# Adversarial Web Search

# Adversarial Web Search

---

## Carlos Castillo

*Yahoo! Research*
*Barcelona 08018*
*Catalunya-Spain*
*chato@yahoo-inc.com*

## Brian D. Davison

*Lehigh University*
*Bethlehem, PA 18015*
*USA*
*davison@cse.lehigh.edu*

**now**

the essence of knowledge

Boston – Delft

# Foundations and Trends® in Information Retrieval

# Foundations and Trends® in Information Retrieval
Volume 4 Issue 5, 2010
## Editorial Board

# Editorial Scope

**Foundations and Trends® in Information Retrieval** will publish survey and tutorial articles in the following topics:

- Applications of IR
- Architectures for IR
- Collaborative filtering and recommender systems
- Cross-lingual and multilingual IR
- Distributed IR and federated search
- Evaluation issues and test collections for IR
- Formal models and language models for IR
- IR on mobile platforms
- Indexing and retrieval of structured documents
- Information categorization and clustering
- Information extraction
- Information filtering and routing

- Metasearch, rank aggregation and data fusion
- Natural language processing for IR
- Performance issues for IR systems, including algorithms, data structures, optimization techniques, and scalability
- Question answering
- Summarization of single documents, multiple documents, and corpora
- Text mining
- Topic detection and tracking
- Usability, interactivity, and visualization issues in IR
- User modelling and user studies for IR
- Web search

now
the essence of knowledge

# Adversarial Web Search

# Carlos Castillo[1] and Brian D. Davison[2]

[1] Yahoo! Research, Diagonal 177, 8th Floor, Barcelona 08018,
  Catalunya-Spain, chato@yahoo-inc.com
[2] Lehigh University, 19 Memorial Drive West, Bethlehem, PA 18015, USA,
  davison@cse.lehigh.edu

## Abstract

Web search engines have become indispensable tools for finding content.
As the popularity of the Web has increased, the efforts to exploit the
Web for commercial, social, or political advantage have grown, making
it harder for search engines to discriminate between truthful signals of
content quality and deceptive attempts to game search engines' rank-
ings. This problem is further complicated by the open nature of the
Web, which allows anyone to write and publish anything, and by the
fact that search engines must analyze ever-growing numbers of Web
pages. Moreover, increasing expectations of users, who over time rely
on Web search for information needs related to more aspects of their
lives, further deepen the need for search engines to develop effective
counter-measures against deception.

In this monograph, we consider the effects of the adversarial rela-
tionship between search systems and those who wish to manipulate
them, a field known as "Adversarial Information Retrieval". We show
that search engine spammers create false content and misleading links
to lure unsuspecting visitors to pages filled with advertisements or mal-
ware. We also examine work over the past decade or so that aims to

discover such spamming activities to get spam pages removed or their effect on the quality of the results reduced.

Research in Adversarial Information Retrieval has been evolving over time, and currently continues both in traditional areas (e.g., link spam) and newer areas, such as click fraud and spam in social media, demonstrating that this conflict is far from over.

# Contents

# 1

## Introduction

Information Retrieval (IR) is a branch of computer science that deals with tasks such as gathering, indexing, filtering, retrieving, and ranking content from a large collection of information-bearing items. It is a field of study that is over 40 years old, and started with the goal of helping users locate information items in carefully curated collections, such as the ones available in libraries. In the mid-1990s, the emergence of the World Wide Web created new research opportunities and challenges for information retrieval. The Web as a whole is larger, less coherent, more distributed and more rapidly changing than the previous document collections in which IR methods were developed [9].

From the perspective of an information retrieval system such as a search engine, the Web is a mixture of two types of content: the "closed Web" and the "open Web" [37]. The closed Web comprises a small number of reputable, high-quality, carefully maintained collections which a search engine can fully trust. The "open Web", on the other hand, includes the vast majority of Web pages, and in which document quality cannot be taken for granted. The openness of the Web has been the key to its rapid growth and success, but the same openness is the most challenging aspect when designing effective Web-scale information retrieval systems.

*Adversarial Information Retrieval* addresses the same tasks as Information Retrieval: gathering, indexing, filtering, retrieving, and ranking information, with the difference that it performs these tasks in collections wherein a subset has been manipulated maliciously [73]. On the Web, the predominant form of such manipulation is "search engine spamming" (also known as *spamdexing* or *Web spam*). Search engine spamming is the malicious attempt to influence the outcome of ranking algorithms, usually aimed at getting an undeservedly high ranking for one or more Web pages [92].

Among the specific topics related to Adversarial Information Retrieval on the Web, we find the following. First, there are several forms of general Web spam including link spam, content spam, cloaking, etc. Second, there are specialized forms of Web spam for particular subsets of the Web, including for instance blog spam (*splogs*), opinion spam, comment spam, referrer spam, etc. Third, there are ways in which a content publisher may attempt to deceive a Web advertiser or advertiser broker/intermediary, including search spam and click spam. Fourth, there are other areas in which the interests of the designers of different Web systems collide, such as in the reverse engineering of ranking methods, the design of content filters for ads or for Web pages, or the development of undetectable automatic crawlers, to name a few.

## 1.1   Search Engine Spam

The Adversarial IR topic that has received the most attention has been search engine spam, described by Fetterly et al. as "Web pages that hold no actual informational value, but are created to lure Web searchers to sites that they would otherwise not visit" [74].

Search engines have become indispensable tools for most users [17]. Web spammers try to deceive search engines into showing a lower-quality result with a high ranking. They exploit, and as a result, weaken, the trust relationship between users and search engines [92], and may damage the search engines' reputation. They also make the search engine incur extra costs when dealing with documents that have little or no relevance for its users; these include network costs for downloading them, disk costs for storing them, and processing costs for

indexing them. Thus, the costs of Web spam are felt both by end-users and those providing a service to them.

Ntoulas et al. [182] measured Web spam across top-level domains (TLDs) by randomly sampling pages from each TLD in a large-scale Web search engine, and then labeling those pages manually. In their samples, 70% of the pages in the `.biz` domain, 35% of the pages in `.us` and 20% of the pages in `.com` were spam. These are uniform random samples, while the top results in search engines are much more likely to be spam as they are the first target of spammers. In a separate study, Eiron et al. [69] ranked 100 million pages using PageRank and found that 11 out of the top 20 achieved such high ranking through link manipulation.

Ignoring Web spam is not an option for search engines. According to Henzinger et al. [98], "Spamming has become so prevalent that every commercial search engine has had to take measures to identify and remove spam. Without such measures, the quality of the rankings suffers severely." In other words, on the "open Web", a naïve application of ranking methods is no longer an option.

## 1.2  Activists, Marketers, Optimizers, and Spammers

The existence of Web spam pages can be seen as a natural consequence of the dominant role of search engines as mediators in information seeking processes [85]. User studies show that search engine users only scan and click the top few results for any given search [87], which means that Web page exposure and visitor traffic are directly correlated with search engine placement. Those who seek visibility need to have pages in the top positions in search engine results pages, and thus have an incentive to try to distort the ranking method.

There are many reasons for seeking visibility on the Web. Some people (activists) spam search engines to further a political message or to help a non-profit achieve its end. This is the case of most link bombs (perhaps better known as *Google bombs*) that spam a particular term or phrase to link it to a particular Web page. A memorable example of this manipulation is the one that affected the query "miserable failure", which during the 2004 presidential election, returned the home page of

George W. Bush as the first result in several Web search engines. This was the result of a coordinated effort by bloggers and Web page authors around the world. We discuss link bombing further in Section 4.

Most search engine spam, however, is created for financial gain. There is a strong economic incentive to find ways to drive traffic to Web sites, as more traffic often translates to more revenue [231]. Singhal [212] estimated the amount of money that typical spammers expected to receive in 2005: a few US dollars per sale for affiliate programs on Amazon or E-Bay, around 6 USD per sale of Viagra, and around 20–40 USD per new member of pornographic sites. Given the small per-sale commissions and the low response rates, a spammer needs to collect millions of page views to remain profitable. Further, some spam pages exist to promote or even install malware [68, 192, 193].

The incentive to drive traffic to Web sites, both for legitimate and illegitimate purposes, has created a whole industry around search engines. The objective of *Search Engine Marketing* (SEM) is to assist marketers in making their Web content visible to users via a search engine.[1] SEM activities are divided by the two principal kinds of information displayed on a search results page: the editorial content and the advertising (or "sponsored search").

Advertising on search engines today is also a ranking process, involving bidding for keywords to match to user queries, the design of the ads themselves, and the design of the "landing pages" to which users are taken after clicking on the ads. An advertiser's goal in sponsored search is to attract more *paid* traffic that "converts" (i.e., buys a product or service, or performs some other action desired by the advertiser), within a given advertising budget.

Sponsored search efforts are fairly self-regulated. First, marketers have to pay the search engine for each click on the ads. Second, the marketer does not simply want to attract traffic to his Web site, but to attract traffic that leads to conversions. Thus, it is in his best interest to bid for keywords that represent the actual contents of his Web site.

