

Self Organizing For Dimensional Reduction in Web Mining

1.C.SADHANA(SP14CAD502)

Research Scholar

Computer Science and Application

St peters University

2.Dr.L.Mary Immaculate Sheela

Professor,

Department of Information Technology,

Pentecost University College,

Ghana.

Abstract

Data mining is a form of extracting data available in the internet. Web mining is a part of data mining. Web mining adopts much of the data mining techniques to discover potentially useful information. Web mining analysis relies on three general sets of information previous usage patterns, degree of shared content and inter memory association link structure corresponding to three subset in web mining namely Web usage mining ,Web content mining, Web structure mining respectively. proposal shares dissimilar goals with many those agents, our approach is automatic that it does not require users explicit input. Moreover, we take a systematic approach to collect and comprehend user activities. We provide a general framework for collecting, mining, and search/query personal usage data, which may be employed by various agents. Web usage mining is used to analyze the behavior of websites users. It involves automatic discovery of user access patterns from one or more web servers. It contains four processing stages including data collection, preprocessing, pattern discovery and analysis. The web content mining refers to the discovery of useful information from web contents which include text, image, audio, video etc. The mining of link structure aims at developing techniques to take advantage of the collective .It includes extraction of structure data from web pages, identification, similarity and integration of data with similar meaning. There are two common tasks involved in web mining they are Clustering and Classification. Neural based approach is used to analysis the performance of the clustering of the number of request. We propose an approach "ENHANCED SELF ORGANIZATION MAP" which is data visualization technique; it reduces the dimensions of data through the use of neural network. In previous study on SOM plot the similarities of data by grouping similar data items together, so they reduces dimension and display similarities SOM organize sample data, which are usually surrounded by similar samples ,similar samples are not always near each other .In ESOM we use users Clustering mining algorithm. ESOM can estimate the center and the number of clustering data set by "dissimilarity computing", it optimizes SOM neural network learning and improve clustering effect.

Introduction

The expansion of the World Wide Web (Web forshort) has resulted in a large amount of data that isnow in general freely available for user access.The different types of data have to be managedand organized that they can be accessedefficiently. Therefore the application of datamining techniques on the Web is now the focus ofan increasing number of researchers. Several data

Mining methods are used to discover the hiddeninformation in the Web. However, Web mining does not only mean applying data miningtechniques to the data stored in the Web. Thealgorithms have to be modified to better suit the demands of the Web. New approaches should beused better fitting to the properties of Web data.Furthermore, not only data mining algorithms, butalso artificial intelligence, information retrievaland natural language processing techniques can beused efficiently. Thus, Web mining has beendeveloped into an autonomous research area. Webmining involves a wide range of applications thataim at discovering and extracting hiddeninformation in data stored on the Web. Another Importantpurpose of Web mining is to provide a mechanismtomake the data access more efficiently andadequately. The third interesting approach is todiscover the information which can be derivedfrom theactivities of users, which are stored in log files for

example for predictive Web caching [1]. Thus, Web mining can be categorized into threedifferent classes based on which part of the Webis to be mined [2], [3], and [4]. These threecategories are Web content mining, Web structuremining and Web usage mining. Web contentmining [7], [6] is the task of discovering usefulinformation available on-line. There are differentkinds of Web content which can provide usefulinformation to users, for example multimedidata, structured (i.e. XML documents), semi structured (i.e. HTML documents) andunstructured data (i.e. plain text). The aim of Webcontent mining is to provide an efficient mechanism to help the users to find theinformation they seek. Web content miningincludes the task of organizing and clustering the documents and providing search engines foraccessing the different documents by keywords,categories, contents. Web structure mining is theprocess of discovering the structure of hyperlinks within the Web. Practically, while Web content mining focuses on the inner-documentinformation, Web structure mining discovers thelink structures at the inter-document level. Theaim is to identify the authoritative and the hubpages for a given subject. Web usage mining isthe task of discovering the activities of the userswhile they are browsing and navigating throughthe Web. The aim of understanding the navigationpreferences of the visitors is to enhance thequality of electronic commerce services(ecommerce), to personalize the Web portals [7]or to improve the Web structure and Web serverperformance [3]. For this reason a model of theusers (User Model - UM) have to be built based on the information gained from the log data.

