

Enhancing The Efficiency of the Website by Mining Web Server Log Files

Jigar H. Jobanputra¹ Banshi D. Soni²

^{1,2}Assistant Professor

¹Department of M.C.A ²Department of BE (I.T)

^{1,2}SRK Institute of Management and Computer Education

Abstract— Web mining - is the application of data mining techniques to discover patterns from the World Wide Web. Web mining can be divided into three different types – Web usage mining, Web content mining and Web structure mining. Web usage mining is the technique that mines the web log. Usage data captures the identity or origin of Web users along with their browsing behavior at a Web site. The web log can be examined by either the client side or server side. In order to gain better efficiency, we can use server web log. It contains many useful information such as name of user, IP Address, number of bytes transferred timestamp, URL etc, so it is very useful in analysis of user behavior and from that it is also useful in the effective maintenance of website. A web server log file is the pool of huge data about client's behavior. In order to increase the efficiency of a website, the website owner should study the user's interest. This can be traced using Web Log Expert or other web mining tools. Web Log Expert is the tool that creates the quick report based on server log file of any website. In this paper we have studied the web log file of a jewellery website and extracted useful patterns of the user interest.

Key words: web mining, web usage mining, web server log mining

I. INTRODUCTION

with the incredible usage of internet now the life of individual become so convenient that consumers prefer e-commerce application/website to shop. So, e-commerce website produces large amount of data. The data has very useful hidden information regarding the consumer's behavior. This information is useful for trading website to find interesting patterns about the consumer behavior which may help to attract & entertain users, implement new marketing strategy, improve the efficiency of website, to compose accessible structure of website. These all useful information of consumers can be mined using web usage mining [1]. Web usage mining is used in the following areas [2]

- Web usage patterns can be used to gather business intelligence to improve consumer attraction, Customer retention, sales, marketing and advertisement, cross sales.
- Web usage mining focuses on techniques that could predict user behavior while the user interacts with the Web.
- Performance and other service quality attributes are crucial to user satisfaction and high quality performance of a web application is expected.
- Web usage mining of patterns provides a key to understanding Web traffic behavior, which can be used to deal with policies on web caching, network transmission, load balancing, or data distribution.

- Web usage and data mining is also useful for detecting intrusion, fraud, and attempted break-ins to the system.

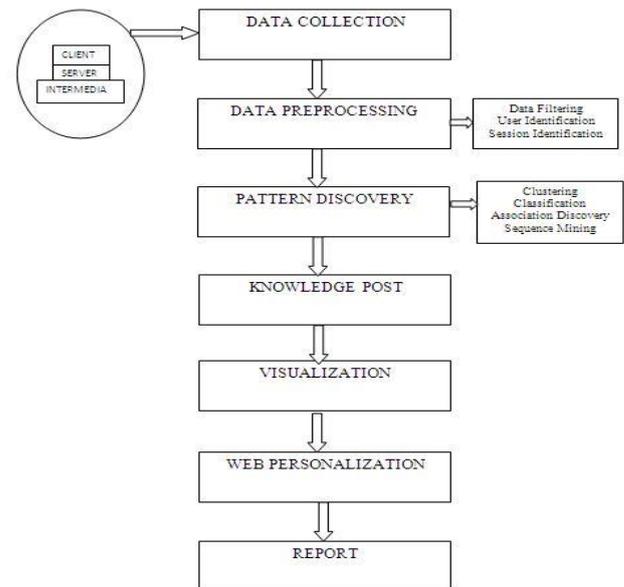


Fig. 1:

A. Steps of web usage mining

- 1) Data collection
The very first step of web usage mining is to collect a data to mine the data is a web log document which consist the log of every visit on the web.
- 2) Data Integration
Coordinate different web logs to a single document.
- 3) Data Pre-processing
Cleaning and organizing information to get ready for pattern extraction.
- 4) Pattern Extraction
Extraction of user patterns which will help to study user behaviour.
- 5) Pattern Analysis and visualization
Pattern Analysis will help to find interesting rules from set of patterns.
- 6) Pattern Application
Apply the patterns in the real world problems.

