# Hyper-Graph Based Documents Categorization On Knowledge From Decision Trees

[1]R.Merjulah

*Address*
[1]Vinayaka Missions University,
*Salem, Tamilnadu, India*

## Abstract

This document has devised a novel representation that compactly captures a Hyper-graph Partitioning and Clustering of the documents based on the weightages. The approach we take integrates data mining and decision making to improve the effectiveness of the approach, we also present a NeC4.5 decision trees. This algorithm is creating the cluster and sub clusters according to the user query. This project is forming sub clustering in the database. Some of the data's in the database may be efficient one, so we are clustering the data's depending upon the ability.

*Keywords:* Hyper-graph agglomerate algorithm, Clustering, Data Mining, NeC4.5 decision trees.

## 1. Introduction

A web document has a lot of features (or heterogeneous features) including content, anchor text, URL, hyperlink, user access log etc. The document classification can be done based on either probability methods or distance measures.

Database is a collection of records stored in a computer in a systematic way. For better retrieval and sorting, each record is usually organized as a set of data elements (facts).Computer program used to manage and query a database is known as a database management system (DBMS).

For a given database there is a structural description of the type of facts held in that database: this description is known as a schema. There are a number of different ways of organizing a schema, that is, of modelling the database structure: these are known as database models (or data models). The model in most common use today is the relational model. Other models such as the hierarchical model and the network model use a more explicit representation of relationships.

A great deal of the internal engineering of a DBMS, however, is independent of the data model, and is concerned with managing factors such as performance, concurrency, integrity, and recovery from hardware failures. In these areas there are large differences between products.

Data mining can be defined as "the nontrivial extraction of implicit, previously unknown, and potentially useful information from data" and "the science of extracting useful information from large data sets or databases". It is usually associated with a business or other organization's need to identify trends.

Data mining involved the process of analysing data to show patterns or relationships; sorting through large amounts of data; and picking out pieces of relative information or patterns that occur e.g., picking out statistical information from some data.

Data mining is often approximated via stepwise regression methods wherein the space of $2^k$ models. This procedure is called all subset or exhaustive regression. Data mining software is one of a number of analytical tools for analysing data. Technically, data mining is the process of finding correlations or patterns among dozens of fields in large relational databases.

Web mining and web usage mining is the application of data mining techniques to discover usage patterns from the web in order to better understand and serve the needs of users or web-based applications. The first web analysis tools simply provided mechanisms to report user activity as recorded in the servers. Using such tools, it was possible to determine such information as the number of accesses to server, the times or time intervals of visits as well as the domain names and the URLs of users of the Web server.

## 2. Proposed Technique

2.1 Deliverables:

Below is the list of deliverables that are expected when the implementation of proposed research work is completed. The proposed system provides an efficient way of multiple features (Keywords example java, html, and oracle) based processing that compactly captures f features based on weight ages.

- Efficient filtering of data
- Takes less computational time.
- Data accuracy can be maintained.
- Advantage of the proposed system is that it is lightweight.

## 2.2 Intellectual Challenges:

- Understanding the complex topic of hyper-graph based documents categorization and evaluating its performance in achieving the performance efficiency.

- An attempt is made to provide sufficient information in the remainder section's of this report and very much confident that reader will be convinced with the abilities and performances of distributed document categorization based on the hyper-graph.

## 2.3 Research Program:

Methodology:
There are three type methodology used:
1. Hyper-graph partition algorithm for clustering
2. NeC4.5-for decision making
3. Multi feature Query Techniques

- Techniques used for data collection

Academic Literature review of published papers to understand the various issues associated with designing and implementation of the protocol .Referring to journals, papers, and articles and to analyse various methods and framework used to model designing processes

## 3. System Implementation

There are four modules in the System implementation.
1. User Query
2. Hyper graph based clustering
3. Neural network processing
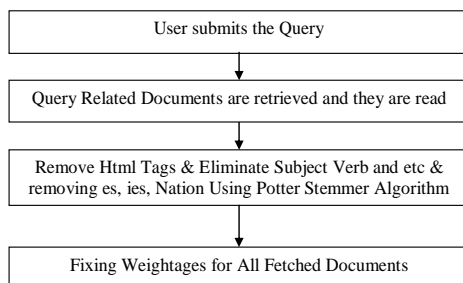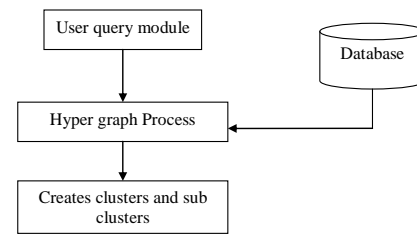4. Efficient information retrieval


