

A new Approached Data in Uncertain Classification Using Support Vector Machine Algorithm

Nijaguna GS¹, Dr. Thippeswamy K²

¹Research Scholar, VTU Belgaum, nijagunags@gmail.com

²Professor & Head, Dept. Of CS&E, VTU PG Centre, Regional Office Mysuru, India.

Received Feb. 19, 2017

Accepted March 11, 2017

ABSTRACT

As we all examine a new algorithm of learning replica in which the experiential data input is corrupted with plenty of noise. Based on the probability of modeling technique, we can derivative a common formulation in statistical data where unnoticed input is replica as a hided mixture basic. Many algorithms exist in literature for users to choose a correct one as per their needs. This research paper gives a concept with the fundamentals of many existing classification of data techniques for uncertain data via KNN approach. We were proficient to proposed evaluation technique that obtains uncertainty input into deliberation. For deterioration problems, the correlation of our technique, Aggravated by this probability model technique and proposed new SVM classification technique that handles input data uncertainty. This technique has a understanding of the geometric perceptive data. Furthermore, two observing demonstration, one with realistic data, was used to demonstrate that the new technique is far better and superior to the existing SVM for problems with noisy input data.

Key Words: *New SVM, fuzzy logics, uncertain clustering of data.*

1. INTRODUCTION

One of the most essential tasks in data mining and machine learning area, classification has been studied for many years. To solve the problem in various aspects more number of effective models and algorithm has been introduced, including support vector machine, rule-based classifier, decision tree, etc. excluding some traditional rule-based algorithms such as discriminative measurements like less confidence threshold and associative classification tries to mine all the frequent patterns from the input data set, taking the user - specified less support threshold.

To select the more number of discriminative patterns Sequential covering technology is employed while covering more number of input Training instances. Based on the mined patterns a test instance is classified after using the associative classification classifier train. CBA is one of algorithms. Which

associative classification algorithm could give better classification accurately than other algorithms on categorical datasets but this approach takes a large amount of running time in both pattern mining and feature selection because most of the mined frequent patterns are not the most discriminative ones and will be neglected after some time.

Several algorithms have been proposed to improve the efficiency of associative classification in recent years, try to mine the large number of discriminative patterns directly during the pattern mining step. Different discriminative measures and different instance covering methods have also been devised. HARMONY is one of the most typical algorithms that use confidence to evaluate the discrimination of patterns. It gives a so- called instance-centric neglecting other methods, associative classification find all the frequent patterns in the input

categorical data satisfying a user-specified less support and other discrimination measures like less information-gain or confidence. This patterns are used later either as training features for support vector machine (SVM) classifier or rules for rule-based classifier, after a feature selection procedure which usually tries to cover many number of input instances with the most discriminative patterns in various ways. To mine the most discriminative patterns directly without costly feature selection several algorithms have been proposed; associative classification could provide better classification accuracy to many datasets. Many studies have been conducted on indecisive data, where fields of indecisive attributes no longer have confident values. Instead probability distribution functions are adopted to represent the possible values and their corresponding probabilities. Noise and measurement limits cause improbability. To solve the classification problem on indecisive data several algorithms have been proposed like by extending traditional rule-based classifier and decision tree to work on indecisive data. In this research , we will propose a novel algorithm which mines discriminative patterns directly and effectively from indecisive data as classification rules, to help train either SVM or rule-based classifier. We will discover patterns directly from the input database, feature selection usually taking a large amount of time could be restricted completely. We will develop Effective techniques for computation of expect confidence of the mined patterns used as the measurement of discrimination will also propose. Numerous studies have been conducted on indecisive data in which fields of indecisive attributes no longer have confident values. To represent the possible

values and their corresponding probabilities probability distribution functions are adopted. The improbability is usually caused by measurement limits, noise or by other possible factors.

To solve the classification problem on indecisive data several algorithms have been proposed recently, for example by decision tree to work on indecisive data and extending traditional rule-based classifier. In this research , we proposed a novel technique this mines discriminative patterns directly and effectively from indecisive data as classification rules, to help train either SVM or rule-based classifier. We discover patterns directly from the input database, feature selection usually taking a large amount of time could be restricted completely. We analysis Effective method for computation of expect confidence of the mined patterns used as the measurement of discrimination are also propose.

2. RELATED WORK

Fabrizio Angiulli in at al [1] nearest neighbor class of a test object is the class that maximizes the probability of given that it's nearest neighbor. The confirmation is that the former thought is a lot more influential than the second in the presence of uncertainty, in that it appropriately models the right semantics of the nearest neighbor verdict rule when applied to the uncertain scenario. An effective and resourceful algorithm to perform uncertain nearest neighbor categorization of a generic (un)certain test object is intended, based on properties that very much reduce the temporal cost connected with nearest neighbor class probability subtraction.

Yongxin Tong in at al[2] behavior a comprehensive learning of every the frequent item set mining algorithms more

than uncertain databases. Since there are two definitions of frequent item sets more than uncertain data, nearly all existing research are categorized into two directions. Though, through our searching, initially clarify that there is a close association among two dissimilar definitions of frequent item sets over uncertain data. Consequently, require not use the existing solution for the subsequent definition and replace them with practical obtainable solution of first meaning.