---

[1] Some practitioners define SEM more narrowly, focusing on the sponsored search side, but from a business perspective, all of these efforts fall under marketing.

Also ad market designers are careful to design systems that provide incentives for advertisers to bid truthfully.

The objective of *Search Engine Optimization* (SEO), on the other hand, is to make the pages of a certain Web site rank higher in the editorial side of search engines, in order to attract more *unpaid* or *organic* traffic to the target Web site.

The efforts of a search engine optimizer, in contrast, are not self-regulating, and in some cases can significantly disrupt search engines, if counter-measures are not taken. For this reason, search engines threaten SEOs that have become spammers with penalties, which may include the demotion or removal from the index of pages that use deceptive practices. The penalties that search engines apply are well known by the SEO community. Boundaries are, of course, fuzzy, as all search engines seem to allow some degree of search engine optimization.

Moran and Hunt [169] advise Web site owners on how to tell search engine spammers from SEOs. A search engine spammer tends to (i) offer a guarantee of top rankings, which no reputable firm can do as there are many variables outside their control; (ii) propose minimal changes to the pages, which indicate that they are likely to create a *link farm* (described in Section 4.3) instead of actually modifying the way the content is presented to users and search engines; and (iii) suggest to use server-level *cloaking* (described in Section 3.5) or other modifications whose typical purpose is to spam.

## 1.3 The Battleground for Search Engine Rankings

In general, search engine results are ranked using a combination of two factors: the *relevance* of the pages to the query, and the *authoritativeness* of the pages themselves, irrespective of the query. These two aspects are sometimes named respectively *dynamic ranking* and *static ranking*, and both have been the subject of extensive studies from the IR community (and discussed in IR textbooks [13, 58, 154]).

Some search engine spammers may be assumed to be knowledgeable about Web information retrieval methods used for ranking pages. Nevertheless, when spammers try to manipulate the rankings of a search engine, they do not know the details about the ranking methods used

by the search engine; for instance they do not know which are the specific features used for computing the ranking. Under those conditions, their best strategy is simply to try to game *any* signal believed to be used for ranking.

In the early days of the Web, search engine spammers manipulated mainly the contents and URLs of the pages, automatically generating millions of pages, including incorporating repetitions or variants of certain keywords in which the spammer was interested. Next, as search engines began to use link-based signals [33, 34, 122, 183], spammers started to create pages interlinked deceptively to generate misleading link-based ranking signals.

As the search engines adapted to the presence of Web spam by using more sophisticated methods, including the usage of machine-learning-based ranking for Web pages [201], more elements of the pages were taken into consideration which pushed spammers to become more sophisticated. Next, the possibility of adding comments to forums and the existence of other world-writable pages such as *wikis* presented new opportunities for spammers as they allowed the insertion of arbitrary links into legitimate pages.

Recently search engines have devised other ways of exploiting the "wisdom of crowds", e.g., through usage data to rank pages, but search engine spammers can also pose as members of the crowds and disrupt rankings as long as they are not detected. Web spam has been evolving over the years, and will continue to evolve to reflect changes in ranking methods used by popular services.

Thus, there are a variety of useful signals for ranking and each of them represents an opportunity for spammers, and in Sections 3–7 we will highlight how spammers have taken advantage of these opportunities to manipulate valuable ranking signals and what work has been done to detect such manipulation.

## 1.4 Previous Surveys and Taxonomies

In 2001, Perkins [189] published one of the earliest taxonomies of Web spam. This taxonomy included content spam, link spam, and cloaking. It also suggested a test for telling spam from non-spam: Spam is "any

attempt to deceive a search engine's relevancy algorithm", non-spam is "anything that would still be done if search engines did not exist, or anything that a search engine has given written permission to do."

In 2005, Gyöngyi and Garcia-Molina [93] proposed a different taxonomy. This taxonomy stressed the difference between boosting techniques and hiding techniques. Boosting techniques are directly aimed at promoting a page or a set of pages by manipulating their contents or links. Hiding techniques, instead, are used by spammers to "cover their tracks", thus preventing the discovery of their boosting techniques.

In 2007, a brief overview of Adversarial IR by Fetterly [73] appeared in *ACM Computing Reviews*. It included a general description of the field, and references to key articles, data sources, and books related to the subject. In the same year Heymann et al. [101] published a survey focused on social media sites, stating that in the case of social media sites, a preventive approach was possible, in addition to detection- and demotion-based approaches. Prevention is possible because in social media sites there is more control over what users can do; for example, CAPTCHAs can be incorporated to prevent automated actions, the rate at which users post content can be limited, and disruptive users can be detected and banned.

Additionally, several Ph.D. and M.Sc. theses have included elements related to Web spam. A partial list of them includes theses in the areas of link spam [95, 149, 160, 208], splogs and spam in blogs [124, 166], content spam [180], Web spam systems in general [45, 232, 236, 251], and search engine optimization [123].

We have left out the closely related subject of e-mail spam. While some methods overlap, particularly in the case of content-based Web-spam detection (which we discuss in Section 3.6), there are substantial differences between the two areas. For a survey on e-mail spam, see, e.g., Cormack [55].

## 1.5  This Survey

In this survey we have tried to be relatively inclusive; this is reflected in citations to about 250 publications, which we consider large for a survey on a young sub-field of study. We also intended to appeal to a

wide audience including developers and practitioners. For this reason, we have chosen to present general descriptions of Web spam techniques and counter-measures, and to be selective with the details.

The rest of this monograph is organized as follows:

> **Section 2** describes general systems for detecting search engine spam, including the choice of a machine learning method, the feature design, the creation of a training set, and evaluation methodologies.
>
> **Section 3** describes content-based spam techniques and how to detect them, as well as malicious mirroring, which is a form of plagiarism for spam purposes.
>
> **Section 4** describes link-based spam techniques and how to detect them, and covers topics such as link alliances and nepotistic linking.
>
> **Section 5** describes methods for propagating trust and distrust on the Web, which can be used for demoting spam pages.
>
> **Section 6** describes click fraud and other ways of distorting Web usage data, including Web search logs; it also deals with the subject of using search logs as part of Web spam detection systems.
>
> **Section 7** describes ways of spamming social media sites and user-generated content in general.

Finally, the discussion in **Section 8** includes future research directions and links to research resources.

# References

[1] B. A, K. Csalogány, and T. Sarlós, "Link-based similarity search to fight Web spam," in *Proceedings of the Second International Workshop on Adversarial Information Retrieval on the Web (AIRWeb)*, 2006.

[2] J. Abernethy and O. Chapelle, "Semi-supervised classification with hyperlinks," in *Proceedings of the ECML/PKDD Graph Labeling Workshop*, September 2007.

[3] J. Abernethy, O. Chapelle, and C. Castillo, "Webspam identification through content and hyperlinks," in *Proceedings of the Fourth International Workshop on Adversarial Information Retrieval on the Web (AIRWEB)*, pp. 41–44, ICPS: ACM Press, April 2008.

[4] J. Abernethy, O. Chapelle, and C. Castillo, "Graph regularization methods for web spam detection," *Machine Learning Journal*, vol. 81, no. 2, pp. 207–225, 2010.

[5] S. Adali, T. Liu, and M. Magdon-Ismail, "Optimal link bombs are uncoordinated," in *Proceedings of the First International Workshop on Adversarial Information Retrieval on the Web (AIRWeb)*, May 2005.

[6] B. T. Adler and L. de Alfaro, "A content-driven reputation system for the Wikipedia," in *Proceedings of the 16th International Conference on World Wide Web (WWW)*, pp. 261–270, New York, NY, USA: ACM, 2007.

[7] E. Amitay, S. Yogev, and E. Yom-Tov, "Serial sharers: Detecting split identities of Web authors," in *Workshop on Plagiarism Analysis, Authorship Identification, And Near-Duplicate Detection*, July 2007.

[8] R. Andersen, C. Borgs, J. Chayes, J. Hopcraft, V. Mirrokni, and S.-H. Teng, "Local computation of PageRank contributions," in *Algorithms and Models*

*for the Web-Graph*, vol. 4863 *of Lecture Notes in Computer Science*, pp. 150–165, Springer, 2007.

[9] A. Arasu, J. Cho, H. Garcia-Molina, A. Paepcke, and S. Raghavan, "Searching the Web," *ACM Transactions on the Internet Technology (TOIT) 1*, vol. 1, pp. 2–43, August 2001.

[10] J. Attenberg and T. Suel, "Cleaning search results using term distance features," in *Proceedings of the Fourth International Workshop on Adversarial Information Retrieval on the Web (AIRWeb)*, pp. 21–24, New York, NY, USA: ACM, 2008.

[11] V. Bacarella, F. Giannotti, M. Nanni, and D. Pedreschi, "Discovery of ads Web hosts through traffic data analysis," in *Proceedings of the 9th ACM SIG-MOD Workshop on Research Issues in Data Mining and Knowledge Discovery (DMKD)*, pp. 76–81, New York, NY, USA: ACM, 2004.

[12] R. Baeza-Yates, C. Castillo, and V. López, "PageRank increase under different collusion topologies," in *First International Workshop on Adversarial Information Retrieval on the Web (AIRWeb)*, pp. 17–24, May 2005.

