

Chapter 21

Web Mining — Concepts, Applications, and Research Directions

Jaideep Srivastava, Prasanna Desikan, Vipin Kumar

Web mining is the application of data mining techniques to extract knowledge from web data, including web documents, hyperlinks between documents, usage logs of web sites, etc. A panel organized at ICTAI 1997 (Srivastava and Mobasher 1997) asked the question “Is there anything distinct about web mining (compared to data mining in general)?” While no definitive conclusions were reached then, the tremendous attention on web mining in the past five years, and a number of significant ideas that have been developed, have certainly answered this question in the affirmative in a big way. In addition, a fairly stable community of researchers interested in the area has been formed, largely through the successful series of WebKDD workshops, which have been held annually in conjunction with the ACM SIGKDD Conference since 1999 (Masand and Spiliopoulou 1999; Kohavi, Spiliopoulou, and Srivastava 2001; Kohavi, Masand, Spiliopoulou, and Srivastava 2001; Masand, Spiliopoulou, Srivastava, and Zaiane 2002), and the web analytics workshops, which have been held in conjunction with the SIAM data mining conference (Ghosh and Srivastava 2001a, b). A good survey of the research in the field (through 1999)

is provided by Kosala and Blockeel (2000) and Madria, Bhowmick, Ng, and Lim (1999).

Two different approaches were taken in initially defining web mining. First was a “process-centric view,” which defined web mining as a sequence of tasks (Etzioni 1996). Second was a “data-centric view,” which defined web mining in terms of the types of web data that was being used in the mining process (Cooley, Srivastava, and Mobasher 1997). The second definition has become more acceptable, as is evident from the approach adopted in most recent papers (Madria, Bhowmick, Ng, and Lim 1999; Borges and Levene 1998; Kosala and Blockeel 2000) that have addressed the issue. In this chapter we follow the data-centric view of web mining which is defined as follows,

Web mining is the application of data mining techniques to extract knowledge from web data, i.e. web content, web structure, and web usage data.

The attention paid to web mining, in research, software industry, and web-based organization, has led to the accumulation of significant experience. It is our goal in this chapter to capture them in a systematic manner, and identify directions for future research.

The rest of this chapter is organized as follows: In section 21.1 we provide a taxonomy of web mining, in section 21.2 we summarize some of the key concepts in the field, and in section 21.3 we describe successful applications of web mining. In section 21.4 we present some directions for future research, and in section 21.5 we conclude the chapter.

21.1 Web Mining Taxonomy

Web mining can be broadly divided into three distinct categories, according to the kinds of data to be mined. Figure 21.1 shows the taxonomy.

21.1.1 Web Content Mining

Web content mining is the process of extracting useful information from the contents of web documents. Content data is the collection of facts a web page is designed to contain. It may consist of text, images, audio, video, or structured records such as lists and tables. Application of text mining to web content has been the most widely researched. Issues addressed in text mining include topic discovery and tracking, extracting association patterns, clustering of web documents and classification of web pages. Research activities on this topic have drawn heavily on techniques developed in other disciplines such as Information Retrieval (IR) and Natural Language Processing (NLP). While

there exists a significant body of work in extracting knowledge from images in the fields of image processing and computer vision, the application of these techniques to web content mining has been limited.

21.1.2 Web Structure Mining

The structure of a typical web graph consists of web pages as nodes, and hyperlinks as edges connecting related pages. Web structure mining is the process of discovering structure information from the web. This can be further divided into two kinds based on the kind of structure information used.

Hyperlinks

A hyperlink is a structural unit that connects a location in a web page to a different location, either within the same web page or on a different web page. A hyperlink that connects to a different part of the same page is called an *intra-document hyperlink*, and a hyperlink that connects two different pages is called an *inter-document hyperlink*. There has been a significant body of work on hyperlink analysis, of which Desikan, Srivastava, Kumar, and Tan (2002) provide an up-to-date survey.

Document Structure

In addition, the content within a Web page can also be organized in a tree-structured format, based on the various HTML and XML tags within the page. Mining efforts here have focused on automatically extracting document object model (DOM) structures out of documents (Wang and Liu 1998; Moh, Lim, and Ng 2000).

21.1.3 Web Usage Mining

Web usage mining is the application of data mining techniques to discover interesting usage patterns from web usage data, in order to understand and better serve the needs of web-based applications (Srivastava, Cooley, Deshpande, and Tan 2000). Usage data captures the identity or origin of web users along with their browsing behavior at a web site. web usage mining itself can be classified further depending on the kind of usage data considered:

Web Server Data

User logs are collected by the web server and typically include IP address, page reference and access time.

