Journal of Chemical and Pharmaceutical Sciences

Analysis of lead in river water quality data using clustering

S. Brintha Rajakumari, C. Nalini

Dept. of CSE, Bharath University, Chennai.

*Corresponding author: E-mail: brintha.ramesh@gmail.com

ABSTRACT

Data Mining is an art and science of extracting hidden information from the large datasets. Water pollution assessment is a critical area of study because it directly affects the human beings. The simple k Mean, EM and cobweb clustering algorithms are used to analyze the lead contaminant in the river water quality dataset. The simple KMean and hierarchical clustering methods takes less time to build the model than the other methods. This paper presents the comparison of different clustering algorithm used for the river water quality dataset and predicts the cluster, which contains the highest value of lead in river water.

KEY WORDS: clustering, Data mining, River Water Quality, Weka tool.

1. INTRODUCTION

Due to the huge amount of data being created from normal business operations, there is a demand to take information in modern data repositories to make effective decisions. Data mining is a data driven approach and applies statistical and machine intelligence tools to extract useful, exact and meaningful patterns from the large volume of data sets. Data mining tools are used every day to solve the real problem in business, engineering and science.

Water quality monitoring helps in drawing the water quality trends and prioritizing pollution control efforts and calculate approximately the environment and the level of pollution control required and effectiveness of pollution control measures. The paper prepared the heavy metal pollution index of the river in India and applied multivariate statistical analysis, factor analysis and cluster analysis to predict the source of heavy metals. The cluster analysis is used for spatial characterization to evaluate the most important parameters that affect the water health of the river.

The analysis of complex data sets with enhanced evaluate the water quality and a variety of environmental issues, including study the spatial and temporal patterns of water quality, chemical group associated with hydrological conditions, amount of pollution sources, the principal component analysis (PCA) and cluster analysis (CA) have been widely used.

2. MATERIALS AND METHODS

Analysis of lead in river water: Weka tool contains lots of algorithms to do the research. We have used WEKA tool to analyze the river water quality data. First, we have upload the lead dataset which contains water quality site name, river name, data of water taken from the river and the numerical value of the lead metal in the river. Totally there are 67 observations in the lead dataset that are all exceeded limit of lead content in the river water. In the dataset 12 data were taken from the period of Feb 2012, 26 data from June 2012, 12 data from Nov 2011, 10 for the period of Oct 2012 and 4 for the period of Mar 2013. The initial lead data sets are loaded into the wake tool is in figure1 and the arff file is in figure.2.



Figure.1. Lead Data in Weka Tool

Clustering methods: The clustering algorithm divides the huge data into a group of similar objects. We have used four different algorithms in this paper to analyze the dataset. They are Simple KMean Clustering Algorithm EM Clustering Algorithm Cobweb Clustering Algorithm Hierarchical Clustering Algorithm

	ISSN: 0974-2115
Journal of Chemical a	nd Pharmaceutical Sciences

No.	1: WQSite Nominal	2: River Nominal	3: Period Nominal	4: Pb Numeric	
1	A.P.Puram	Chittar	01-02-2012	11.77	-
2	A.P.Puram	Chittar	01-06-2012	13.24	ίΠ
3	Addoor	Gurupur	01-02-2012	12.63	
4	Agra(P.G.)	Yamuna	01-06-2012	24.58	
5	Anandpur	Baitarni	01-11-2011	14.3	1.1
6	Andhiyar	Hamp	01-10-2012	15.35	17
7	Ankinghat	Ganga	01-02-2012	10.79	
8	Ankinghat	Ganga	01-06-2012	12.33	
9	Ankinghat	Ganga	01-10-2012	10.11	
10	Balrampur	Rapti	01-06-2012	10.66	
11	Balrampur	Rapti	01-10-2012	21.64	
12	Bareilly	Ramga	01-11-2011	10.1	
13	Bareilly	Ramga	01-06-2012	22.29	
14	Basantpur	Mahanadi	01-10-2012	12.51	
15	Basti	Kwano	01-02-2012	14.13	
16	Basti	Kwano	01-06-2012	11.79	
17	Bhitaura	Ganga	01-06-2012	11.34	
18	Birdghat	Rapti	01-02-2012	22.75	
19	Birdghat	Rapti	01-06-2012	10.3	
20	Dabri	Ramga	01-06-2012	12.72	
21	DelhiRlyB	Yamuna	01-11-2011	13.35	
22	DelhiRlyB	Yamuna	01-02-2012	18.31	
23	DelhiRlyB	Yamuna	01-06-2012	27.45	
24	Dholai	Rukni	01-11-2011	13.59	

