

Insights of Web Mining

Tejas S Bhaise

Department of Computer Engineering
SSBT's College of Engineering and Technology
Bambhori, Jalgaon - 425 001, Maharashtra, India

Abstract: Mining means extraction the data .data mining means the extraction of data from the huge data file. Web data mining belongs to data mining i.e.(subpart) which deals with huge content over the internet, today there is the rapid growth of information or data on internet and user needs more time to find a particular data in a short period for that web mining technique is used.

I.Introduction

The goal of web mining is to extract the information or data from the internet. The data is in the form of picture, audio, video, text, reports, mathematical form. For extraction of data different techniques are used in web mining .this paper gather all techniques, algorithms, issues about web mining

II.Web Mining

The main goal of web mining is to extract useful structured data and according to user needs. Internet has huge content of data and all the data interconnected. Government, education field, management sector commercial advertisement, news is the areas of knowledge service. Resource finding, preprocessing and information selection, generalization and analysis are a subtask of web mining [3]

The process in which we extract the data either from online and offline text data available on the web involves in resource finding

The automatic selection and preprocessing of specific data from fetched web resources involved in information selection and preprocessing

General patterns at an individual, as well as multiple sites, discover automatically in generalization

Authorization and perception plays important role in mind patterns to analysis

III.Web mining Classification

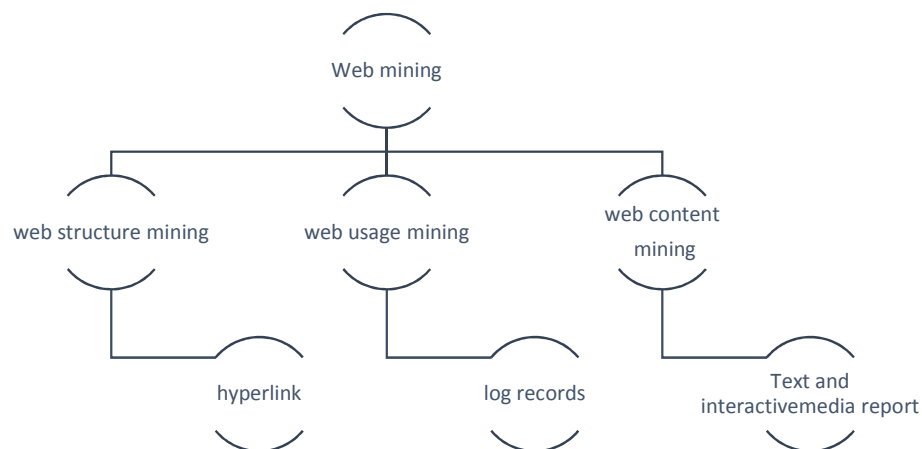


Figure-web mining Classification

1. Web content mining

It is the technique extracting the data from the internet into more structured form. Web content mining includes multimedia reports and text. Web content mining do not always provide structure data it may unstructured or interpret document

1.1 Web content mining techniques-

1	Structured	1)Technique in which it represents the host page on the internet 2)structured mining use to mine the structured data from internet 3)It is easy to extract data comparatively unstructured mining
2	Unstructured	1)On internet mostly data is in unstructured form so it need text and data mining approaches 2)the exact content extraction from internet is major part of text mining 3)text mining is covered by web content mining
3	Semi structured	1)the presentation of semi structured data over the internet is in the form of tags such as HTML,XML 2)semi structured is partially form of structured data
4	Multimedia	1) Multimedia data like images, videos, audio etc. are available on internet 2)this mining technique is use for retrievals interesting multimedia data from internet

Figure-table of web content mining techniques [4]

1.2 Web content mining algorithms-

Support vector machine-The conversion of original training data into a higher dimension done by uses of non-definite mapping.it finds for the definite optimal separating hyperplane with appropriate non definite mapping to an adequately high aspect, among the new aspect. Hyper plane partitioned by two classes from data. [7]

Neural network- This is another algorithm for web content mining. It contains an input layer and output layer, one and more than one hidden layer .the basic unit is neuron. The attributes measured for each training tuple is corresponded by inputs to the network. For making input layer the input supply all together. Supply concurrently to the hidden layer and it will be weighted. A number of hidden layers are inconsistent even it only one. To emit network's prediction, weighted output of last hidden layer are inputs to units making up the output layer. For input and output unit none of the weights cycle back as network is feed-forward. [7]

2. Web structure mining

Discover structure data from web reports called web structure mining.it is classified into two levels one is document (intra-page) and other is hyperlink. The hyperlink grant point to pages with authority of the same topic of the page containing the link.

2.1 Web structure mining techniques-

1	Link cardinality	to find corresponding websites and analyze them, page distribution
2	Link strength	It related to link weights
3	Link type	prediction of link type between two or more than two item
4	Link based cluster analysis	1)Data is distributed or grouped together 2)similar and dissimilar data grouped in single group and separate group respectively
5	Link based classification	1)anticipate web page categories like HTML,text,tags, 2)association between web pages

Figure-table of web structure mining techniques [4]

2.2 Web structure mining algorithms-

Page rank algorithm- Its name itself indicate that, it is based on the ranking of pages.it is one of the relevant algorithm of web structure mining. Page rank of a web page is established by a number of inbound links joined to it. The page is considered on the basis of links joined to that page or links taken away from most popular or usable web pages. The calculation of ranking based on the following formula.

$$PR(J) = (1 - d) + d \left[\frac{PR(K1)}{C(K1)} + \dots + \frac{PR(Kn)}{C(Kn)} \right]$$

Here,

$PR(J)$ = Page Rank of J (web page)

