**IJESRT**

# INTERNATIONAL JOURNAL OF ENGINEERING SCIENCES & RESEARCH TECHNOLOGY

## PRODUCT CATEGORIES INTEGRATION BY USING DATA MINING

**Miss. Dhaygude Tejashri Mohan \*, Patil Bharati R., Prof. Mr. Bere S.S.**
*Student, Dept. of Computer Engineering, DGOI, COE, Daund, Pune ,Maharastra,India.
 Student, Dept. of Computer Engineering, DGOI, COE, Daund, Pune ,Maharastra,India.
 Assistant Professor, Dept. of Computer Engineering , DGOI, COE, Daund, Pune ,Maharastra,India.

## ABSTRACT

The major important task of online ecommerce based web portals and commerce search engine based application is the integration of products into its own categories. Categorization of products from the data provider into the master taxonomy and whereas make use of the data provider taxonomy information becomes major problem. To avoid this problem classify the products based on their textual based classifier and taxonomy-aware step that adjusts the results of a textual based classifier to ensure that products that are close together in the provider taxonomy remain close in the master taxonomy. If large data sets are used then integration of product become difficult so hierarchical clustering are used by using this system become efficient and scalable.
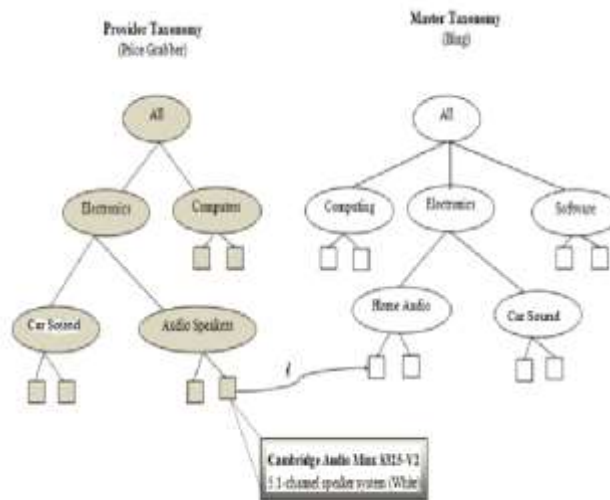
**KEYWORDS**: Categorization, Data Mining, Integration, Classification, Taxonomies.

## INTRODUCTION

There has been increased number of web portal give a user knowledge centered on online shopping. These web portals include several commercial sites and commerce search Engines. Hence data integration task is important for these commercial portals. The data integration task faced by these web portals is the integration of data coming from numerous data providers into a particular product catalog. This process is known as product categorization.All web portals maintain their own master taxonomy for organizing products and it is used for both online shopping and searching purposes. When a new product arrives from dissimilar providers, it should automatically categorize the products into the master taxonomy according to their structure. The provider taxonomy may be different from the master taxonomy, but in most cases, there is still a powerful signal coming from the provider classification. Intuitively, products that are in nearby categories in the provider taxonomy, should be classified into nearby categories in the master taxonomy.A taxonomy-aware processing step that adjusts the results of a text-based classifier to ensure that products that are close together in the provider taxonomy remain close in the master taxonomy. If large data sets are used then integration of product become difficult so hierarchical clustering are used by using this system become efficient and scalable.To illustrate this point, consider the example in the provider taxonomy is an excerpt from the taxonomy used by Shoppers Stop, and the master taxonomy is an excerpt from the taxonomy used by TACI searching site Shopping. Now, given a product tagged with a category from provider taxonomy(Price Grabber), we want to categorize it in the TACI searching site Master taxonomy(Bing). Suppose we are given the product "Audio Speakers" from the category Electronics Category in the Provider taxonomy and we use a text-based classifier to categorize this product into the TACI searching site taxonomy .suppose we are given the product "Cambridge Audio Minux S325-V2", it is unclear whether this product should be classified intoElectronics/car sound or Electronics/Audio speakers in provider(Price Grabber)Taxonomy. If we know that most of the products in the electronics category. If we use a text-based classifier to categorize this product into the Bing taxonomy,it is unclear whether this product should be classified into Electronics/Home Audio or Electronics/Car sound.Provider taxonomy are categorized  to Electronics category in Master Taxonomy, then we can conclude that most likely the new product should also be classified in Electronics category.However, for most products in the Electronics Category from the provider taxonomy, the classifier is actually unable to decide if they should be classified in Electronics. Therefore, we cannot use the categorization of the products in the same category as the new product to guide as for the correct decision. In this case, the taxonomy information can help to determine the correct categorization. As we go up the taxonomy tree of the Provider taxonomy, we observe that many more products in