Web Usage Mining

The aim of Web usage mining is to discover patterns of user activities in order to better serve the needs of the users for example by dynamic link handling, by page recommendation etc. The aim of a Web site or Web portal is to supply the user the information which is useful for him. There is a great competition between the different commercial portals and Web sites because every user means eventually money (through advertisements, etc.). Thus the goal of each owner of a portal is to make his site more attractive for the user. For this reason the response time of each single site has to be kept below 2s. Moreover some extras have to be provided such as supplying dynamic content or links or recommending pages for the user that are possible of interest of the given user. Clustering of the user activities stored in different types of log files is a key issue in the Web community. There are three

types of log files that can be used for Web usage mining [4]. Log files are stored on the server side, on the client side and on the proxy servers. By having more than one place for storing the information of navigation patterns of the users makes the mining process more difficult. Really reliable results could be obtained only if one has data from all three types of log file. The reason for this is that the server side does not contain records of those Web page accesses that are cached on the proxy servers or on the client side. Besides the logfile on the server, that on the proxy server provides additional information. However, the page requests stored in the client side are missing. Yet, it is problematic to collect all the information from the client side. Thus, most of the algorithms work based only on the Server Side data. Web usage mining consists of three main steps:

(i) Pre-processing (ii) Pattern discovery (iii) Pattern analysis

In the pre-processing phase the data have to be collected from the different places it is stored (client side, server side, proxy servers). After identifying the users, the click-streams of each user has to be split into sessions. In general the timeout for determining a session is set to 30 minutes [5]. The pattern discovery phase means applying data mining techniques on the pre-processed log data. It can be frequent pattern mining, association rule mining or clustering. In this paper we are dealing only with the task of clustering web usage log. In web usage mining there are two types of clusters to be discovered: usage clusters and page clusters. The aim of clustering users is to establish groups of users having similar browsing behaviour. The users can be clustered based on several information. In the one hand, the user can be requested filling out a form regarding their interests, for example when registration on the web portal. The clustering of the users can be accomplished based on the forms. On the other hand, the clustering can be made based on the information gained from the log data collected during the user was navigating through the portal. Different types of user data can be collected using these methods, for example (i) characteristics of the user (age, gender, etc.), (ii) preferences and interests of the user, (iii) user's behaviour pattern. The aim of clustering web pages is to have groups of pages that have similar content. This information can be useful for search engines or for applications that create dynamic index pages. The last step of the whole web usage mining process is to analyze the patterns found during the pattern discovery step. Web Usage Mining tries to understand the patterns detected in the previous step. The most common techniques is data visualization applying filters, zooms, etc.

Web Usage Mining (WUM) consists in the analysis of the way a web site is browsed by its users so as to improve it (in a very broad sense). The practical goals include to a few; improving the performances of web servers with intelligent caching and proxy that anticipate user requests; improving the structure of a site based on user typical navigation, for instance by creating automatically site map and bypassing links; introducing recommendations, such as the proposed by on-line bookstores, based on recognition of typical browsing sessions or of buying patterns.

Self Organized Map

The Self-Organizing Map (SOM) [17] was developed by professor Kohonen. It is one of the most popular neural network models. It belongs to the category of competitive learning networks based on unsupervised learning, which means that no human intervention is needed during the learning and that little need to be known about the characteristics of the input data. SOM is used for clustering the data without knowing the class. The SOM can be used to detect features inherent to the problem and thus has also been called SOFM, the Self-Organizing Feature Map, Provides a topology preserving mapping from the high dimensional space to map units. Map units, or neurons, usually form a two-dimensional lattice and thus the mapping is a mapping from high dimensional space onto a plane. The property of topology preserving means that the mapping preserves the relative distance between the points. Points that are near each other in the input space are mapped to nearby map units in the SOM. The SOM can thus serve as a cluster analyzing tool of high dimensional data.

Self Organising Algorithm

1. Select output layer network topology. Initialize current neighbourhood distance, $D(0)$, to a positive value.
2. Initialize weights from inputs to outputs to small random values.
3. Let $t=0$
4. While computational bounds are not exceeded do ($t \leq 1$).
 - i) Select an input sample t_i, k .
 - ii) Compute the square of the Euclidean distance of t_i, k . From weight vectors (w_j) associated with each output node. $t_i, k - w_j, k(t)$
 - iii) Select output node j^* that has weight vector with minimum value from step 2.
 - iv) Update weights to all nodes within a topological distance given by $D(t)$ from j^* , using the weight update rule: $w_j(t+1) = w_j(t) + n(t)(t_i - w_j(t))$
 - v) Increment t .
5. End while.

