# Novel Technique Uncertain Data Classification Using Support Vector Machine

**Mohit Jain[*],  Asha Khatri**
CS, Medicaps Institute of Technology, RGPV, Indore,
Madhya Pradesh, India

*Abstract— we examine a novel learning replica in which the experiential input is corrupted with noise. Based on a probability modeling technique, we derivative a common statistical formulation where unnoticed input is replica as a hidden mixture basic. We were proficient to proposed evaluation technique that obtains input uncertainty into deliberation. For deterioration problems, the correlation of our technique. Aggravated by this probability model technique and proposed novel SVM classification technique that handles input uncertainty. This technique has a perceptive geometric understanding. Furthermore, two observed illustration, one with real data, were used to demonstrate that the novel technique is superior to the typical SVM for problems with noisy input data.*

*Keywords— novel SVM, fuzzy c-means, uncertain data clustering.*

## I.    INTRODUCTION

One of the most essential tasks in data mining and machine learning area, classification has been studied for many years. Many effective models and algorithm have been proposed to solve the problem in different aspects, including decision tree, rule-based classifier, support vector machine, etc. Unlike some traditional rule-based algorithms like associative classification tries to mine the complete set of frequent patterns from the input dataset, given the user-specified minimum support threshold and/or discriminative measurements like minimum confidence threshold. Sequential covering technology is further employed to select the most discriminative patterns while covering most input Training instances. A test instance is classified later using classifier trained based on the mined patterns. CBA is one of the most classical associative classification algorithms. that associative  classification algorithm could provide better classification accuracy than other algorithms on categorical datasets. However, this approach takes a great amount of running time in both pattern mining and feature selection, since most of the mined frequent patterns are not the most discriminative ones and will be dropped later. To improve the efficiency of associative classification, several algorithms have been prop used in recent years to try to mine the most discriminative patterns directly during the pattern mining step. Different discriminative measures and different instance covering technologies have also been devised. One of the most typical algorithms is HARMONY .which uses confidence to evaluate the discrimination of patterns. It employs a so-called instance-centric Unlike other methods, associative classification tries to find all the frequent patterns existing in the input categorical data satisfying a user-specified minimum support and/or other discrimination measures like minimum confidence or information-gain. Those patterns are used later either as rules for rule-based classifier or training features for support vector machine (SVM) classifier, after a feature selection procedure which usually tries to cover as many as the input instances with the most discriminative patterns in various manners. Several algorithms have also been proposed to mine the most discriminative patterns directly without costly feature selection; associative classification could provide better classification accuracy over many datasets. Recently, 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. The improbability is usually caused by noise, measurement limits, or other possible factors. Several algorithms have been proposed to solve the classification problem on indecisive data recently, for example 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 features/rules, to help train either SVM or rule-based classifier.  We will discover patterns directly from the input database, feature selection usually taking a great amount of time could be avoided completely. We will develop Effective method for computation of expect confidence of the mined patterns used as the measurement of discrimination will also propose.

Recently, numerous 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. The improbability is usually caused by noise, measurement limits, or other possible factors. Several algorithms have been proposed to solve the classification problem on indecisive data recently, for example by extending traditional rule-based classifier and decision tree to work on indecisive data. In this research , we proposed a novel technique which mines discriminative patterns directly and effectively from indecisive data as classification features/rules, to help train either SVM or rule-based classifier.  We discover patterns directly from the input database, feature selection usually taking a great amount of time could be avoided completely. We analysis Effective method for computation of expect confidence of the mined patterns used as the measurement of discrimination are also propose.

## II. RELATED WORK

Fabrizio Angiulli in at al[1]nearest neighbor class of a test object is the class that maximize the probability of given that its 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 itemset mining algorithms more than uncertain databases. Since there are two definitions of frequent itemsets more than uncertain data, nearly all existing research are categorize into two directions. Though, through our searching, initially clarify that there is a close association among two dissimilar definitions of frequent itemsets 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 at al[3] develop an efficient algorithm to straight mine discriminative k-ee subtrees, 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 at 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 at al[5]proposed efficient algorithm, StreamGen, to mine frequent itemset 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 technique are also proposed to accelerate the mining process.

