

## Enhanced Non-Parametric Summarization Using Concept Evolution

Balaji.R<sup>1</sup>, Jayanthi.S<sup>2</sup>, Dinesh Raj.R<sup>3</sup>, Dhivya.E<sup>4</sup>

<sup>1</sup>Department of ME-CSE, Srinivasan Engineering College

<sup>2</sup>Assistant Professor/CSE, Srinivasan Engineering College

<sup>3</sup>Department of ME-CSE, Srinivasan Engineering College

<sup>4</sup>Department of ME-AE, Anna University-Regional Centre

---

**Abstract**— To abstract a concept from the raw data user may choose any classification algorithm of processor interest, or choose one that appears to be good at learning the current data. Information is commonly used and can achieve reasonable classification accuracy in general. A trigger detection algorithm finds instances, across which the underlying concept has changed and the prediction model should be modified. It is especially important when concept shifts. A classification methodology is used here with two parameters window size and error threshold. The beginning of the window is always a misclassified instance. When the window is full and its error rate exceeds the error threshold, the beginning instance is taken as a trigger; otherwise, the beginning of the window is slid to the next misclassified instance (if there is any) and previous instances are dropped from the window. The temporarily package holds potential novel class instances. These instances are analyzed periodically in order to detect novel class, which is explained in the next paragraph, needs to be cleared periodically to remove instances that no longer contribute to novel class detection. Besides, instances in that have reached classification deadline.

**Keywords**— *Data Streams, Concept-Drift, Novel Class, Ensemble Classification.*

---

### I. INTRODUCTION

Data Mining is an important part of Knowledge discovery Process that we can analyze an enormous set of data and get hidden and useful knowledge. Classification and clustering is the major technique in data mining to solve machine learning process. Data Stream Mining is the process of extracting knowledge structures from continuous, rapid data records. A data stream is an ordered sequence of instances that in many applications of data stream mining can be read only once or a small number of times using limited computing and storage capabilities. Examples of data streams include computer network traffic, ATM transactions, web searches, and sensor data. Data stream mining can be considered a subfield of data mining, machine learning, and knowledge discovery. Machine learning techniques can be used to learn this prediction task from labeled examples in an automated fashion. Often, concepts from the field of incremental learning, a generalization of Incremental heuristic search are applied to cope with structural changes and real-time demands. In many applications, especially operating within non-stationary environments, the distribution underlying the instances or the rules underlying their labeling may change over time, i.e. the goal of the prediction, the class to be predicted or the target value to be predicted, may change over time.

Concept drift is statistical properties of the target variable, which the model is trying to predict, change over time in unforeseen ways. Concept-drift happens in the medical data stream when the original concept of the data changes over time i.e. When detecting novel class in the process. While learning a combination of classifiers, ensemble classification is needed to categorize unique information and decision tree method for entire classification. Thus, the classification model can able to efficient continuously so that it reflects the most recent concept.

Most of the existing solutions assume that the total number of classes in the data stream is fixed. But in real world data stream classification problems, such as intrusion detection, text classification and fault detection, novel classes may appear at any time in the stream (e.g. a new intrusion). Traditional data stream classification techniques would be unable to detect the novel class until the classification models are trained with labeled instances of the novel class. Thus, all novel class instances will go undetected (i.e., misclassified) until the novel class is manually detected by experts, and training data with the instances of that class is made available to the learning algorithm. They address this concept-evolution problem and provide a solution that handles all three problems, namely, infinite length, concept-drift, and concept-evolution.

## II. RELATED WORK

Initially synthetic data has been collected and saving into the repository. Each and every attributes has been extracting and move on to pre-processing work. In specified grouping process, values are clustered and storing into the database. Classification based on structured and unstructured values. Novel class has been detect and classified in single step. The overall values can be storing into the database.

