

Pattern Classification using Web Mining

Krunal Jaha¹ Krishna Patel² Krishna Punwar³ Mr. Romil Patel⁴ Mr. Jignesh Tilva⁵

^{1,2,3}B.E. Student ^{4,5}Assistant Professor

^{1,2,3,4,5}Department of Information Technology

^{1,2,3,4,5}Sigma Institute of Engineering, Vadodara, India

Abstract— It is the application of data mining techniques on the web data to solve the problems to extracting useful information. As the information in the internet increases, the search engines lack the efficiency of providing relevant and required information. This project proposes an approach for the web content mining using the algorithm. The aim of our project is the pattern classification of dataset and analysis the data on E-commerce business or organization websites.

Key words: Pattern Classification, Web Mining, C 4.5 Algorithm

I. INTRODUCTION

Web mining is type of the data mining techniques to discover some interesting patterns from the web data. Classify interesting pattern from using some classification algorithm and techniques. [5] In this paper we read some papers which are related to the pattern classification from the web data. There are three types of web mining. Web content mining, web structure mining, and web usage mining.

Classification algorithms are used for the complex and real web lo data.

Pattern classification is the technique that discover the pattern from the dataset which we used in it. Here we used the amazon dataset for discovering the pattern from the dataset and we classified the dataset. We find the pattern that how many customer uses Cash on delivery to pay for the product and how many customer pay by credit/debit card for particular product.

It helps the e-commerce companies for productivity flow, e-business depends on the information and the data which we conclude to take right decision.

Our system will improve the business of any e-commerce websites or e-commerce companies.

II. PATTERN CLASSIFICATION USING WEB MINING

Classification is the technique to classify the dataset and find the particular pattern from the classification technique which we apply for the system. After identifying the dataset there are various fields like IP addresses, use session, pin code, product id, and buy id and payment method. There are various kinds of pattern techniques to display the result and there are also different types of the classification algorithm for classification. When dataset is cleaned than after apply a c4.5 classification algorithm and then after we found particular pattern which we needed from the data. We used classification analysis, data items are classified according to predefined categories. [1]

In our work there are web log data which we divided in particular session and divided by the zip code or pin code. We classify the IP addresses and divided by the payment method either COD (cash on Delivery) or payment by card (credit card, debit card). From this type we easily understand the dataset and we work on the proper pattern.

III. PROPOSED SYSTEM

In proposed methodology for classification of web data in order some predefine our criteria. In this mode we present the steps of the system. In our proposed system there are some types in the system which are data modification, data integration, data cleaning, apply classification algorithm and then last one is the pattern analysis. Data cleaning and the apply classification algorithm is the main part of the system. Data cleaning is the process that we remove the inconsist data and noisy data than also we clean the field that we don't need. After all the process pattern analysis is the term and technique to describe the pattern and shows from it.

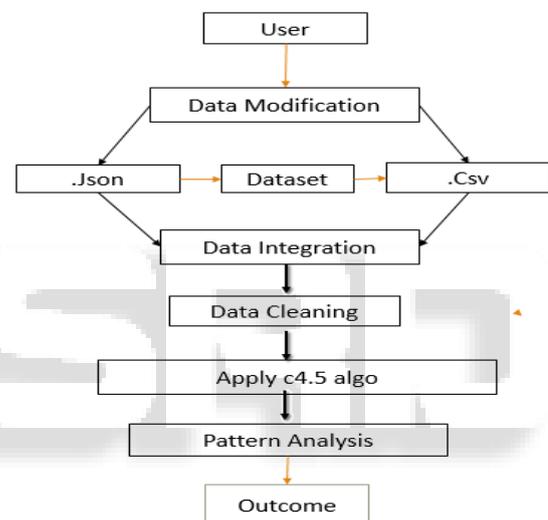


Fig. 1: Proposed System

A. User

User is the person the uses the system and work on the dataset and follow the particular pattern or system which we created.

B. Data Modification

Data modification is term that there are lots of data in the dataset and there are data which is in the .json format so we convert the dataset from .json to .csv format.

C. Data Integration

Data integration means that the data has been integrate with another dataset like we merge the dataset into one single dataset and we work on it.

D. Data Cleaning

Data cleaning is the term and the technique that remove the inconsist and the noisy data from the dataset. It's also called as remove the error from the dataset. We work on the dataset that there are different types of fields in it so we can remove the fields which we don't need from it. We create our tool for the data cleaning process and remove the error from the data.

