



OUTLIER DETECTION AND CLUSTERING IMPROVEMENT USING A NOVEL APPROACH

¹**Author: Jagannath Mahajan, research scholar at Department of Computer Science at SSSUTMS, Sehore-MP**

²**Author: Dr. Jitendra Sheethlani, Professor at Department of Computer Science at SSSUTMS, Sehore-MP**

Abstract: Outlier detection is utilised in a variety of applications, including clustering-based illness onset identification, gene expression analysis, computer network intrusion detection, financial fraud detection, and human behaviour analysis. Because of their poor accuracy and lack of a comprehensive strategy, existing methods for detecting outliers are ineffective. Small clusters are usually considered outliers, and most algorithms provide an outlier score to each data object. Due to high computational complexity and misidentification of regular data objects as outliers, these techniques have disadvantages. Using a modified k-means clustering technique, we provide a unique unsupervised approach for detecting outliers in this study. To increase clustering accuracy, outliers are eliminated from the dataset. By comparing our method to existing approaches and benchmark performance, we are able to prove that it is effective.

Keywords: Outlier, Clustering, k-means clustering, Unsupervised

I. INTRODUCTION

Anomaly location is a significant information examination task. The principal objective of exception location is to recognize peculiar or unusual information from a given dataset. This is an intriguing region of information mining research as it includes finding new and uncommon examples from a dataset. Exception identification has been broadly concentrated in measurements and AI. It is otherwise called inconsistency location, curiosity recognition, deviation discovery and exemption mining [1]. Anomaly is characterized by numerous specialists in different ways in view of the application area. One generally acknowledged meaning of anomaly is given by Hawkins [2]. As indicated by Hawkins, 'An anomaly is a perception which veers off so much from different perceptions as to stir doubts that it was produced by an alternate instrument'. Exceptions are reviewed as significant on the grounds that they show huge however, uncommon occasions, and can incite basic moves to be made in a wide scope of utilization spaces. For instance, a surprising traffic design in an organization could show that a PC is hacked and sending information to unapproved objections, a bizarre way of behaving in Visa exchanges could show fake exercises, an exception in a MRI picture might demonstrate the presence of a threatening cancer. Anomaly identification has been generally applied to endless application spaces. We examine probably the main applications to propel its utilization. Interruption Detection: Computer interruption incorporates hacking what's more, spreading of infection and worm across organizations to penetrate a neighbourhood or remote machine, or cause harm utilizing Distributed Refusal of Service (DDoS) assaults. Notwithstanding, interruptions comprise just a little level of the absolute organization and PC utilization that are viewed as typical use. Anomaly location can

be utilized to distinguish noxious exercises of projects as well as programmers [3] from network traffic information and PC exercises.

Extortion Detection: Most weak areas of deceitful exercises are in unapproved charge card use, cell bill, pointless protection guarantee and stock trade insider exchanging. Taken charge cards are utilized in a surprising manner than the typical example. The use example of a taken Visa is looked at against the customary utilization information of the genuine proprietor what's more, consequently anomalies are identified from the charge card exchange information [4]. Criminal rings of unlawful protection petitioners and suppliers control the case handling framework for unapproved claims. Following such exercises help the organization to keep away from monetary misfortunes. Brain network based procedures have been effectively applied to identify such exceptions. Insider exchanging is a crime the securities exchange, where benefits are made by inside data before it is unveiled. Clinical and Public Health: Anomalous records can be reproduced because of patient condition or instrumental mistake or recording mistakes. An off-base test report could have genuine repercussions. Then again, anomaly recognition in this area is a vital apparatus that might possibly save living souls [5] by identifying issue right on time from test results also, pictures. Exception discovery has additionally been applied to a few other spaces including Image Processing, Sensor Networks, Astronomy, Biology, Speech Recognition and some more. In this paper, we acquaint an original unaided methodology with identify exceptions utilizing an adjusted k-implies calculation. To work on the grouping exactness, we eliminate the recognized exceptions from the groups. We contrast our method and existing procedures furthermore, benchmark execution. We likewise tried different things with irregular circumstances to assess whether our methodology recognizes exception by some coincidence or not. Trial investigation gives aintensive comprehension of the exhibition of our strategy which outflanks existing techniques on a few measures. The rest of the paper is coordinated as follows. Area II portrays the connected work. In Section III, essential k-implies calculation is momentarily presented as foundation. Area IV makes sense of our methodology. Area V incorporates broad exploratory outcomes that approves proposed approach for anomaly location, further developing grouping proficiency, precision as a classifier utilizing different datasets from UCI Machine Learning Archive [6] alongside correlation with different methods. Segment VI closes the paper.

