

The information ecology of social media and online communities*

Tim Finin, Anupam Joshi, Pranam Kolari, Akshay Java, Anubhav Kale and Amit Karandikar
University of Maryland, Baltimore County, Baltimore MD 21250

Abstract

Social media systems such as weblogs, photo- and link-sharing sites, wikis and on-line forums are currently thought to produce up to one third of new Web content. One thing that sets these “Web 2.0” sites apart from traditional Web pages and resources is that they are intertwined with other forms of networked data. Their standard hyperlinks are enriched by social networks, comments, trackbacks, advertisements, tags, RDF data and metadata. We describe recent work on building systems that use models of the Blogosphere to recognize spam blogs, find opinions on topics, identify communities of interest, derive trust relationships, and detect influential bloggers.

Introduction

Web-based social media systems such as blogs, wikis, media-sharing sites and message forums have become an important new way to publish information, engage in discussions and form communities on the Internet. Their reach and impact is significant, with tens of millions of people providing content on a regular basis around the world. Recent estimates suggest that social media systems are responsible for as much as one third of new Web content. Corporations, traditional media companies, governments and NGOs are working to understand how to adapt to them and use them effectively. Citizens, both young and old, are also discovering how social media technology can improve their lives and give them more voice in the world. We must better understand the information ecology of these new publication methods in order to make them and the information they provide more useful, trustworthy and reliable.

We are developing a model of information flow, influence and trust on the Blogosphere and are exploring how this model can be used to answer questions like the following: how can blog communities be identified

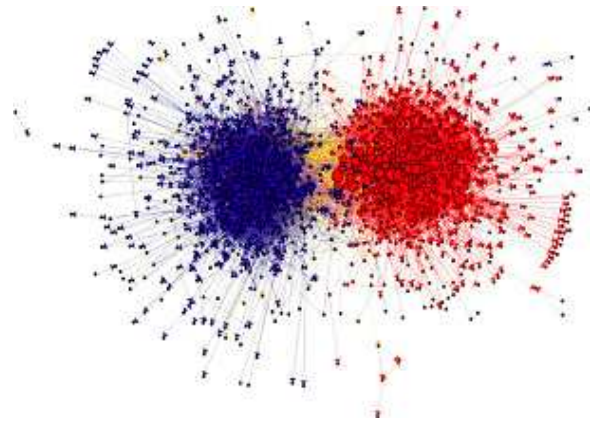


Figure 1: Modeling influence and information flow on the Blogosphere and other social media systems requires attention to many factors, including link structure, sentiment analysis, readership data, conversational structure, topic classification, and temporal analysis (figure from (Adamic & Glance 2005)).

based on a combination of topic, bias, and underlying beliefs; which authors and blogs are most influential within a given community; from where do particular beliefs or ideas originate and how do they spread; what are the most trustworthy sources of information about a particular topic; and what opinions and beliefs characterize a community and how do these opinions change.

The Blogosphere is part of the Web and therefore shares most of its general characteristics. It differs, however, in ways that impact how it should be modeled, analyzed and exploited. The common model for the general World Wide Web is as a directed graph of Web pages with undifferentiated links between pages. The Blogosphere has a much richer network structure in that there are more *types* of nodes which have more *kinds* of relations between them. For example, the people who contribute to blogs and author blog posts form a social network with their peers, which can be induced by the links between blogs. The blogs themselves form a graph, with direct links to other blogs through *blogrolls* and indirect links through their posts. Blog posts are linked to their host blogs and typically to other blog

*Submitted to AI Magazine, special issue on networks, Fall 2007. Corresponding author: Tim Finin, finin@umbc.edu, 410-455-3522, fax: 410-455-3969). Copyright © 2007, American Association for Artificial Intelligence (www.aaai.org). All rights reserved.

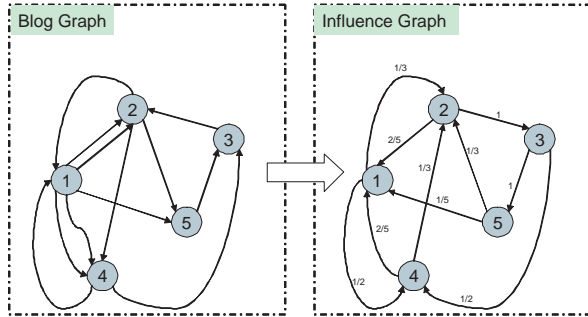


Figure 2: A graph of blogs and Web links between them can be converted into an *influence graph*. A link from u to v indicates that u is influenced by v . The edges in the influence graph are reverse of the blog graph to indicate this influence. Multiple edges indicate stronger influence and are weighed higher

posts and Web resources as part of their content. A typical blog post has a set of comments that link back to people and blogs associated with them. Finally, the blogosphere trackback protocol generates implicit links between blog posts. Still more detail can be added by taking into account post tags and categories, syndication feeds, and semi-structured metadata in the form of XML and RDF content.

In the rest of this article, we discuss our ongoing research in modeling the Blogosphere and extracting useful information from it. We begin by describing an overarching task of discovering which blogs and bloggers are most influential within a community or about a topic. Pursuing this task uncovers a number of problems that must be addressed, three of which we describe in more detail. The first is recognizing spam in the form of spam blogs (splogs) and spam comments. The second is developing more effective techniques to recognize the social structure of blog communities. The final one involves devising a better abstract model for the underlying blog network structure and how it evolves.

Modeling influence in the Blogosphere

The Blogosphere provides an interesting opportunity to study online social interactions including spread of information, opinion formation and influence. Through original content and sometimes via commentary on topics of current interest, bloggers influence each other and their audience. We are working to study and characterize these social interaction by modeling the Blogosphere and providing novel algorithms for analyzing social media content. Figure 2 shows a hypothetical blog graph and its corresponding flow of information in the *influence graph*.

Studies on influence in social networks and collaboration graphs have typically focused on the task of identifying key individuals who play an important role in propagating information. This is similar to finding

authoritative pages on the Web. Epidemic-based models like linear threshold and cascade models (Kempe, Kleinberg, & Tardos 2003; 2005; Leskovec *et al.* 2007) have been used to find a small set of individuals who are most influential in social network. However, influence on the Web is often a function of topic. For example, Engadget’s¹ influence is in the domain of consumer electronics and Daily Kos² in politics. A post in the former is unlikely to be very effective in influencing opinions on political issues even though Engadget is one of the most popular blogs on the Web.