Dissimilarities of WUM

A difficulty included by the proposed representation is that N (the number of sessions) is in general quite high (more than 16000 in the proposed application).

2 Related Algorithms

2.1 SOM

Teuvo Kohonen [4] introduced the SOM network that reduced the dimensions of data through the use of self-organizing neural networks. The SOM network produces a map of usually one or two dimensions which plot the similarities of the data by grouping similar data items together. This mapping process reduces the problem dimensions. The SOM network integrates dimensions reducing and clustering in one network. Figure 1 shows the mapping from a one-dimensional input to a two-dimensional array.

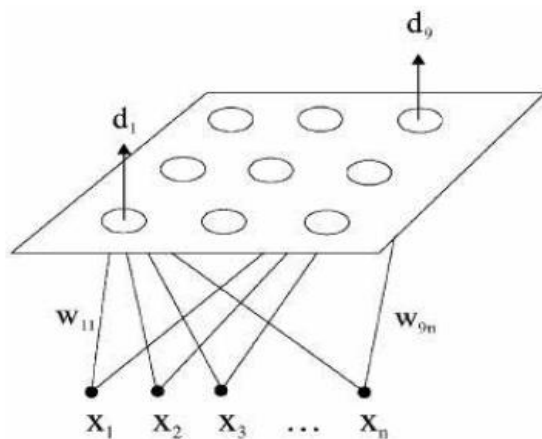


Figure 1: The Mapping from a one-dimensional input to two-dimensional array [11]. The SOM network organizes itself by competing representation of the samples. Neurons are also allowed to change themselves in hoping to win the next competition. This selection and learning process makes the weights to organize themselves into a map representing similarities. The algorithm of the SOM network is shown as follows:

1. Initialize Map
2. Set $t = 0$ and repeat the following steps until $t > 1$
 Randomly select a sample
 Get best matching unit
 Scale neighbors
 Increase t by a small amount
3. End for

The first step in constructing a SOM is to initialize the weight vectors. From there the algorithms select a sample vector randomly and search the map of weight vectors to find the weight that can represent the sample best. Since each weight vector has a location, it also has neighbouring weights that are close to it. The chosen weight is rewarded to perform better than a randomly selected sample vector. In addition to this reward, the neighbours of the weight are also rewarded. From this step we increase t some small amount because the number of neighbours and how much each weight can learn decreases over the time. This whole process is then repeated a large number of times, usually at least 1000 times. The main advantage of using the SOM network is that SOM automatically (self-organizing) clusters documents. The SOM network also can be applied to a large scale of data.

2.2 K-Means

The k -means algorithm was introduced by J. MacQueen, and it had been one of the most popular clustering Algorithms. This clustering algorithm represents each of k clusters C_j by the mean (weighted average) c_j of its point (called centroid). It initially selects clusters such that points are mutually farthest apart. Next, it examines each point and assigns it to one of the clusters depending on the minimum distance. The centroid's position is recalculated every time a point is added to the cluster. This process will be repeated until all the points are grouped into the k clusters. However, this algorithm does not work well if there are large differences in the data set. The equation for k -means algorithm is in Equation 1 and 2.

$$u_{ij} \in U_{c \times n}; C_i = \frac{1}{n_i} \sum_{j=1}^n X_j \tag{1}$$

$$\min J = \sum_{i=1}^c \sum_{j=1}^n u_{ij} \|X_j - C_i\|^2 \tag{2}$$

In equations (1) and (2), X_j represents each point j 's co-ordinates and u_{ij} represents the hypothetical belonging of point j into cluster i (i.e., $u_{ij} = 1$ if j belongs to cluster i ; $u_{ij} = 0$ if j belongs to any other cluster different from i)

3 SOM-based Web page clustering

Overview of the SOM Algorithm

We have a spatially *continuous input space*, in which our input vectors live. The aim is to map from this to a low dimensional spatially *discrete output space*, the topology of which is formed by arranging a set of neurons in a grid. Our SOM provides such a nonlinear transformation called a *feature map*.