II. RELATED WORK

In this part, we audit the current grouping based exception location draws near. Grouping based anomaly location approaches don't need pre-labeled information and can recognize anomalies alongside bunching. Grouping calculations like ROCK, DBSCAN and BIRCH spotlight just on grouping information, but they have exemption dealing with limits. A large portion of the bunching techniques are created to enhance grouping process, yet all at once not the exception recognition capacity. Furthermore, these methodologies do not perform well when utilized with high layered datasets. Svetlona et al. [7] introduced an anomaly evacuation bunching calculation (ORC) that gives anomaly discovery and information bunching all the while. Their proposed calculation has two stages. At first the k-implies bunching is applied and afterward a distance factor, o_i for every one of the information point is determined by taking the proportion of a guide's distance toward the centroid and the most extreme separation from centroid to some other point. A limit T is set under 1 to check for exceptions. If remote variable for any point is more noteworthy than the edge then it is considered as an anomaly and eliminated from the dataset. Their trial information incorporates engineered information and some guide pictures. Mean Absolute Error (MAE) is utilized to assess calculation execution. The boundary T esteem is reliant upon the dataset which might cause cluster showing in heterogeneous huge scope datasets.

$$o_i = \frac{\|x_i - C_{p_i}\|}{d_{max}} \quad (1)$$

Another definition for cluster-based local outliers was suggested by He et al. [8]. According to their definition, all data points in a cluster are deemed outliers rather than a normal distribution. Figure 1 shows a single point. The smaller groups C1, C2 and C3 are both regarded to be outliers. They made use of some numbers, parameters, i.e., to distinguish between Small Cluster (SC) and Large Cluster (LC). These are what the clustering approach is based on, parameters, however it's unclear how to specify the values for a variety of datasets. The SQUEEZER algorithm was employed to cluster data since it delivers good clustering quality and can deal with data with a lot of dimensions. Then there's the FindCBLOF. Each individual record's outlier factor is determined using an algorithm in the data set. For each record, CBLOF(t).

$$CBLOF(t) = \begin{cases} |C_i| * \min(d(t, C_j)) & \text{where } t \in C_i, C_i \in SC \text{ and} \\ & C_j \in LC \text{ for } j=1 \text{ to } b \\ |C_i| * (d(t, C_i)) & \text{where } t \in C_i \text{ and } C_i \in LC \end{cases} \quad (2)$$