The other related dimension of influence is readership. With the large number of niches existing on the Blogosphere, a blog that is relatively low ranked can be highly influential in this small community of interest. In addition, influence can be subjective and based on the interest of the users. Thus by analyzing the readership of a blog we gain insights into the community that is likely to be influenced by it.

We have implemented a system called Feeds That Matter³ (Java *et al.* 2007b) that aggregates subscription information across thousands of Bloglines users to automatically categorize blogs into different topics. Bloglines⁴ is a popular *feed reader* service that lets users to manage subscriptions and monitor a number of feeds for any unread posts. Bloglines provides a feature allows users to share their subscriptions. We conduct a study of the publicly listed OPML⁵ feeds from 83,204 users consisting of a total of 2,786,687 subscriptions of which 496,879 are unique. A Blogline user’s feeds are typically organized into named folders, such as *Podcasts* or *Politics* and the folder structure is maintained in the OPML representation. The folder names can be used as an approximation of the topic that a user associated with a feed. By clustering related folders, we can induce an intuitive set of topics for feeds and blogs. Figure 3 shows a tag cloud of popular topics aggregated from readership information. Finally, we rank the feeds relevant to each of the topics generated. In our approach, we say that a feed is topically relevant and authoritative if many users have categorized it under similar folder names. For example, Table 1 shows the top political blogs ranked using readership-based influence metrics.

An important component in understanding influence is to detect the sentiment and opinions expressed in blog posts. An aggregated opinion over many users is a predictor for an interesting trend in a community. Sufficient adoption of this trend could lead to a “tipping point” and consequently influence the rest of the community. The BlogVox system (Java *et al.* 2007a; Martineau *et al.* 2007; Balijepalli 2007) retrieves opinionated blog posts specified by ad hoc queries identi-

¹<http://engadget.com/>

²<http://dailykos.com/>

³<http://ftm.umbc.edu/>

⁴<http://bloglines.com/>

⁵OPML, or Outline Processor Markup Language, is an XML format commonly used to share lists of Web feed URLs.



Figure 3: The tag cloud generated from the top 200 folders before and after merging related folders. The size of the word is scaled to indicate how many users use the folder name.

| | |
|----|---|
| 1 | http://www.talkingpointsmemo.com |
| 2 | http://www.dailykos.com |
| 3 | http://atrios.blogspot.com |
| 4 | http://www.washingtonmonthly.com |
| 5 | http://www.wonkette.com |
| 6 | http://instapundit.com |
| 7 | http://www.juancole.com |
| 8 | http://powerlineblog.com |
| 9 | http://americablog.blogspot.com |
| 10 | http://www.crooksandliars.com |

Table 1: The Feeds That Matter for ‘Politics’ shows the top political blogs ranked using readership-based influence metrics.

fying an entity or topic of interest (e.g., “March of the Penguins”). After retrieving posts relevant to a topic query, the system processes them to produce a set of independent features estimating the likelihood that a post expresses an opinion about the topic. These are combined using an SVM-based system and integrated with the relevancy score to rank the results.

Since blog posts are often informally written, poorly structured, rife with spelling and grammatical errors, and feature non-traditional content they are difficult to process with standard language analysis tools. Performing linguistic analysis on blogs is plagued by two additional problems: (i) the presence of spam blogs and spam comments and (ii) extraneous non-content including blog rolls, link rolls, advertisements and sidebars. In the next section we describe techniques designed to eliminate spam content from a blog index. This is a vital task before any useful analytics can be supported on social media content.

In the following sections we also introduce a technique we call “link polarity”. We represent each edge in the influence graph with a vector of topic and corresponding weights indicating either positive or negative

sentiment associated with the link for a Web resource. Thus if a blog A links to a blog B with a negative sentiment for a topic T, influencing B would have little effect on A. Opinions are also manifested as biases. A community of ipod fanatics, for example, needs little or no convincing that it is a good product. Thus, attempting to influencing an opinion leader in such already positively biased communities will have less impact. Using link polarity and trust propagation we have demonstrated how like-minded blogs can be discovered and the potential of using this technique for more generic problems such as detecting trustworthy nodes in web graphs (Kale *et al.* 2007).

Existing models of influence have considered a static view of the network. The Blogosphere, on the other hand, is extremely dynamic and “buzzy”. New topics emerge and blogs constantly rise and fall in popularity. By considering influence as a temporal phenomenon, we can find key individuals that are early adopters or “buzz generators” for a topic. We propose an abstract model of the Blogosphere that provides a systematic approach to modeling the evolution of the link structure and communities. Thus in order to model influence on the Blogosphere, we need to consider topic, readership, community structure, sentiment and time.

In the following sections, we provide a detailed description of various issues that need to be handled in order to model influence. Detecting influence and understanding its role in how people perceive and adopt a product or service provides a powerful tool for marketing, advertising and business intelligence. This requires new algorithms that build on social network analysis, community detection and opinion extraction.

Detecting blog spam

As with other forms of communication, spam has become a serious problem in blogs and social media, both for users and for systems that harvest, index and analyze generated content. Two forms of spam are common in blogs: spam blogs (also known as splogs) where the entire blog and hosted posts are machine generated, and spam comments where authentic posts feature machine generated comments. Though splogs continue to be a problem for web search engines and are considered a special case of web spam, they present a new set of challenges for blog analytics. Given the context of this paper and the intricacies of indexing blogs (Mishne 2007) we limit our discussion to splogs.

Blog search engines index new blog posts by processing pings from update ping servers, intermediary systems that aggregate notifications from updated blogs. Scores of spam pages infiltrate at these ping servers increasing computational requirements, corrupting results, and eventually reducing user satisfaction. We estimate that more than 50% of all pings are from spam sources (Kolari, Java, & Finin 2006). Two kinds of spam content sources are prevalent:



Figure 4: This example of a splog contains plagiarized content(ii), promotes other spam pages by linking to them (iii) and (i) hosts high paying advertisements automatically selected to match the splog’s content.