[13] R. Baeza-Yates and B. Ribeiro-Neto, "Modern Information Retrieval," Addison Wesley, May 1999.

[14] J. Bar-Ilan, "Web links and search engine ranking: The case of Google and the query "jew"," *Journal of the American Society for Information Science and Technology*, vol. 57, no. 12, pp. 1581–1589, 2006.

[15] J. Bar-Ilan, "Google bombing from a time perspective," *Journal of Computer-Mediated Communication*, vol. 12, no. 3, 2007.

[16] Z. Bar-Yossef, I. Keidar, and U. Schonfeld, "Do not crawl in the DUST: Different URLs with similar text," *ACM Transactions on the Web*, vol. 3, no. 1, pp. 1–31, 2009.

[17] J. Battelle, *The Search: How Google and Its Rivals Rewrote the Rules of Business and Transformed Our Culture*. New York: Portfolio, 2005.

[18] L. Becchetti, C. Castillo, D. Donato, R. Baeza-Yates, and S. Leonardi, "Link analysis for Web spam detection," *ACM Transactions on the Web*, vol. 2, no. 1, pp. 1–42, February 2008.

[19] L. Becchetti, C. Castillo, D. Donato, S. Leonardi, and R. Baeza-Yates, "Using rank propagation and probabilistic counting for link-based spam detection," in *Proceedings of the Workshop on Web Mining and Web Usage Analysis (WebKDD)*, ACM Press, August 2006.

[20] A. A. Benczúr, I. Bíró, K. Csalogány, and M. Uher, "Detecting nepotistic links by language model disagreement," in *Proceedings of the 15th International Conference on World Wide Web (WWW)*, pp. 939–940, ACM Press, 2006.

[21] A. A. Benczúr, K. Csalogány, T. Sarlós, and M. Uher, "SpamRank: Fully automatic link spam detection," in *Proceedings of the First International Workshop on Adversarial Information Retrieval on the Web (AIRWeb)*, May 2005.

[22] F. Benevenuto, T. Rodrigues, V. Almeida, J. Almeida, C. Zhang, and K. Ross, "Identifying video spammers in online social networks," in *Proceedings of the Fourth International Workshop on Adversarial Information Retrieval on the Web (AIRWeb)*, pp. 45–52, New York, NY, USA: ACM, 2008.

[23] P. Berkhin, "A survey on PageRank computing," *Internet Mathematics*, vol. 2, no. 2, pp. 73–120, 2005.

[24] K. Berlt, E. S. de Moura, C. M. André, N. Ziviani, and T. Couto, "A hypergraph model for computing page reputation on Web collections," in *Proceedings of the Simpósio Brasileiro de Banco de Dados (SBBD)*, pp. 35–49, October 2007.

[25] K. Bharat and M. R. Henzinger, "Improved algorithms for topic distillation in hyperlinked environments," in *Proceedings of the 21st International ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 104–111, August 1998.

[26] J. Bian, Y. Liu, E. Agichtein, and H. Zha, "A few bad votes too many?: Towards robust ranking in social media," in *Proceedings of the Fourth International Workshop on Adversarial Information Retrieval on the Web (AIRWeb)*, pp. 53–60, New York, NY, USA: ACM, 2008.

[27] A. Bifet, C. Castillo, P.-A. Chirita, and I. Weber, "An analysis of factors used in search engine ranking," in *Proceedings of the First International Workshop on Adversarial Information Retrieval (AIRWeb)*, May 2005.

[28] I. Bíró, D. Siklósi, J. Szabó, and A. Benczúr, "Linked latent dirichlet allocation in Web spam filtering," in *Proceedings of the 5th International Workshop on Adversarial Information Retrieval on the Web (AIRWeb)*, pp. 37–40, ACM Press, 2009.

[29] I. Bíró, J. Szabó, and A. A. Benczúr, "Latent dirichlet allocation in Web spam filtering," in *Proceedings of the 4th International Workshop on Adversarial Information Retrieval on the Web (AIRWeb)*, pp. 29–32, New York, NY, USA: ACM, 2008.

[30] D. M. Blei, A. Y. Ng, and M. I. Jordan, "Latent dirichlet allocation," *Journal of Machine Learning Research*, vol. 3, pp. 993–1022, 2003.

[31] A. Borodin, G. O. Roberts, J. S. Rosenthal, and P. Tsaparas, "Finding authorities and hubs from link structures on the World Wide Web," in *Proceedings of the 10th International Conference on World Wide Web (WWW)*, pp. 415–429, 2001.

[32] O. P. Boykin and V. Roychowdhury, "Personal Email networks: an effective anti-spam tool," Condensed Matter cond-mat/0402143, 2004.

[33] S. Brin, R. Motwani, L. Page, and T. Winograd, "What can you do with a Web in your pocket?," *Data Engineering Bulletin*, vol. 21, no. 2, pp. 37–47, 1998.

[34] S. Brin and L. Page, "The anatomy of a large-scale hypertextual Web search engine," in *Proceedings of the 7th International Conference on the World Wide Web*, pp. 107–117, April 1998.

[35] A. Brod and R. Shivakumar, "Advantageous semi-collusion," *The Journal of Industrial Economics*, vol. 47, no. 2, pp. 221–230, 1999.

[36] A. Z. Broder, S. C. Glassman, M. S. Manasse, and G. Zweig, "Syntactic clustering of the Web," *Computer Networks and ISDN Systems*, vol. 29, no. 8–13, pp. 1157–1166, September 1997.

[37] T. A. Brooks, "Web search: How the Web has changed information retrieval," *Information Research*, vol. 8, no. 3, April 2003.

[38] G. Buehrer, J. W. Stokes, and K. Chellapilla, "A large-scale study of auto-
     mated Web search traffic," in *Proceedings of the 4th International Workshop
     on Adversarial Information Retrieval on the Web (AIRWeb)*, pp. 1–8, New
     York, NY, USA: ACM, 2008.

[39] S. Büttcher, C. L. A. Clarke, and B. Lushman, "Term proximity scoring for
     ad-hoc retrieval on very large text collections," in *Proceedings of the 29th
     ACM Annual SIGIR Conference on Research and Development in Information
     Retrieval*, pp. 621–622, New York, NY, USA: ACM Press, 2006.

[40] C. Castillo, "Effective Web Crawling," PhD thesis, University of Chile, 2004.

[41] C. Castillo, C. Corsi, D. Donato, P. Ferragina, and A. Gionis, "Query log
     mining for detecting polysemy and spam," in *Proceedings of the KDD Work-
     shop on Web Mining and Web Usage Analysis (WEBKDD)*, Springer: LNCS,
     August 2008.

[42] C. Castillo, C. Corsi, D. Donato, P. Ferragina, and A. Gionis, "Query-log
     mining for detecting spam," in *Proceedings of the 4th International Workshop
     on Adversarial Information Retrieval on the Web (AIRWeb)*, ICPS: ACM
     Press, April 2008.

[43] C. Castillo, D. Donato, L. Becchetti, P. Boldi, S. Leonardi, M. Santini, and
     S. Vigna, "A reference collection for Web spam," *SIGIR Forum*, vol. 40, no. 2,
     pp. 11–24, December 2006.

[44] C. Castillo, D. Donato, A. Gionis, V. Murdock, and F. Silvestri, "Know your
     neighbors: Web spam detection using the web topology," in *Proceedings of the
     30th Annual International ACM SIGIR Conference on Research and Devel-
     opment in Information Retrieval*, ACM, July 2007.

[45] J. Caverlee, "Tamper-Resilient Methods for Web-Based Open Systems," PhD
     thesis, College of Computing, Georgia Institute of Technology, August 2007.

[46] J. Caverlee and L. Liu, "Countering Web spam with credibility-based link
     analysis," in *Proceedings of the Twenty-Sixth Annual ACM Symposium on
     Principles of Distributed Computing (PODC)*, pp. 157–166, New York, NY,
     USA: ACM, 2007.

[47] J. Caverlee, L. Liu, and S. Webb, "Socialtrust: Tamper-resilient trust estab-
     lishment in online communities," in *Proceedings of the 8th ACM/IEEE-CS
     Joint Conference on Digital Libraries (JCDL)*, pp. 104–114, 2008.

[48] J. Caverlee, L. Liu, and S. Webb, "Towards robust trust establishment in
     Web-based social networks with SocialTrust," in *Proceedings of the 17th Inter-
     national World Wide Web Conference (WWW)*, pp. 1163–1164, ACM, 2008.

[49] J. Caverlee, S. Webb, and L. Liu, "Spam-resilient Web rankings via influence
     throttling," in *Proceedings of the IEEE International Parallel and Distributed
     Processing Symposium (IPDPS)*, pp. 1–10, 2007.

[50] K. Chellapilla and D. M. Chickering, "Improving cloaking detection using
     search query popularity and monetizability," in *Proceedings of the 2nd Interna-
     tional Workshop on Adversarial Information Retrieval on the Web (AIRWeb)*,
     pp. 17–24, August 2006.