B. Web Logs

A Web log file [3] records activity information when a Web user submits a request to a Web Server. The web log file is the main source of raw data which we shall refer to as log file. The web log file can be stored in three different places: 1). Web servers 2). Web proxy servers 3). Client browsers

C. Web servers

a web server generally stores the most complete and accurate usage data.

D. Web proxy server logs

A proxy server takes the HTTP requests from users and passes them to a web server then returns to the users the results passed by the web server. Web proxy server can also stores the log files.

E. Client side logs

clients can download the special software that records the web usage HTTP cookies could also be used for this purpose. These are a piece of information generated by a web server and stored in user's computer.

Web log file of jewelers' website is considered for the purpose of analysis.

F. Description of log file

A file produced by a Web server to record activities on the Web server. It usually has the following features[5]:

- The log file is text file. Its records are identical in format.
- Each record in the log file represents a single HTTP request.
- A log file record contains essential information about a request: the client side host name or IP address, the date and time of the request, the requested file name, the HTTP response status and size, the referring URL, and the browser information.
- A browser may shoot multiple HTTP requests to Web server to display a single Web page. This is because a Web page not only needs the main HTML document; it may also need additional files, like images and JavaScript files.
- The main HTML document and additional files all require HTTP requests.
- Each Web server has its own log file format.
- If your Web site is hosted by an ISP (Internet Service Provider), they may not keep the log files for you, because log files can be very huge if the site is very busy. Instead, they only give you statistics reports generated from the logs files.
- Remote IP address or domain name: An IP address is a 32-bit host address defined by the Internet Protocol; a domain name is used to determine a unique Internet address for any host on the internet. One IP address is usually defined for one domain name.
- Authuser: Username and password if the server requires user authentication
- Entering and exiting date and time.
- Modes of request: GET, POST or HEAD method of CGI (Common Gateway Interface).
- Status: The HTTP status code returned to the client, e.g., 200 is "ok" and 404 are "not found".
- Bytes: The content-length of the document transferred.
- Remote log and agent log.
- Remote URL.
- "request:." The request line exactly as it came from the client.
- Requested URL

G. Types of Log Analyzer Tools

There are many web log analyzer tools are available and some of them are open source. Here the list of some free available web log analyzer tool [4].

- Web log Expert
- Deep log Analyzer
- Google Analytics
- AWstats
- W3Perl
- Visitors
- Analog
- Webtrax
- Weblog Expert Lite

H. Example with snapshot analysis

In this study we have analyzed the web log file of web server with the help of weblog Expert analyzer tool. The sample log files consist the data form 8th dec, 2015 to 16th dec, 2015. The web log file have stored 1393 KB of data. The weblog Analyzer tools are required as they help extensively in analyzing the information about visitors, hits, page views, bandwidth, browsers used, operating systems etc, which can be utilized by system administrator and web designer to increase the effectiveness of the website.

I. General Statistics

In this section we get general information regarding the website.

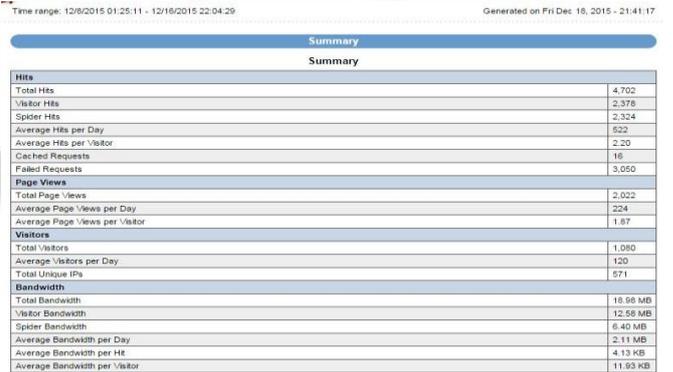


Fig. 2:

J. Activity Statistics

in this section the activities performed by user will be displayed.