- **Non-Blogs** are pages that attempt to increase the visibility of hosted and linked-to content, by feigning to be blogs to leverage higher trust and quicker indexing by web search engines. An example of one such non-blog is an Amazon affiliate (third-party vendor of products) book-selling site that pings an update ping server.
- **Spam blogs** constitute the second kind of spam. These are blogs created using splog creation tools (Finin 2006), and are either fully or partly machine generated. Splogs have two often overlapping motives. The first is the creation of blogs containing gibberish or hijacked content from other blogs and news sources with the sole purpose of hosting profitable context based advertisements. The second is the creation of blogs which realize link farms intended to increase the ranking of affiliate sites (blogs or non-blog web-pages). One such splog is shown in figure 4.

Detecting Splogs

Over the past year (Kolari 2007) we have developed techniques to detect spam blogs as they fit the overall architecture (figure 5), arrived at through our discussions with practitioners. Our existing and continuing work has explored all aspects of this architecture. We discuss highlights of our effort based on splog detection using blog home-pages with local and relational features. Interested readers are referred to (Kolari *et al.* 2006; Kolar, Finin, & Joshi 2006) for further details.

Results reported in the rest of this section are based on a seed data set of 700 positive (splogs) and 700 negative (authentic blog) labeled examples containing the entire HTML content of each blog home-page. All of the models are based on SVMs (Boser, Guyon, & Vapnik 1992), which are known to perform well in clas-

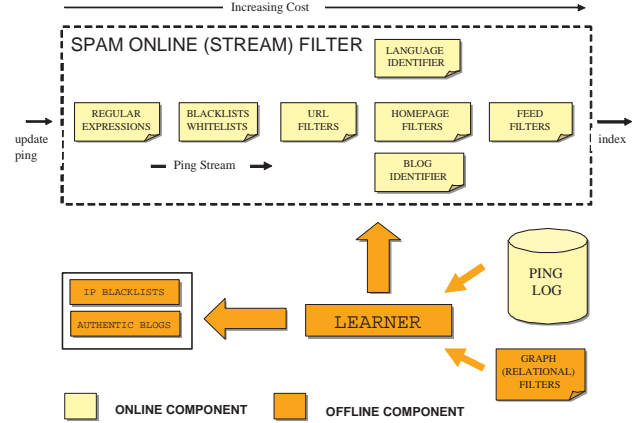


Figure 5: In our online spam detection system, pings from existing ping servers are aggregated and pass through a stream based filtering system. Sub-filters are ordered based on their cost of filtering with each deciding whether to reject a ping as spam or pass it through. Relational techniques are used in the offline component. The system uses all these detection techniques together to adapt and co-evolve through a learning component.

sification tasks (Joachims 1998). We use linear kernel with top features chosen using mutual information, and models evaluated using one-fold cross validation. We view detection techniques as local and relational, based on feature types used.

Local Features

A blog’s local features can be quite effective for splog detection. A *local feature* is one that is completely determined by the contents of a single web page, i.e. it does not require following links or consulting other data sources. A local model built using only these features can provide a quick assessment of the authenticity of blogs. We have experimented with many such models, and our results are summarized in Figure 6.

(i) **Words.** To verify their utility, we created bag-of-words for the samples based on their textual content. We also analyzed discriminating features by ordering features based on weights assigned to them by the linear kernel. It turns out that the model was built around features which the human eye would have typically overlooked. Blogs often contain content that expresses personal opinions, so words like “I”, “We”, “my”, “what” appear commonly on authentic blog posts. To this effect, the bag-of-words model is built on an interesting “blog content genre”. In general, such a content genre is not seen on the Web, which partly explains why spam detection using local textual content is less effective there.

(ii) **Word N-Grams.** An alternative methodology to using textual content for classification is the bag-of-word-N-Grams, where N adjacent words are used

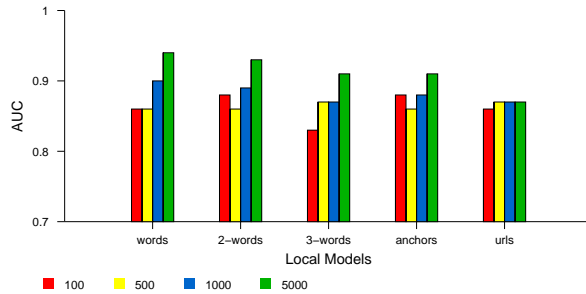


Figure 6: The performance of local models, as measured by the standard, *area under the curve* metric, varies for different feature types and sizes.

as a feature. We evaluated both bag-of-word-2-Grams and bag-of-word-3-Grams, which turned out to be almost as effective as bag-of-words. Interesting discriminative features were observed in this experiment. For instance, text like “comments-off” (comments are usually turned-off in splogs), “new-york” (a high paying advertising term), “in-uncategorized” (spammers do not bother to specify categories for blog posts) are features common to splogs, whereas text like “2-comments”, “1-comment”, “i-have”, “to-my” were some features common to authentic blogs. Similar features ranked highly in the 3-word gram model.

(iii) Tokenized Anchors. Anchor text is the text that appears in an HTML link (i.e., between the `<a . . . >` and `` tags.) and is a common link-spamming technique around profitable contexts. We used a bag-of-anchors feature, where anchor text on a page, with multiple word anchors split into individual words, is used. Note that anchor text is frequently used for web page classification, but typically to classifying the target page rather than the one hosting the link. We observed that “comment” and “flickr” were among the highly ranked features for authentic blogs.

(iv) Tokenized URLs. Intuitively, both local and outgoing URLs can be used as effective attributes for splog detection. This is motivated by the fact that many URL tokens in splogs are in profitable contexts. We term these features as bag-of-urls, arrived at by tokenizing URLs using “/” and “.”. Results indicate this can be a useful approach complementing other techniques.

Relational Features

A global model is one that uses some non-local features, i.e., features requiring data beyond the content of Web page under test. We have investigated the use of link distributions to see if splogs can be identified once they place themselves on the blog (web) hyper-link graph. The intuition is that that authentic blogs are very unlikely to link to splogs and that splogs frequently do link to other splogs. We have evaluated this approach

by extending our seed dataset with labeled in-links and out-links, to achieve AUC values of close to 0.85. Interested readers are referred to (Kolari *et al.* 2006; Kolari 2007) for further details.