[51] K. Chellapilla and A. Maykov, "A taxonomy of JavaScript redirection spam,"
     in *Proceedings of the 3rd International Workshop on Adversarial Information
     Retrieval on the Web (AIRWeb)*, pp. 81–88, New York, NY, USA: ACM Press,
     2007.

[52]  A. Cheng and E. Friedman, "Manipulability of PageRank under sybil strate-gies," in *Proceedings of the First Workshop on the Economics of Networked Systems (NetEcon06)*, 2006.

[53]  Y.-J. Chung, M. Toyoda, and M. Kitsuregawa, "A study of link farm distri-bution and evolution using a time series of Web snapshots," in *Proceedings of the 5th International Workshop on Adversarial Information Retrieval on the Web (AIRWeb)*, pp. 9–16, ACM Press, 2009.

[54]  A. Clausen, "The cost of attack of PageRank," in *Proceedings of the Inter-national Conference on Agents, Web Technologies and Internet Commerce (IAWTIC)*, July 2004.

[55]  G. V. Cormack, "Email spam filtering: A systematic review," *Foundations and Trends in Information Retrieval 1*, pp. 335–455, 2008.

[56]  G. V. Cormack, M. Smucker, and C. L. Clarke, "Efficient and effective spam filtering and re-ranking for large web datasets," Unpublished draft, available from http://durum0.uwaterloo.ca/clueweb09spam/spamhunt.pdf, retrieved 10, April 2010.

[57]  N. Craswell, O. Zoeter, M. Taylor, and B. Ramsey, "An experimental compar-ison of click position-bias models," in *Proceedings of the First International Conference on Web Search and Data Mining (WSDM)*, pp. 87–94, New York, NY, USA: ACM, 2008.

[58]  B. Croft, D. Metzler, and T. Strohman, *Search Engines: Information Retrieval in Practice*. Addison Wesley, 2009.

[59]  A. L. C. da Costa-Carvalho, P.-A. Chirita, E. S. de Moura, P. Calado, and W. Nejdl, "Site level noise removal for search engines," in *Proceedings of the 15th International Conference on World Wide Web (WWW)*, pp. 73–82, New York, NY, USA: ACM Press, 2006.

[60]  N. Dai, B. Davison, and X. Qi, "Looking into the past to better classify Web spam," in *Proceedings of the 5th International Workshop on Adversarial Infor-mation Retrieval on the Web (AIRWeb)*, pp. 1–8, ACM Press, 2009.

[61]  N. Daswani and M. Stoppelman, "The anatomy of Clickbot.A," in *Proceedings of the USENIX HOTBOTS Workshop*, April 2007.

[62]  B. D. Davison, "Recognizing nepotistic links on the Web," in *Artificial Intel-ligence for Web Search*, pp. 23–28, AAAI Press, July 2000.

[63]  B. D. Davison, M. Najork, and T. Converse, "Adversarial information retrieval on the Web (AIRWeb 2006)," *SIGIR Forum*, vol. 40, no. 2, pp. 27–30, 2006.

[64]  J. Douceur, "The sybil attack," in *Proceedings of the First International Peer To Peer Systems Workshop (IPTPS)*, pp. 251–260, Springer: Vol. 2429 of *Lecture Notes in Computer Science*, January 2002.

[65]  I. Drost and T. Scheffer, "Thwarting the nigritude ultramarine: Learning to identify link spam," in *Proceedings of the 16th European Conference on Machine Learning (ECML)*, pp. 233–243, Vol. 3720 of *Lecture Notes in Arti-ficial Intelligence*, 2005.

[66]  Y. Du, Y. Shi, and X. Zhao, "Using spam farm to boost PageRank," in *Proceedings of the 3rd International Workshop on Adversarial Information Retrieval on the Web (AIRWeb)*, pp. 29–36, New York, NY, USA: ACM, 2007.

[67] O. Duskin and D. G. Feitelson, "Distinguishing humans from robots in Web search logs: Preliminary results using query rates and intervals," in *Proceedings of the WSDM Workshop on Web Search Click Data (WSCD)*, pp. 15–19, New York, NY, USA: ACM, 2009.

[68] M. Egele, C. Kruegel, and E. Kirda, "Removing web spam links from search engine results," *Journal in Computer Virology*, In press. Published online 22 August, 2009.

[69] N. Eiron, K. S. Curley, and J. A. Tomlin, "Ranking the Web frontier," in *Proceedings of the 13th International Conference on World Wide Web*, pp. 309–318, New York, NY, USA: ACM Press, 2004.

[70] E. Enge, "Matt cutts interviewed by eric enge," Article online at http://www.stonetemple.com/articles/interview-matt-cutts-012510.shtml and retrieved on 11 April 2010, March 2010.

[71] M. Erdélyi, A. A. Benczúr, J. Masanes, and D. Siklósi, "Web spam filtering in internet archives," in *Proceedings of the 5th International Workshop on Adversarial Information Retrieval on the Web (AIRWeb)*, pp. 17–20, ACM Press, 2009.

[72] J. Feigenbaum, S. Kannan, M. A. McGregor, S. Suri, and J. Zhang, "On graph problems in a semi-streaming model," in *Proceedings of the 31st International Colloquium on Automata, Languages and Programming (ICALP)*, pp. 531–543, Springer: Vol. 3142 of *LNCS*, 2004.

[73] D. Fetterly, "Adversarial Information Retrieval: the manipulation of Web content," *ACM Computing Reviews*, July 2007.

[74] D. Fetterly, M. Manasse, and M. Najork, "Spam, damn spam, and statistics: Using statistical analysis to locate spam Web pages," in *Proceedings of the Seventh Workshop on the Web and databases (WebDB)*, pp. 1–6, June 2004.

[75] D. Fetterly, M. Manasse, and M. Najork, "Detecting phrase-level duplication on the World Wide Web," in *Proceedings of the 28th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 170–177, New York, NY, USA: ACM, 2005.

[76] R. Fishkin, "Lessons learned building an index of the WWW," Retrieved 15 June 2009 from http://www.seomoz.org/blog/lessons-learned-building-an-index-of-the-www, April 2009.

[77] S. Fox, M. Madden, and A. Smith, "Digital footprints," Pew Internet and American Life report. Retrieved 15 June 2009 from http://www.pewinternet.org/Reports/2007/Digital-Footprints.aspx, December 2007.

[78] Q. Gan and T. Suel, "Improving Web spam classifiers using link structure," in *Proceedings of the 3rd International Workshop on Adversarial Information Retrieval on the Web (AIRWeb)*, pp. 17–20, New York, NY, USA: ACM, 2007.

[79] D. Gibson, R. Kumar, and A. Tomkins, "Discovering large dense subgraphs in massive graphs," in *Proceedings of the 31st International Conference on Very Large Data Bases (VLDB)*, pp. 721–732, VLDB Endowment, 2005.

[80] Y. Gil and D. Artz, "Towards content trust of Web resources," in *Proceedings of the 15th International Conference on World Wide Web (WWW)*, pp. 565–574, New York, NY, USA: ACM, 2006.

[81] A. Gkanogiannis and T. Kalamboukis, "An algorithm for text categorization," in *Proceedings of the 31st Annual International ACM SIGIR Conference on*

*Research and Development in Information Retrieval*, pp. 869–870, New York, NY, USA: ACM, 2008.

[82] A. Gkanogiannis and T. Kalamboukis, "A novel supervised learning algorithm and its use for spam detection in social bookmarking systems," in *Proceedings of the ECML/PKDD Discovery Challenge*, 2008.

[83] H. L. Gomes, B. R. Almeida, A. M. L. Bettencourt, V. Almeida, and M. J. Almeida, "Comparative graph theoretical characterization of networks of spam and legitimate email," April 2005.

[84] M. Goodstein and V. Vassilevska, "A two player game to combat WebSpam," Technical Report, Carnegie Mellon University, 2007.

[85] M. Gori and I. Witten, "The bubble of Web visibility," *Communications of the ACM*, vol. 48, no. 3, pp. 115–117, March 2005.

[86] P. Gramme and J.-F. Chevalier, "Rank for spam detection — ECML discovery challenge," in *Proceedings of the ECML/PKDD Discovery Challenge*, 2008.

[87] A. L. Granka, T. Joachims, and G. Gay, "Eye-tracking analysis of user behavior in www search," in *Proceedings of the 27th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 478–479, New York, NY, USA: ACM, 2004.

[88] J. Grappone and G. Couzin, *Search Engine Optimization: An Hour a Day.* Wiley, 2006.

[89] R. Guha, R. Kumar, P. Raghavan, and A. Tomkins, "Propagation of trust and distrust," in *Proceedings of the 13th International Conference on World Wide Web (WWW)*, pp. 403–412, New York, NY, USA: ACM Press, 2004.

[90] D. M. Guinness, H. Zeng, Li, D. Narayanan, and M. Bhaowal, "Investigations into trust for collaborative information. repositories: A Wikipedia case study," in *Proceedings of workshop on Models of Trust for the Web (MTW06)*, May 2006.