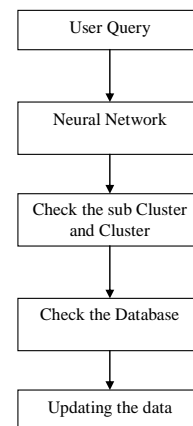
Fig. 1 User Query



Fig. 2 Hyper-graph Based Clustering


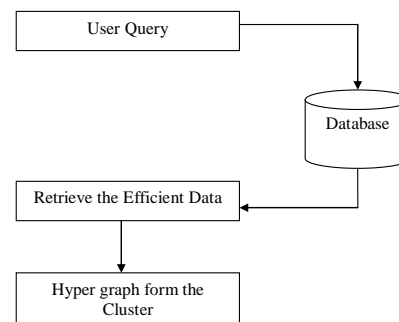
Fig. 3 Neural Network Process



Fig. 4 Efficient Information Retrieval

## 4. Literature Review

Extensive research in data mining has been done on discovering distributional knowledge about the underlying data. Models such as Bayesian models, decision trees, support vector machines, and association rules have been applied to various industrial applications such as customer relationship management (CRM).

IJCSI International Journal of Computer Science Issues, Vol. 9, Issue 2, No 1, March 2012
ISSN (Online): 1694-0814
www.IJCSI.org

317

The data sets are often cost sensitive and unbalanced. If we predict a valuable customer who will be an attritor as loyal, the cost is usually higher than the case when we classify a loyal customer as an attritor. Similarly, in direct marketing, it costs more to classify a willing customer as a reluctant one. Such information is usually given by a cost matrix, where the objective is to minimize the total cost.

Describes a new hyper-graph formulation for document categorization, where hyper clique patterns, strongly affiliated documents in this case, are used as hyper edges. Compared to frequent item sets, the objects in a hyper clique patterns have a guaranteed level of global pair wise similarity to one another as measured by the cosine or Jaccard similarity measure.

In similarity-based clustering, similarity combination is a common method to exploit the different features of an object. Proposed a new document similarity function based on both term similarity and hyperlink similarity factors. Conventional approaches to document categorization are generally based on either probability methods or distance measures. There are number of problems in using these methods for document clustering. First it is easy to define an appropriate distance measure in such a space. Second, the total number of words in the whole document dataset can be very large.

Due to the inherent sparsity of the high dimensional space, close points do not necessarily belong to the same document class, which is the basic assumption of many distance based clustering methods such as K-means.to overcome difficulties of high dimensional clustering, one promising approach is the Association Rule Hyper-graph Partitioning (ARHP) algorithm. It is used for clustering related items (e.g., documents) in transaction (e.g., Words) based databases.

Since clustering is mainly based on reasonable similarity, there may be advantages of using the hyperclique based clustering algorithm. In this paper, we investigate using hyperclique patterns as an alternative to frequent item sets and present an analytic comparison between Hyperclique based hyper-graph Partitioning (HYPA) and ARHP.

A major challenge in document clustering is the extremely high dimensionality. For example, the vocabulary for a document set can easily be thousands of words. On the other hand, each document often contains a small fraction of words in the vocabulary. These features require special handlings. Another requirement is hierarchical clustering where clustered document can be browsed according to the increasing specificity of topics.

## 4. Analysis and design

### 4.1 System Study:

Generally system involved the process of gathering facts, analysing and identifying the problem using it to improve the system implementation.

- Need for the proposed system

To overcome difficulties of high dimensional clustering, one promising approach is the Association Rule Hyper-graph Partitioning (ARHP) algorithm. It is used for clustering related items (e.g., documents) in transaction (e.g., words) based databases. The proposed system provides an efficient way of multiple feature based processing that compactly captures f features based on weightages after retrieving the data we form sub-clustering based on the keywords. For example the Java keyword for sub-cluster is RMI, applet etc.. After forming the sub-clustering using Nec4.5 algorithm is used to form decision trees, for fastest data retrieving and the loss of data missing is almost reduced.

## 5. Project Evaluation

The main objective of this is to evaluate the project functionality, tools used and the product developed. For the above discussed area of investigation I have found the various issues that affect the existence and present level of performance of Document categorization using hyper-graph. As its one major business that had great future.

I have discussed in detail about the implementation of the new system which overcomes the limitations, Difference From Previous Works In business marketing, most of the marketing planning activities have been done in a human-heavy process, which is carried out by hand. An important computational aspect is to segment a customer group into subgroup, often in terms a binary decision variable such as customer's willingness to buy a product or not to buy. Where the aim is to generate a ranking function for the customers sorted on their likelihood to buy a product, so that a "gain chart" of "lift chart" can be created for human analysis. Machine learning and data mining research has contributed to the business practice by addressing some new issues in marketing.