Data cleaning is the process to remove error data and inconsistent data. There are some missing values which can be removed from the data cleaning process and then after it reduces the size of the dataset. [1]

IV. C4.5 CLASSIFICATION ALGORITHM

Input: training dataset T; attributes S.
Output: decision tree
1: if T is NULL then
2: return failure
3: end if
4: if S is NULL then
5: return Tree as a single node S
6: end if
7: if |S|=1
8: return Tree as a single node S
9: end if
10: set Tree
11: for a ∈ S do
12: set Info (a, T) = 0, and Split Info (a, T) = 0
13: compute Entropy (a)
14: for v ∈ values (a, T) do
15: set Ta, v as the subset of T with attribute a= v
16: Info a, T) += - |Ta, v| / |Ta| Entropy (av)
17: Split Info (a, T) += - |Ta, v| / |Ta| log |Ta, v| / |Ta|
18: end for
19: Gain (a, T) = Entropy (a) – Info (a, T)
20: Gain Ratio (a, T) = Gain (a, T) / Split Info (a, T)
21: end for
22: set abest = argmax {Gain Ratio (a, T)}
23: attach abest into Tree
24: for v ∈ values (abest, T) do
25: call C4.5 (Ta, v)
26: end for
27: return Tree.

This algorithm used for the classification. It generates the tree at the end of it. [4]

Here Entropy calculation,

$$Entropy(S) = \sum_{i=1}^c - p_i \log_2 p_i$$

Where, pi – proportion of S, [2]

Information Gain Calculation,

$$Gain(S,A) = Entropy(S) - \sum_{v \in Value(A)} \frac{|S_v|}{|S|} Entropy(S_v)$$

Where, values (A) – set of all possible values of A

Sv – subset of s. [2]

V. PROBLEM STATEMENT

Recently lots of users or people uses the internet they surfing over it and find some interesting things like images and videos. Online shopping is one of the most important part of the internet that most of people buy or sell products online and they done their payments through Cash on Delivery and Credit card or Debit card. Most of people like to buy products online but there are some issues like every users not feel safe for online payment through credit card and debit card they preferred Cash on Delivery. But some products when they buy, COD is not available for such locations, and users not buy that products. So we have

dataset of the users preferred COD or Credit card/Debit card and which places or how much cities the ratio of the payment methods. So we can find the pattern and it might be helpful for such e-commerce websites. Our aim to pattern classification using the web mining data.

VI. EXPERIMENT RESULT

The above C4.5 classification algorithm is implemented in the java using NetBeans 8.0.2 and database work in MySQL.

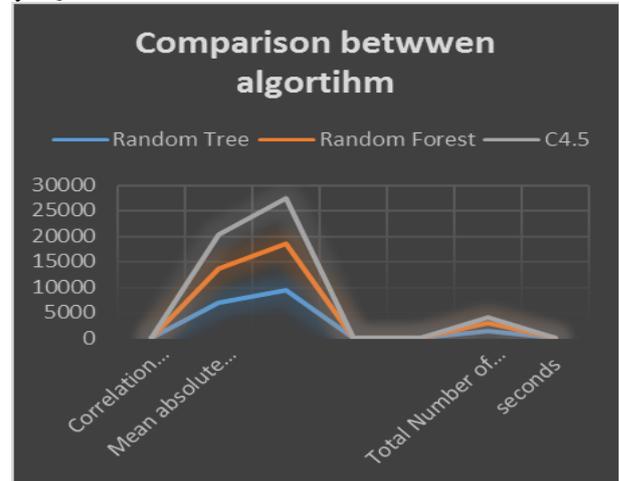


Fig. 2: Comparison Between Other Algorithms

We compare our algorithm with other two algorithms to prove that our algorithm is much better than the other algorithm and it saves the time. We also work on the different clustering and classification algorithm for better result.

Classification Algorithm	Random Tree	Random Forest
Correlation coefficient	0.0588	0.2161
Mean absolute error	6790.8325	6704.9141
Root mean squared error	9315.7272	9111.1842
Relative absolute error	100.22%	98.96%
Root relative squared error	101.09%	98.87%
Total Number of Instances	1364	1364
Seconds	1.66	0.03