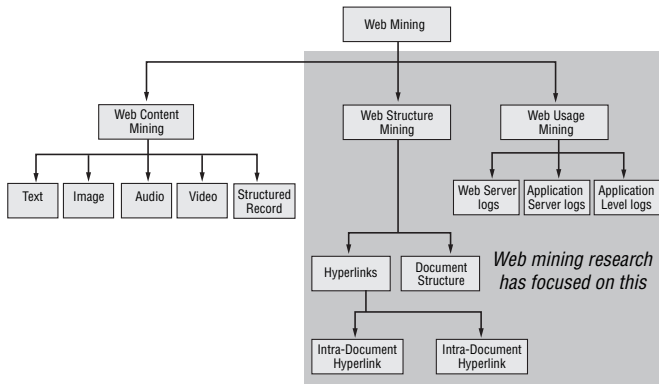


Figure 21.1: Web mining Taxonomy

Application Server Data

Commercial application servers such as Weblogic,^{1,2} StoryServer,³ have significant features to enable E-commerce applications to be built on top of them with little effort. A key feature is the ability to track various kinds of business events and log them in application server logs.

Application Level Data

New kinds of events can be defined in an application, and logging can be turned on for them — generating histories of these events.

It must be noted, however, that many end applications require a combination of one or more of the techniques applied in the above the categories.

21.2 Key Concepts

In this section we briefly describe the new concepts introduced by the web mining research community.

21.2.1 Ranking Metrics—for Page Quality and Relevance

Searching the web involves two main steps: *Extracting the pages relevant to a query* and *ranking them according to their quality*. Ranking is important as it

¹<http://www.bea.com/products/weblogic/server/index.shtml>

²<http://www.bvportal.com/>.

³http://www.cio.com/sponsors/110199_vignette_story2.html.

helps the user look for “quality” pages that are relevant to the query. Different metrics have been proposed to rank web pages according to their quality. We briefly discuss two of the prominent ones.

PageRank

PageRank is a metric for ranking hypertext documents based on their quality. Page, Brin, Motwani, and Winograd (1998) developed this metric for the popular search engine Google⁴ (Brin and Page 1998). The key idea is that a page has a high rank if it is pointed to by many highly ranked pages. So, the rank of a page depends upon the ranks of the pages pointing to it. This process is done iteratively until the rank of all pages are determined. The rank of a page p can be written as:

$$PR(p) = d/n + (1 - d) \sum_{(q,p) \in G} \left(\frac{PR(q)}{Outdegree(q)} \right)$$

Here, n is the number of nodes in the graph and $OutDegree(q)$ is the number of hyperlinks on page q . Intuitively, the approach can be viewed as a stochastic analysis of a random walk on the web graph. The first term in the right hand side of the equation is the probability that a random web surfer arrives at a page p by typing the URL or from a bookmark; or may have a particular page as his/her homepage. Here d is the probability that the surfer chooses a URL directly, rather than traversing a link⁵ and $1 - d$ is the probability that a person arrives at a page by traversing a link. The second term in the right hand side of the equation is the probability of arriving at a page by traversing a link.

Hubs and Authorities

Hubs and authorities can be viewed as “fans” and “centers” in a bipartite core of a web graph, where the “fans” represent the hubs and the “centers” represent the authorities. The hub and authority scores computed for each web page indicate the extent to which the web page serves as a hub pointing to good authority pages or as an authority on a topic pointed to by good hubs. The scores are computed for a set of pages related to a topic using an iterative procedure called HITS (Kleinberg 1999). First a query is submitted to a search engine and a set of relevant documents is retrieved. This set, called the “root set,” is then expanded by including web pages that point to those in the “root set” and are pointed by those in the “root set.” This new set is called the “base set.” An adjacency matrix, A is formed such that if there exists at least one

⁴<http://www.google.com>.

⁵The parameter d , called the dampening factor, is usually set between 0.1 and 0.2 (Brin and Page 1998).

hyperlink from page i to page j , then $A_{i,j} = 1$, otherwise $A_{i,j} = 0$. HITS algorithm is then used to compute the hub and authority scores for these set of pages.

There have been modifications and improvements to the basic page rank and hubs and authorities approaches such as SALSA (Lempel and Moran 2000), topic sensitive page rank, (Haveliwala 2002) and web page reputations (Mendelzon and Rafiei 2000). These different hyperlink based metrics have been discussed by Desikan, Srivastava, Kumar, and Tan (2002).