[91] Z. Gyöngyi and H. Garcia-Molina, "Link spam alliances," in *Proceedings of the 31st International Conference on Very Large Data Bases (VLDB)*, pp. 517–528, 2005.

[92] Z. Gyöngyi and H. Garcia-Molina, "Spam: It's not just for inboxes anymore," *IEEE Computer Magazine*, vol. 38, no. 10, pp. 28–34, 2005.

[93] Z. Gyöngyi and H. Garcia-Molina, "Web spam taxonomy," in *Proceedings of the First International Workshop on Adversarial Information Retrieval on the Web (AIRWeb)*, pp. 39–47, May 2005.

[94] Z. Gyöngyi, H. Garcia-Molina, and J. Pedersen, "Combating Web spam with TrustRank," in *Proceedings of the 30th International Conference on Very Large Data Bases (VLDB)*, pp. 576–587, Morgan Kaufmann, August 2004.

[95] Z. I. Gyöngyi, "Applications of Web link analysis," PhD thesis, Stanford University, Adviser: Hector Garcia-Molina, 2008.

[96] Z. P. Gyöngyi, P. Berkhin, H. Garcia-Molina, and J. Pedersen, "Link spam detection based on mass estimation," in *Proceedings of the 32nd International Conference on Very Large Databases (VLDB)*, pp. 439–450, 2006.

[97] H. T. Haveliwala, "Topic-sensitive PageRank," in *Proceedings of the Eleventh World Wide Web Conference (WWW)*, pp. 517–526, ACM Press, May 2002.

[98] R. M. Henzinger, R. Motwani, and C. Silverstein, "Challenges in Web search engines," *SIGIR Forum*, vol. 37, no. 2, 2002.

 [99] R. Herbrich, T. Graepel, and K. Obermayer, "Large margin rank bound-
aries for ordinal regression," in *Advances in Large Margin Classifiers*, (Smola,
Bartlett, Schoelkopf, and Schuurmans, eds.), pp. 115–132, Cambridge, MA:
MIT Press, 2000.

[100] A. Heydon and M. Najork, "Mercator: A scalable, extensible web crawler,"
*World Wide Web*, vol. 2, no. 4, pp. 219–229, 1999.

[101] P. Heymann, G. Koutrika, and H. Garcia-Molina, "Fighting spam on social
Web sites: A survey of approaches and future challenges," *IEEE Internet Com-
puting*, vol. 11, no. 6, pp. 36–45, 2007.

[102] J. Hopcroft and D. Sheldon, "Manipulation-resistant reputations using hitting
time," in *Proceedings of the Workshop on Algorithms and Models for the Web-
Graph (WAW)*, pp. 68–81, Springer: Vol. 2863 of *Lecture Notes in Computer
Science*, December 2007. Also appears in *Internet Mathematics 5, 5:71–90,
2009.*

[103] J. Hopcroft and D. Sheldon, "Network reputation games," Techical Report,
Cornell University, October 2008.

[104] M. Hu, E.-P. Lim, A. Sun, W. H. Lauw, and B.-Q. Vuong, "Measur-
ing article quality in Wikipedia: Models and evaluation," in *Proceedings of
the Sixteenth ACM Conference on Information and Knowledge Management
(CIKM)*, pp. 243–252, New York, NY, USA: ACM, 2007.

[105] S. Hutcheon, "Google pardons BMW website," in *Sydney Morning Herald*,
Retrieved 18 June 2009 from http://www.smh.com.au/news/breaking/google-
pardons-bmw-website/2006/02/09/1139379597733.html, February 2006.

[106] N. Immorlica, K. Jain, and M. Mahdian, "Game-theoretic aspects of design-
ing hyperlink structures," in *Proceedings of the 2nd Workshop on Internet
and Network Economics (WINE)*, pp. 150–161, Vol. 4286, Springer LNCS,
December 2006.

[107] N. Immorlica, K. Jain, M. Mahdian, and K. Talwar, "Click fraud resistant
methods for learning click-through rates," in *Proceedings of the Workshop
on Internet and Network Economics (WINE)*, pp. 34–45, Springer, Berlin:
Vol. 3828 of *Lecture Notes in Computer Science*, 2005.

[108] M. Jamali and M. Ester, "Trustwalker: A random walk model for combining
trust-based and item-based recommendation," in *Proceedings of the 15th ACM
SIGKDD International Conference on Knowledge Discovery and Data Mining*,
pp. 397–406, New York, NY, USA: ACM, 2009.

[109] J. B. Jansen, "Click fraud," *Computer*, vol. 40, no. 7, pp. 85–86, 2007.

[110] Q. Jiang, L. Zhang, Y. Zhu, and Y. Zhang, "Larger is better: Seed selection
in link-based anti-spamming algorithms," in *Proceedings of the 17th Interna-
tional Conference on World Wide Web (WWW)*, pp. 1065–1066, New York,
NY, USA: ACM, 2008.

[111] N. Jindal and B. Liu, "Analyzing and detecting review spam," in *Proceedings
of the 7th IEEE International Conference on Data Mining (ICDM)*, pp. 547–
552, 2007.

[112] N. Jindal and B. Liu, "Review spam detection," in *Proceedings of the 16th
International Conference on World Wide Web (WWW)*, pp. 1189–1190, New
York, NY, USA: ACM Press, 2007.

[113] N. Jindal and B. Liu, "Opinion spam and analysis," in *Proceedings of the International Conference on Web Search and Data Mining (WSDM)*, pp. 219–230, New York, NY, USA: ACM, 2008.

[114] T. Joachims, "Optimizing search engines using clickthrough data," in *Proceedings of the ACM Conference on Knowledge Discovery and Data Mining (KDD)*, pp. 133–142, New York, NY: ACM Press, 2002.

[115] T. Joachims, L. Granka, B. Pan, H. Hembrooke, and G. Gay, "Accurately interpreting clickthrough data as implicit feedback," in *Proceedings of the 28th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 154–161, New York, NY, USA: ACM Press, 2005.

[116] T. Jones, D. Hawking, and R. Sankaranarayana, "A framework for measuring the impact of Web spam," in *Proceedings of 12th Australasian Document Computing Symposium (ADCS)*, December 2007.

[117] T. Jones, D. Hawking, R. Sankaranarayana, and N. Craswell, "Nullification test collections for Web spam and SEO," in *Proceedings of the 5th International Workshop on Adversarial Information Retrieval on the Web (AIRWeb)*, pp. 53–60, New York, NY, USA: ACM, April 2009.

[118] T. Z. Jr., "Gaming the search engine, in a political season," in *New York Times*, November 2006.

[119] D. S. Kamvar, T. M. Schlosser, and H. Garcia-Molina, "The Eigentrust algorithm for reputation management in P2P networks," in *Proceedings of the 12th International Conference on World Wide Web (WWW)*, pp. 640–651, New York, NY, USA: ACM Press, 2003.

[120] G. Karypis and V. Kumar, "Multilevel k-way partitioning scheme for irregular graphs," *Journal of Parallel and Distributed Computation*, vol. 48, no. 1, pp. 96–129, 1998.

[121] T. Katayama, T. Utsuro, Y. Sato, T. Yoshinaka, Y. Kawada, and T. Fukuhara, "An empirical study on selective sampling in active learning for Splog detection," in *Proceedings of the 5th International Workshop on Adversarial Information Retrieval on the Web (AIRWeb)*, pp. 29–36, ACM Press, 2009.

[122] M. J. Kleinberg, "Authoritative sources in a hyperlinked environment," *Journal of the ACM*, vol. 46, no. 5, pp. 604–632, 1999.

[123] J. Köhne, "Optimizing a large dynamically generated Website for search engine crawling and ranking," Master's thesis, Technical University of Delft, 2006.

[124] P. Kolari, "Detecting Spam Blogs: An Adaptive Online Approach," PhD thesis, Department of Computer Science and Electrical Engineering, University of Maryland-Baltimore County, 2007.

[125] P. Kolari, T. Finin, A. Java, and A. Joshi, "Splog blog dataset," Techical Report, UMBC ebiquity, 2006.

[126] P. Kolari, T. Finin, A. Java, and A. Joshi, "Towards spam detection at ping servers," in *Proceedings of the International Conference on Weblogs and Social Media (ICWSM)*, AAAI Press. Demo, March 2007.

[127] P. Kolari, A. Java, and T. Finin, "Characterizing the Splogosphere," in *Proceedings of the 3rd Annual Workshop on the Weblogging Ecosystem*, 2006.

[128] P. Kolari, A. Java, T. Finin, T. Oates, and A. Joshi, "Detecting spam blogs: A machine learning approach," in *Proceedings of the National Conference on Artificial Intelligence (AAAI)*, July 2006.

[129] A. Korolova, K. Kenthapadi, N. Mishra, and A. Ntoulas, "Releasing search queries and clicks privately," in *Proceedings of the 18th International Conference on World Wide Web (WWW)*, pp. 171–180, New York, NY, USA: ACM, 2009.

[130] M. Koster, "A standard for robot exclusion," http://www.robotstxt.org/wc/robots.html, 1996.