K. Access Statistics

This section will display the statistics regarding the accessed pages by users. This section is very useful to know the user interest.

Page	Hits	Incomplete Requests	Visitors	Bandwidth (KB)
1 http://umajewellers.in/	730	0	674	6,972
2 http://umajewellers.in/Photogallery.aspx	252	0	210	1,915
3 http://umajewellers.in/index.aspx	14	0	12	135
4 http://umajewellers.in/ContactUs.aspx	16	0	12	148
5 http://umajewellers.in/Default.aspx	10	0	8	48
6 http://umajewellers.in/Traditional.aspx	8	0	6	95
7 http://umajewellers.in/Style.aspx	8	0	6	62
8 http://umajewellers.in/Beauty.aspx	8	0	6	96
9 http://umajewellers.in/Index.aspx	4	0	2	38
Total	1,050	0	N/A	9,512

Fig. 4:
Most Requested Directories

Directory	Hits	Incomplete Requests	Visitors	Bandwidth (KB)
1 http://umajewellers.in/	1,054	0	706	9,512
2 http://umajewellers.in/images/	54	0	40	1,308
3 http://umajewellers.in/css/	110	0	12	575
Total	1,218	0	N/A	11,396

Fig. 5:

L. Visitors

This section will display the host information along with their hits and bandwidth.

Host	Hits	Visitors	Bandwidth (KB)
1 212.54.240.254	336	52	1,933
2 192.99.33.150	40	20	219
3 87.235.214.63	24	12	131
4 62.210.35.155	45	12	253
5 62.210.141.4	24	12	12
6 193.27.127.171	22	10	104
7 192.76.161.153	40	10	52
8 193.154.165.43	39	10	165
9 85.105.16.214	20	10	155
10 94.121.144.43	16	8	87
11 38.87.42.212	24	8	143
12 121.205.198.216	22	8	28
13 188.248.48.70	16	8	87
14 117.26.162.89	24	8	31
15 62.210.203.121	12	6	65
16 172.248.82.168	8	6	41
17 62.210.181.15	6	6	1
18 102.43.11.169	20	6	25
19 157.165.159.220	6	6	41
20 201.254.166.90	6	4	3
21 212.132.28.132	12	6	65
22 212.83.151.12	8	4	43
23 114.17.242.33	6	4	41
24 216.176.100.161	4	4	21
25 134.248.82.107	8	4	43
26 212.83.151.6	8	4	43
27 216.155.218.58	6	4	41
28 112.111.163.75	4	4	5
29 188.214.65.171	4	4	22
30 107.172.218.41	8	4	43
31 188.214.14.99	4	4	22
32 193.201.226.74	8	4	44
33 81.264.166.227	8	4	3
34 193.154.165.158	8	4	43
35 62.210.170.175	12	4	74
36 173.209.148.117	4	4	5
37 2134.15.291	4	4	5
38 172.248.71.130	4	4	22
39 27.151.162.21	4	4	5
40 175.44.19.205	4	4	5
41 2134.75.30	4	4	22
42 223.83.103.54	4	4	17
43 172.242.214.121	4	4	39
44 46.165.141.199	4	4	39
45 2.175.152.140	8	4	5
46 617624.176	8	4	43
47 21.209.35.155	8	4	43
48 178.137.18.18	8	4	43
49 42.79.137.83	4	4	0
50 42.102.83.141	6	4	24
Subtotal	902	356	4,565
Total	2,378	1,086	12,882

Fig. 5:

M. Browsers

This section will display the information regarding the browser used to visit the web site. And also operating system on client machine.