However, all the above research works are aimed at either finding a segmentation of the customer database, or deciding to take a predefined action for every customer based on that customer's current status. None of them addressed the issue discovering actions that might be taken from a customer database. To the best of our knowledge, ours is the first such work in machine learning and business marketing that addressed this action-discovery issue.
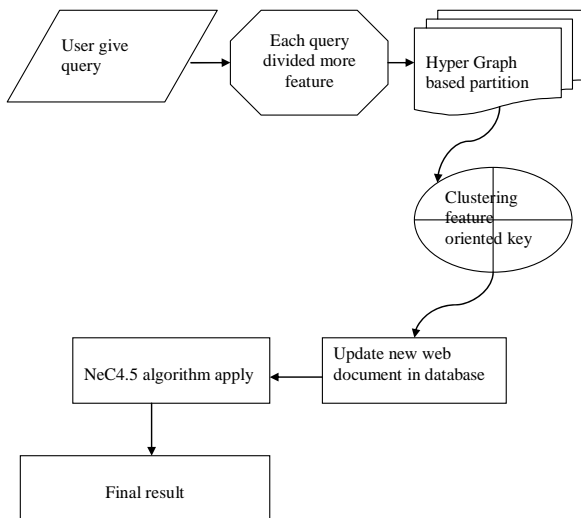
IJCSI International Journal of Computer Science Issues, Vol. 9, Issue 2, No 1, March 2012
ISSN (Online): 1694-0814
www.IJCSI.org

318

Fig. 5 Overall Block Diagram

## 6. Conclusion

We have proposed Multi type Features Co-selection for Clustering (MFCC), a novel algorithm to exploit different types of features to perform Web document clustering. We use the intermediate clustering result in one feature space as additional information to enhance the feature selection in other spaces. Consequently, the better feature set co-selected by heterogeneous features will produce better clusters in each space. After that, the better intermediate result will further improve co-selection in the next iteration. Finally, feature co-selection is implemented iteratively and can be well integrated into an iterative clustering algorithm. In this paper, we present a novel technique to take these results as input and product a set of actions that can be applied to transform customers from undesirable classes to desirable ones.

We presented a new approach. Hyperclique based hyper-graph Partitioning (HYPA), for document clustering. Since similarity among vertices in the hyper-edge is the basic assumption for hyper-graph partitioning, the hyperclique is a better candidate for hyper-edge than the frequent item set. Indeed, our experiments on real world document datasets validated the advantages of HYPA against ARHP in terms of various external clustering criteria.

## References

[1] D.S. Hochbaum, "Approximation Algorithms for Np-Hard Problems," chapter 3, p.136. PWS Publishing Company, 1995.

[2] Tian-Ming Hul, Ji Ouyang, Chao Qu, Sam Yuan Sung "Hyper-graph Based Document Categorization: Frequent Itemsets Vs Hypercliques", Dongguan University Of Technology, Dongguan, Guangdong 523808, China South Texas College, Mcallen, Tx78501, Usa E-Mail: Tmhu05@gmail.com Proceedings Of Sixth International Conference On Machine Learning Cybernetics, Hong Kong, 19-22 August 2007

[3] Ng, R.T., Han, J. Efficient and Effective clustering methods for spatial data mining. In: Bocca, J.B., Jarke, M., Zaniolo, c., eds. Proceedings of the 20th International Conference on Very Large Data Bases. Santiago: Morgan Kaufmann, 1994.144-155.

[4] Fayyad, U., Piatetsky-Shapiro, G., Smyth, P. Knowledge discovery and data mining: towards a unifying frame work. In: Simoudis, E., Han, J., Fayyad, U.M., eds. Proceedings of the 2nd International Conference on Knowledge Discovery and Data Mining. Portland Oregon: AAAI Press, 1996.82-

[5] B.C.M. Fung, K. Wang, and M.Ester. Large hierarchical document clustering using frequent itemsets. In Proc. 3rd SIAM Int. Conf. On Data Mining, 2003.

[6] Ester, M., Kriegel, H.P., Sander, J., et al. A density based algorithm for discovering clusters in large spatial databases with noise. In: Simoudis, E., Han, J., Fayyad, U.M., eds. Proceedings of the 2nd International Conference on Knowledge Discovery and Data Mining. Portland, Oregon: AAAI Press, 1996.226-231.