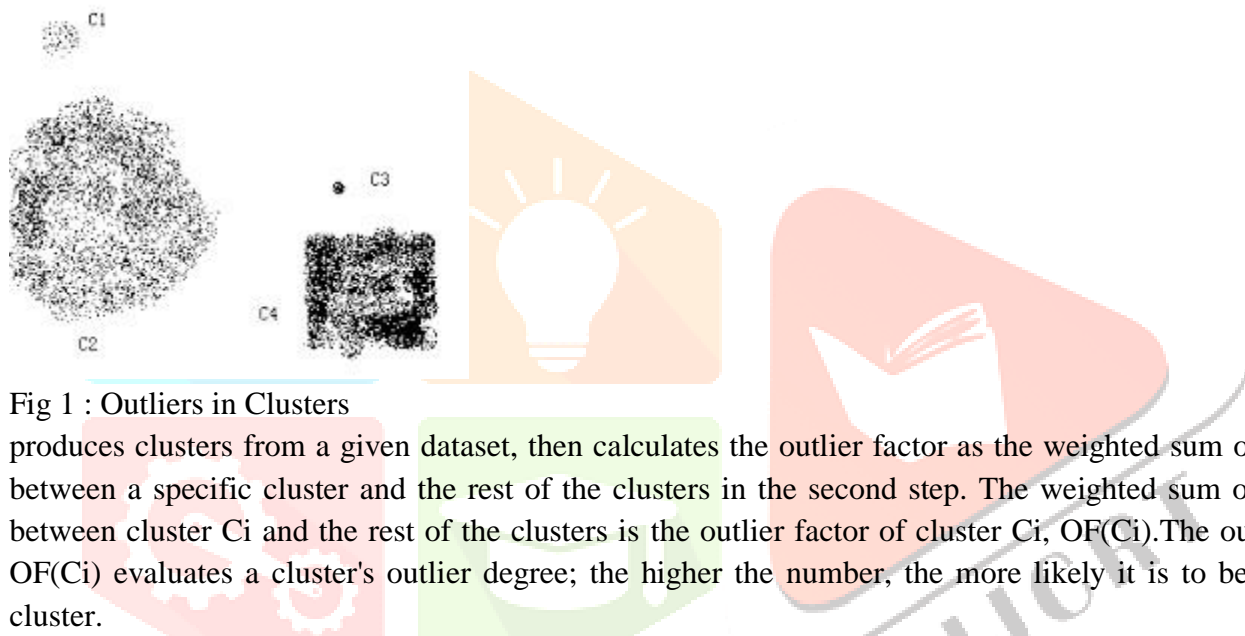


Fig 1 : Outliers in Clusters

produces clusters from a given dataset, then calculates the outlier factor as the weighted sum of distances between a specific cluster and the rest of the clusters in the second step. The weighted sum of distances between cluster C_i and the rest of the clusters is the outlier factor of cluster C_i , $OF(C_i)$. The outlier factor $OF(C_i)$ evaluates a cluster's outlier degree; the higher the number, the more likely it is to be an outlier cluster.

$$OF(C_i) = \sum_{j \neq i} |C_j| * d(C_i, C_j) \quad (3)$$

Outlier clusters are defined as minimum b clusters that meet the following requirements. To assess performance, they looked at the detection rate and false alarm rate.

$$\frac{\sum_{i=1}^b |C_i|}{|D|} \geq \epsilon (0 < \epsilon < 1) \quad (4)$$

Jiang et al. [10] introduced a two-stage bunching method to recognize exceptions. To begin with, they utilized an altered k -implies calculation to make bunches. On the off chance that the focuses in a similar group are not sufficiently close, the bunch can be parted into two more modest groups what's more, consolidated when a given edge surpasses. In the second step, they develop a base traversing tree with the group focuses and eliminate the longest edge. The more modest sub trees are considered as anomalies. Their procedure thinks about a whole group as an anomaly, which may not be pertinent for some datasets and increment False Positive rate. Yoon et al. [11] utilized k -implies bunching on the total dataset and to track down the appropriate worth of k , utilized Cubic Clustering Criterion (CCC). CCC is a method used to appraise the quantity of bunches assessed by Monte Carlo strategy. When the grouping is done, the space master looks for outside and inward exceptions in the bunching results. Outside anomalies are the information focuses situated at a more prominent distance

than different groups and inner exceptions are the data of interest remotely situated inside a group, displayed in Figure 2. In the event that the expulsion of exceptions make significant groups then the method stops. Their methodology has been applied uniquely for programming estimation information and a space master is required for translation.