The stages of the SOM algorithm can be summarised as follows:

1. *Initialization* – Choose random values for the initial weight vectors w_j .
2. *Sampling* – Draw a sample training input vector x from the input space.
3. *Matching* – Find the winning neuron $I(x)$ with weight vector closest to input vector.
4. *Updating* – Apply the weight update equation $Dw_{ji} = \eta (x_j - w_{ji})$.
5. *Continuation* – keep returning to step 2 until the feature map stops changing.

Dissimilarity/Similarity metric: Similarity is expressed in terms of a distance function, which is typically metric: $d(i, j)$

There is a separate "quality" function that measures the "goodness" of a cluster.

The definitions of distance functions are usually very different for interval-scaled, boolean, categorical, ordinal and ratio variables.

It is hard to define "similar enough" or "good enough" the answer is typically highly subjective.

Similarity and Dissimilarity Between Objects

Distances are normally used to measure the similarity or dissimilarity between two data objects

Some popular ones include: *Minkowski distance*:

where $i = (x_{i1}, x_{i2}, \dots, x_{ip})$ and $j = (x_{j1}, x_{j2}, \dots, x_{jp})$ are two p -dimensional data objects, and q is a positive integer

If $q = 1$, d is Manhattan distance

$$d(i, j) = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + \dots + |x_{ip} - x_{jp}|$$

If $q = 2$, d is Euclidean distance:

Properties

- $d(i, j) \geq 0$
- $d(i, i) = 0$
- $d(i, j) = d(j, i)$
- $d(i, j) \leq d(i, k) + d(k, j)$

Also one can use weighted distance, parametric Pearson product moment correlation, or other dissimilarity measures.

Dissimilarity between Binary Variables

gender is a symmetric attribute

the remaining attributes are asymmetric binary

let the values Y and P be set to 1, and the value N be set to 0

3.2 Data Pre-processing

There are several pre-processing tasks to be done before executing the data mining algorithms on the Webserver logs. These processes include data formatting, user identification, session identification, and transaction identification. The original server logs are formatted and grouped into meaningful transactions before being processed by theming system. We describe each of these processes in the following paragraphs. Data formatting The access log is saved to keep record of every request made by the users. Since our main purpose is to facilitate more effective and efficient navigation, we only want to keep the log entries with information

relevant to our purpose of organizing the Web pages. Some irrelevant log entries are deleted from the log file. Sometimes a user requests a page that does not exist. This will create an error entry in the log. Since we are organizing the existing Web URLs, we are not interested in this kind of error entries, and hence these error entries shall be deleted. A user's request to view a particular page often results in several log entries because the page consists of several materials such as graphics or small applets. However, we are only interested in, and hence only keep, what the user explicitly requests because we intend to design a system that is user-oriented.

User identification The task of identifying unique users is greatly complicated by the existence of local caches, corporate firewalls, and proxy servers. Therefore, some heuristics are commonly used to help identify unique users. We use the machine's IP addresses to identify unique users.

User-session identification For logs that span a long period of time, it is very likely that different users may use the same machine to access the server Web sites. Therefore, we differentiate the entries into different user-sessions through a session timeout. That is, if two time stamps between page requests exceeds a certain limit, we assume the pages are requested by two different user-sessions, even though the IP address is the same.