[131] Z. Kou, "Stacked graphical learning," PhD thesis, School of Computer Science, Carnegie Mellon University, 2007.

[132] G. Koutrika, A. F. Effendi, Z. Gyöngyi, P. Heymann, and H. Garcia-Molina, "Combating spam in tagging systems," in *Proceedings of the 3rd International Workshop on Adversarial Information Retrieval on the Web (AIRWeb)*, pp. 57–64, New York, NY, USA: ACM Press, 2007.

[133] G. Koutrika, A. F. Effendi, Z. Gyöngyi, P. Heymann, and H. Garcia-Molina, "Combating spam in tagging systems: An evaluation," *ACM Transactions on the Web*, vol. 2, no. 4, pp. 1–34, 2008.

[134] B. Krause, H. A. Schimitz, and G. Stumme, "The anti-social tagger — detecting spam in social bookmarking systems," in *Proceedings of the Fourth International Workshop on Adversarial Information Retrieval on the Web (AIRWeb)*, April 2008.

[135] V. Krishnan and R. Raj, "Web spam detection with anti-TrustRank," in *Proceedings of the 2nd International Workshop on Adversarial Information Retrieval on the Web (AIRWeb)*, pp. 37–40, 2006.

[136] R. Kumar, J. Novak, B. Pang, and A. Tomkins, "On anonymizing query logs via token-based hashing," in *Proceedings of the 16th International Conference on World Wide Web (WWW)*, pp. 629–638, New York, NY, USA: ACM Press, 2007.

[137] J. Kupke and M. Ohye, "Specify your canonical," Retrieved 18 June 2009 from http://googlewebmastercentral.blogspot.com/2009/02/specify-your-canonical.html, February 2009.

[138] N. A. Langville and D. C. Meyer, "Deeper inside PageRank," *Internet Mathematics*, vol. 1, no. 3, pp. 335–380, 2003.

[139] N. A. Langville and D. C. Meyer, *Google's PageRank and Beyond: The Science of Search Engine Rankings*. Princeton, NJ: Princeton University Press, 2006.

[140] H.-T. Lee, D. Leonard, X. Wang, and D. Loguinov, "Irlbot: Scaling to 6 billion pages and beyond," *ACM Transactions on the Web*, vol. 3, no. 3, pp. 1–34, 2009.

[141] R. Lempel and S. Moran, "The stochastic approach for link-structure analysis (SALSA) and the TKC effect," *Computer Networks*, vol. 33, no. 1–6, pp. 387–401, 2000.

[142] R. Levien and A. Aiken, "Attack-resistant trust metrics for public key certification," in *Proceedings of the 7th USENIX Security Symposium*, pp. 229–242, 1998.

[143] J. Lewis, "Google bombs," in *LA Weekly*, Retrieved June 1, 2009 from http://www.laweekly.com/2003-12-25/news/google-bombs, December 2003.

[144] G. Liang, "Surveying Internet usage and impact in five Chinese cities," Report of the Research Center for Social Development, Chinese Academy of Social Sciences. Retrieved 15 June 2009 from http://news.bbc.co.uk/1/shared/bsp/hi/pdfs/10_02_06_china.pdf, November 2005.

[145] M. Lifantsev, "Voting model for ranking Web pages," in *Proceedings of the International Conference on Internet Computing (IC)*, (P. Graham and M. Maheswaran, eds.), pp. 143–148, CSREA Press, June 2000.

[146] J.-L. Lin, "Detection of cloaked Web spam by using tag-based methods," *Expert Systems with Applications*, vol. 36, no. 4, pp. 7493–7499, May 2009.

[147] Y.-R. Lin, H. Sundaram, Y. Chi, J. Tatemura, and L. B. Tseng, "Splog detection using self-similarity analysis on blog temporal dynamics," in *Proceedings of the 3rd International Workshop on Adversarial Information Retrieval on the Web (AIRWeb)*, pp. 1–8, New York, NY, USA: ACM Press, 2007.

[148] Y.-R. Lin, H. Sundaram, Y. Chi, J. Tatemura, and L. B. Tseng, "Detecting splogs via temporal dynamics using self-similarity analysis," *ACM Transations on the Web*, vol. 2, no. 1, pp. 1–35, 2008.

[149] T. Liu, "Analyzing the importance of group structure in the Google PageRank algorithm," Master's thesis, Rensselaer Polytechnic Institute, November 2004.

[150] Y. Liu, R. Cen, M. Zhang, S. Ma, and L. Ru, "Identifying Web spam with user behavior analysis," in *Proceedings of the 4th International Workshop on Adversarial Information Retrieval on the Web (AIRWeb)*, pp. 9–16, New York, NY, USA: ACM, 2008.

[151] Y. Liu, B. Gao, T.-Y. Liu, Y. Zhang, Z. Ma, S. He, and H. Li, "BrowseRank: Letting web users vote for page importance," in *Proceedings of the 31st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 451–458, New York, NY, USA: ACM, 2008.

[152] J. Ma, K. L. Saul, S. Savage, and M. G. Voelker, "Beyond blacklists: Learning to detect malicious web sites from suspicious urls," in *Proceedings of the 15th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 1245–1254, New York, NY, USA: ACM, 2009.

[153] C. C. Mann, "How click fraud could swallow the internet," *Wired*, vol. 14, no. 1, January 2006.

[154] D. C. Manning, P. Raghavan, and H. Schütze, *Introduction to Information Retrieval*. Cambridge University Press, 2008.

[155] B. Markines, C. Cattuto, and F. Menczer, "Social spam detection," in *Proceedings of the 5th International Workshop on Adversarial Information Retrieval on the Web (AIRWeb)*, pp. 41–48, New York, NY, USA: ACM, 2009.

[156] K. Marks and T. Celik, "Microformats: The rel=nofollow attribute," Techical Report, Technorati, 2005. Online at http://microformats.org/wiki/rel-nofollow. Last accessed 29 January 2009.

[157] K. Marks and T. Celik, "Microformats: Vote links," Technical Report, Technorati, 2005. Online at http:// microformats.org/wiki/vote-links. Last accessed 29 January 2009.

[158] S. Marti and H. Garcia-Molina, "Taxonomy of trust: Categorizing P2P reputation systems," *Computer Networks*, vol. 50, no. 4, pp. 472–484, March 2006.

[159] J. Martinez-Romo and L. Araujo, "Web spam identification through language model analysis," in *Proceedings of the 5th International Workshop on Adversarial Information Retrieval on the Web (AIRWeb)*, pp. 21–28, ACM Press, 2009.

[160] K. Mason, "Detecting Colluders in PageRank: Finding Slow Mixing States in a Markov Chain," PhD thesis, Department of Engineering Economic Systems and Operations Research, Stanford University, September 2005.

[161] P. Massa and C. Hayes, "Page-reRank: Using trusted links to re-rank authority," in *Proceedings of the IEEE/WIC/ACM International Conference on Web Intelligence (WI)*, pp. 614–617, 2005.

[162] A. Mathes, "Filler Friday: Google Bombing," Retrieved June 1, 2009 from http://uber.nu/2001/04/06/, April 2001.

[163] T. McNichol, "Engineering Google results to make a point," *New York Times*, January 2004.

[164] T. P. Metaxas and J. Destefano, "Web spam, propaganda and trust," in *Proceedings of the First International Workshop on Adversarial Information Retrieval on the Web (AIRWeb)*, May 2005.

[165] A. Metwally, D. Agrawal, and E. A. Abbadi, "Detectives: Detecting coalition hit inflation attacks in advertising networks streams," in *Proceedings of the 16th International Conference on World Wide Web (WWW)*, pp. 241–250, New York, NY, USA: ACM Press, 2007.

[166] A. G. Mishne, "Applied Text Analytics for Blogs," PhD thesis, University of Amsterdam, April 2007.

[167] G. Mishne, D. Carmel, and R. Lempel, "Blocking blog spam with language model disagreement," in *Proceedings of the 1st International Workshop on Adversarial Information Retrieval on the Web (AIRWeb)*, May 2005.

[168] T. Moore and R. Clayton, "Evil searching: Compromise and recompromise of internet hosts for phishing," in *Financial Cryptography and Data Security*, pp. 256–272, Springer, 2009.

[169] M. Moran and B. Hunt, *Search Engine Marketing, Inc.* Upper Saddle River, NJ: IBM Press, 2006.

[170] A. Moshchuk, T. Bragin, D. S. Gribble, and M. H. Levy, "A crawler-based study of spyware on the web," in *Proceedings of the Network and Distributed System Security Symposium (NDSS)*, pp. 17–33, February 2006.

[171] R. Moulton and K. Carattini, "A quick word about Googlebombs," Retrieved June 1, 2009 from http://googlewebmastercentral.blogspot.com/2007/01/quick-word-about-googlebombs.html, January 2007.

[172] L. Mui, M. Mohtashemi, and A. Halberstadt, "A computational model of trust and reputation," in *Proceedings of the 35th Hawaii International Conference on System Science (HICSS)*, 2002.

[173] M. Najork, "System and method for identifying cloaked web servers," U.S. Patent 6,910,077 (issued June 2005), 2002.