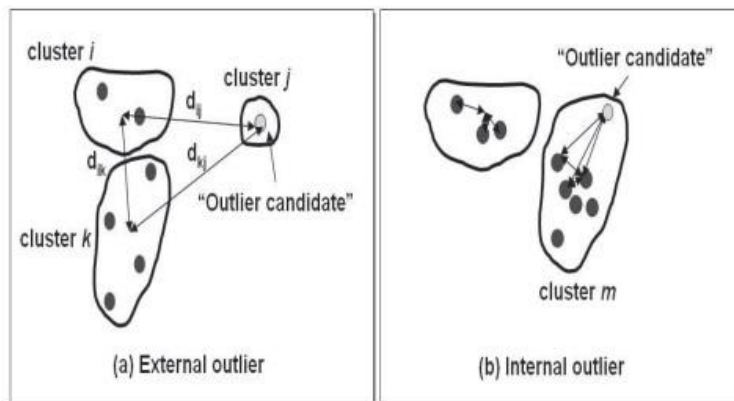


Fig 2: Outliers on the Outside and Inside

k-means (basic) Algorithm

1. Choose k locations to serve as the initial centroids.
2. Do it again.
3. Assign each point to its nearest centroid to create k clusters.
4. Recalculate each cluster's centroid.
5. Until the centroids remain unchanged.

Fundamental k-means calculation shows the essential strides of k-means bunching. In k-means bunching each item is doled out to definitively one of k groups. The calculation takes an info boundary, k , which is known deduced by information master and parts a bunch of n objects into k -groups. When the grouping has been done, the subsequent intra-group similitude is high be that as it may, the between bunch comparability is low. The calculation functions as follows: right away, k -beginning centroids are picked haphazardly from the arrangement of n objects. Each item is allocated to its nearest group in view of its Euclidian distance to the bunch centroid. The arrangement of focuses relegated to a centroid is viewed as a group. This is the means by which k -bunches are shaped. Then, the centroid of each group is refreshed in view of the mean of the articles appointed to it. This interaction is rehashed, so that each point is allocated to the closest group in light of the progressions in the position of the group centroid. The interaction stops when no object changes the group or the centroids quit moving. To gauge the nature of a bunching, k-means calculation utilizes the Sum of Squared Error (SSE) and Total Sum of Squares (SST). Euclidian distance between each article and the centroid of the bunch to which it has a place, addresses a blunder, what's more, from this the complete amount of squared blunders is figured. SST is the squared all out amount of distances between the mean of the dataset and every one of the places in dataset. Table I characterizes the documentations utilized in the rest of the paper. Given two various arrangements of bunches created by k-means, the grouping which has a lower SSE/SST, is viewed as better. SSE and SST are officially characterized as

$$SSE = \sum_{i=1}^k \sum_{C_i} dist(c_i, x)^2 \quad \text{where } \forall x \in C_i \quad (5)$$

$$SST = \sum dist(C_m, x)^2 \quad \text{where } \forall x \text{ in Dataset } D \quad (6)$$

IV. EXCEPTION DISCOVERY AND CLUSTERING APPROACH

This segment makes sense of how an exception is distinguished and how the essential k-implies grouping calculation [12] can be adjusted to recognize anomalies. First we characterize exception utilizing proposed strategy and afterward clear up the calculation for identify exceptions what's more, at the same time further develop bunching productivity.

TABLE I

Symbol	Description
x	A vector representing an object
Capital C_i	The i th cluster
Small c_i	Centroid of cluster C_i
K or k	The number of clusters
C_m	The mean of all points

A unique outlier based on clusters G

The distance between a location and its centroid is used to designate an outlier. An outlier is a data point that is a fixed multiple of the mean distances of all other data points from the centroid. 'An item o in a group of n objects is an outlier if the distance between o and the centroid is higher than p times the mean of the distances between the centroid and other objects,' according to the formal definition. p is always bigger than 1 in this case. Take, for example, six items in a cluster with a centroid in Table II (5,6). Figure 3 shows the sample data points from Table II displayed to show an outlier (9,9). Figure 3 demonstrates that.