## III. PROPOSED METHODOLOGY

Beside with the development in knowledge, huge quantity of data are moreover getting produce and are accumulate in digital form. Throughout the production of data, errors or uncertainties move stealthily in with or devoid. The produce data is accumulate in a database and can be use 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 an alternative, a related probability value with every item in both records. Conventional data mining technique cannot be practical straight on uncertain databases. This direct to the need of propose narrative techniques that will be capable to handle the uncertain databases. As there is a group of uncertainty in data to be mine. When user searches for anything it is completely uncertain that what is to be searched. This approach works for this uncertainty, i.e. indecisive data. This approach will directly mine the different patterns based on probability function because indecisive data fields' attributes have no longer confident values. The proposed approach mines the most discriminative patterns directly and effectively on indecisive data .This approach will be less time consuming as it directly mines the patterns the time consumed in pattern mining and feature selection is reduced.

Costly sequential covering technology is also replaced by instance strategy to assure probability of each instance cover by at least one pattern. The probability should be higher than threshold. In previously done work there were a lot of work on finding discrimination among patterns, but they all are time consuming, as they have to first mine the complete set of frequent patterns using some association classification technique. Association classification uses some minimum support or discriminative measurement like minimum confidence.
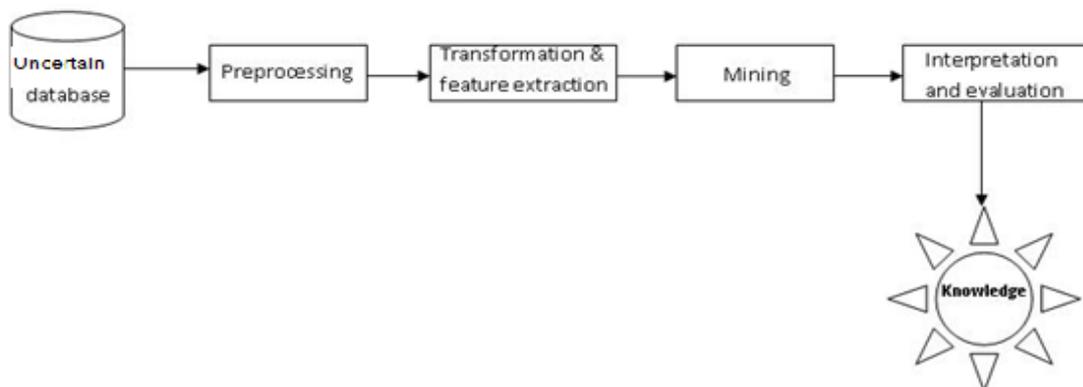


Figure 1: uncertain data classification technique

Important research curiosity in data uncertainty managing has been increasing in the past a small number of years. In common, data uncertainty is categorized into two types, that is existential uncertainty and value uncertainty. Initial category believes the uncertainty of a tuple's continuation in a database and the subsequent type deal with probable values of an objective. The greater part of the works listening carefully on study uncertain data management for straightforward database queries, in its place of comparatively more difficult data mining problems. Instigate of the several classification algorithms proposed previous, construction classifiers based on uncertainty has remain a challenge. There are a only some simple technique developed for behavior missing or noisy data values such as in , which might as well be use for conduct uncertainty. The technique of clustering has as well been well considered in data mining explore.

Also, merely a little learns on clustering for uncertain data have been description. We used SVM, association rules based on fuzzy clustering; a cluster is symbolized by a fuzzy separation of a set of objects. Every object can fit in to additional than one cluster, every with a dissimilar degree. Dissimilar fuzzy clustering technique has been functional on normal data or fuzzy data to create fuzzy clusters amongst which fuzzy c-means is the majority widely used technique. In the fuzzy replica, substance belong to each (fuzzy) set with a number of degree of relationship where as in uncertain databases, objects might or may not fit in to every set. There is an uncertainty or chance of whether an object belongs to every set or not. In the text, uncertainty is classically generated unnaturally using dissimilar technique. One technique is to prefer an uncertainty model for a specified dataset and ran-domly put in the uncertainties to every data point. In this variety of technique, uncertainty information or noise is produce using a probability density function for every dimension of the data and these values are close to the data points. the majority of the density functions used are Gaussian and uniform. One more way of computing uncertainty is to make uncertainty from random probability density function and attaching the probability values to the data points. only some other works have produce probabilities in such a way that the sum of probabilities diagonally dissimilar dimensions of a specified data point . In reality, the noise strength not go after the Gaussian or uniform distribution. In our work there is no such assumption laid on the uncertainty noise values, as they are straight capture from real-world data. We illustrate the dissimilar technique of uncertain data categorization and uncertain data clustering. For nearly every one of these papers it is the inference that the enhanced the uncertainty cans replica errors, the superior the accuracy of the classifiers that they can construct with their algorithms. Consequently, for research purposes, they produce the uncertainty information by guess the error replica. While, in this research there is no requirement of decide the accurate error function to initiate uncertainty as the uncertainty that is confine is in real-world data. as well, these works are mainly speculative, where request areas are optional but not experiment and are supposed. Whereas, our technique are functional in a developed real-world call for such as text data classification and clustering.