Future Challenges

Though current techniques work well, the problem of spam detection is an adversarial challenge. In our continuing efforts we are working towards better addressing concept drift and leveraging community and relational features. The problem of spam in social media is now extending well beyond blogs and is quite common in popular social tools like Myspace and Facebook. The nature of these social tools demand additional emphasis on relational techniques, a direction we are exploring as well.

Recognizing Blogosphere communities

Underlying most forms of social media is the concept of a community of people. Identifying these communities, and their possible sub-communities, continues to be a ubiquitous and important task. Most work on community detection (Gibson, Kleinberg, & Raghavan 1998) is based on the analysis of the networks associated with social system.

We have addressed the problem of community detection in the Blogosphere by modeling trust and influence (Kale *et al.* 2007; Kale 2007). Our approach uses the link structure of blog graph to associate sentiments with the links connecting blogs. Such links are manifested as URL that blogger *a* uses in his blog post to refer to blogger *b*’s post. We call this sentiment as *link polarity* and the sign and magnitude of this value is based on the sentiment of text surrounding the link. These polar edges are evidence of bias/trust/distrust between respective blogs. We then use trust propagation models to “spread” the polarity values from a subset of nodes to all possible pairs of nodes. We have evaluated this technique of using trust propagation on polar links in the domain of political blogs by predicting the “like-mindedness” of blogs oriented toward either a Democratic or Republican position. In order to determine a blog’s bias, we compute its trust/distrust score from a seed set of influential blogs (discussed later) and use a hand-labeled dataset to validate our results. More generally, we address the problem of detecting all such nodes that a given node would trust even if it is not directly connected to them.

Link Polarity

The term “link polarity” represents the opinion of the source blog about the destination blog. In order to determine the sentiment based on links, we analyze section of text around the link in the source blog post to determine the sentiment of source blogger about the destination blogger. The text neighboring the link provides direct meaningful insight into blogger *a*’s opinion about

blogger b . Hence, we consider a window of x characters (x is variable parameter for our experimental validations) before and after the link. Note that this set of $2x$ characters does not include html tags.

For our requirements, we do not need to employ complex natural language processing techniques since bloggers typically convey their bias about the post/blog pointed by the link in a straightforward manner. Hence, we use a manually created lexicon of positive and negative oriented words and match the token words from the set of $2x$ characters against this corpus to determine the polarity. Since bloggers frequently use negation of sentimental words to indicate bias about another blog-post (“What B says is not bad”), our corpus includes simple bi-gram patterns of the form “not positive/negative word”.

We adopted the following formula for calculating the link polarity between two posts:

$$\text{Polarity} = (Np - Nn) / (Np + Nn).$$

Np : Number of positively oriented words

Nn : Number of negatively oriented words

Notice that our formula incorporates zero polarity links automatically. The term in the denominator ensures that the polarity is weighed according to the number of words matched against the lexicon. We use summation as the aggregation technique for computing the polarity between two blogs. For our experiments, we choose a domain with a low probability of “off-the-topic” posts within a single blog, hence the notion of summing post-post polarity values to yield a blog-blog polarity value holds.

Trust Propagation

Since blog graphs are not always be densely connected, we will not have the trust scores between many pairs of nodes. Hence, we have investigated techniques for inferring a trust relationship for nodes where one is not explicitly known. Guha et al. (Guha et al. 2004) have used a framework to spread trust in a network bootstrapped by a known set of trusted nodes. Their approach uses a “belief matrix” to represent the initial set of beliefs in the graph. This matrix is generated through a combination of known trust and distrust among a subset of nodes. This matrix is then iteratively modified by using “atomic propagations”. The “atomic propagation” step incorporates direct propagation, co-citation, transpose trust and trust coupling as described in Fig 7. Finally “rounding” technique is applied on the final matrix to produce absolute values of trust between all pair of nodes.

In order to form clusters after the step of trust propagation, we take the approach of averaging trust score for all blog nodes from a predefined set of “influential” nodes belonging to each community. A positive trust score indicates that the blog node belongs to the community influenced by the trusted node of that community. Specifically, we selected top three influential (using the number of inlinks as the measure) Democratic

and Republican bloggers. A positive trust score for a blog from top three Democratic blogs indicates that it belongs to the Democratic cluster and a negative score indicates that it is a Republican blogger.

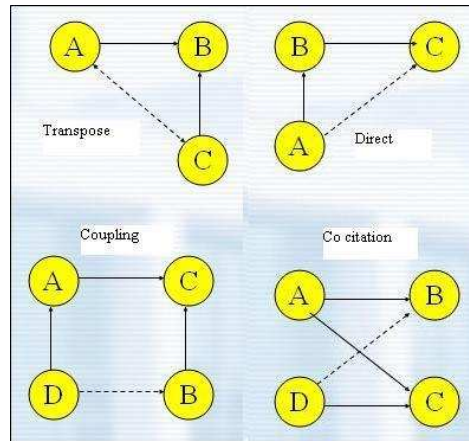


Figure 7: These four graphs represent our atomic propagation patterns. The solid and dotted arrows represent known and inferred trust scores, respectively. The first, for example, indicates that if A and C trust B, then A and C are likely to trust each other.

Experiments

We choose political blogs as our domain; one of the major goals of the experiments was to validate that our proposed approach can correctly classify the blogs into two sets: Republican and Democratic. Through some manual analysis of the political blogs, we observed that the link density among political blogs is reasonably high and hence we could deduce the effectiveness of our approach by running our algorithms over fairly small number of blogs.

Guha’s work argues that “one step distrust” provides the best trust propagation results in their domain of experiments. They propose the notion of “trust and distrust” between two nodes in the graph where the same set of two nodes can trust or distrust each other. The “one step distrust” approach uses the “trust matrix” as the belief matrix. However, we believe that in our domain the initial belief matrix should incorporate both trust and distrust (positive and negative polarities from blog A to blog B). Hence, we use the difference between the trust and distrust matrices as our initial belief matrix. We experimented with various values of the “alpha vector” (the vector used to define the fractional weights for atomic propagation parameters) to confirm that Guha’s conclusion of using the values they proposed $\{0.4, 0.4, 0.1, 0.1\}$ yields best results. Further, Guha et al. recommend performing “atomic propagations” approximately 20 times to get best results; we took the approach of iteratively applying atomic propagations till convergence and our experiments indeed indicate a value close to 20.