21.2.2 Robot Detection and Filtering—Separating Human and Nonhuman Web Behavior

Web robots are software programs that automatically traverse the hyperlink structure of the web to locate and retrieve information. The importance of separating robot behavior from human behavior prior to building user behavior models has been illustrated by Kohavi (2001). First, e-commerce retailers are particularly concerned about the unauthorized deployment of robots for gathering business intelligence at their web sites. Second, web robots tend to consume considerable network bandwidth at the expense of other users. Sessions due to web robots also make it difficult to perform click-stream analysis effectively on the web data. Conventional techniques for detecting web robots are based on identifying the IP address and user agent of the web clients. While these techniques are applicable to many well-known robots, they are not sufficient to detect camouflaged and previously unknown robots. Tan and Kumar (2002) proposed a classification based approach that uses the navigational patterns in click-stream data to determine if it is due to a robot. Experimental results have shown that highly accurate classification models can be built using this approach. Furthermore, these models are able to discover many camouflaged and previously unidentified robots.

21.2.3 Information Scent—Applying Foraging Theory to Browsing Behavior

Information scent is a concept that uses the snippets of information present around the links in a page as a “scent” to evaluate the quality of content of the page it points to, and the cost of accessing such a page (Chi, Pirolli, Chen, and Pitkow 2001). The key idea is to model a user at a given page as “foraging” for information, and following a link with a stronger “scent.” The “scent” of a path depends on how likely it is to lead the user to relevant information, and is determined by a network flow algorithm called spreading activation. The snippets, graphics, and other information around a link are called “proximal cues.”

The user's desired information need is expressed as a weighted keyword vector. The similarity between the proximal cues and the user's information need is computed as "proximal scent." With the proximal cues from all the links and the user's information need vector, a "proximal scent matrix" is generated. Each element in the matrix reflects the extent of similarity between the link's proximal cues and the user's information need. If enough information is not available around the link, a "distal scent" is computed with the information about the link described by the contents of the pages it points to. The proximal scent and the distal scent are then combined to give the scent matrix. The probability that a user would follow a link is then decided by the scent or the value of the element in the scent matrix.

21.2.4 User Profiles — Understanding How Users Behave

The web has taken user profiling to new levels. For example, in a "brick-and-mortar" store, data collection happens only at the checkout counter, usually called the "point-of-sale." This provides information only about the final outcome of a complex human decision making process, with no direct information about the process itself. In an on-line store, the complete click-stream is recorded, which provides a detailed record of every action taken by the user, providing a much more detailed insight into the decision making process. Adding such behavioral information to other kinds of information about users, for example demographic, psychographic, and so on, allows a comprehensive user profile to be built, which can be used for many different purposes (Masand, Spiliopoulou, Srivastava, and Zaiane 2002).

While most organizations build profiles of user behavior limited to visits to their own sites, there are successful examples of building web-wide behavioral profiles such as Alexa Research⁶ and DoubleClick⁷. These approaches require browser cookies of some sort, and can provide a fairly detailed view of a user's browsing behavior across the web.

21.2.5 Interestingness Measures — When Multiple Sources Provide Conflicting Evidence

One of the significant impacts of publishing on the web has been the close interaction now possible between authors and their readers. In the preweb era, a reader's level of interest in published material had to be inferred from indirect measures such as buying and borrowing, library checkout and renewal, opinion surveys, and in rare cases feedback on the content. For material published on the web it is possible to track the click-stream of a reader to observe the exact

⁶<http://www.alexa.com>.

⁷<http://www.doubleclick.com/>.

path taken through on-line published material. We can measure times spent on each page, the specific link taken to arrive at a page and to leave it, etc. Much more accurate inferences about readers' interest in content can be drawn from these observations. Mining the user click-stream for user behavior, and using it to adapt the "look-and-feel" of a site to a reader's needs was first proposed by Perkowitiz and Etzioni (1999).

While the usage data of any portion of a web site can be analyzed, the most significant, and thus "interesting," is the one where the usage pattern differs significantly from the link structure. This is so because the readers' behavior, reflected by web usage, is very different from what the author would like it to be, reflected by the structure created by the author. Treating knowledge extracted from structure data and usage data as evidence from independent sources, and combining them in an evidential reasoning framework to develop measures for interestingness has been proposed by several authors (Padmanabhan and Tuzhilin 1998, Cooley 2000).