Figure.2. ARFF file of Lead Data

3. RESULT AND DISCUSSION

Weka tool is used to cluster the river water quality data [8] is in figure 3. In the diagram 4, X axis represent the period of river water taken to monitor the status of the lead metal and y represent the value of the lead.



Figure.3. Clustering lead data



Figure.4. Visualization of lead data using EM

Time taken to build a model for training set is 0.03 seconds in EM clustering algorithm and the incorrectly classified instance is 42.0 and the percentage is 62.6866 %. The visualization of lead data using an EM clustering algorithm is in figure 4. Time taken to build a model for full training data is 0.03 seconds in Cobweb clustering algorithm and the incorrectly clustered instances is 61.0 and percentage 91.0448 %. The visualization of the cobweb clustering algorithm is in the figure.5.



Figure.5. Visualization of lead data in Cobweb

Plot : 17:48:59	-	SimpleKMeans	(pbdata-weka.filters.unsupervised.attribute.Remove-R1)
Instance: 47			
Instance_number	÷	46.0	
WQSite	;	Moradabad	
River	÷	Ramganga	
Period	÷	01-06-2012	
Pb	:	48.92	
Cluster	÷	cluster0	

Figure.6. The highest value of lead using SimpleKeans

The KMean clustering algorithm takes Time to build a model for full training data is 0.02 seconds and incorrectly clustered instances is 47.0 and the percentage is 70.1493 %. Similarly Time taken to build model the full training data is 0.02 seconds and incorrectly clustered instances is 38.0 and the percentage is 56.7164 %. The figure 6 and 6 contains the highest lead river contamination of the river Ramganga in June 2012 and the lead amount is 48.92 is found in the cluster 1.

Table.1. Time and mediteet classified instance					
Algorithm	Time in seconds	Incorrect classified instance			
Cobweb	0.03	91.04%			
EM	0.03	62.68%			
Hierarchial Cluster	0.02	56.71%			
Simple K Means	0.02	70.14%			

Table.1.Time and Incorrect classified Instance

Comparision of Algorithms



Figure.8.Comparison of algorithms

We have compared the time taken to build the model using training dataset. From that analysis the SimpleKMean and Hierarchical Clustering method takes less time than the Cobweb and EM clustering methods.

4. CONCLUSION

Any country river is the main source of water. Degradation of river water quality affects the health status of human beings. Lead metal causes reduce hemoglobin and kidney damage. We have applied data mining technique to predict the value of lead in Indian River water quality dataset. We have used cobweb, EM, simpleKMean and Hierarchical clustering algorithm to group the data and compared the time taken to cluster the lead data and also displayed the highest lead in river water quality dataset.

REFERENCES

Achudhan M, Prem Jayakumar M, Mathematical modeling and control of an electrically-heated catalyst, International Journal of Applied Engineering Research, 9 (23), 2014, 23013.

Bhardwaj RM, Water Quality Monitoring In India-Achievements And Constraints, Iwg-Env, International Work Session on Water Statistics, Vienna, 2005, 20-22.

Gopalakrishnan K, Sundeep Aanand J, Udayakumar R, Electrical properties of doped azopolyester, Middle - East Journal of Scientific Research, 20 (11), 2014, 1402-1412.

Gopinath S, Sundararaj M, Elangovan S, Rathakrishnan E, Mixing characteristics of elliptical and rectangular subsonic jets with swirling co-flow, International Journal of Turbo and Jet Engines, 32 (1), 2015, 73-83.