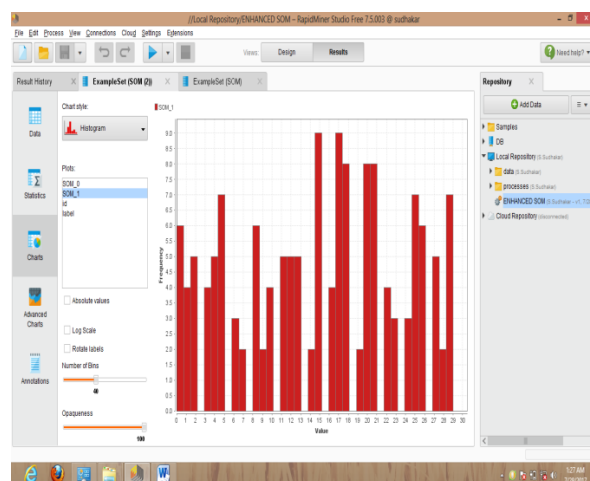
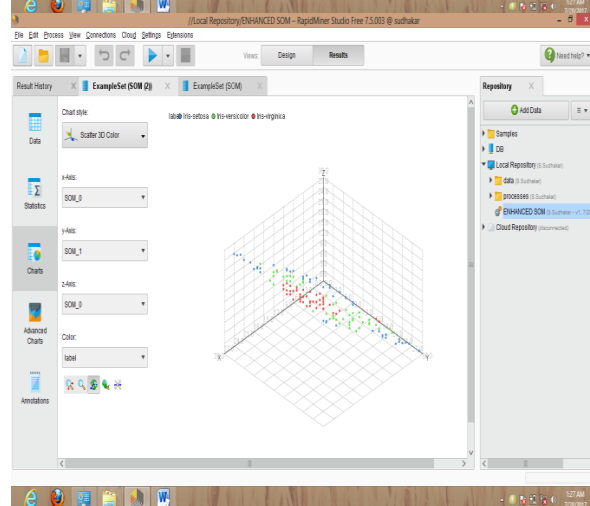
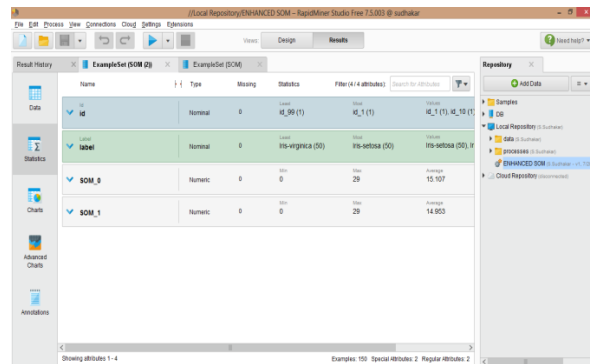
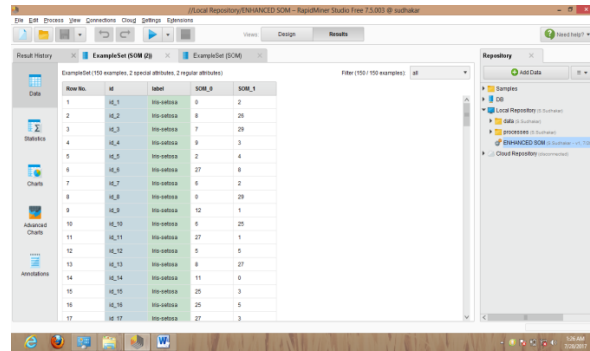
Transaction identification The transactions are identified using maximal forward references. Each time a backward reference is made, a transaction is identified. A new forward reference indicates the next transaction for that session.

3.3 Web Page Mapping

K-Means Clustering After the user sessions and transactions are identified, we make a two-dimensional array in which each row is arranged for a transaction and each column is for a URL. Initially, the URLs that appear in a transaction are set to one in the corresponding row, and rest values are set to zero. Initially, k transactions are selected at random for the k clusters. Then the means of the k clusters will be calculated. Afterwards, the distance between every transaction and the k clusters is calculated using the means of the k clusters. A transaction will be grouped into the cluster to which the distance is the shortest. For each of these k clusters, we sum up the values of each column and calculate its new mean. The mean values are used as the weights for the groups, which are used to indicate the similarity between groups. The algorithm will be repeated until the weights become stable. SOM The k groups of transactions and the set of unique URLs are the input to the SOM network. The input is represented by a two-dimensional m by k matrix, where m is the number of unique URLs and k is the number of transaction groups.

4 Experimental Results

We used Web log file for October, 2006 from the our test data. The data size is about 30MB with about 300,000 entries. Table 1, 2 and 3 show the example of user identifications, session identifications, and transaction identifications. The number of unique URLs generated by preprocessing is 188. We used a fixed value of 20 as the number of clusters, so the input to the SOM network is a 188 by 20 array. We have tested different parameters for the SOM network as follows: α varies from 0.2 to 0.9 and!