21.2.6 Preprocessing—Making Web Data Suitable for Mining

In the panel discussion referred to earlier (Srivastava and Mobasher 1997), preprocessing of web data to make it suitable for mining was identified as one of the key issues for web mining. A significant amount of work has been done in this area for web usage data, including user identification and session creation (Cooley, Mobasher, and Srivastava 1999), robot detection and filtering (Tan and Kumar 2002), and extracting usage path patterns (Spiliopoulou 1999). Cooley's Ph.D. dissertation (Cooley 2000) provides a comprehensive overview of the work in web usage data preprocessing.

Preprocessing of web structure data, especially link information, has been carried out for some applications, the most notable being Google style web search (Brin and Page 1998). An up-to-date survey of structure preprocessing is provided by Desikan, Srivastava, Kumar, and Tan (2002).

21.2.7 Identifying Web Communities of Information Sources

The web has had tremendous success in building communities of users and information sources. Identifying such communities is useful for many purposes. Gibson, Kleinberg, and Raghavan (1998) identified web communities as "a core of central authoritative pages linked together by hub pages. Their approach was extended by Ravi Kumar and colleagues (Kumar, Raghavan, Rajagopalan, and Tomkins 1999) to discover emerging web communities while crawling. A different approach to this problem was taken by Flake, Lawrence,

and Giles (2000) who applied the “maximum-flow minimum cut model” (Jr and Fulkerson 1956) to the web graph for identifying “web communities.” Imafuji and Kitsuregawa (2002) compare HITS and the maximum flow approaches and discuss the strengths and weakness of the two methods. Reddy and Kitsuregawa (2002) propose a dense bipartite graph method, a relaxation to the complete bipartite method followed by HITS approach, to find web communities. A related concept of “friends and neighbors” was introduced by Adamic and Adar (2003). They identified a group of individuals with similar interests, who in the cyber-world would form a “community.” Two people are termed “friends” if the similarity between their web pages is high. Similarity is measured using features such as text, out-links, in-links and mailing lists.

21.2.8 Online Bibliometrics

With the web having become the fastest growing and most up to date source of information, the research community has found it extremely useful to have online repositories of publications. Lawrence observed (Lawrence 2001) that having articles online makes them more easily accessible and hence more often cited than articles that are offline. Such online repositories not only keep the researchers updated on work carried out at different centers, but also makes the interaction and exchange of information much easier.

With such information stored in the web, it becomes easier to point to the most frequent papers that are cited for a topic and also related papers that have been published earlier or later than a given paper. This helps in understanding the state of the art in a particular field, helping researchers to explore new areas. Fundamental web mining techniques are applied to improve the search and categorization of research papers, and citing related articles. Some of the prominent digital libraries are Science Citation Index (SCI),⁸ the Association for Computing Machinery’s ACM portal,⁹ the Scientific Literature Digital Library (CiteSeer),¹⁰ and the DBLP Bibliography.¹¹

21.2.9 Visualization of the World Wide Web

Mining web data provides a lot of information, which can be better understood with visualization tools. This makes concepts clearer than is possible with pure textual representation. Hence, there is a need to develop tools that provide a graphical interface that aids in visualizing results of web mining.

⁸<http://www.isinet.com/isi/products/citation/sci/>.

⁹<http://portal.acm.org/portal.cfm>.

¹⁰<http://citeseer.nj.nec.com/cs>

¹¹<http://www.informatik.uni-trier.de/ley/db/>.

Analyzing the web log data with visualization tools has evoked a lot of interest in the research community. Chi, Pitkow, Mackinlay, Pirolli, Gossweiler, and Card (1998) developed a web ecology and evolution visualization (WEEV) tool to understand the relationship between web content, web structure and web usage over a period of time. The site hierarchy is represented in a circular form called the “Disk Tree” and the evolution of the web is viewed as a “Time Tube.” Cadez, Heckerman, Meek, Smyth, and White (2000) present a tool called WebCANVAS that displays clusters of users with similar navigation behavior. Prasetyo, Pramudiono, Takahashi, Toyoda, and Kitsuregawa developed Naviz, an interactive web log visualization tool that is designed to display the user browsing pattern on the web site at a global level, and then display each browsing path on the pattern displayed earlier in an incremental manner. The support of each traversal is represented by the thickness of the edge between the pages. Such a tool is very useful in analyzing user behavior and improving web sites.