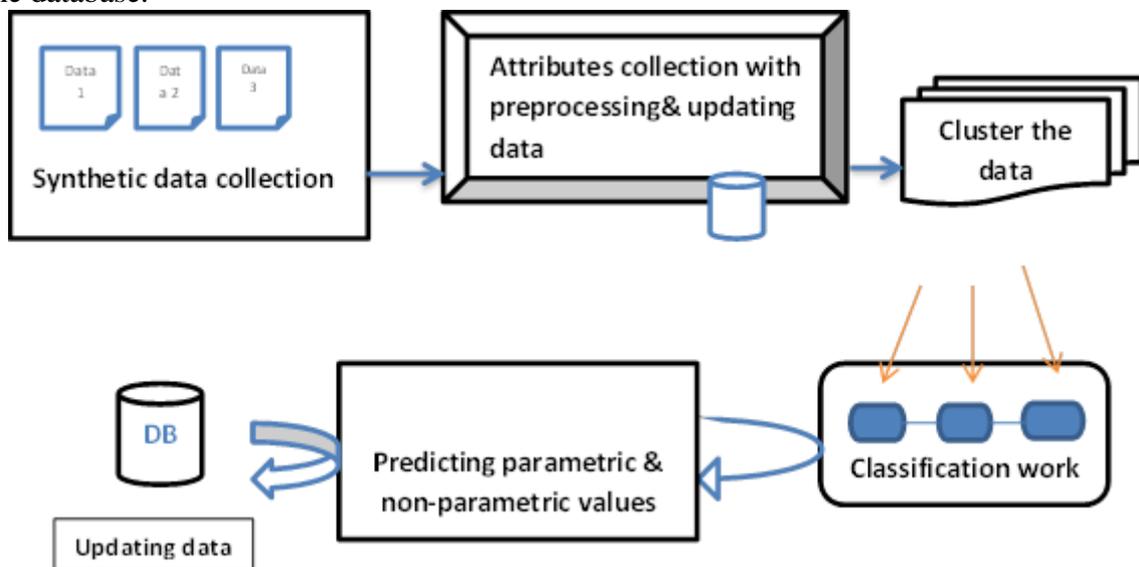


Figure 1 system architecture

The ideal classification scenario is to detect the changes when they come, and retrain the classifier automatically to suit the new distributions. Most methods for novelty detection rely on some form of modeling of the probability distributions of the data and monitoring the likelihood of the new-coming observation. Use direct detection of changes in data. After a drift is detecting older data are removed and a classifier is updated. However, the most widely used such direct approaches (called triggers) are based on observing decrease in the classification accuracy, which requires access to labeled stream of examples.

Data gathering methods are often loosely controlled, resulting in out-of-range values, impossible data combinations, missing values. Data collection will be the initial work in data mining progress. Input data can be txt format or .csv file format. By getting the input file information into back end storage for the future usage. Representation and quality of data is first and foremost before running an analysis. After collecting attributes of entire group required data stream extraction can be proceed for clustering. Grouping the content made in clustering process by using k means clustering technique. In the name of pre-processing flow move on to collecting the non-parametric data i.e. incomplete data can be detect and progress for the completion work.

$k$ - Means clustering aims to partition  $n$  observations into  $k$  clusters in which each observation belongs to the cluster with the nearest mean. In centroid-based clustering when the number of clusters is fixed to  $k$ , it formally optimizing the process. Under the clustering process  $k$  cluster centers and assign the objects to the nearest cluster center, such that the squared distances from the cluster are minimized. Data can be drift by number of classification process with data updating. Classification involves deriving a function that will separate data into categories, or classes, characterized by a distinct set of features. An evolutionary ensemble classifier has given lead to categorize the attributes with its specific constrains. From the result of process will update the data stream by labeled sort & unlabeled sort for the step by step classification process.

Collected information has been gathered and ready to process for reclassification work. By taking different entity for ordering the elements in second time, each attribute have taking and arranging according to the cluster values grouped. Changing the primary classification attribute, new datum has been taking chance to updating the process. Depends on clustered values, structured values has been ordered. Adding the novel class there is no need to develop the clustering work from the beginning; we can update it from the last segments. Both grouping and validation work can be achieved with this new class entrance. Proof has been collect from redirecting into the data frame.

By categorizing the primal and secondary information from the details the major work balancing for the maintenance. Maintaining the biological information defines that how the admin will using that optimized data in their working process. Maintenance can be achieved by showing the entire classification process to view the hospitalized content heavily. From that proposing the full communication in between the medical attributes. This will lead to show our improved performance.