TABLE II

X	Y	Distance
4	8	2.23
5	5	1.0
5	6	0
4	7	1.41
6	5	1.41
9	9	5.0
Mean Distance		1.84

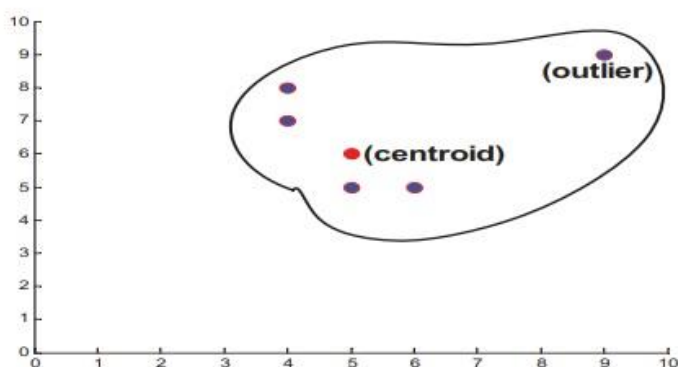


Fig 3: Example of Outlier

Algorithm ODC

Input: $D(A_1, A_2, \dots, A_n), k, p$ // The dataset, no. of group and limit

Yield: Clustered Data, Outliers and SSE/SST.

Start

1. Pick a worth of k .
2. Select k items haphazardly and use them as beginning arrangement of centroids
3. Ascertain the distances between k centroids and every one of the articles in dataset D
4. Ascertain the mean distances (M_d) between k centroids and all the objects in dataset D
5. Relegate each item to the bunch for which it is closest centroid furthermore, work out SSE/SST.
6. for each article x in dataset D
7. In the event that distance $(x, c_k) > p * (M_d)$
8. Consider x as an anomaly and eliminate from dataset D and work out SSE/SST.
9. end
10. Recalculate the centroids.
11. Rehash stages 3-10 until objects quit evolving groups.

End

its nearest centroid, then, at that point, it ascertains the SSE/SST of the grouping to lessen mistake. Then, at that point, it recognizes the anomalies as per the proposed meaning of an exception (i.e., $p > 1$). On the off chance that any anomaly is recognized, it is taken out from the dataset and put away independently as exceptions. Then, the centroids are recalculated. This interaction stops when the items don't change their situation starting with one bunch then onto the next. The worth of p is resolved tentatively. For all the datasets we have tried different things with, no. of group, $k = 3$ and edge, $p = 2$ was utilized.

V. TEST ANALYSIS

In this part we approve our calculation with reality datasets got from UCI Machine Learning Repository [6]. Our calculation was executed in MATLAB and all the tests were directed on Windows 7 64-bit rendition with center i7 processor and 8GB DDR3 RAM. The trial segment is isolated into four sections. In the initial segment we analyze our methodology against existing procedures. Second part incorporates the exploratory after effects of grouping improvement on various benchmark datasets. In the third part we think about the characterization exactness for anomaly recognition. In the last part we contrast our calculation and irregular determination situation. At long last, we give a concise thought on computational intricacy. A. Recognizable proof of Rare Classes In this investigation we have utilized the Lymphography dataset from UCI AI Repository [6]. This dataset has 148 cases with 18 credits. The dataset has 142 cases in the normal class (ordinary) and just 6 occasions in the interesting class (strange). Here, the goal is to test if the proposed anomaly location procedure can recognize the couple of intriguing class cases present in the dataset. The probabilistic translation of the term Recall is that, it is the likelihood that significant record is recovered in a search. With regards to our examination, we are utilizing review bend to show that the anomalies recognized from the benchmark dataset are recovered accurately and among those anomalies the

Algorithm for tracking down Frequent Outliers

1. for $I = 1$ to 10
 2. Run ODC calculation.
 3. $S(i) =$ Detected Outliers
 4. end
 4. $S = S(1) \cup S(2) \cup \dots \cup S(10)$
 5. Sort S as indicated by recurrence of identification.
- End
-

uncommon class cases are obviously distinguished. As k-implies calculation [12] introduces the centroids haphazardly, hence, the beginning centroids will change in various runs of the calculation. Thus, we ran our calculation multiple times and created a bunch of applicant anomalies identified by joining the exceptions from all of the ten executions.