## IV.    CONCLUSION

 In this research we have proposed a novel classifier that handles uncertainty nearby in the data. The proposed classifier uses an associative categorization technique to classify class of an object in a given text data. This technique takes benefit of using latent frequent patterns. The classifier as well utilize the standardize .recurrent item sets detain every one the foremost associations among items in dataset. Resourceful item set mining algorithms subsist. These classifiers in nature handle missing values and as they simply deal with statistically important associations which create the classification to be robust.  General performance learns have exposed such classifiers to be normally added accurate. As well, unlike most classifiers which rely on together positive and negative class sets for training, our proposed classifier depends simply on the positive class data for training.

**REFERENCES**
[1]     Fabrizio Angiulli ,Fabio Fassetti   DIMES, University of Calabria, "Nearest Neighbor-Based Classification of Uncertain Data"Italy ACM Transactions on Knowledge Discovery from Data (TKDD) TKDD Homepage archive Volume 7 Issue 1, March 2013  Article No. 1.
[2]     Yongxin Tong, Lei Chen, Yurong Cheng , Philip S. Yu "Mining Frequent Itemsets over Uncertain Databases" August 27th 31st 2012, Istanbul, Turkey. Proceedings of the VLDB Endowment, Vol. 5, No. 11.
[3]     Sangkyum Kim, Hyungsul Kim, Tim Weninger, Jiawei Han, Hyun Duk Kim," Authorship Classification: A Discriminative Syntactic Tree Mining Approach " SIGIR '11 July 24-28 2011, Beijing, China.
[4]     Ibrahim Ozkan, Burhan Türkşen," MiniMax ε-Stable Cluster Validity Index for Type-2 Fuzziness" 978-1-4244-7858-3/10/-2010 IEEE.
[5]     Chuancong Gao, Jianyong Wang," Efficient Itemset Generator Discovery over a Stream Sliding Window" CIKM'09, November 2–6, 2009, Hong Kong, China.
[6]      Metanat HooshSadat and R. Osmar Zaiane. An associative classifier for uncertain datasets. In Proceedings of the 16th Pacific-Asia Conference on Knowledge Discovery and Data Mining, PAKDD'12, 2012.
[7]     Michael Chau, Reynold Cheng, and Ben Kao. Uncertain data mining: A new research direction. In Proceedings of the Workshop on the Sciences of the Artificial, 2005.
[8]     Charu C. Aggarwal, Yan Li, JianyongWang, and JingWang. Frequent pattern mining with uncertain data. In Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining, KDD '09, pages 29–38, 2009.
[9]     Jiye Liang," Feature selection for large-scale data sets in GrC" IEEE International Conference on Granular Computing-2012.
[10]    Ashfaqur Rahman, Daniel V. Smith, Greg Timms," Multiple Classifier System for Automated Quality Assessment of Marine Sensor Data" IEEE ISSNIP 2013.
[11]    Xiaojing Shen, Yunmin Zhu Yingting Luo, Jiazhou He," Minimized Euclidean Error Data Association for Multi-Target and Multisensor Uncertain Dynamic Systems"
[12]    X. Shen, Y. Zhu, E. Song, and Y. Luo, "Minimizing Euclidian state estimation error for linear uncertain dynamic systems based on multisensory and multi-algorithm fusion," IEEE Transactions on Information Theory, vol. 57, pp. 7131–7146, October 2011.