model of spam from training sets that include spam and non-spam examples. WEBSPAM—UK collections from 2006 to 2007 are very popular training sets (Castillo et al., 2006). Many projects have used these data sets to test web spam detectors. However, most detectors were verified and proved effective only on data coming from a single year (c.f. Mahmoudi et al., 2010; Algur and Pendari, 2012). One may conclude that the presented solutions are temporary.

---

* Corresponding author.
E-mail addresses: mluckner@mini.pw.edu.pl (M. Luckner), michal.gad@eo.pl (M. Gad), psobkowiak@acm.org (P. Sobkowiak).
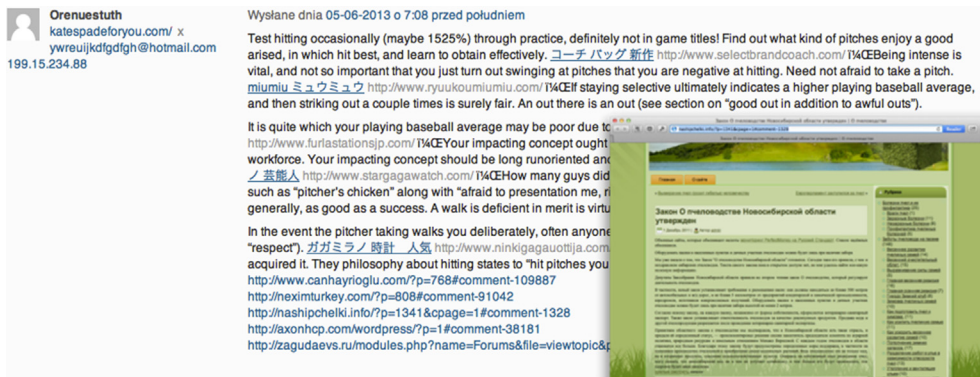
Fig. 1 shows several web spam injections into web forum comments. Fig. 1(a) presents a typical example of web spam comment with several injected links. This spam was created to promote link farms and provide credibility to the spammer website. The second example of text spam on a web blogging platform Fig. 1(c) refers to the image. Any human can recognise the linked page as spam, but an automaton needs special techniques of image spam detection (Gao et al., 2009; Wakade et al., 2013). The link contained in the last example of pingback spam on Wordpress platform Fig. 1(b) directs users to the Youtube spam movie. Because of computation costs, video spam is very hard to detect and approaches to detect video spam are still being developed (Luz et al., 2012).
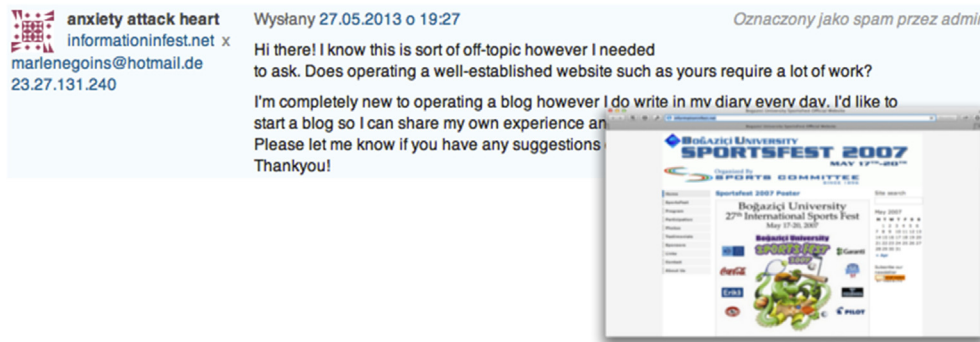
The aim of this study was to create a web spam detector that works over years. We selected several features based on lexical items and distinguished web spam from non-spam context. We believe that such features could be common for the WEBSPAM–UK2006 and WEBSPAM–UK2007 data sets. To test this hypothesis we created three web spam detectors: trained and tested using 2006 data; trained and tested using 2007 data; trained using 2006 data, but tested using 2007 data. The results of the first and the second detector formed points of reference for the third one. We expected that the accuracy obtained by the last detector would be similar to the accuracy obtained using 2006 data. The results obtained using 2007 data allowed us to estimate the reduction of the accuracy obtained by the third detector.
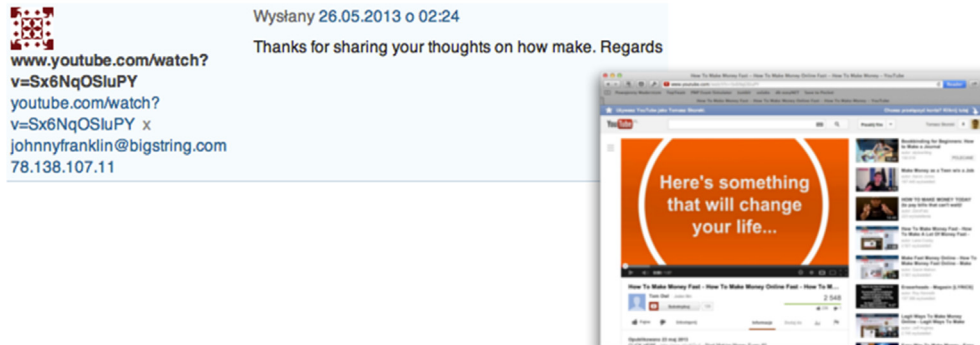
To estimate the influence of the selected features, we repeated the test using only commonly used features. Our hypothesis was that the obtained results would be worse than the ones achieved in the first test.



(a)

(b)

(c)

**Fig. 1 – The examples of web spam comments with linked content. The links under signatures lead to the web sites presented in the bottom right corners.**