# Outlier Based Fraud Detection System

**Prerana M More, Mrudul K Rushipathak, Shweta G Bairagi, Trupti S Nalnikar**
Department of Computer Engineering
KKWIEER, Nashik
Email id: prerana.more2012@gmail.com, mrudul40@gmail.com , shwetabairagi12@gmail.com, trupti.nalnikar@gmail.com

## Abstract

*Data mining has the vital task of Outlier detection which aims to detect an outlier from given datasets. The analysis or detection of outlier data is referred to as Outlier Mining. In Data mining, outlier detection is the identification of unusual or distant data records that might be require further investigation or analysis. This paper provides the data driven methods for various fraud detection systems based on literature review, fraudulent activities or cases and comparative research. Outlier detection is the technique which discovers such type of data from the given data set. Several techniques of outlier detection have been introduced which requires input parameter from the user. The goal of this proposed work is to partition the input data set into the number of clusters using K-NN algorithm. Then the clusters are given as an input to the outlier detection methods namely cluster based outlier algorithm and Local Outlier Factor Algorithm. The Performance evaluation of this algorithm confirms that our approach of finding local outliers can be practically implemented.*

**Keywords:** *Data Clustering, K-NN, Outlier Detection, ODA, LOF*

## INTRODUCTION

A number of surveys, review articles and books cover outlier detection techniques in the field of data mining. In this paper we make an attempt to bring together many different outlier detection techniques in a structured and generic description.

Outlier detection has been used for long time to find and remove identical observations from data. Outliers are generated due to system faults, errors in data entry, fraudulent behaviour, data scaled inappropriately or simply through natural deviations in populations. Their detection can identify system faults and fraud before they escalate with potentially identical consequences. It can identify errors and remove their significant effect on the data set and as such to purify the data for processing. The earlier outlier detection methods were complexed but now, different techniques are used, related to Statistics. In this paper, we introduce

different algorithms for outlier detection. We identify the motivations of it and differentiate among their advantages and disadvantages in a comparative review.

The LOF algorithm is an outlier detection method which is unsupervised method and it calculates the density deviation of a given data point with respect to its neighbors in the database. It notes the data as outlier samples that have a relatively lower density than their neighbors in the database.

In this paper basic methodologies currently used for solving this problem are considered, and their advantages and disadvantages are discussed. It is based on methods of fuzzy set theory and the use of kernel functions and possesses a number of advantages compared to the existing methods. The performance of the algorithm used is studied by the example of outlier detection arising in credit card
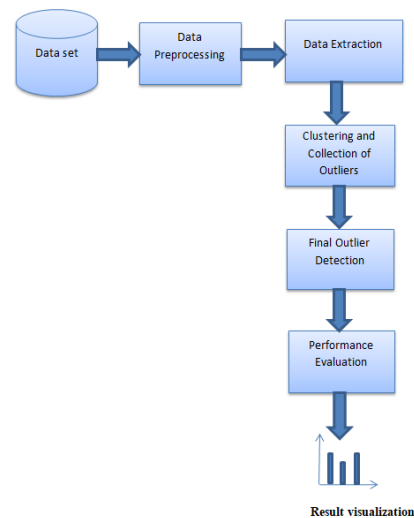
systems, the so-called fraud detection systems.

## Literature Survey

a. Kadam, Pund and Aggarwal went through various approaches to detect outliers including the cluster-based approach.

b. Aparna and Nair proposed the CHB-K-Means algorithm by using a weighted attribute matrix to detect outliers.

c. Jiang et al. introduced a two-phase clustering algorithm for outlier detection. In the first phase, the k-means algorithm is changed to partition the data in such a way that a data point is assigned to be a new cluster center only if the data point is far away from all clusters. In the second phase, a minimum spanning tree is constructed based on the cluster centers obtained from the first phase.

d. Zhang et al. (2009) proposed a measure called Local Distance-based Outlier Factor (LDOF) to measurethe outliers present in scattered data.

e. He et al. declared the concept of cluster based local outlier and also designed a measure called clusterbased local outlier factor (CBLOF) .

## Proposed Methodology

The proposed system works in three phase's viz. preprocessing, clustering and outlier detection. Before performing clustering input dataset is preprocessed and then the clustering is performed using k-nn algorithm. After clustering, each cluster out of k clusters is given as an input to the outlier detection methods. There are different methods that are used for detecting an outlier. Output provided is set of outliers from input dataset. System architectural diagram is explained below.



Result visualization

## Step 1 :Pre-processing

Input : Open the file or URL to be processed.

Dataset consists of numeric data from various sources. Dataset should be in xml or csv format. The data can be generated by different transactions, inventory management systems and DBMS.

User need to mention attributed which are unwanted or don't play any role in outlier detection as well as mark class attribute from imported file.

Save the file to proceed further.

Processing : System will now remove undesirable attributes and prepare the file for further processing.