Test Dataset Our test dataset consists of a blog graph created from the link structure of Buzzmetrics (Nielsen-Buzzmetric) dataset. The dataset consists of about 14 million weblog posts from three million weblogs collected by Nielsen BuzzMetrics for May 2006. The data is annotated with 1.7 million blog-blog links (buz).

Reference Dataset Adamic and Glance (Adamic & Glance 2005) provided us with a reference dataset of 1490 blogs with a label of *Democratic* and *Republican* for each blog. Their data on political leaning is based on analysis of blog directories and manual labeling and has a timeframe of 2004 presidential elections.

Our test dataset from Buzzmetrics did not provide a classified set of political blogs. Hence, for our experiments we used a snapshot of Buzzmetrics that had a complete overlap with our reference dataset to validate the classification results. The snapshot contained 297 blogs, 1309 blog-blog links and 7052 post-post links. The reference dataset labeled 132 blogs as Republicans and 165 blogs as Democrats (there did not exist any *neutral* labels).

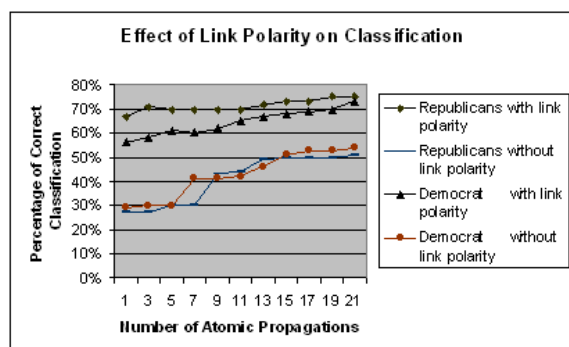


Figure 8: Our experiments show that using polar links for classification yields better results than plain link structure.

Effect of Link Polarity The results in Fig 8 indicate a clear improvement on classifying Republican and Democratic blogs by applying polar weights to links followed by trust propagation. We get a “cold-start” for Democratic blogs and we observe that the overall results are better for Republican blogs than Democratic blogs. The results being better for Republican blogs can be attributed to the observations from (Adamic & Glance 2005) that Republican blogs typically have a higher connectivity than Democratic blogs in the political blogosphere.

We believe that the idea of *polar links* is quite useful and can be applied to multiple domains. The main contribution of this work lies in applying trust propagation models over polar links. We demonstrated one such application in the domain of political blogosphere where we used natural language processing to deduce the link polarity. We would like to emphasize that the

specific techniques to generate polar links is orthogonal to our main contribution and our approach can be easily adapted to different domains for modeling trust and detecting pre-defined communities.

A generative model for the Blogosphere

The blog analysis techniques described in the earlier sections are data-intensive. They require large amounts of blog data that must be obtained by crawling the Blogosphere. Moreover, if the collection is not comprehensive, attention needs to be paid to ensure a representative sample. Once collected, the data must be cleaned to remove spurious spam blogs, or splogs (Kolari *et al.* 2006) and preprocessed in various ways. (Leskovec *et al.* 2007; Shi, Tseng, & Adamic 2007). To overcome the similar difficulty in Web analysis, various graph models have been proposed for the structural and statistical analysis, including the BA model (Barabasi & Albert 1999) and Pennock model (Pennock *et al.* 2002). However, these models are not suitable for generating the blog graphs. While the blog networks resemble many properties of Web graphs, the dynamic nature of the Blogosphere and the evolution of the link structure due to blog readership and social interactions is not well expressed by the existing models.

There are several motivations for developing a generative model for the Blogosphere. First, we hope that such a model might help us understand various aspects of the Blogosphere at an abstract level. Secondly, noticing that portion of the blog graph deviates from the Blogosphere graph can signal that something is amiss. Spam blogs, for example, often form communities whose structural properties are very unlike those of naturally occurring blogs. Third, a generative model can be used to create artificial datasets of varying sizes that simulate portions of the Blogosphere which is often useful for testing and comparing algorithms and systems. For example, testing a model with hidden variables that measures a blog’s influence can benefit from the simulated Blogosphere with different blog graph structures.

Modeling blogger characteristics

Our generative model to construct blog graphs is based on the general characteristics of the bloggers as observed in a recent PEW Internet Survey (Lenhart & Fox 2006) which can be summarized as follows:

- Blog writers are enthusiastic blog readers.
- Most bloggers post infrequently.
- Blog readership can be inferred through blogrolls, a list of links to related or friends’ blogs. Active bloggers are more likely to have a blogroll and follow it regularly.

Our model uses the elements of the existing preferential attachment (Barabasi & Albert 1999) and random attachment models (Chung & Lu 2006). Each blogger is

assumed to be reading, writing or being idle according to the preferential selection of the bloggers. This helps to capture the linking pattern arising in the Blogosphere through *local interactions*. Local interactions refers to the interaction of the bloggers among the other blogs that are generally connected to them either by an inlink or an outlink. We have studied the properties including the degree distributions, degree correlation, clustering coefficient, average degree, reciprocity and the distribution of connected components. To the best of our knowledge, there exist no general models to generate the blog and post networks which possess the properties observed in the real world blogs. Table gives a quick comparison of the properties of the existing Web models and shows the need for a model for Blogosphere.

Defining blog and post networks

Our model of the Blogosphere includes two related networks, one for blogs and one for their posts. The *blog network* (9a) is defined as a network of blogs obtained by collapsing all directed post links between blog posts into directed edges between blogs. Blog networks give a macroscopic view of the blogosphere and help to infer a social network structure, under the assumption that blogs that are “friends” link each other more often. The *post network* (Leskovec *et al.* 2007) (9b) is formed by ignoring the posts’ parent blogs and focus on the link structure among posts only. Each post also has a timestamp of the post associated with it. Post networks give a microscopic view of the blogosphere with details like which post linked to which other post and at what time.