the Car sound&Audio Speakers in Provider taxonomy(Price Grabber) are classified to the Electronics category, as opposed to the Computer category. Using this information we can conclude that most likely the products from the Electronics Category should be mapped to laptop rather than Computer. First use the text based classifier to adjust the results of taxonomy information. The text based classifier representation makes use of the taxonomy construction of the master and provider taxonomies in order to attain associations amongst the disparate categories in the taxonomy. Text based classifier the labeling problem occurs related to a diversity of optimization problems such as the metric labeling problem or structured prediction problems .The text based classifier model makes use of the taxonomy structure of the master and provider taxonomies in order to achieve relationships among the dissimilar categories in the taxonomy with the purpose of relying on the category membership information of the products. To overcome these problems we proposed hierierchical clustering and supervised learning. The major steps of the work as follows:

1. First originate the taxonomy-aware catalog integration difficulty as a structured prediction difficulty. In this method the approach that leverages the structure of the taxonomies in order to enhance catalog integration. .

2. Second describe the taxonomy aware classification process with two steps: In first step product are classified under base classification step, then use taxonomies aware processing steps.

3. During the taxonomy aware classification step the optimization problem or label classification problem have been overcomed with scalable algorithm for the taxonomyaware classification process.

4. Tuning the parameters of the k, Ө, ɾ is important for the performance of the system. supervised learning algorithm makes best result for classification result.



*Fig.1 A simple catalog integration example.*

## RELATED WORK

Agrawal and Srikant [1]. method scales to large data sets (like ours), introduce the problem of pervasive in web portal environment. The current technology for automating this process consists of building a classifier that uses the categorization of documents in the master catalog to construct a model for predicting the category of unknown documents. But many of the data sources have their own categorization, and the accuracy of classification can be improved by factoring in the implicit information in these source categorizations. Our solution insight is that numerous of the data sources have their have possession of categorization and classification accuracy can be improved by factoring in the implicit information in these source categorizations. It makes use of source category information, but treats the basis and target taxonomy as flat.

A. Fraser et. al and P. Ravikumar et. Al[2] provide the catalog integration problem as an optimization problem and this problem is stimulated by the metric labeling problem. The metric labeling problem aims to discover the optimal labeling of a number of objects consequently that they reduce an assignment and a separation cost.

Sarawagi et al.[3]provide a cross training model is established for document classification occurrence of multiple label sets.A document classifier is original trained using documents with preassigned labels or classes picked from a set of labels it is named as taxonomy or catalog. Once the classifier is trained, it is offered test documents for which it must guess the best labels.

D. Zhang and W.S. Lee [4] have also developed approaches to catalog integration by using boosting and transductive learning[5] anddiscuss a straightforward approach to automating the process of catalog integration would be to learn a classifier that can classify objects from the source taxonomy into categories of the master taxonomy. The key vision is that the availability of the source taxonomy data could be helpful to build better classifiers for the master taxonomy if their categorizations have some semantic overlap.

Nandi and Bernstein [6] proposed an approach for matching taxonomies based on query term distributions.Primary it perform the mapping at the taxonomy level,mapping category from the source to the target, while we achieve the mapping at the occurrence level by categorizing personality product instances to the target taxonomy.

C. Chekuri et. al [7] consider the method of finding a label at minimum cost where the cost of a labeling is determined by the pairwise relations between the objects is considered. A distance function on labels; the distance function is assumed to be a metric is used. Each object also incurs an assignment cost that is label, and vertex dependent.The problem captures many classification problems that arise in computer vision and related fields. The solution to the problem is obtained from a general formulationand that formulation allows us to extend the ideas to obtain the first non-trivial approximation for the truncated quadratic distance function.

G. Ifrim et. al [8] discuss a Bayesian logistic regression approach that uses a Laplace prior to avoid over fitting and produces sparse predictive models for text data. This approach is applied to a range of document classification problems and show that it produces compact predictive models at least as effective as those produced by other classifiers.Lasso logistic regression provides state-of-the-art text categorization effectiveness while producing sparse and thus efficient models. This approach is also useful in other high dimensional data analysis problems.