21.3 Prominent Applications

Excitement about the web in the past few years has led to the web applications being developed at a much faster rate in the industry than research in web related technologies. Many of these are based on the use of web mining concepts, even though the organizations that developed these applications, and invented the corresponding technologies, did not consider it as such. We describe some of the most successful applications in this section. Clearly, realizing that these applications use web mining is largely a retrospective exercise. For each application category discussed below, we have selected a prominent representative, purely for exemplary purposes. This in no way implies that all the techniques described were developed by that organization alone. On the contrary, in most cases the successful techniques were developed by a rapid “copy and improve” approach to each other’s ideas.

21.3.1 Personalized Customer Experience in B2C E-commerce—Amazon.com

Early on in the life of Amazon.com,¹² its visionary CEO Jeff Bezos observed,

“In a traditional (brick-and-mortar) store, the main effort is in getting a customer to the store. Once a customer is in the store they are likely to make a purchase — since the cost of going to another store is high — and thus the marketing budget (focused on getting

¹²<http://www.amazon.com>.

the customer to the store) is in general much higher than the in-store customer experience budget (which keeps the customer in the store). In the case of an on-line store, getting in or out requires exactly one click, and thus the main focus must be on customer experience in the store.”¹³

This fundamental observation has been the driving force behind Amazon’s comprehensive approach to personalized customer experience, based on the mantra “a personalized store for every customer” (Morphy 2001). A host of web mining techniques, such as associations between pages visited and click-path analysis are used to improve the customer’s experience during a “store visit.” Knowledge gained from web mining is the key intelligence behind Amazon’s features such as “instant recommendations,” “purchase circles,” “wish-lists,” etc.

21.3.2 Web Search—Google

Google¹⁴ is one of the most popular and widely used search engines. It provides users access to information from over 2 billion web pages that it has indexed on its server. The quality and quickness of the search facility makes it the most successful search engine. Earlier search engines concentrated on web content alone to return the relevant pages to a query. Google was the first to introduce the importance of the link structure in mining information from the web. PageRank, which measures the importance of a page, is the underlying technology in all Google search products, and uses structural information of the web graph to return high quality results.

The Google toolbar is another service provided by Google that seeks to make search easier and informative by providing additional features such as highlighting the query words on the returned web pages. The full version of the toolbar, if installed, also sends the click-stream information of the user to Google. The usage statistics thus obtained are used by Google to enhance the quality of its results. Google also provides advanced search capabilities to search images and find pages that have been updated within a specific date range. Built on top of Netscape’s Open Directory project, Google’s web directory provides a fast and easy way to search within a certain topic or related topics.

The advertising program introduced by Google targets users by providing advertisements that are relevant to a search query. This does not bother users with irrelevant ads and has increased the clicks for the advertising companies

¹³The truth of this fundamental insight has been borne out by the phenomenon of “shopping cart abandonment,” which happens frequently in on-line stores, but practically never in a brick-and-mortar one.

¹⁴<http://www.google.com>.

by four to five times. According to BtoB, a leading national marketing publication, Google was named a top 10 advertising property in the Media Power 50 that recognizes the most powerful and targeted business-to-business advertising outlets¹⁵.

One of the latest services offered by Google is Google News¹⁶. It integrates news from the online versions of all newspapers and organizes them categorically to make it easier for users to read “the most relevant news.” It seeks to provide latest information by constantly retrieving pages from news site worldwide that are being updated on a regular basis. The key feature of this news page, like any other Google service, is that it integrates information from various web news sources through purely algorithmic means, and thus does not introduce any human bias or effort. However, the publishing industry is not very convinced about a fully automated approach to news distillation (Springer 2002).

21.3.3 Web-Wide Tracking—DoubleClick

“Web-wide tracking,” i.e. tracking an individual across all sites he visits, is an intriguing and controversial technology. It can provide an understanding of an individual’s lifestyle and habits to a level that is unprecedented, which is clearly of tremendous interest to marketers. A successful example of this is DoubleClick Inc.’s DART ad management technology¹⁷. DoubleClick serves advertisements, which can be targeted on demographic or behavioral attributes, to the end-user on behalf of the client, i.e. the web site using DoubleClick’s service. Sites that use DoubleClick’s service are part of The DoubleClick Network and the browsing behavior of a user can be tracked across all sites in the network, using a cookie. This makes DoubleClick’s ad targeting to be based on very sophisticated criteria. Alexa Research¹⁸ has recruited a panel of more than 500,000 users, who have voluntarily agreed to have their every click tracked, in return for some freebies. This is achieved through having a browser bar that can be downloaded by the panelist from Alexa’s website, which gets attached to the browser and sends Alexa a complete click-stream of the panelist’s web usage. Alexa was purchased by Amazon for its tracking technology.