[174] N. Neubauer, R. Wetzker, and K. Obermayer, "Tag spam creates large non-giant connected components," in *Proceedings of the 5th International Workshop on Adversarial Information Retrieval on the Web (AIRWeb)*, pp. 49–52, ACM Press, 2009.

[175] B. News, "Miserable failure' links to Bush: George W Bush has been Google bombed," http://news.bbc.co.uk/2/hi/americas/3298443.stm, December 2003.

[176] L. Nie, D. B. Davison, and B. Wu, "Incorporating trust into Web authority," Technical Report LU-CSE-07-002, Department of Computer Science and Engineering, Lehigh University, 2007.

[177] L. Nie, B. Wu, and D. B. Davison, "A cautious surfer for PageRank," in *Proceedings of the 16th International Conference on World Wide Web (WWW)*, pp. 1119–1120, New York, NY, USA: ACM Press, 2007.

[178] L. Nie, B. Wu, and D. B. Davison, "Winnowing wheat from the chaff: Propagating trust to sift spam from the Web," in *Proceedings of the 30th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 869–870, July 2007.

[179] Y. Niu, Y.-M. Wang, H. Chen, M. Ma, and F. Hsu, "A quantitative study of forum spamming using context-based analysis," in *Proceedings of the 14th Annual Network and Distributed System Security Symposium (NDSS)*, pp. 79–92, February 2007.

[180] A. Ntoulas, "Crawling and searching the Hidden Web," PhD thesis, University of California at Los Angeles, Los Angeles, CA, USA, 2006. Adviser-Cho, Junghoo.

[181] A. Ntoulas, J. Cho, and C. Olston, "What's new on the Web?: The evolution of the Web from a search engine perspective," in *Proceedings of the 13th International Conference on World Wide Web*, pp. 1–12, New York, NY, USA: ACM Press, 2004.

[182] A. Ntoulas, M. Najork, M. Manasse, and D. Fetterly, "Detecting spam Web pages through content analysis," in *Proceedings of the 15th International Conference on World Wide Web (WWW)*, pp. 83–92, May 2006.

[183] L. Page, S. Brin, R. Motwani, and T. Winograd, "The PageRank citation ranking: Bringing order to the Web," Techincal Report, Stanford University, 1998. Available from http://dbpubs.stanford.edu/pub/1999-66.

[184] R. C. Palmer, B. P. Gibbons, and C. Faloutsos, "ANF: A fast and scalable tool for data mining in massive graphs," in *Proceedings of the Eighth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 81–90, New York, NY, USA: ACM Press, 2002.

[185] K. Park, V. S. Pai, K.-W. Lee, and S. Calo, "Securing Web service by automatic robot detection," in *Proceedings of the USENIX Annual Technical Conference*, pp. 255–260, 2006.

[186] J.-X. Parreira, D. Donato, C. Castillo, and G. Weikum, "Computing trusted authority scores in peer-to-peer networks," in *Proceedings of the 3rd International Workshop on Adversarial Information Retrieval on the Web (AIRWeb)*, pp. 73–80, ACM Press, May 2007.

[187] L. A. Penenberg, "Click fraud threatens Web," Wired, October 2004.

[188] L. A. Penenberg, "Legal showdown in search fracas," in *Wired*, Retrieved 18 June 2009 from http://www.wired.com/culture/lifestyle/news/2005/09/68799, September 2005.

[189] A. Perkins, "The classification of search engine spam," Available online at http://www.silverdisc.co.uk/articles/spam-classification/, September 2001.

[190] J. Piskorski, M. Sydow, and D. Weiss, "Exploring linguistic features for Web spam detection: A preliminary study," in *Proceedings of the Fourth International Workshop on Adversarial Information Retrieval on the Web (AIRWeb)*, pp. 25–28, New York, NY, USA: ACM, 2008.

[191] G. Price, "Google and Google bombing now included New Oxford American Dictionary," Retrieved June 1, 2009 from http://blog.searchenginewatch.com/blog/050516-184202, May 2005.

[192] N. Provos, P. Mavrommatis, A. M. Rajab, and F. Monrose, "All your iFRAMEs point to us," in *Proceedings of the 17th USENIX Security Symposium*, pp. 1–15, Berkeley, CA, USA: USENIX Association, 2008.

[193] N. Provos, D. McNamee, P. Mavrommatis, K. Wang, and N. Modadugu, "The ghost in the browser analysis of web-based malware," in *Proceedings of the First Workshop on Hot Topics in Understanding Botnets (HotBots)*, 2007.

[194] X. Qi and D. B. Davison, "Knowing a web page by the company it keeps," in *Proceedings of the 15th ACM International Conference on Information and Knowledge Management (CIKM)*, pp. 228–237, New York, NY: ACM Press, November 2006.

[195] X. Qi and D. B. Davison, "Web page classification: Features and algorithms," *ACM Computing Surveys*, vol. 41, no. 2, February 2009.

[196] X. Qi, L. Nie, and D. B. Davison, "Measuring similarity to detect qualified links," in *Proceedings of the 3rd International Workshop on Adversarial Information Retrieval on the Web (AIRWeb)*, pp. 49–56, May 2007.

[197] R. J. Quinlan, *C4.5: Programs for Machine Learning*. San Mateo, CA: Morgan Kaufmann, 1993.

[198] S. R. Rainwater, "Nigritude ultramarine FAQ," Retrieved June 1, 2009 from http://www.nigritudeultramarines.com/, 2005.

[199] Y. Rasolofo and J. Savoy, "Term proximity scoring for keyword-based retrieval systems," in *ECIR*, pp. 207–218, 2003.

[200] P. Resnick, K. Kuwabara, R. Zeckhauser, and E. Friedman, "Reputation systems," *Communications of the ACM*, vol. 43, no. 12, pp. 45–48, 2000.

[201] M. Richardson, A. Prakash, and E. Brill, "Beyond PageRank: Machine learning for static ranking," in *Proceedings of the 15th International Conference on World Wide Web (WWW)*, pp. 707–715, New York, NY, USA: ACM Press, May 2006.

[202] G. Roberts and J. Rosenthal, "Downweighting tightly knit communities in World Wide Web rankings," *Advances and Applications in Statistics (ADAS)*, vol. 3, pp. 199–216, 2003.

[203] S. Robertson, S. Walker, M. M. Beaulieu, M. Gatford, and A. Payne, "Okapi at TREC-4," in *NIST Special Publication 500-236: The Fourth Text REtrieval Conference (TREC-4)*, pp. 73–96, 1995.

[204] G. Salton, A. Wong, and S. C. Yang, "A vector space model for automatic indexing," *Communications of the ACM*, vol. 18, no. 11, pp. 613–620, November 1975.

[205] Y. Sato, T. Utsuro, Y. Murakami, T. Fukuhara, H. Nakagawa, Y. Kawada, and N. Kando, "Analysing features of Japanese splogs and characteristics of keywords," in *Proceedings of the Fourth International Workshop on Adversarial*

*Information Retrieval on the Web (AIRWeb)*, pp. 33–40, New York, NY, USA: ACM, 2008.

[206] C. Schmidt, "Page hijack: The 302 exploit, redirects and Google," http://clsc.net/research/google-302-page-hijack.htm, 2005.

[207] F. Sebastiani, "Machine learning in automated text categorization," *ACM Computing Surveys*, vol. 34, no. 1, pp. 1–47, 2002.

[208] D. Sheldon, "Manipulation of PageRank and Collective Hidden Markov Models," PhD thesis, Cornell University, 2009.

[209] G. Shen, B. Gao, T.-Y. Liu, G. Feng, S. Song, and H. Li, "Detecting link spam using temporal information," in *Proceedings of the Sixth IEEE International Conference on Data Mining (ICDM)*, December 2006.

[210] N. Shrivastava, A. Majumder, and R. Rastogi, "Mining (social) network graphs to detect random link attacks," in *Proceedings of the International Conference on Data Engineering (ICDE)*, IEEE CS Press, April 2008.

[211] F. Silvestri, "Mining query logs: Turning search usage data into knowledge," *Foundations and Trends in Information Retrieval*, vol. 3, 2009.

[212] A. Singhal, "Challenges in running a commercial search engine," Keynote presentation at SIGIR 2005, August 2005.

[213] M. Sirivianos, X. Yang, and K. Kim, "FaceTrust: Assessing the credibility of online personas via social networks," Technical Report, Duke University, 2009. Retrieved 15 June 2009 from http://www.cs.duke.edu/~msirivia/publications/facetrust-tech-report.pdf.

[214] M. Sobek, "PR0 — Google's PageRank 0 penalty," http://pr.efactory.de/e-pr0.shtml, 2002.

[215] A. Stassopoulou and D. M. Dikaiakos, "Web robot detection: A probabilistic reasoning approach," *Computer Networks*, vol. 53, no. 3, pp. 265–278, February 2009.

[216] A.-J. Su, C. Y. Hu, A. Kuzmanovic, and C.-K. Koh, "How to improve your Google ranking: Myths and reality," in *Proceedings of the IEEE/WIC/ACM International Conference on Web Intelligence and Intelligent Agent Technology (WI-IAT)*, pp. 50–57, Vol. 1, IEEE, 2010.