Then the exceptions are arranged by the quantity of times a specific exception shows up in S, these are called incessant anomalies. The calculation for identifying regular anomalies is obviously expressed. Successive exceptions have been looked at against the top proportion anomalies in the FindCBLOF approach [8] and distinguished anomalies in ORC approach [7]. Top proportion exceptions are the quantity of exceptions determined as top-k anomalies to that of the records in the dataset. Utilizing the proposed strategy 45 exceptions were distinguished that showed up at least a time or two in S. From this arranged rundown of 45 anomalies, all of the 6 uncommon class cases of the dataset have been seen as in only the initial 28 generally successive exceptions, as displayed in Figure 4. This outcome plainly illustrates that not just the proposed strategy accomplishes 100 percent review in observing all the interesting class occurrences (6) however it likewise accomplishes this in less applicant anomalies (28), contrasted with FindCBLOF (30) also, ORC (40). With just 28 exceptions the proposed approach recognizes 6 out of 6 (100 percent review) uncommon classes contrasted with review of FindCBLOF (4 out of 6, 67%) and ORC (5 out of 6, 84%).

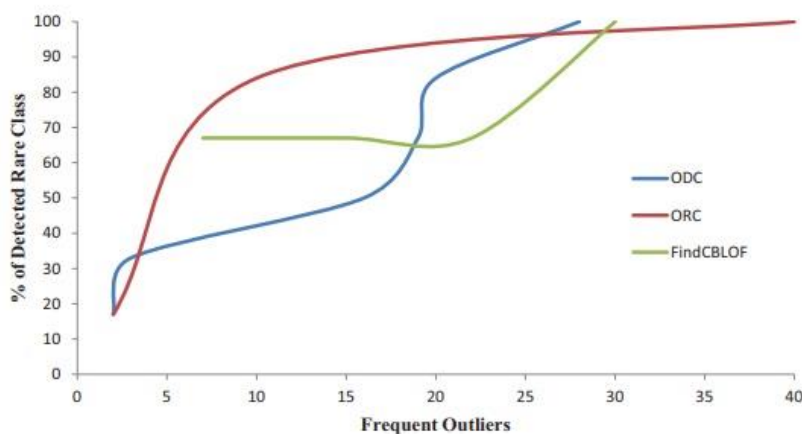


Fig 4: Anomaly Detection Techniques on Lymphography Data are Compared.

B. Grouping Improvement

One method for assessing an exception identification strategy would be to gauge how much exactness is expanded by the anomaly recognition and expulsion process. Here, we give the test after effects of SSE/SST on five benchmark datasets; Toughening (798 occasions with 38 characteristics), Lymphography (148 cases with 18 ascribes), Iris (150 occasions with 4 ascribes), Glass (214 cases with 10 credits), Yeast (1484 cases with 8 credits). As referenced in Section III, lower upsides of SSE/SST show better grouping [12]. We determined the worth of SSE/SST on three distinct situations, (1) preceding anomaly evacuation, (2) after exception expulsion, (3) later the eliminating same number of identified anomalies arbitrarily.

For instance, in Glass information, the worth of SSE/SST previously anomaly discovery was 0.92. The proposed calculation identified 18 exceptions and eliminated them, thus, the SSE/SST esteem dropped down to 0.67. Then, 18 occurrences were taken out haphazardly from the dataset coming about a SSE/SST worth of 0.89. From Figure 5 plainly our ODC approach recognizes the interesting cases accurately and consequently further develops grouping exactness by limiting the worth of SSE/SST.

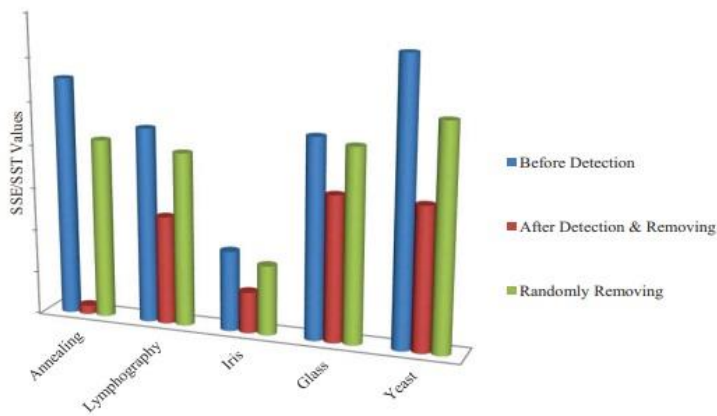


Fig 5: SSE/SST comparison in many settings.