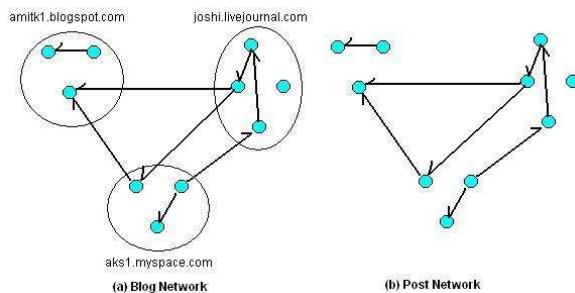


Figure 9: The graph representation for the Blogosphere includes both a blog network and post network.

Design considerations

In designing our model we relied on our experience in analyzing, using and generating blog content as well as an investigation into previous models for social media graphs. We describe how we addressed some of these observations in our model in the following paragraphs.

Linking of new blogs: The new blogger (blog node⁶) may join the existing network by linking to a popular

⁶Blog and blogger are synonymous and the exact meaning is evident from the context

blog (a blog node with high in-degree) or may not link at all. Many Web models using continuous growth (Chung & Lu 2006) also use similar techniques for addition of the new node.

Linking in blogosphere: Generally active bloggers read several posts and tend to link to some of the posts that they read recently but the only “observable behavior” is the creation of a link to the read post (destination). We model this behavior by having the blogger keep track of the recently read blog posts and link to them.

Linking to a post: Leskovec *et al.* (Leskovec *et al.* 2007) observed that any post gathers most of its inlinks within 24 hours of post time. We approximate this behavior by linking to recent posts (within a fixed window) of the visited blog when our blogger visits any blog node.

Blogger neighborhood: Active bloggers tend to subscribe to the well known blogs of interest and read the subscriptions regularly forming blog readership (Lenhart & Fox 2006) (e.g., Bloglines⁷). Hence we see that the blogger interactions are largely concentrated in the blog neighborhood, i.e. nodes connected by either inlinks or outlinks.

Use of emerging tools in Blogosphere: New tools in Blogosphere help to discover popular blogs. For example, blog search engines such as Technorati, social bookmarking and ranking systems like Del.icio.us and Digg, blog classification systems like Feeds that Matter, and trend discovery systems like BlogPulse. The availability and use of these systems mean that blog post reads are not totally random but biased towards the popularity of the blogs which are initially not known to the blogger.

Conversations through comments and trackbacks: The exchange of links among bloggers through comments and trackbacks leads to *higher reciprocity* (i.e., through reciprocal links) in the Blogosphere than the random networks. Bloggers tend to link to the blogs to which they have linked in the past either through comments, trackbacks or general readership. We expect these local interactions to provide for a *higher clustering coefficient* (as observed in Blogosphere) than the random networks.

Activity in the Blogosphere: Not all bloggers are “active” (either reading or writing) at all times. Only a small portion of Blogosphere is *active* with the remainder identified as *idle*. This activity can be approximated by observing the number of links created within a time span. We use a *super linear growth function* to model the activity as defined by Leskovec *et al.* (Leskovec, Kleinberg, & Faloutsos 2007). The outlinks from a blog can be considered as the measure of an *active blog writer* because an active writer will naturally look for more interesting sources to link to. The reverse may not be true that the blogger who reads a lot also writes

⁷<http://www.bloglines.com>

| Property | ER model | BA model | Blogosphere | Our Simulation |
|--------------------------|------------------------|-------------------------|---------------------|-------------------|
| Type | undirected | undirected | directed | directed |
| Degree distribution | poisson | power law | power law | power law |
| Slope [inlinks,outlinks] | N/A | [2.08,-] | [1.66-1.8,1.6-1.75] | [1.7-2.1,1.5-1.6] |
| Avg. degree | constant (for given p) | constant (adds m edges) | increases | increases |
| Component distribution | N/A (undirected) | N/A (undirected) | Power law | Power law |
| Correlation coefficient | - | 1 (fully preferential) | 0.024 (WWE) | 0.1 |
| Avg clustering coeff. | 0.00017 | 0.00018 | 0.0235 (WWE) | 0.0242 |
| Reciprocity | N/A (undirected) | N/A (undirected) | 0.6 (WWE) | 0.6 |

Table 2: This table shows various graph properties for two popular network models (ER and BA), an empirical student of a Blogosphere sample, and our simulated model.

more.

Experiments and Results

The preferential attachment model as proposed by Barabasi (Barabasi & Albert 1999) obtains the power law degree distributions in an *undirected network*. However, this model is not defined for a *directed network*. The model proposed by Pennock et al. (Pennock et al. 2002) captures the random behavior but does not capture the local interactions among nodes in the graph. We use the “alpha preferential attachment” model proposed by Chung et al. (Chung & Lu 2006) to obtain power law degree distributions in a directed graph. We have modified this model to reflect local interaction among the bloggers by using preferential attachment among neighboring nodes. The details of our algorithm can be found in (Karandikar 2007).

Part of our evaluation is done by comparing the distinguishing properties of the real blog graphs with the results of our simulation. These properties were verified against two large blog datasets available for researchers namely WWE⁸ 2006 and ICWSM⁹ 2007.

Figure 10 shows the power law curve for inlink distribution observed in the blog network of the “simulated” blogosphere. Similarly, figure 11 compares the distribution of the strongly connected components (SCC) in the simulation and the blogosphere.

Figure 12 shows the scatter plot for in-degree and out-degree correlations in the simulated Blogosphere. The plot shows the low degree correlation as observed in the real Blogosphere by Leskovec et al (Leskovec et al. 2007).

Being able to generate synthetic blog data is useful for testing and evaluating algorithms for analyzing and extracting information from the Blogosphere. The utility, however, depends on the generated graphs being similar to the actual ones in key properties. We have created a model that accounts for several key features that influence how the Blogosphere grows and its result-

| Blog network | ICWSM | WWE | Simulation |
|-------------------------|---------|-----------|------------|
| Total blogs | 159,036 | 650,660 | 650,000 |
| Blog-blog links | 435,675 | 1,893,187 | 1,451,069 |
| Unique links | 245,840 | 648,566 | 1,158,803 |
| Average degree | 5.47 | 5.73 | 4.47 |
| Indegree distribution | -2.07 | -2.0 | -1.71 |
| Outdegree distribution | -1.51 | -1.6 | -1.76 |
| Degree correlation | 0.056 | 0.002 | 0.10 |
| Diameter | 14 | 12 | 6 |
| Largest WCC size | 96,806 | 263,515 | 617,044 |
| Largest SCC size | 4,787 | 4,614 | 72,303 |
| Clustering coefficients | 0.04429 | 0.0235 | 0.0242 |
| Percent Reciprocity | 3.03 | 0.6838 | 0.6902 |

Table 3: This table compares a simulated blog network with two datasets based on blogs harvested from the Web using a number of standard graph metrics.

ing structure. By selecting appropriate parameters for our model, we can generate graphs that more closely approximate observed blogosphere network properties than previous network models.