Glue [9], introduce the machine learning to learn how to map between ontologies; and Iliads , a system which makes use of machine learning and logical inference techniques to output alignments.

Ming Ji et. al [10] explain a Simple algorithm for semi-supervised learning that on one hand is easy to implement, and on the other hand is guaranteed to improve the generalization performance of supervised learning under appropriate assumptions. This is learned from the labeled examples the best prediction function that can be used for the parameter calibration and it can also be used to incrementally re-train the base classifier.

## CONCLUSION
In this paper we presented a well-organized and efficient and scalable approach to catalog integration that is based on the use of source category and taxonomy structure information. We also showed that this approach leads to substantial gains in accuracy with respect to existing classifiers. Supervised learning algorithm were used for retrain the base classifier during the product calibration step, they can also be used for other problems. The output of the parameter result as chosen might be second-hand as a feature for item identical,while would like to match elements classified under the master taxonomy to incoming offers from the providers.

## ACKNOWLEDGEMENTS

## REFERENCES:
1. R. Agrawal and R. Srikant, "On Integrating Catalogs," Proc.10th Int'l Conf. World Wide Web(WWW),pp.603-612,2001.
2. P. Ravikumar and J. Lafferty, "Quadratic Programming Relaxations for Metric Labeling and Markov Random Field Map Estimation," Proc. 23rd Int'l Conf. Machine Learning (ICML), pp. 737-744,2006.
3. S. Sarawagi, S. Chakrabarti, and S. Godbole, "Cross-Training: Learning Probabilistic Mappings between Topics," Proc. Ninth ACM SIGKDD Int'l Conf. Knowledge Discovery and Datamining (KDD),2003.
4. D. Zhang and W.S. Lee, "Web Taxonomy Integration through Co-Bootstrapping," Proc. 27th Ann.Int'l ACM SIGIR Conf. Research and Development in Information Retrieval, pp. 410-417, 2004.

5. D. Zhang, X. Wang, and Y. Dong, "Web Taxonomy Integration Using Spectral Graph Transducer,"Proc.ER Workshop, pp. 300- 312, 2004.

6. A. Nandi and P.A. Bernstein, "Hamster: Using Search Clicklogs for Schema and Taxonomy Matching,"Proc. VLDB Endowment, vol. 2, no. 1, pp. 181-192, 2009.

7. C. Chekuri, S. Khanna, J.S. Naor, and L.Zosin, "Approximation Algorithms for the Metric Labeling Problem via a New Linear Programming Formulation," Proc. 12th Ann. ACM-SIAM Symp. Discrete Algorithms (SODA), pp. 109-118, 2001

8. Georgiana Ifrim, Gökhan Bakir, Gerhard Weikum, "Fast logistic regression for text categorization with variable-length ngrams" KDD '08 Proceedings of the 14thACM SIGKDD international conference on Knowledge discovery and data mining Pages 354-362.

9. A. Doan, J. Madhavan, R. Dhamankar, P. Domingos, and A. Halevy,"Learning to Match Ontologies on the Semantic Web," The VLDB J.,vol. 12, no. 4, pp. 303-319, 2003.

10. Ming Ji, Tianbao Yang, Binbin Lin, Jiawei Han, "A Simple Algorithm for Semisupervised Learning with Improved Generalization Error Bound" Proc. 29th Int'l Conf. on Machine Learning,Edinburgh, Scotland,UK, 2012.

## AUTHOR BIBLIOGRAPHY

| | |
|---|---|
|  | **Miss.DhaygudeTejashri Mohan**<br>Received his B.E. degree in Computer Engineering (Distinction) in the year 2013 from Pune University. He is currently working toward the M.E. Degree in Computer Engineering from the University of Pune,Pune. |
|  | **Patil Bharati R.**<br>Received his B.E. degree in Information Technology (Distinction) in the year 2013 from Shivaji University. He is currently working toward the M.E. Degree in Computer Engineering from the University of Pune,Pune. |
|  | **Prof. Bere S.S.**<br>received his B.E. degree in Technology (First-class) in the year 2008 from Pune University and M.Tech Degree (Distinction) in Computer Engineering in2013 from JNTU University. He has 07 years of teaching experience at undergraduate and postgraduate level. Currently he is working as Assistant Professor and HOD in Department of Computer Engineering of DGOI, FOE, swami-chincholi, Daund, Pune University. His research paper has been published in IJTITCC, IJISET year 2014. His research interests are Digital Image processing. |