C. Classifier Evaluation

We assess the accuracy of outlier-based anomaly identification in terms of True Positive (TP), True Negative (TN), False Positive (FP), and False Negative (FN) in this experimental study (FN). On benchmark datasets, the results are compared to FindCBLOF and ORC approaches. The labels true and false relate to whether the classifier's prediction matches to the external judgement or ground truth. The following formula is used to calculate accuracy:

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} \quad (7)$$

The results of an experiment conducted using the Lymphography dataset to determine the accuracy of rare/anomalous class recognition are shown in Table III. Figure 6 illustrates that the suggested strategy (66 percent) outperforms both ORC (50 percent) and FindCBLOF (50 percent) based on testing findings (63 percent).

To compute Euc, we use True Positive Rate (T Prate) and False Positive Rate (FPrate) to evaluate the classification accuracy of the proposed approach. The distance between a classifier and the ideal classifier on the Receiver Operating Curve graph is measured by Euc. The formal definition of Euc is as follows:

$$TP_{rate} = \frac{TP}{TP + FN} \quad \text{and} \quad FP_{rate} = \frac{FP}{FP + TN} \quad (8)$$

TABLE III

	ODC	FindCBLOF	ORC
TP	6	6	6
FP	22	24	34
TN	92	88	68
FN	28	30	40

Euc = 0 would be the best feasible classifier. When FPrate = 0 and T Prate = 1, the Euc has the smallest feasible value of 0. When FPrate = 1 and T Prate = 0, the largest possible value is 2. Euc values that are lower indicate a good classifier. We evaluated the Euc values for each of the three ways using the value from Table III, and it is clear from Figure 7 that the suggested strategy has the best classification accuracy.

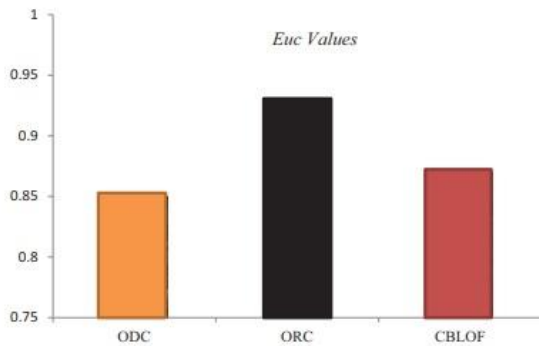


Fig 7. Comparison of classification accuracy.

D. Contrast with Random Selection

In this experiment, we compare the suggested technique against an algorithm that chooses the same number of outliers at random to see if it detects outliers by chance. We apply our method to spot common outliers and look for uncommon occurrences. We also identify the same amount of outliers at random and look for uncommon occurrences. For example, ODC discovers 28 outliers in the Lymphography dataset, and among

Algorithm for Random Selection

Begin

1. Execute the ODC algorithm.
2. The number of Frequent Outliers is denoted by the letter O.
3. Look for uncommon occurrences in O.
4. R = Pick O random occurrences from the dataset.
5. In R End, look for uncommon occurrences.

End

We obtain 100% rare cases when we randomly choose 28 outliers from the dataset, but we only get 17% uncommon instances when we randomly select 28 outliers from the dataset. This comparison, shown in Figure 8, indicates the efficacy of our technique for detecting uncommon class occurrences.

E. Computational Complexity

Time intricacy of our calculation is $O(I*k*m*(n-o))$. Here I is the emphasizes expected for union, m is the number of traits, n is the quantity of items, o is the quantity of exceptions and k is the quantity of groups. Space intricacy is $O((n-o) + k)*m$ as it stores just information items and centroids. The ORC approach additionally utilizes k -implies strategy, yet it identifies anomaly furthermore, eliminates anomaly subsequent to grouping the information by k -implies, so it requires more cycle to join than our methodology. The FindCBLOF approach utilized SQUEEZER calculation to bunch information and allocate distance variable to every one of the information object, consequently they require two specific calculations (SQUEEZER also, FindCBLOF) for exception identification so clearly, their methodology appreciates more computational intricacy than our own.

