



Introduction

The aim of this study is to test data mining techniques on 'Big data'. In this regards this project was done using the proxy server log entries of an institution to identify user groups based on their Internet access patterns. The proxy log entries specify the detail of each request for a resource on the Internet and the outcome of their request. Big data is a term used to describe the exponential growth and availability of data, both structured and unstructured. These data sets are so large or traditional processing complex that data applications are inadequate.

Motivation

Extract meaningful data from large quantity of proxy log entries and study on applying data mining techniques

Big data analysis is an emerging field and many organisations are exploring its potential. This study would provide a better understanding of this area and would show ways to analyse the data collected from day to day activities of an organization.

Objective

- \checkmark Identify useful web usage patterns using proxy access logs and to enhance the performance of a proxy server to provide a better service to **Users**
- \checkmark Providing a case study for the application of Big data analytics and data mining techniques

Major issues

Proxy log files contains vast number of missing values and blank entries, making it difficult to read correctly. Several methods were tested to handle missing and erroneous values. As part of this research work several tools and methods were studied and utilised such as Apache Hadoop framework, Java libraries, H2O predictive analytic tool and Matlab. One major problem in data analysis is handling numerical and categorical data together. Several approaches to handle different data types were tested and k-means and k-mediods clustering techniques were applied to cluster the data.

Web User Access Pattern Analysis on Proxy Log Data

A. Ann Sinthusha (annsinthu21@gmail.com) & E. Y. A. Charles Department of Computer Science, Faculty of Science, University of Jaffna



According to the institutions LAN design, each client IP addresses are replaced with an easily understandable label – virtual LAN ID, to which that IP address belongs.

Date a	nd 1	Time			
Before		Aft		Miss bec	
25/May/2014:04:02:13	25	May	2014	Slot-4	bee
he selected date_tin late and time. Then nonth, year, and day and the time value a ime slot labels	ne v from info re c	values of the da rmation onverte	are divi te valu are ex d as o	ided as le date, atracted ne hour	Thre of o Elbo
Respo	nse	Time			
Before			After		
25172		Above_AVG			
Aumerical values eplaced with the lo lata points having average value and points having the va value	of i bel the Belo alue	Respons Above value w_AVG below	e Tim _AVG abov for th the a	e are for the ve the e data verage	
Rep	oly S	ize			
Before	After				The
304		Below_AVG			Ma
he numerical value eplaced with the locate lata points having overage value and points having the value	es Ibel the Belo alue	of Rep Above value w_AVG below	oly Siz _AVG abov for th the a	te are for the ve the e data verage	tec big onl Pro reto
Request URL					
Before			After		
http://185.8.105.27/din.asp	(?	185	.8.105.27	,	
The hostnames are ex	(trac	ted fron	n each	URLs.	clu dc
Squi	d Sta	atus			
Poforo			Aftor		

Before	After
DIRECT/185.8.105.27	DIRECT

Squid status appeared as 'None' when it accessed from the squid and 'Direct/ <IP address of the site>' when it accessed from the hosted server location. From these data only the 'Direct' and 'None' have been considered as the attribute value

ata. 7, No. 3 ✓ en.wikipedia.org



Removal of Error & Missing Values ing values and error data which obtained ause of the advertisement's pop-ups, have n replaced with "null" string.

Clustering

ee clusters were identified form a sample set data using the Matlab tool K-Medoids the ow metric. Further analysys is necessary to erstand the meaning of the clusters.



Discussion & Conclusion

ere are some tools like H2O and Apache ahout provide facilities to apply data mining chniques like clustering and classification on data. But these all limited to numerical data ly and failed to identify any cluster.

posed pre-processing methods were able to ain the knowledge

medioid was able to identify three clusters on sample data set.

Future Scope

ocess the whole data to obtain meaningful usters rather than processing a sample set of

References

✓ V.Chitraa and Dr. Antony Selvdoss Davamani , 2010. A Survey on Preprocessing Methods for Web Usage Data. (IJCSIS) International Journal of Computer Science and Information Security, Vol.

✓ hadoop.apache.org ✓ 0xdata.com