Clearly web-wide tracking is a very powerful idea. However, the invasion of privacy it causes has not gone unnoticed, and both Alexa/Amazon and DoubleClick have faced very visible lawsuits.¹⁹ Microsoft’s Passport²⁰ technology

¹⁵<http://www.google.com/press/pressrel/b2b.html>

¹⁶<http://news.google.com>

¹⁷<http://www.doubleclick.com/dartinfo/>

¹⁸<http://www.alexa.com>.

¹⁹See <http://www.wired.com/news/business/0,1367,36434,00.html>.

²⁰<http://www.microsoft.com/netservices/passport/>.

also falls into this category. The value of this technology in applications such as cyber-threat analysis and homeland defense is quite clear, and it might be only a matter of time before these organizations are asked to provide information to law enforcement agencies.

21.3.4 Understanding Web Communities—AOL

One of the biggest successes of America Online (AOL)²¹ has been its sizeable and loyal customer base. A large portion of this customer base participates in various AOL communities, which are collections of users with similar interests. In addition to providing a forum for each such community to interact amongst themselves, AOL provides them with useful information and services. Over time these communities have grown to be well-visited water-holes for AOL users with shared interests. Applying web mining to the data collected from community interactions provides AOL with a very good understanding of its communities, which it has used for targeted marketing through advertisements and e-mail solicitation. Recently, it has started the concept of “community sponsorship,” whereby an organization, say Nike, may sponsor a community called “Young Athletic TwentySomethings.” In return, consumer survey and new product development experts of the sponsoring organization get to participate in the community, perhaps without the knowledge of other participants. The idea is to treat the community as a highly specialized focus group, understand its needs and opinions on new and existing products, and also test strategies for influencing opinions.

21.3.5 Understanding Auction Behavior—eBay

As individuals in a society where we have many more things than we need, the allure of exchanging our useless stuff for some cash, no matter how small, is quite powerful. This is evident from the success of flea markets, garage sales and estate sales. The genius of eBay’s²² founders was to create an infrastructure that gave this urge a global reach, with the convenience of doing it from one’s home PC. In addition, it popularized auctions as a product selling and buying mechanism and provides the thrill of gambling without the trouble of having to go to Las Vegas. All of this has made eBay as one of the most successful businesses of the internet era. Unfortunately, the anonymity of the web has also created a significant problem for eBay auctions, as it is impossible to distinguish real bids from fake ones. eBay is now using web mining techniques to analyze bidding behavior to determine if a bid is fraudulent (Colet

²¹ See <http://www.aol.com>.

²² <http://www.ebay.com>.

2002). Recent efforts are geared towards understanding participants' bidding behaviors/patterns to create a more efficient auction market.

21.3.6 Personalized Portal for the Web—MyYahoo

Yahoo²³ was the first to introduce the concept of a “personalized portal,” i.e. a web site designed to have the look-and-feel and content personalized to the needs of an individual end-user. This has been an extremely popular concept and has led to the creation of other personalized portals such as Yodlee²⁴ for private information like bank and brokerage accounts. Mining MyYahoo usage logs provides Yahoo valuable insight into an individual's web usage habits, enabling Yahoo to provide personalized content, which in turn has led to the tremendous popularity of the Yahoo web site.²⁵

21.3.7 CiteSeer—Digital Library and Autonomous Citation Indexing

NEC ResearchIndex, also known as CiteSeer²⁶ (Bollacker, Lawrence, and Giles 1998) is one of the most popular online bibliographic indices related to computer science. The key contribution of the CiteSeer repository is its “Autonomous Citation Indexing” (ACI) (Lawrence, Giles, and Bollacker 1999). Citation indexing makes it possible to extract information about related articles. Automating such a process reduces a lot of human effort, and makes it more effective and faster.

CiteSeer works by crawling the web and downloading research related papers. Information about citations and the related context is stored for each of these documents. The entire text and information about the document is stored in different formats. Information about documents that are similar at a sentence level (percentage of sentences that match between the documents), at a text level or related due to cocitation is also given. Citation statistics for documents are computed that enable the user to look at the most cited or popular documents in the related field. They also maintain a directory for computer science related papers, to make search based on categories easier. These documents are ordered by the number of citations.

²³<http://www.yahoo.com>.

²⁴See <http://www.yodlee.com>.

²⁵Yahoo has been consistently ranked as one of the top web properties for a number of years (See <http://www.jmm.com/xp/jmm/press/mediaMetrixTop50.xml>).