VI. CONCLUSION

Anomaly recognition is a significant assignment for KDD applications. In this paper another calculation has been proposed to recognize exceptions and group information at the same time utilizing segment based calculation. This proposed method is an alteration of the well known k -implies calculation. Exploratory outcomes illustrate that it beats existing methods for exception location as well as grouping precision on benchmark datasets. In future, hypothetical premise will be laid out on the anomaly definition what's more, new grouping methods will be proposed to work on the exactness considerably further.

REFERENCES

- [1] V. Chandola, A. Banerjee, and V. Kumar, "Anomaly detection: A survey," *ACM Comput. Surv.*, vol. 41, no. 3, pp. 15:1–15:58, Jul. 2009. [Online]. Available:
- [2] D. Hawkins, *Identification of Outliers*. London: Chapman and Hall, 1980.
- [3] D. J. Marchette, *Computer Intrusion Detection and Network Monitoring: A Statistical Viewpoint*, V. Nair, M. Jordan, S. L. Lauritzen, and J. Lawless, Eds. Secaucus, NJ, USA: Springer-Verlag New York, Inc., 2001.
- [4] R. J. Bolton and D. J. H., "Unsupervised profiling methods for fraud detection," in *Proc. Credit Scoring and Credit Control VII*, 2001, pp. 5–7.
- [5] J. Lin, E. Keogh, A. Fu, and H. Van Herle, "Approximations to magic: Finding unusual medical time series," in *Proceedings of the 18th IEEE Symposium on Computer-Based Medical Systems*, ser. CBMS '05. Washington, DC, USA: IEEE Computer Society, 2005, pp. 329–334. [Online]. Available: <http://dx.doi.org/10.1109/CBMS.2005.34>
- [6] A. Frank and A. Asuncion, "UCI machine learning repository," 2010. [Online]. Available: <http://archive.ics.uci.edu/ml>
- [7] V. Hautamaki, S. Cherednichenko, I. Karkainen, T. Kinnunen, and P. Franti, "Improving k-means by Outlier Removal," in *Proc. 14th Scandinavian Conference on Image Analysis (SCIA'05)*, 2005, pp. 978–987.
- [8] Z. He, X. Xu, and S. Deng, "Discovering cluster based local outliers," *Pattern Recognition Letters*, vol. 2003, pp. 9–10, 2003.
- [9] S.-y. Jiang and Q.-b. An, "Clustering-based outlier detection method," in *Proceedings of the 2008 Fifth International Conference on Fuzzy Systems and Knowledge Discovery - Volume 02*, ser. FSKD '08. Washington, DC, USA: IEEE Computer Society, 2008, pp. 429–433. [Online]. Available: <http://dx.doi.org/10.1109/FSKD.2008.244>
- [10] M. F. Jaing, S. S. Tseng, and C. M. Su, "Two-phase clustering process for outliers detection," *Pattern Recogn. Lett.*, vol. 22, no. 6-7, pp. 691–700, May 2001. [Online]. Available: [http://dx.doi.org/10.1016/S0167-8655\(00\)00131-8](http://dx.doi.org/10.1016/S0167-8655(00)00131-8)
- [11] K.-A. Yoon, O.-S. Kwon, and D.-H. Bae, "An approach to outlier detection of software measurement data using the k-means clustering method," in *Empirical Software Engineering and Measurement, 2007. ESEM 2007. First International Symposium on*, sept. 2007, pp. 443–445.
- [12] J. B. MacQueen, "Some methods for classification and analysis of multivariate observations," in *Proc. of the fifth Berkeley Symposium on Mathematical Statistics and Probability*, L. M. L. Cam and J. Neyman, Eds., vol. 1. University of California Press, 1967, pp. 281–297.