Conclusion

Social media systems are increasingly important on the Web and account for a significant fraction of new content. The various kinds of social media are alike in that they all have rich underlying network structures that provide metadata and context that can help when extracting information from their content. We have described some initial results from ongoing work that is focused on extracting, modeling and exploiting this structural information from the underlying networks.

As the Web continues to evolve, we expect that the ways people interact with it, as content consumers as well as content providers, will also change. The result, however, will continue to represent an interesting and extravagant mixture of underlying networks – networks of individuals, groups, documents, opinions, beliefs, advertisements, and scams. These interwoven networks present new opportunities and challenges for extracting information and knowledge from them.

⁸a dataset developed for the Workshop on the Weblogging Ecosystem held at the 2006 World Wide Web Conference

⁹International Conference on Weblogs and Social Media: <http://www.icwsm.org/data.html>

| Post network | ICWSM | WWE | Simulation |
|-------------------------|-----------|-----------|------------|
| Total posts | 1,035,361 | 1,527,348 | 1,380,341 |
| Post-post links | 1,354,610 | 1,863,979 | 1,451,069 |
| Unique links | 458,950 | 1,195,072 | 1,442,525 |
| Avg post outlinks | 1.30 | 1.22 | 1.051 |
| Average degree | 2.62 | 2.44 | 2.10 |
| Indegree distribution | -1.26 | -2.6 | -2.54 |
| Outdegree distribution | -1.03 | -2.04 | -2.04 |
| Degree correlation | -0.113 | -0.035 | -0.006 |
| Diameter | 20 | 24 | 12 |
| Largest WCC size | 134,883 | 262,919 | 1,068,755 |
| Largest SCC size | 14 | 13 | 3 |
| Clustering coefficients | 0.0026 | 0.00135 | 0.00011 |
| Percent Reciprocity | 0.029 | 0.021 | 0.01 |

Table 4: Comparison of post network properties of datasets and simulation

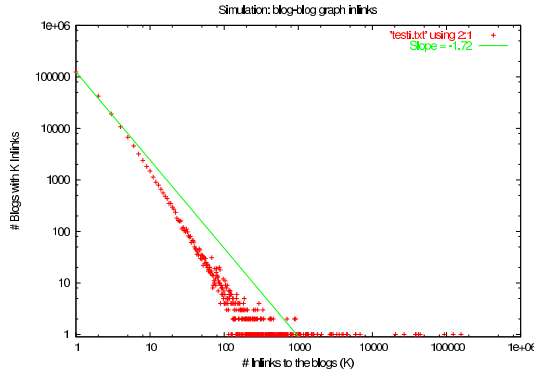


Figure 10: Simulation: Blog inlinks distribution

Acknowledgements

We thank Dr. James Mayfield, Justin Martineau and Sandeep Balijepalli for their contributions to the work on sentiment detection. Partial support for this research was provided by the National Science Foundation (awards ITR-IIS-0325172 and NSF-ITR-IDM-0219649) and I.B.M.

References

- Adamic, L. A., and Glance, N. 2005. The political blogosphere and the 2004 U.S. election: divided they blog. In *LinkKDD '05: Proceedings of the 3rd international workshop on Link discovery*, 36–43. New York, NY, USA: ACM Press.
- Balijepalli, S. 2007. Blogvox2: A modular domain independent sentiment analysis system. Master's thesis, University of Maryland, Baltimore County.
- Barabasi, A.-L., and Albert, R. 1999. Emergence of scaling in random networks. *Science* 286:509.
- Boser, B. E.; Guyon, I. M.; and Vapnik, V. N. 1992.

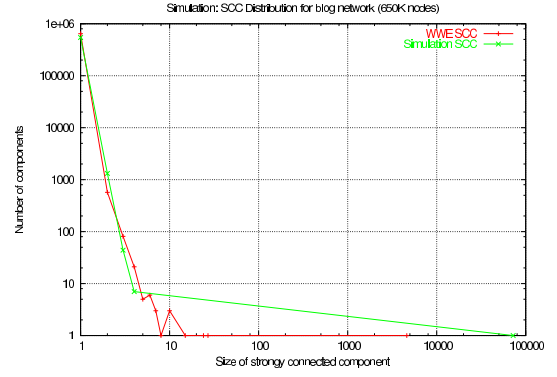


Figure 11: Distribution of SCC in blog network (Simulation and Blogosphere)

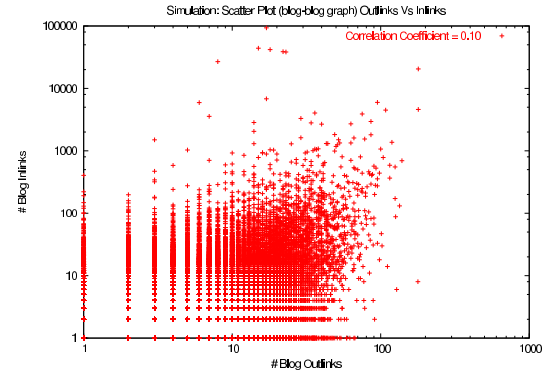


Figure 12: Simulated blog network scatter plot: Outdegree vs. Indegree

A training algorithm for optimal margin classifiers. In *COLT '92: Proceedings of the fifth annual workshop on Computational learning theory*, 144–152. New York: ACM Press.

Buzzmetrics dataset.

Chung, F., and Lu, L. 2006. *Complex Graphs and Networks (Cbms Regional Conference Series in Mathematics)*. Boston, MA, USA: American Mathematical Society.