²⁶See <http://citeseer.nj.nec.com/cs>.

21.4 Research Directions

Although we are going through an inevitable phase of irrational despair following a phase of irrational exuberance about the commercial potential of the web, the adoption and usage of the web continues to grow unabated.²⁷ As the web and its usage grows, it will continue to generate ever more content, structure, and usage data, and the value of web mining will keep increasing. Outlined here are some research directions that must be pursued to ensure that we continue to develop web mining technologies that will enable this value to be realized.

21.4.1 Web Metrics and Measurements

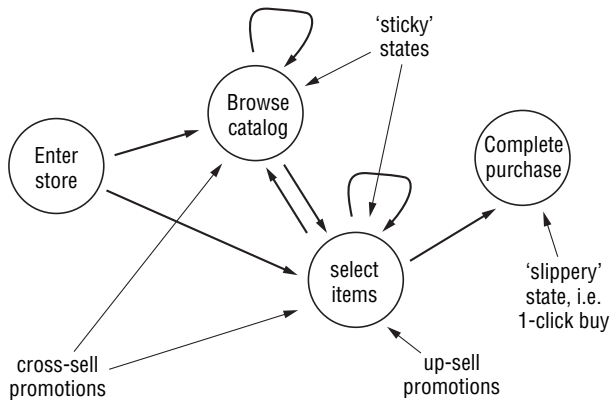
From an experimental human behaviorist's viewpoint, the web is the perfect experimental apparatus. Not only does it provide the ability of measuring human behavior at a micro level, it eliminates the bias of the subjects knowing that they are participating in an experiment, and allows the number of participants to be many orders of magnitude larger than conventional studies. However, we have not yet begun to appreciate the true impact of this revolutionary experimental apparatus for human behavior studies. The web Lab of Amazon²⁸ is one of the early efforts in this direction. It is regularly used to measure the user impact of various proposed changes, on operational metrics such as site visits and visit/buy ratios, as well as on financial metrics such as revenue and profit, before a deployment decision is made. For example, during Spring 2000 a 48 hour long experiment on the live site was carried out, involving over one million user sessions, before the decision to change Amazon's logo was made. Research needs to be done in developing the right set of web metrics, and their measurement procedures, so that various web phenomena can be studied.

21.4.2 Process Mining

Mining of market basket data, collected at the point-of-sale in any store, has been one of the visible successes of data mining. However, this data provides only the end result of the process, and that too decisions that ended up in product purchase. Click-stream data provides the opportunity for a detailed look at the decision making process itself, and knowledge extracted from it can be used for optimizing, influencing the process, etc. (Ong and Keong 2003). Underhill (2000) has conclusively proven the value of process information in

²⁷See, for example, <http://thewhir.com/marketwatch/ser053102.cfm>.

²⁸See <http://www.amazon.com>.



Overall goal:

- Maximize probability of reaching final state
- Maximize expected sales from each visit

Figure 21.2: Shopping pipeline modeled as state transition diagram.

understanding users' behavior in traditional shops. Research needs to be carried out in (1) extracting process models from usage data, (2) understanding how different parts of the process model impact various web metrics of interest, and (3) how the process models change in response to various changes that are made, i.e. changing stimuli to the user. Figure 21.2 shows an approach of modeling online shopping as a state transition diagram.

21.4.3 Temporal Evolution of the Web

Society's interaction with the web is changing the web as well as the way people interact with each other. While storing the history all of this interaction in one place is clearly too staggering a task, at least the changes to the web are being recorded by the pioneering internet archive project.²⁹ Research needs to be carried out in extracting temporal models of how web content, web structures, web communities, authorities, hubs, etc. evolve over time. Large organizations generally archive usage data from their web sites. With these sources of data available, there is a large scope of research to develop techniques for analyzing of how the web evolves over time.

²⁹See <http://www.archive.org/>.

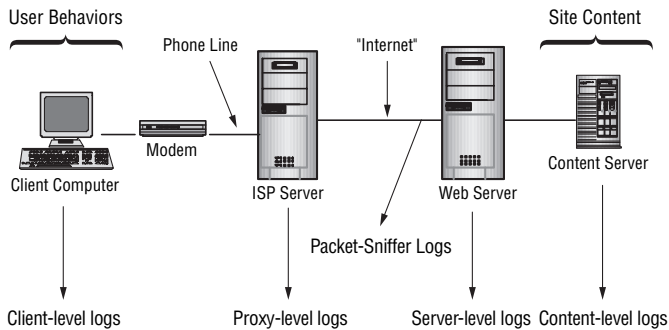


Figure 21.3: High level architecture of different web logs.