[217] M. K. Svore, Q. Wu, C. J. C. Burges, and A. Raman, "Improving Web spam classification using rank-time features," in *Proceedings of the Third International Workshop on Adversarial Information Retrieval on the Web (AIRWeb)*, pp. 9–16, New York, NY, USA: ACM Press, 2007.

[218] P.-N. Tan and V. Kumar, "Discovery of Web robot sessions based on their navigational patterns," *Data Mining and Knowledge Discovery*, vol. 6, no. 1, pp. 9–35, 2002.

[219] E. Tardos and T. Wexler, "Network formation games and the potential function method," in *Algorithmic Game Theory*, (N. Nisan, T. Roughgarden, E. Tardos, and V. Vazirani, eds.), Cambridge University Press, 2007.

[220] C. Tatum, "Deconstructing Google bombs: A breach of symbolic power or just a goofy prank?," *First Monday*, vol. 10, no. 10, October 2005.

[221] A. Thomason, "Blog spam: A review," in *Proceedings of Conference on Email and Anti-Spam (CEAS)*, August 2007.

[222] A. Toffler, "The Third Wave," Bantam Books, 1980.

[223] N. Tran, B. Min, J. Li, and L. Submaranian, "Sybil-resilient online content voting," in *Proceedings of the 6th Symposium on Networked System Design and Implementation (NSDI)*, 2009.

[224] T. Urvoy, E. Chauveau, P. Filoche, and T. Lavergne, "Tracking Web spam with HTML style similarities," *ACM Transactions on the Web*, vol. 2, no. 1, 2008.

[225] T. Urvoy, T. Lavergne, and P. Filoche, "Tracking Web spam with hidden style similarity," in *Proceedings of the Second International Workshop on Adversarial Information Retrieval on the Web (AIRWeb)*, August 2006.

[226] N. Vidyasaga, "India's secret army of online ad 'clickers'," *The Times of India*, May 2004.

[227] L. A. von Ahn and L. Dabbish, "Labeling images with a computer game," in *Proceedings of the Conference on Human Factors in Computing Systems (CHI)*, pp. 319–326, ACM Press, 2004.

[228] K. Walsh and G. E. Sirer, "Fighting peer-to-peer spam and decoys with object reputation," in *Proceedings of the ACM SIGCOMM Workshop on Economics of Peer-to-Peer systems (P2PECON)*, pp. 138–143, New York, NY, USA: ACM Press, 2005.

[229] W. Wang, G. Zeng, M. Sun, H. Gu, and Q. Zhang, "EviRank: An evidence based content trust model for Web spam detection," in *Proceedings of Workshop on Emerging Trends of Web Technologies and Applications WAIM/APWeb*, pp. 299–307, 2007.

[230] Y.-M. Wang, D. Beck, X. Jiang, and R. Roussev, "Automated web patrol with Strider HoneyMonkeys: Finding web sites that exploit browser vulnerabilities," in *Proceedings of the Network and Distributed System Security Symposium (NDSS)*, February 2006.

[231] Y.-M. Wang, M. Ma, Y. Niu, and H. Chen, "Spam double-funnel: Connecting Web spammers with advertisers," in *Proceedings of the 16th International Conference on World Wide Web (WWW)*, pp. 291–300, New York, NY, USA: ACM Press, 2007.

[232] S. Webb, "Automatic Identification and Removal of Low Quality Online Information," PhD thesis, College of Computing, Georgia Institute of Technology, December 2008.

[233] S. Webb, J. Caverlee, and C. Pu, "Introducing the Webb spam corpus: Using email spam to identify Web spam automatically," in *Proceedings of the Third Conference on Email and Anti-Spam (CEAS)*, July 2006.

[234] S. Webb, J. Caverlee, and C. Pu, "Predicting Web spam with HTTP session information," in *Proceeding of the 17th ACM Conference on Information and Knowledge Management (CIKM)*, pp. 339–348, New York, NY, USA: ACM, 2008.

[235] H. I. Witten and E. Frank, *Data Mining: Practical Machine Learning Tools and Techniques with Java Implementations*. Morgan Kaufmann, 1999.

[236] B. Wu, "Finding and Fighting Search Engine Spam," PhD thesis, Department of Computer Science and Engineering, Lehigh University, March 2007.

[237] B. Wu and K. Chellapilla, "Extracting link spam using biased random walks from spam seed sets," in *Proceedings of the 3rd International Workshop on*

*Adversarial Information Retrieval on the Web (AIRWeb)*, pp. 37–44, New York, NY, USA: ACM Press, 2007.

[238] B. Wu and D. B. Davison, "Cloaking and redirection: A preliminary study," in *Proceedings of the 1st International Workshop on Adversarial Information Retrieval on the Web (AIRWeb)*, 2005.

[239] B. Wu and D. B. Davison, "Detecting semantic cloaking on the Web," in *Proceedings of the 15th International World Wide Web Conference (WWW)*, pp. 819–828, ACM Press, 2006.

[240] B. Wu and D. B. Davison, "Undue influence: Eliminating the impact of link plagiarism on Web search rankings," in *Proceedings of The 21st ACM Symposium on Applied Computing (SAC)*, pp. 1099–1104, April 2006.

[241] B. Wu and D. B. Davison, "Identifying link farm spam pages," in *Special interest tracks and posters of the 14th International Conference on World Wide Web (WWW)*, pp. 820–829, New York, NY, USA: ACM Press, 2005.

[242] B. Wu, V. Goel, and D. B. Davison, "Propagating trust and distrust to demote Web spam," in *Workshop on Models of Trust for the Web (MTW)*, May 2006.

[243] B. Wu, V. Goel, and D. B. Davison, "Topical TrustRank: Using topicality to combat Web spam," in *Proceedings of the 15th International World Wide Web Conference (WWW)*, pp. 63–71, ACM Press, May 2006.

[244] L. Xiong and L. Liu, "PeerTrust: Supporting reputation-based trust for peer-to-peer electronic communities," *IEEE Transactions on Knowledge and Data Engineering*, vol. 16, no. 7, pp. 843–857, 2004.

[245] H. Yu, M. Kaminsky, B. P. Gibbons, and A. Flaxman, "SybilGuard: Defending against sybil attacks via social networks," in *Proceedings of the ACM Conference on Applications, Technologies, Architectures, and Protocols for Computer Communications (SIGCOMM)*, pp. 267–278, New York, NY, USA: ACM Press, 2006.

[246] K. Yusuke, W. Atsumu, K. Takashi, B. B. Bahadur, and T. Toyoo, "On a referrer spam blocking scheme using Bayesian filter," *Joho Shori Gakkai Shinpojiumu Ronbunshu*, vol. 1, no. 13, pp. 319–324, In Japanese, 2005.

[247] H. Zhang, A. Goel, R. Govindan, K. Mason, and B. V. Roy, "Making eigenvector-based reputation systems robust to collusion," in *Proceedings of the Third Workshop on Web Graphs (WAW)*, pp. 92–104, Springer: Vol. 3243 of *Lecture Notes in Computer Science*, October 2004.

[248] L. Zhang, Y. Zhang, Y. Zhang, and X. Li, "Exploring both content and link quality for anti-spamming," in *Proceedings of the Sixth IEEE International Conference on Computer and Information Technology (CIT)*, IEEE Computer Society, 2006.

[249] Y. Zhang and A. Moffat, "Some observations on user search behavior," *Australian Journal of Intelligent Information Processing Systems*, vol. 9, no. 2, pp. 1–8, 2006.

[250] L. Zhao, Q. Jiang, and Y. Zhang, "From good to bad ones: Making spam detection easier," in *Proceedings of the Eighth IEEE International Conference on Computer and Information Technology Workshops (CIT)*, IEEE Computer Society, 2008.

[251] B. Zhou, "Mining page farms and its application in link spam detection," Master's thesis, Simon Fraser University, 2007.

[252] B. Zhou and J. Pei, "Sketching landscapes of page farms," in *Proceedings of the 7th SIAM International Conference on Data Mining (SDM)*, SIAM, April 2007.

[253] B. Zhou, J. Pei, and Z. Tang, "A spamicity approach to Web spam detection," in *Proceedings of the SIAM International Conference on Data Mining (SDM)*, April 2008.

[254] D. Zhou, C. J. C. Burges, and T. Tao, "Transductive link spam detection," in *Proceedings of the 3rd International Workshop on Adversarial Information Retrieval on the Web (AIRWeb)*, pp. 21–28, New York, NY, USA: ACM Press, 2007.

[255] C.-N. Ziegler and G. Lausen, "Propagation models for trust and distrust in social networks," *Information Systems Frontiers*, vol. 7, no. 4–5, pp. 337–358, December 2005.

[256] R. P. Zimmermann, *The Official PGP User's Guide*. Cambridge, MA: MIT Press, 1995.

[257] J. Zittrain, *The Future of the Internet — And How to Stop It*. Yale University Press, April 2008.