This modules enables users get to know graphically the value of entire progress. Process will get to evaluate the entire classification work by concluding the better performances easily. Reducing time on current working instances will show the integrity of project. Time consumption and accuracy has been screening with comparing the other classification development as improvement. By initialize the time value at the starting stage of classification will be ended up on final progress. Evaluation work will be based on calculating time in between the process. Finally graph can be estimate with the help of this time.

### **III. SCOPE OF THE PAPER**

Proposed a cross-validation-based framework to choose data and compare sensible choices. Exhaustively select the training data by comparing all the sensible choice. This paper proposed work based on decision tree ensemble to “sift through” old data and combine with new data to construct the optimal model for evolving concept. The basic idea is to train a number of random and uncorrelated decision trees. Each decision tree is constructed by randomly selecting available features. A simpler approach always uses data from a fixed number of periods. The unselective use of old data definitely helps improve model accuracy. Loss functions are used jointly in order to minimize the expected loss. The structure of the tree is uncorrelated. Their only correlation is on the training data itself. It Can't Detect New Class's in the Dataset. Time Constraints not considered in the existing methods. There is no formal definition of data sufficiency.

This paper Proposed to mine high-order models in evolving data. This paper advocates exploiting historically trained models to improve the qualities of stream classifiers for evolving data. This approach does not require users to tune any parameters to achieve a satisfying result on stream since data streams have infinite length. We report significant improvement of classification accuracy. This approach is not only desirable, but also feasible. Prediction for data streams is however not a trivial

task. Methodology is an effective and efficient solution to prediction for data streams. This solves the fundamental problem of model over-fitting in classifying data streams. Concept's classifier is used to predict the record, and the overall classification result is updated.

A trigger detection algorithm finds instances, across which the underlying concept has changed and the prediction model should be modified. A sliding window methodology is used here with two parameters window size and error threshold. No significant effort has been devoted to foreseeing a bigger picture. The problem of predicting the oncoming concept can be far more complicated. A classification algorithm is used to abstract a concept from the raw data. A user may choose any classification algorithm of his/her interest, or choose one that appears to be good at learning the current data. A trigger detection algorithm finds instances, across which the underlying concept has changed and the prediction model should be modified. Commonly used and can achieve a reasonable classification accuracy in general. Detecting the concept change and accordingly modifying the prediction model for oncoming instances. It does not explore their associations and hence cannot be proactive. Approaches have not well organized nor made good use of the history of data streams. Concept shift does not hold the required assumption.

Infinite length & concept drift is the main glitches in data streams. The true label of a data point can be accessed as soon as it has been classified by the classification model and updated directly by the branded occurrence. For classification, k-NN used to weight the contributions of the neighbors, so that it donates more to the average than the more distant ones. Traditional novelty detection techniques assume that there is only one normal class and any instance that does not fit to the normal class is an anomaly class instance. Therefore, they are unable to distinguish among different types of irregularity. Traditional stream classification techniques make unreasonable expectations about the obtainability of labeled data. Most assume that the true label of a data point can be accessed as soon as it has been classified by the classification model. Thus, according to their statement, the existing model can be updated immediately using the branded instance. Some ensemble methods for stream data mining have been proposed to require relatively simpler operations to update the current concept than their single model counterparts, and also handle concept-drift efficiently. Several nonparametric approaches available, such as parson window method, k-nearest neighbor (k-NN)-based approach, kernel-based method, and rule-based approach. A cluster-based novel concept detection technique is applicable to data streams. Though, a "single-class" novelty detection technique, where authors assume that there is only one "normal" class and all other classes are novel. Existing novelty detection techniques only consider whether a test point is significantly different from the normal data.

Outliers are the by-product of intermediate computation steps in our algorithm. Ensemble and decision tree methods use multiple learning algorithms to obtain better predictive performance than could be obtained from any of the constituent learning algorithms. A simple and easy way to classify a given data set through a certain number of clusters (k –means clusters) fixed a significance data. Our technique is a nonparametric approach, and therefore, it is not restricted to any specific data distribution. Arrival of novel classes in the stream causes the classifiers in the ensemble to have different sets of class labels.