$PR(Ki)$ = Rank of pages Ki link to X

$C(Ki)$ = No of outward links to Ki

d = Damped Factor (value ranges from 0-1, normally value is 0.85) [6]

Weighted page rank algorithm- Term weighted indicates 'value'. This algorithm based on rank value of page according to their importance (outgoing and incoming links) rather than dividing it uniformly.

$$W^{in}(v, u) = wt. of link (v, u)$$

$$W^{out}(v, u) = wt. of link (v, u)$$

$$W_{(v,u)}^{in} = \frac{Iu}{\sum_{p \in R(v)} Ip} \quad \dots 1)$$

Here,

$R(v)$ = associate page list of page 'v' and Iu and Ip = No. of up-links of page u, p.

$$w_{(u,v)}^{out} = \frac{O_u}{\sum_{p \in R(v)} O_p} \quad \dots 2)$$

Here,

$R(v)$ = associate page list of page 'v' and $O(u)$ and $O(p)$ = No. of down-links of page u, p.

Afterwards, tally of importance of web pages the formula is,

$$WPR(u) = (1 - d) + d \sum_{v \in B(u)} PR(v) w_{(u,v)}^{in} w_{(u,v)}^{out} [6]$$

3. Web usage mining

The mining in which elicitation of useful usage patterns from web data. Especially server log, web log records, browser log

3.1 Web usage mining techniques-

1	Pattern analysis	In this technique all relevant rules are extracted and irrelevant rules are separated.
2	Pattern discovery	1) This technique used to find interesting patterns. 2) statistical analysis, association rule, clustering method, classification and sequential sessions are the techniques in pattern discovery
3	Data preprocessing	1) Conversion of noisy data into an pure errorless data 2) user and session identification and data cleaning are the techniques in data preprocessing

Figure-table of web usage mining techniques [4]

3.2 Web usage mining algorithms-

Fuzzy Cmean- Two or more than two clusters relate to same group for distribution of data. A membership grade assign by algorithm to, the highest page having the highest membership and the lowest pages having the lowest membership. The user uses a web page in future according to it. it is clustering basis algorithm. Membership calculated by a distance between the cluster and data point. [5]

Sequential pattern- To find data in web access log their technique is used .It depends on timestamp. It is trying to discover which item is followed by another set of item. to access weblog it is required that each activity contain the field and that field denotes the period of time for which we are mining sequential pattern. [8]

IV.Applications

- **E-Learning:** It is generally based on web usage. Machine learning techniques and web usage intensify web-based learning surrounding.
- **Digital libraries:** It is crucial application, to provide relevant data all over the world.
- **E-Government:** For better social service different Gov. Systems interact with people. The main aspect is to join people in E system, use of technology to deliver service electronically, eye on people needs by providing better information.

- Electronic commerce: The main thing in E-commerce is to recognize customer need. It can improve the scope of service.
- E-Politics: It contains all the data about regional Gov., member of houses and parties. It is important to clarify political transparency and democracy
- Security and Crime Investigation: Cybercrime is a big issue in a world of internet. it is used for protection of user system. Cyber terrorism like internet fraud, hacking, fake websites, virus spreading. Clustering and classification like that techniques are used.
- Electronic Business: It is for improvement of marketing, sale support and customer support by web mining.

V.Challenges and Issues

There are certain challenges and Issues regarding web mining, like for finding the data on web people generally use browser, while surfing on web they usually put a word query and the query reply is the list of pages ranked based on their closeness of query, Trouble to *finding relevant information* to get the result. People want the content in a particular format and it is available on the internet in another format, this is one of the challenge to gather information in the same format to *Personalization of the data*. When already a huge data available on internet and we have to create new data from available data (data-triggered process) this issue regarding to *create new relevant knowledge from available content on internet*. [2]

VI.Conclusion and Future Work

As the rapid growth of data on internet web mining is an important technique. This paper highlights algorithms, techniques, challenges and applications. Process Mining, Web Mining, and Privacy, fraud and threat analysis, web services performance optimization are the areas of future work [1].

References

1. Tulasi Gayatri Devi, Aparna Ks, A Survey On Web Mining: Overview, Techniques, Tools And Applications, International Journal For Research In Applied Science And Engineering Technology(Ijrasnet) Volume 4 Issue I, January 2016
2. R.Munilatha, K.Venkataramana, A Study On Issues And Techniques Of Web Mining, International Journal of Computer Science and Mobile Computing, Vol.3 Issue.5, May-2014
3. Mr.T.Suresh Kumar, T.Rajamadhangi, A Survey Of Web Mining Algorithms, International Journal Of Engineering Science Invention Research And Development, Vol. II Issue VIII, February 2016
4. Muhammad Jawad Hamid Mughal, Data Mining: Web Data Mining Techniques, Tools And Algorithms: An Overview, (IJACSA) International Journal of Advanced Computer Science and Applications, Vol. 9, No. 6, 2018
5. P. Sampath and Prabhavathy m., web page access prediction using fuzzy clustering by local approximation memberships (flame) algorithm, vol. 10, no. 7, April 2015
6. Sanjay and dharmender Kumar, A Review Paper on Page Ranking Algorithms, International Journal of Advanced Research in Computer Engineering & Technology (IJARCET) Volume 4 Issue 6, June 2015

7. R.Malarvizhi,K.Saraswathi,Web Content Mining Techniques Tools and Algorithms – A Comprehensive Study, International Journal of Computer Trends and Technology(IJCTT) – volume 4 Issue 8–August 2013
8. S.Jagan,Dr.S.P.Rajagopalan,A Survey On Web Personalization Of Web Usage mining, International Research Journal of Engineering and Technology (IRJET)Volume: 02 Issue: 01 | March-2015