## Step 2 : Outlier Detection Algorithm

Input : User has to click on Start Algorithm

Processing : Data Records which may contain outliers are detected and differentiated from remaining records
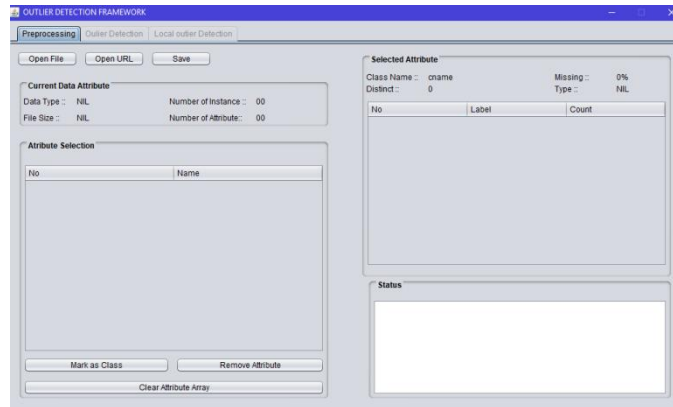
## Step 3 : LOF Algorithm

Input : Records which may contain outliers are considered for execution of Local Outlier Factor algorithm.

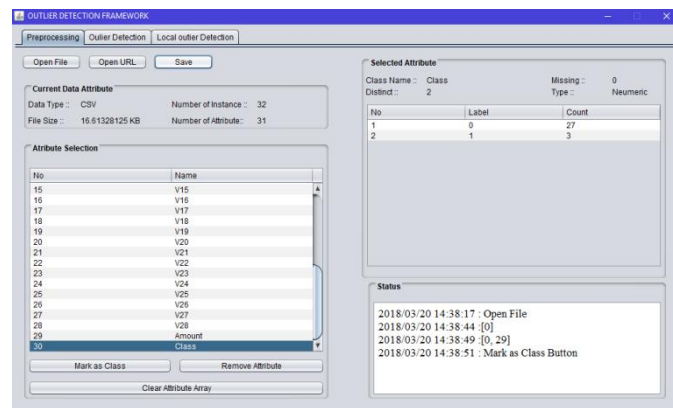User has to mention Number of Parameters for KNN algorithm

Processing: KNN algorithm will be run followed by LOF to detect outliers

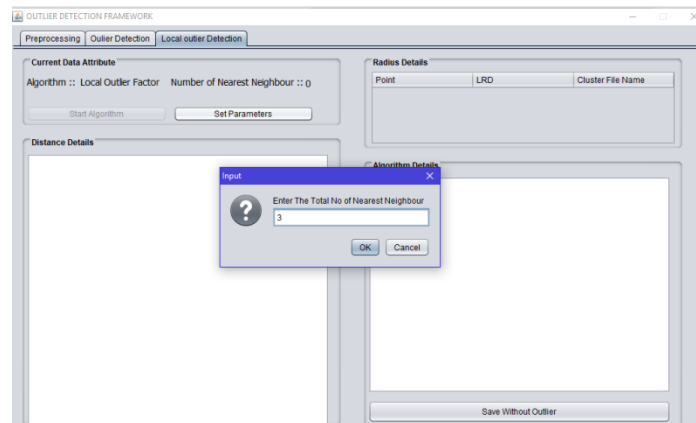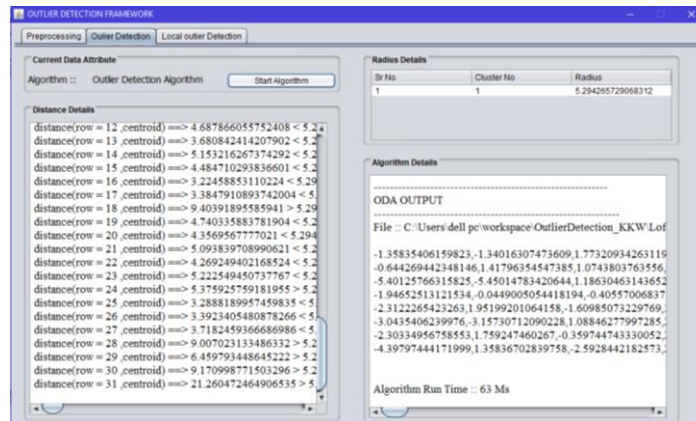Detected outliers will be displayed followed by time in ms