Finin, T. 2006. Splog software from hell. [Online; accessed 31-August-2006; <http://ebiquity.umbc.edu/blogger/splog-software-from-hell/>].

Gibson, D.; Kleinberg, J.; and Raghavan, P. 1998. Inferring web communities from link topology. In *HYPERTEXT '98: Proceedings of the ninth ACM conference on Hypertext and hypermedia*, 225–234. New York: ACM Press.

Guha, R.; Kumar, R.; Raghavan, P.; and Tomkins, A.

2004. Propagation of trust and distrust. In *WWW '04: Proceedings of the 13th international conference on World Wide Web*, 403–412. New York, NY, USA: ACM Press.
- Java, A.; Kolari, P.; Finin, T.; Joshi, A.; Martineau, J.; and Mayfield, J. 2007a. The BlogVox Opinion Retrieval System. In *Proceedings of the Fifteenth Text REtrieval Conference (TREC 2006)*.
- Java, A.; Kolari, P.; Finin, T.; Joshi, A.; and Oates, T. 2007b. Feeds That Matter: A Study of Bloglines Subscriptions. In *Proceedings of the International Conference on Weblogs and Social Media (ICWSM 2007)*. Computer Science and Electrical Engineering. To Appear.
- Joachims, T. 1998. Text categorization with support vector machines: Learning with many relevant features. In *ECML '98: Proceedings of the 10th European Conference on Machine Learning*, 137–142. London, UK: Springer-Verlag.
- Kale, A.; Karandikar, A.; Kolari, P.; Java, A.; Joshi, A.; and Finin, T. 2007. Modeling Trust and Influence in the Blogosphere Using Link Polarity. In *Proceedings of the International Conference on Weblogs and Social Media (ICWSM 2007)*. Short Paper.
- Kale, A. 2007. Modeling trust and influence in blogosphere using link polarity. Master's thesis, University of Maryland, Baltimore County.
- Karandikar, A. 2007. Generative Model To Construct Blog and Post Networks In Blogosphere. Master's thesis, University of Maryland, Baltimore County.
- Kempe, D.; Kleinberg, J. M.; and Tardos, É. 2003. Maximizing the spread of influence through a social network. In *proceedings of the Ninth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 137–146.
- Kempe, D.; Kleinberg, J. M.; and Tardos, É. 2005. Influential nodes in a diffusion model for social networks. In *Proceedings of the 32nd International Colloquium on Automata, Languages and Programming*, 1127–1138.
- Kolari, P.; Java, A.; Finin, T.; Oates, T.; and Joshi, A. 2006. Detecting Spam blogs: A machine learning approach. In *Proceedings of the 21st National Conference on Artificial Intelligence*. AAAI Press.
- Kolari, P.; Finin, T.; and Joshi, A. 2006. SVMs for the blogosphere: Blog identification and splog detection. In *AAAI Spring Symposium on Computational Approaches to Analyzing Weblogs*. AAAI Press.
- Kolari, P.; Java, A.; and Finin, T. 2006. Characterizing the splogosphere. In *WWW 2006, 3rd Annual Workshop on the Weblogging Ecosystem: Aggregation, Analysis and Dynamics*.
- Kolari, P. 2007. *Detecting Spam Blogs: An Adaptive Online Approach*. Ph.D. Dissertation, University of Maryland, Baltimore County.
- Lenhart, A., and Fox, S. 2006. Bloggers: A portrait of the internet's new storytellers.
- Leskovec, J.; McGlohon, M.; Faloutsos, C.; Glance, N.; and Hurst, M. 2007. Cascading behavior in large blog graphs. In *SIAM International Conference on Data Mining (SDM 2007)*.
- Leskovec, J.; Kleinberg, J.; and Faloutsos, C. 2007. Graph evolution: Densification and shrinking diameters. *ACM Transactions on knowledge discovery from data* 1:1.
- Martineau, J.; Java, A.; Kolari, P.; Finin, T.; Joshi, A.; and Mayfield, J. 2007. Blogvox: Learning sentiment classifiers. In *Proceedings of the 22nd National Conference on Artificial Intelligence*, 1888–1889. AAAI Press. Student Abstract.
- Mishne, G. 2007. *Applied Text Analytics for Blogs*. Ph.D. Dissertation, University of Amsterdam.
- NielsenBuzzmetric. <http://www.nielsenbuzzmetrics.com>.
- Pennock, D. M.; Flake, G. W.; Lawrence, S.; Glover, E. J.; and Giles, C. L. 2002. Winners don't take all: Characterizing the competition for links on the web. *Proceedings of the National Academy of Sciences* 99(8):5207–5211.
- Shi, X.; Tseng, B.; and Adamic, L. 2007. Looking at the blogosphere topology through different lenses. In *Proceedings of the International Conference on Weblogs and Social Media (ICWSM 2007)*.

Biographical sketches

Tim Finin is a Professor of Computer Science and Electrical Engineering at the University of Maryland Baltimore County (UMBC). He has over 35 years of experience in the applications of AI to problems in information systems, intelligent interfaces and robotics. He holds degrees from MIT and the University of Illinois and has held positions at Unisys, the University of Pennsylvania, and the MIT AI Laboratory.



Anupam Joshi is a UMBC Professor with research interests in the broad area of networked computing and intelligent systems. He currently serves on the editorial board of the International Journal of the Semantic Web and Information.



Pranam Kolari is a member of the technical staff at Yahoo! Applied Research. He received a Ph.D. in Computer Science. His dissertation was focused on spam blog detection, with tools developed in use both by academia and industry. He has active research interest in internal corporate blogs, the Semantic Web and blog analytics.



Akshay Java is a UMBC PhD student. His dissertation is on identifying influence and opinions in social media. His research interests include blog analytics, information retrieval, natural language processing and the Semantic Web.



Anubhav Kale received a M.S. degree in Computer Science from UMBC in May 2007. His thesis research demonstrated the effectiveness of detecting sentiment associated with links between blog posts and using this to enhance blog community recognition algorithms. He is currently a software engineer at Microsoft.



Amit Karandikar received a M.S. degree in Computer Science from UMBC in May 2007. His thesis research produced a generative model for the blogosphere which modeled both the reading and writing activities of bloggers. He is currently a software engineer at Microsoft.