- [9] M. Nakagawa and B. Mobasher. A hybrid web personalization model based on site connectivity. In *Proceedings of the International Workshop on Web Knowledge Discovery and Data Mining*, pages 59–70, 2003.
- [10] O. Nasraoui and C. Petenes. Combining web usage mining and fuzzy inference for website personalization. In *Proceedings of the International Workshop on Web Knowledge Discovery and Data Mining*, pages 37–46, 2003.
- [11] Kate A. Smith and Alan Ng. Web page clustering using a self-organizing map of user navigation patterns. *Decision Support Systems*, 35(2):245–256, 2003.
- [12] Jaideep Srivastava, Robert Cooley, Mukund Deshpande, and Pang-Ning Tan. Web usage mining: Discovery and applications of usage patterns from web data. *SIGKDD Explorations*, 1(2):12–23, 2000.
- [13] Zhong Su, Qiang Yang, Hong-Jiang Zhang, Xiaowei Xu, and Yu-Hen Hu. Correlation-based document clustering using web logs. In *34th Hawaii International Conference On System Sciences*, pages 5022–5027, Hawaii, 2001. IEEE Computer Society.
- [14] A. Ypma and T. Heskes. Categorization of web pages and user clustering with mixtures of hidden Markov models. In *Proceedings of the International Workshop on Web Knowledge Discovery and Data Mining*, Edmonton, Canada, 2002.
- [15] S. Sharma, M. Varshney, “An Efficient approach for web log mining using ART”, *International Conference on Education and Management Technology*, 2010 (ICEMT 2010).
- [16] Zhang Y., X. Yu, and J. Hou, “Web communities: Analysis and construction,” *Berlin Heidelberg*, 2006.
- [17] S. Chakrabarti, M. van den Berg, and B. Dom, “Focused crawling: A new approach to topic-specific web resource discovery,” presented at the *8th World Wide Web Conf.*, Toronto, ON, Canada, May 1999.
- [18] N. Tyagi, A. Solanki and S. Tyagi, “An algorithmic approach to data preprocessing in web usage mining”, *Int. journal of information technology and knowledge management*, July-December 2010, Volume 2, No. 2, pp. 279-283.
- [19] R. Cooley, B. Mobasher, and J. Srivastava. “Web mining: Information and pattern discovery on the World Wide Web”, *Technical Report TR 97-027*, University of Minnesota, Dept. of Computer Science, Minneapolis, 1997.
20. Aguilar J.S, Ruiz R, Riquelme J.C and Giráldez R, (2001) SNN: A Supervised Clustering Algorithm, in: *14th Int. Conf. on Industrial & Engineering Applications of Artificial Intelligence & Expert Systems (IEA/AIE 2001): Lecture Notes in Artificial Intelligence*, Springer-Verlag, Budapest, Hungary, June 4–7, 2001, pp. 207–216.
21. Cooley R. (2000) Web Usage Mining: Discovery and Application of Interesting Patterns from Web Data. PhD thesis, University of Minnesota, May 2000.
22. Graham J and Starzyk J, (2008) A Hybrid self Organizing Neural Gas Network, *IEEE World Conference on Computational Intelligence (WCCI'08) June 1-6*
23. Jirayusakul A and Auwatanamongkol S (2000) A supervised growing neural gas algorithm for cluster analysis.
24. Kohonen T, (1995). Self-Organizing Maps, Berlin, Germany: Springer,
25. Martinet Mz, Berkovich S and Schulten K, (1993). Neural-gas network for vector quantization and its application to time series prediction, *IEEE Trans. Neural Networks* 4(1993)558-569.
26. Pedrycz W and Vukovich G, (2004) Fuzzy clustering with supervision, *Pattern Recognition* 37(7), 1339–1349.
27. Qu Y and Xu S, (2004) Supervised cluster analysis for microarray data based on multivariate Gaussian mixture, *Bioinformatics* 20(12), 1905–1913.
28. Slonim N and Tishby N, (1999) Agglomerative information bottleneck, in: *Proceedings of the 13th Neural Information Processing Systems*, (NIPS).
29. Sonali Muddalwar and Shashank Kavar, (2012) Applying Artificial neural network in Web Usage Mining, *International Journal of Computer Science and Management Research* Vol 1 Issue 4 November 2012.
30. Yu F, Sandhu K, and Shih M. (2000) A generalization-based approach to clustering of web usage sessions. In *Proc. of the 1999 KDD*