Table 1: Classification Algorithm Result

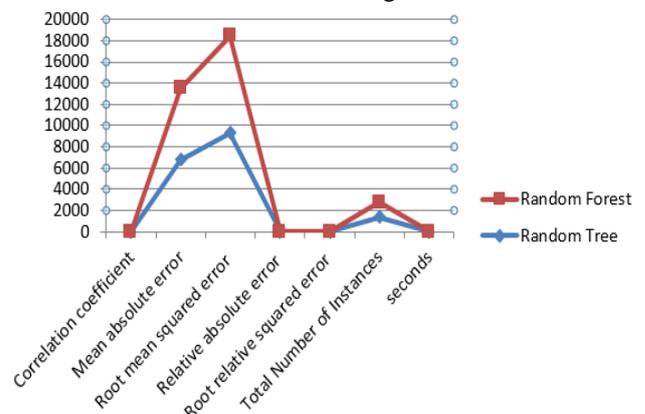


Fig. 3: Displaying Classification Result

Clustering Algorithm	simple k means	Make Density Based Cluster
0	983	945
1	381	419
seconds	0.02	0.03

Table 2: Clustering Result

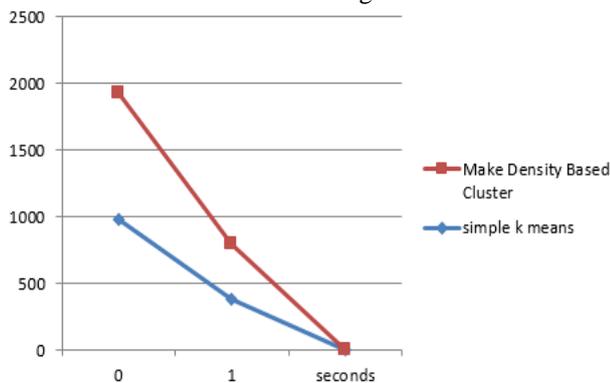


Fig. 4: Displaying Result Of Clustering Algorithm

VII. CONCLUSION & FUTURE WORK

In the report to study about Pattern classification using web mining. Also study C 4.5 classification algorithm. In this project we are working on E-commerce Dataset, payment method and zip code that where the cod is placed and where card payment made. In this project we will done pattern classification using the new algorithm which named as c4.5 classification algorithm. This system is very useful for the E-commerce companies. In future we will implement the system. Pattern classification is in tried to the implementation with new ideas. We have already implemented new approach of classification. We will tried to use new classification algorithm and uses for the system for the better result and better classification.

REFERENCES

- [1] Er. Romil V Patel and Dheeraj Kumar Singh, "Pattern classification based on Web Usage Mining Using Neural Network Technique", IJCA, vol.71-No.21, June 2013.
- [2] Anurag Upadhyay, suneet shukla and sudsanshu kumar, "Empirical comparison by data mining classification algorithm (C 4.5 & C 5.0) for thyroid cancer data set", International journal of computer science & communication network, vol.3(1), 64-68.
- [3] Rutvija Pandya and Jayati Pandya, "C 5.0 algorithm to improved decision tree feature selection and reduced error pruning", IJCA, vol. 117 – NO.16, May 2015.
- [4] A. S. Galathiya, A. P. Ganatra and C. K. Bhensdadia, "Improved Decision Tree Induction Algorithm with Feature Selection, Cross Validation, Model Complexity and Reduced Error Pruning", IJCSIT, vol.3(2), 2012,3427-3431
- [5] Monika Yadav and Mr. Pradeep Mittal, "Web Mining – An Introduction", IJARCSSE, vol. 3, 3 March 2013.
- [6] Priyanka sharma and manavjeet kaur, "Classification in pattern recognition: A Review", IJARCSSE, vol. 3, 4 April 2013.
- [7] Sonali sonksusare and Mr. Jayesh Surana, "A survey on different techniques for data classification and

information extraction from the websites", IJARCSSE, vol. 3, 11 November 2013.

- [8] Govind murari Upadhyay and kanika dhingra, "web content mining: its technique and usage", IJARCSSE, vol. 3, 11 November 2013.
- [9] Lalani, A.S., "Data mining of web access logs", School of Computer Science and Information Technology. Royal Melbourne Institute of Technology. Melbourne, Victoria, Australia, 2003.
- [10] J. Srivastava, R. Cooley, M. Deshpande and P-N. Tan (2000). "Web Usage Mining: Discovery and Applications of usage patterns from Web Data", SIGKDD Explorations, Vol 1, Issue 2.
- [11] <http://en.wikipedia.org/wiki/classification>
- [12] Marisa S. Viveros, John P Nearhos, Michael J. Rothman," Applying Data Mining Techniques to a Health Insurance Information System", 2003.
- [13] Rulequest Research, "data mining tools see5 and c5.0", <http://www.rulequest.com/see5-info.html>, 1997-2004.
- [14] Thair Nu Phyu, "Survey of Classification Techniques in Data Mining", International MultiConference of Engineers and Computer Scientists 2009 Vol I IMECS 2009, March 18 - 20, 2009, Hong Kong.