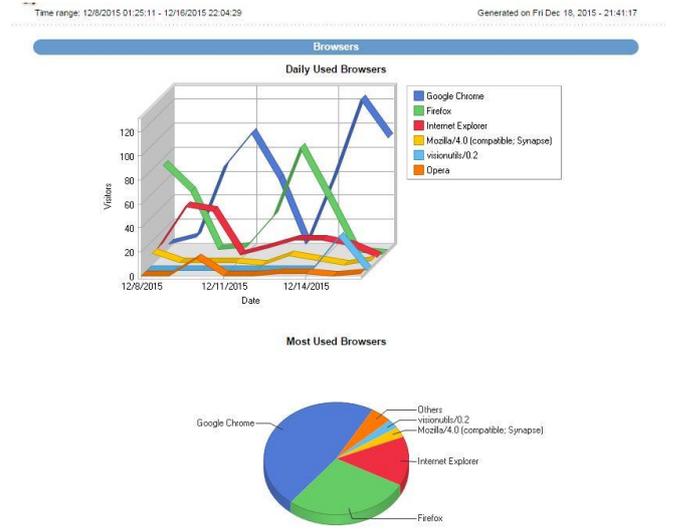


Fig. 6:

Browser	Hits	Visitors	% of Total Visitors
1 Google Chrome	1,228	516	47.51%
2 Firefox	536	304	27.99%
3 Internet Explorer	314	154	14.18%
4 Mozilla/4.0 (compatible; Synapse)	32	32	2.95%
5 visionalib/0.2	32	28	2.58%
6 Opera	32	20	1.84%
7 Others	18	10	0.92%
8 Mozilla/5.0	6	6	0.55%
9 Food Research	8	4	0.37%
10 Wada.vn Vietnamese Search	2	2	0.18%
11 !&t71 - (http://t71.com/)	2	2	0.18%
12 Seolentriquiry.com bot	2	2	0.18%
13 Mozilla/5.0 (X11; Linux x86_64; AppleWebKit/534.34 (KHTML, like Gecko) Qt/4.8.1 Safari/534.34	110	2	0.18%
14 Mozilla/5.0 (compatible; sptbot/4.2; http://OpenLink-Profiler.org/bot)	40	2	0.18%
15 Mozilla/4.0 (compatible; Win32; WinHttp; WinHttpRequest/5)	16	2	0.18%
Total	2,378	1,086	100.00%

Fig. 7:

Operating System	Hits	Visitors	% of Total Visitors
1 Windows 8	796	356	32.96%
2 Windows 7	744	300	27.78%
3 Windows 8.1	458	256	23.70%
4 Others	158	86	7.96%
5 Windows XP	74	50	4.63%
6 Windows / Vista	22	12	1.11%
7 Mac OS	12	8	0.74%
8 Linux	116	8	0.74%
9 Windows 2000	2	2	0.19%
10 Windows 10	4	2	0.19%
11 Windows Server 2003	2	0	0.00%
Total	2,378	1,086	100.00%

Fig. 8:

II. CONCLUSION

The internet is very huge and common platform to attract the consumers for a product or service. When users visit the website they leave their footsteps in forms of web log file. And this web log will help the website owner or designer to understand the consumer interest or behavior. Web site owner can improve the efficiency of website using these reports and can attract more consumers.

REFERENCES

- [1] Yanduo Zhao, "The Review of Web Mining in E-Commerce", ICCIS, 2013, 2013 Fifth International Conference on Computational and Information Sciences (ICCIS), 2013 Fifth International Conference on Computational and Information Sciences (ICCIS) 2013, pp. 571-574, doi:10.1109/ICCIS.2013.158
- [2] Naresh Barsagade "Web usage Mining and Pattern Discovery "a Survey Paper" –Dec 8 2003.
- [3] Jaideep Srivastava, Robert Cooley, Mukund Deshpande, Pang Ning Tan "Web usage mining: Discovery and Applications of usage patterns from web data" SIGKDD Explorations- vol-1, issue-2 Jan 2000 pages 12-33

- [4] S.Padmaja, Dr.Ananthi Sheshasaayee “WEB SERVER LOGS TO ANALYZING USER BEHAVIOR USING LOG ANALYZER TOOL” (PAGE: 514-525)
- [5] <http://www.w3.org/Daemon/User/Config/>
- [6] <http://www.weblogexpert.com/>