21.4.4 Web Services Performance Optimization

As services over the web continue to grow (Katz 2002), there will be a continuing need to make them robust, scalable and efficient. Web mining can be applied to better understand the behavior of these services, and the knowledge extracted can be useful for various kinds of optimizations. The successful application of web mining for predictive prefetching of pages by a browser has been demonstrated in Pandey, Srivastava, and Shekhar (2001). It is necessary to do analysis of the web logs for web services performance optimization as shown in figure 21.3. Research is needed in developing web mining techniques to improve various other aspects of web services.

21.4.5 Fraud and Threat Analysis

The anonymity provided by the web has led to a significant increase in attempted fraud, from unauthorized use of individual credit cards to hacking into credit card databases for blackmail purposes (Scarponi 2000). Yet another example is auction fraud, which has been increasing on popular sites like eBay. Since all these frauds are being perpetrated through the internet, web mining is the perfect analysis technique for detecting and preventing them. Research issues include developing techniques to recognize known frauds, characterize them and recognize emerging frauds. The issues in cyber threat analysis and intrusion detection are quite similar in nature (Lazarevic Dokas, Ertöz, Kumar, Srivastava, and Tan 2002).

21.4.6 Web Mining and Privacy

While there are many benefits to be gained from web mining, a clear drawback is the potential for severe violations of privacy. Public attitude towards privacy

seems to be almost schizophrenic, i.e. people say one thing and do quite the opposite. For example, famous cases like those involving Amazon³⁰ and Doubleclick³¹ seem to indicate that people value their privacy, while experience at major e-commerce portals shows that over 97% of all people accept cookies with no problems, and most of them actually like the personalization features that are provided based on it. Spiekerman, Grossklags, and Berendt (2001) have demonstrated that people were willing to provide fairly personal information about themselves, which was completely irrelevant to the task at hand, if provided the right stimulus to do so. Furthermore, explicitly bringing attention to information privacy policies had practically no effect. One explanation of this seemingly contradictory attitude towards privacy may be that we have a bi-modal view of privacy, namely that “I’d be willing to share information about myself as long as I get some (tangible or intangible) benefits from it, and as long as there is an implicit guarantee that the information will not be abused.” The research issue generated by this attitude is the need to develop approaches, methodologies and tools that can be used to verify and validate that a web service is indeed using user’s information in a manner consistent with its stated policies.

21.5 Conclusions

As the web and its usage continues to grow, so too grows the opportunity to analyze web data and extract all manner of useful knowledge from it. The past five years have seen the emergence of web mining as a rapidly growing area, due to the efforts of the research community as well as various organizations that are practicing it. In this chapter we have briefly described the key computer science contributions made by the field, a number of prominent applications, and outlined some areas of future research. Our hope is that this overview provides a starting point for fruitful discussion.

Acknowledgements

The ideas presented here have emerged in discussions with a number of people over the past few years — far too numerous to list. However, special mention must be made of Robert Cooley, Mukund Deshpande, Joydeep Ghosh, Ronny Kohavi, Ee-Peng Lim, Brij Masand, Bamshad Mobasher, Ajay Pandey, Myra Spiliopoulou, Pang-Ning Tan, Terry Woodfield, and Masaru Kitsuregawa discussions with all of whom have helped develop the ideas presented herein.

³⁰<http://www.ecommercetimes.com/perl/story/2467.html>.

³¹<http://www.wired.com/news/business/0,1367,36434,00.html>.

This work was supported in part by the Army High Performance Computing Research Center contract number DAAD19-01-2-0014. The ideas and opinions expressed herein do not necessarily reflect the position or policy of the government (either stated or implied) and no official endorsement should be inferred. The AHPCRC and the Minnesota Super-Computing Institute provided access to computing facilities.

Jaideep Srivastava is a professor of computer science and engineering at the University of Minnesota, specializing in databases, data mining, and multi-media computing. He can be reached at www.cs.umn.edu/srivasta.

Prasanna Desikan is Ph.D student in the Department of Computer Science and Engineering at the University of Minnesota specializing in link analysis and web mining. He can be reached at www.cs.umn.edu/desikan

Vipin Kumar is currently director of the Army High Performance Computing Research Center and a professor of computer science and engineering at the University of Minnesota specializing in high-performance computing and data mining. He can be reached at www.cs.umn.edu/kumar