Sangkyum Kim in et al[3] develop an efficient algorithm to straight mine discriminative k-ee sub trees, which are not binary but numeric acceptable features, in one iteration. Through complete experiments on a variety of datasets. Exhibit the utility of projected framework to give an effective explanation for the authorship classification problem.

Ibrahim Ozkan in et al[4] it is rational to propose that the level of the fuzziness is a extremely powerful parameter and surely helps us to appreciate both the relation among the data vectors and the overall structure of the data itself.

Chuancong Gao in et al[5] proposed efficient algorithm, Stream Gen, to mine frequent item set generators in excess of sliding windows on stream data. It accept the FP-Tree structure to succinctly store the transactions of the obtainable window, and devise a narrative details tree structure to keep every the mined generator and their edge to the non- generators. In the interim, a number of optimization techniques are also proposed to accelerate the mining process.

3. PROPOSED METHODOLOGY

Beside with the development in knowledge, more quantity of data are moreover getting produce and are accumulate in the form of digital. Throughout the production of data, uncertainties move stealthily with in or

devoid. The produce data is accumulate in a database it can be used to mine the imperative patterns and leaning from the data. Uncertain databases enclose records with items whose occurrence in those is not completely certain. There is as alternative, a related probability value with whole item in both records.

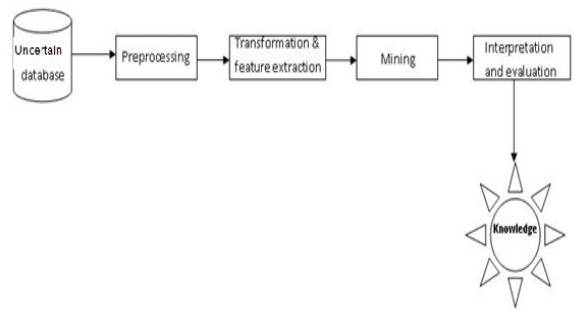


Figure 1: uncertain data classification technique

Conventional data mining technique can't be practical straight on uncertain databases. This direct to the need of propose the narrative techniques that will be capable to handle the un preferred databases. As there is the group of uncertainty in data to be mine. When the user searches about anything it is completely uncertain that what's to be searched. This approach will works for this uncertainty, i.e. indecisive data. This approach will be directly mine the different patterns is based on probability function because of indecisive data fields' to attributes have no longer confident values. The proposed approach mines will be the most discriminative patterns directly and effectively on the indecisive data .This approach will be the less time consuming as it is directly mines be the patterns the time is consumed in pattern mining and feature of selection is reduced.

The uncertain objects contain arbitrary shapes of the uncertain regions. In addition to that, pruning rules are connected with the

occurrence level of uncertain objects in the expensive manner because each uncertain object contains more and large number of occurrences.

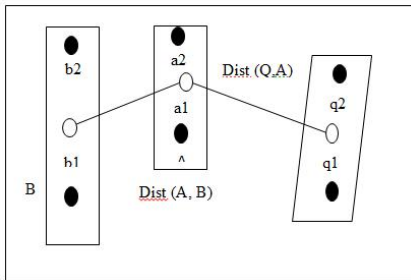


Figure 1. Probabilistic RNN Query

Uncertainty is used in many web applications like information extraction, information integration and web data mining. In uncertain database, probabilistic threshold queries are studied where all results satisfy the queries with possibilities equal to or larger than the threshold values.

least of one pattern. The probability should be more than threshold. In previously done work there will be a lot of work in finding discrimination mentioned patterns, but they all are the time consuming, as they have to be first mine to complete set of frequent patterns using some of association classification technique. Association classification uses some of minimum support or the discriminative measurement like minimum confidence.

Important research curiosity in the data uncertainty managing has to be increasing in the past a less number of years. Data uncertainty is categorized in two types that are existential uncertainty and value uncertainty. Initial category will believes the uncertainty of a tuple's continuation in the database and the subsequent type of deal with probable values of a objective.

The greater part will be works listening carefully study of uncertain data management for the straightforward database queries, in its place of comparatively more difficult than the data mining problems. Instigate of several classification algorithms proposed in previous, construction classifiers based on the uncertainty has remain a challenge. There are a only some of simple technique to be developed for behavior missing or a noisy data values such as, which might as well be the use for conduct uncertainty. The technique of clustering has well been considered in data mining explore.

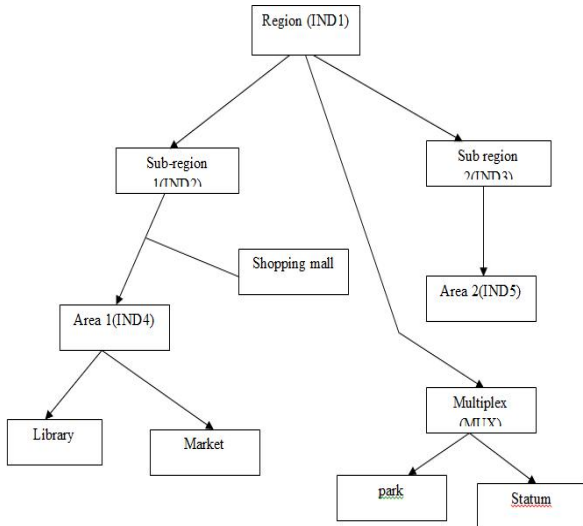


Figure 2. Probabilistic XML Data Tree

As the flexibility of XML data model allocates a natural demonstration of uncertain data, uncertain XML data management contains significant problem.