## Results



After selecting a file Remove unwanted attributes &mark some attributes as class and save the file
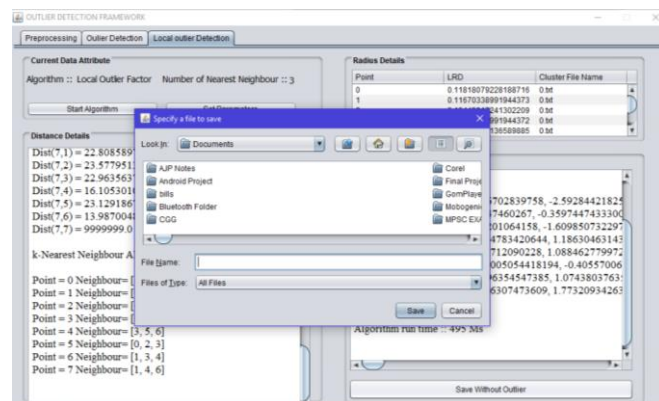


In order to preprocess dataset we have to select & import file through Open file or Open URL
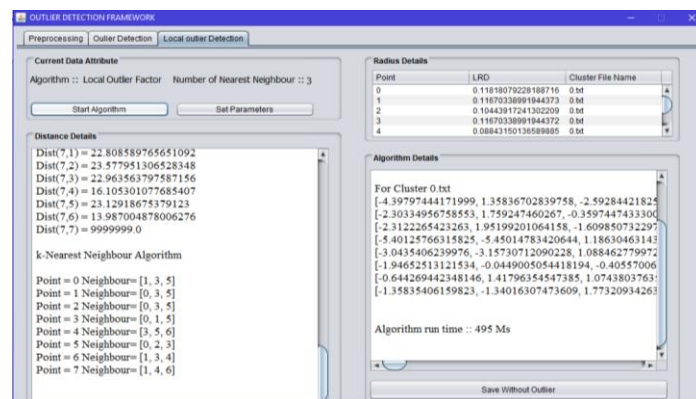


After clicking on Outlier detection tab/button ODA algorithm will run and entire dataset will be checked and suspicious records will be extracted for further processing

Set parameters i.e. Nearest neighbor for outlier detection algorithm and click on start algorithm



Set parameters i.e. Nearest neighbor for outlier detection algorithm and click on start algorithm



It displays outliers. We can now save file without outliers

This system recognizes the outliers from given dataset successfully. In order to produce precise and accurate details it first pre-processes the data and then perform Outlier Detection Algorithm i.e. ODA and then it displays outliers and processing time in mini seconds.

**CONCLUSION**

Outlier detection is a significant problem. It has direct applications in a wide variety of domains. A salient observation with outlier detection is that it was not a precisely formulated problem. Several approaches have been proposed by many researchers to target a particular application domain. To detect cluster based outliers first input data set is clustered into number of clusters using K-nn algorithm and then outlier is detected from each cluster by applying a cluster based outlier detection algorithm and the

Outlier Detection Algorithm. The ODA algorithm gives better accuracy as compared to cluster based outlier detection algorithm for same set of datasets.

We introduce the notion of the local outlier factor LOF, which captures exactly this relative degree of isolation. We show that our definition of LOF enjoys many desirable properties. There are two directions for ongoing work. First is how to describe or explain why the identified local outliers are exceptional. This is specifically meant for high-dimensional data or datasets, as a local outliers may be outlying only on some, but not on all dimensions. The second one is to further improve the performance of LOF computation. On the one hand, such an algorithm may provide more detailed information about the local outliers, e.g., by analyzing the clusters relative to which they are outlying. Also, computation may be shared between LOF processing and clustering. The shared computations may contain k-nn queries and reachability distances.

**REFERENCES**
1. Aggarwal, C.C., 2015. Data Mining: The Textbook. Springer, New York, NY.Aggarwal, C.C., Reddy, C.K. (Eds.), 2013. Data Clustering: Algorithms andApplications. CRC Press, Boca Raton, FL, USA.
2. Ahmed, M., Naser, A., 2013. A novel approach for outlier detection and clusteringimprovement, in: Proceedings of the 8th IEEE Conference on IndustrialElectronics and Applications (ICIEA), pp. 577–582.
3. Aparna, K., Nair, M.K., 2016. Computational Intelligence in Data Mining.Springer. volume 2. chapter Effect of Outlier Detection on Clustering Accuracyand Computation Time of CHB K-Means Algorithm. pp. 25–35.
4. Chawla, S., Gionis, A., 2013. *k*-means: A unified approach to clustering and outlier detection. SIAM. chapter 20. pp. 189–197.
   Dave, R., Krishnapuram, R., 1997. Robust clustering methods: a unified view.IEEE Transactions on Fuzzy Systems 5, 270–293.
5. Duan, L., Xu, L., Liu, Y., Lee, J., 2009. Cluster-based outlier detection. Annalsof Operations Research 168, 151–168.
6. Fahad, A., Alshatri, N., Tari, Z., Alamri, A., Y.Zomaya, A., Khalil, I., Foufou,S., Bouras, A., 2014. A survey of clustering algorithms for big data: Taxonomyamp; empirical analysis. IEEE Transactions on Emerging Topics inComputing PP, 1–1