Ilayaraja K, Ambica A, Spatial distribution of groundwater quality between injambakkam-thiruvanmyiur areas, south east coast of India, Nature Environment and Pollution Technology, 14 (4), 2015, 771-776.

Juahir H, Zain SM, Yusoff MK, Hanidza TIT, Samah MAA, Toriman ME and Mokhtar, Spatial water quality assessment of Langat River Basin (Malaysia) using environmetric techniques. Env. Monit. Ass, 173 (1-4), 2010, 625-641.

Kazi TG, Arain MB, Jamali MK, Jalbani, N, Afridi, HI, Sarfraz, RA., et al. Assessment of water quality of polluted lake using multivariate statistical techniques: A case study, Ecotox Environ Safe, 72, 2009, 301-309.

Kerana Hanirex D, Kaliyamurthie KP, Kumaravel A, Analysis of improved tdtr algorithm for mining frequent itemsets using dengue virus type 1 dataset: A combined approach, International Journal of Pharma and Bio Sciences, 6 (2), 2015, 288-295.

Kumar Manoj, Pratap Kumar Padhy and Shibani Chaudhury, Study of Heavy Metal Contamination of the River Water through Index Analysis Approach and Environmetrics, Bull. Environ. Pharmacol. Life Sci, 1 (10), 2012, 7-15.

Lingeswaran K, Prasad Karamcheti SS, Gopikrishnan M, Ramu G, Preparation and characterization of chemical bath deposited cds thin film for solar cell, Middle - East Journal of Scientific Research, 20 (7), 2014, 812-814.

Omo-Irabor OO, Olobaniyi SB, Oduyemi K, Akunna, J. Surface and ground water quality assessment using mutivariate analytical methods: a case study of the Western Niger Delta, Nigeria, Phys Chem Earth, 33, 2008, 663-673.

Ouyang Y, Evaluation of river water quality monitoring stations by principal component analysis. Water Res, 39, 2005, 2621-2635.

Premkumar S, Ramu G, Gunasekaran S, Baskar D, Solar industrial process heating associated with thermal energy storage for feed water heating, Middle - East Journal of Scientific Research, 20 (11), 2014, 1686-1688.

www.jchps.com

Journal of Chemical and Pharmaceutical Sciences

Status of Trace and Toxic Metals in Indian Rivers, River Directorate River Data Directorate Planning Organisation Planning & Development Organisation, 2014.

Sukhdev Kundu, Categorization of Pollution Load in surface Water System using Multivariate Techniques, Advances in Bioresearch, 3 (1), 2012, 64 - 68.

Sundar Raj M, Saravanan T, Srinivasan V, Design of silicon-carbide based cascaded multilevel inverter, Middle - East Journal of Scientific Research, 20 (12), 2014, 1785-1791.

Thooyamani KP, Khanaa V, Udayakumar R, Application of pattern recognition for farsi license plate recognition, Middle - East Journal of Scientific Research, 18 (12), 2013, 1768-1774.

Thooyamani KP, Khanaa V, Udayakumar R, Efficiently measuring denial of service attacks using appropriate metrics, Middle - East Journal of Scientific Research, 20 (12), 2014, 2464-2470.

Thooyamani KP, Khanaa V, Udayakumar R, Partial encryption and partial inference control based disclosure in effective cost cloud, Middle - East Journal of Scientific Research, 20 (12), 2014, 2456-2459.

Thooyamani KP, Khanaa V, Udayakumar R, Wide area wireless networks-IETF, Middle - East Journal of Scientific Research, 20 (12), 2014, 2042-2046.

Thooyamani, K.P., Khanaa, V., Udayakumar, R., Using integrated circuits with low power multi bit flip-flops in different approch, Middle - East Journal of Scientific Research, 20 (12), 2014, 2586-2593.

Thooyamani, KP, Khanaa V, Udayakumar R, Virtual instrumentation based process of agriculture by automation, Middle - East Journal of Scientific Research, 20 (12), 2014, 2604-2612.

Udayakumar R, Kaliyamurthie KP, Khanaa, Thooyamani KP, Data mining a boon: Predictive system for university topper women in academia, World Applied Sciences Journal, 29 (14), 2014, 86-90.