The precision of our outlier detection technique is easier to the overall performance. Getting much better results than state-of-the-art stream classification by time consumption. Our process not only able to detect the intrusion, but also deduce a new kind of intrusion. By reducing the number of classifiers useful to improve the speed of classification & also updating record easily.

#### IV. CONCLUSION

Classification and clustering summarize the challenging in new class detection and class updating on concept evolution work. By enhancing the classification process will increase the accuracy on the side of application. Updating the content at each step gives possibility to check the data at any time. The main motto of novel class detection process justified at the stage of developing work. Performance has been induced on the centralized cluster vale categorization. Validating the clustered elements has led to reduce the time of verifying work. Hence a mechanism that performs well in general is very useful. As a result, the problem of the intractable amount of streaming data is solved since concepts are much more compact than raw data while still retain the essential information.

#### REFERENCES

- [1]S.J. Roberts, "Extreme Value Statistics for Novelty Detection in Biomedical Signal Processing,"Proc. Int'l Conf. Advances in Medical Signal and Information Processing,pp. 166-172, 2000.
- [2]M. Scholz and R. Klinkenberg, "An Ensemble Classifier for Drifting Concepts,"Proc. Second Int'l Workshop Knowledge Dis-covey in Data Streams (IWKDDS),pp. 53-64, Oct. 2005.
- [3]E.J. Spinosa, A.P. de Leon, F. de Carvalho, and J. Gama, "Cluster-Based Novel Concept Detection in Data Streams Applied to Intrusion Detection in Computer Networks," Proc. 2008 ACM Symp. Applied Computing,pp. 976-980, 2008.
- [4]S. Subramaniam, T. Palpanas, D. Papadopoulos, V. Kalogeraki, and D. Gunopulos, "Online Outlier Detection in Sensor Data Using Non-Parametric Models,"Proc. Int'l Conf. Very Large Data Bases (VLDB),pp. 187-198, 2006.
- [5]G. Tandon and P. Chan, "Weighting versus Pruning in Rule Validation for Detecting Network and Host Anomalies,"Proc. ACM SIGKDD,pp. 697-706, 2007.
- [6]K. Tumer and J. Ghosh, "Error Correlation and Error Reduction in Ensemble Classifiers,"ConnectionScience,vol. 8, no. 304, pp. 385-403, 1996.
- [7]H. Wang, W. Fan, P.S. Yu, and J. Han, "Mining Concept-Drifting Data Streams Using Ensemble Classifiers,"Proc. ACM SIGKDD, pp. 226-235, Aug. 2003.
- [8]D. yanYeung and C. Chow, "Parzen-Window Network Intrusion Detectors," Proc. Int'l Conf. Pattern Recognition, pp. 385-388, 2002.
- [9]Y. Yang, X. Wu, and X. Zhu, "Combining Proactive and Reactive Predictions for Data Streams,"Proc. ACM SIGKDD,pp. 710-715, 2005.
- [10]Y. Yang, J. Zhang, J. Carbonell, and C. Jin, "Topic-Conditioned Novelty Detection,"Proc. ACM SIGKDD,pp. 688-693, 2002.
- [11]X. Zhu, "Semi-Supervised Learning Literature Survey," Technical Report TR 1530, Univ. of Wisconsin Madison, July 2008
- [12]P. Mahoney and M.V. Chan, "Learning Rules for Anomaly Detection of Hostile Network Traffic,"Proc. IEEE Int'l Conf. Data Mining (ICDM),pp. 601-604, 2003.
- [13] R. Elwell and R. Polikar,(Oct.2011), "Incremental learning of concept drift in nonstationary environments," IEEE Trans. Neural Netw., vol. 22, no. 10, pp. 1517-1531.
- [14] W. Fan,(2004), "Systematic data selection to mine concept-drifting data streams," in Proc. 10th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining, pp. 128-137.
- [15] W. Fan, Y. A. Huang, H. Wang, and P. S. Yu, (Apr. 2004), "Active mining of data streams," in Proc. 4th SIAM Int. Conf. Data Mining, pp. 457-461.