Costly sequential covering technology is been replaced by instance strategy to assure the probability of each of instance cover by at

4. PERFORMANCE ANALYSIS OF SPACE QUERY CLASSIFIER INDEXING FOR MINING UNCERTAIN DATA

In order to compare the space query classifier indexing for mining uncertain data using different techniques, number of queries is taken to perform the experiment. Various

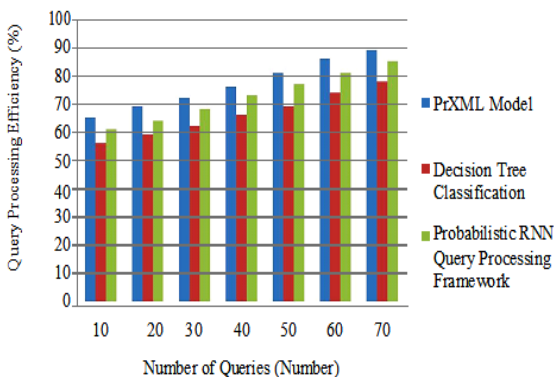
parameters are used for query classification of uncertain data.

Finally, complete content and organizational editing before formatting. Please take note of the following items when proofreading spelling and grammar:

Query processing efficiency is defined based on the queries addressed to the total number of queries by the user with different interval of time periods. It is measured in terms of percentage (%).

Table 1.
Tabulation of query processing efficiency for query classifier indexing for mining uncertain data

Number of	Query Processing Efficiency (%)		
Queries (Number)	<i>PrXML Model</i>	<i>Decision Tree Classification</i>	<i>Probabilistic RNN Query Processing Framework</i>
10	65	56	61
20	69	59	64
30	72	62	68
40	76	66	73
50	81	69	77
60	86	74	81
70	89	78	85



Execution time is defined as the difference between starting time and ending time of query classification of uncertain data. It is

measured in terms of millisecond (ms). Memory Consumption comparison takes place on existing Probabilistic XML Model (PrXML), Decision Tree Classification and Probabilistic Reverse Nearest Neighbor (RNN) Query Processing Framework.

7. CONCLUSION

In terms of the performance the algorithms is developed so far for the precise data in the different data mining techniques like classification, clustering and the association rule mining, we get satisfactory of results but the uncertain data will be provides the completely different scenario and most of algorithms give the different results when applied on data. In this paper we have studied about few techniques of K nearest neighbor classification algorithms on uncertain data. Uncertain data mining is the area of interest for the researchers and more work is required to handle the unrequired data in better way. We further plan to go into the detail of specified classification technique for uncertain data to demonstrate the accurate results as possible with a certain data.

6 REFERENCES

1. Zhou Y., Youwen L., Shixiong X., An Improved KNN Text Classification Algorithm Based on Clustering, Journal of Computers, 2009, 4(3): 230-237.
2. Romero C., Ventura S., Espejo P.G., and Hervas C., Data Mining Algorithms to Classify Students, Proceedings of the 1st Int'l conference on educational data mining, Canada, 2008, pp: 8-17.
3. Zhang J., Mani I., kNN Approach to Unbalanced Data Distributions: A Case Study involving Information Extraction, In Proceedings of The Twentieth International Conference on Machine Learning (ICML-2003), Workshop on Learning from Imbalanced Data Sets II, August 21, 2003.

4. Z.G. Liu, Q. Pan, J. Dezert. "A new belief-based K-nearest neighbor classification method," *Pattern.Recogn.*, vol. 46, No. 3, pp. 834-844, March, 2013
5. Fabrizio Angiulli, Fabio Fassetti," Nearest Neighbor- Based Classification of Uncertain Data," *Journal ACM Transactions on Knowledge Discovery from Data (TKDD)*, Vol. 7,No.1, March 2013.
6. Jianping Gou, Zhang Yi, Lan Du and Taisong Xiong," A Local Mean-Based k-Nearest Centroid Neighbor Classifier," *The computer journal*, Vol.54, No. 1, January 2012.
7. Destercke S, A k-nearest neighbours method based on imprecise probabilities. *Soft Compute* 16(5):833-844, 2012.
8. Reynold Cheng , Lei Chen , Jinchuan Chen, XikeXie, Evaluating probability threshold k-nearest-neighbor queries over uncertain data, *Proceedings of the 12th International Conference on Extending Database Technology: Advances in Database Technology*, March 24-26, 2009, Saint Petersburg, Russia.

Happiness is not something ready-made. It comes from your own actions.

~ Dalai Lama