



## International Journal of Control Theory and Applications

ISSN : 0974-5572

© International Science Press

Volume 9 • Number 46 • 2016

# Quantification of the Projected Outliers Using the Distance Based and Gaussian Mixture Model

Kamal Malik<sup>a</sup>, Harsh Sadawarti<sup>b</sup> and G.S. Kalra<sup>c</sup>

<sup>a</sup>Research Scholar, IKGPTU, Jalandhar, Punjab India

<sup>b</sup>Principal cum Director, RIMTIET (Affiliated to Punjab Technical University), Punjab, India

<sup>c</sup>Assistant Professor, Lovely Professional University, Jalandhar, Punjab. Email: gursharanjeetkalra@yahoo.com

**Abstract:** Outliers are the data objects that do not confirm to the normal behaviour and usually deviates from the remaining data objects may be due to some outlying property which distinguishes them from the whole dataset. Usually, the detection of outliers is followed by the clustering of the dataset which sometimes ignores the prominence of outliers. In this paper, we have tried to detect the outliers and pruned the clustering elements initially so that the outliers can be prominently highlighted and secondly Gaussian Mixture model is used to cluster the higher dimensional data so as to detect the projected outliers, We have proposed an algorithm which effectively prunes the similar data objects from the large datasets and its experimental results compare the neighbouring points and shows the better performance than the existing methods.

**Keywords:** Clusters, Pruning, Distance-Based, Projected Outliers.

## 1. INTRODUCTION

Outlier detection is one of the very important aspects of the data mining. Outliers are the data objects or the points that do not comply with the normal behaviour of the datasets. The important applications of the data mining include data cleaning, fraud detection stock market analysis, intrusion detection marketing, network sensors etc. To find the suspicious or erroneous pattern is to find out the outliers. There are many approaches used by the researchers for the outlier detection which may be classified as supervised that includes the exploitation of the training data set of the normal and abnormal objects, semi supervised which includes only normal examples and unsupervised which searches the unlabelled data set to detect the outliers without giving the proper reason for having their outlying property. The problem of detecting the outliers has been extensively studied in the statistics community<sup>1</sup> Distance-Based techniques<sup>2</sup> are very popular for relating each pair of objects in the data set. There are so many methods for defining the outliers from the perspective of distance based outliers methods are based on the concepts of local neighbourhood of  $k$ -nearest neighbours (KNN) of the data points<sup>3</sup>. The notion of the distance based outlier methods is that the user might not have the idea about the underlying assumptions of the data and usually generalizes the concepts from the distribution methods. Moreover, distance based methods are

usually very easy to scale up for the higher dimensional data to detect the projected outliers. There are several metrics that are used to measure the outlierness of the data points. In this paper, we have tried to find the various outliers and analyse them using the clustering and distance function. The notion behind it is first to prune the points of nearest neighbourhood of the centroid i.e., the points which are of similar properties are collected in a different dataset and are pruned and the remaining points which are outliers are taken into consideration. So, without considering the points of the cluster, only the outliers are taken into the account which reduces the computations and complexity as a whole. Moreover, the ODF i.e., outlier defining factor is used to provide the proper ranking or the measure of outlierness of an outlier.

The main aim of this paper includes two very important points:-

1. The detection of the subset  $S$  from the input population which are homogenous in nature that is the points which are quite similar to each other.
2. Detecting the outliers which crosses the threshold of their inner radius and are at much distance from the centroid.
3. Thirdly, to quantify the outliers using the outlierness defining factor, so as to detect the deviations and differences from the cluster.

## **2. DISTANCE –BASED OUTLIER DETECTION**

Distance based outlier techniques were first of all introduced by knorr and Ng<sup>4,5</sup>. According to them - An object  $p$  in a data set  $DS$  is a DB ( $q, dist$ ) - an outlier if at least fraction of the objects in  $DS$  lie at the greater distance from  $p$ , it can generalize the several statistical tests. Then Ramaswamy et. al<sup>6</sup>, proposed the extension of the above method as they proposed a notion that all the outlier points are ranked based on the outlier score. Moreover, Anguelli and Pizzuti<sup>3</sup> proposed the way of ranking the outliers by considering the whole neighbourhood of the objects. In this case, the points are ranked on the sum of the distances from the  $k$ -nearest neighbours, rather than considering individual distance from centroid. Then breuing et. al.,<sup>7</sup> also proposed the local outlier factor to indicate the degree of outlierness in each outlier. They were first to quantify the outliers and used the term local outlier factor because only neighbourhood of each object is taken into account. It was a density method and has the stronger capability in a sense that a data object is gathered by how many members of its neighbourhood that decides its density. More the data object is denser, lesser the probability of being its outlier. Then, Zhang<sup>12</sup>, proposed the local distance based outlier detection method which is known as ldof function, which has its overall complexity as  $O(N^2)$ , where  $N$  is the no. of the points in the data set. Moreover, there are many clustering algorithms like DBSCAN<sup>8</sup> CURE<sup>9</sup> BIRCH<sup>10</sup> etc, to detect the outliers, but the limitations of these clustering algorithms is that they may optimize the clustering of various data objects but they do not optimize the detection of the outliers which is our prime priority. In this paper, we have used the pruning based algorithm PLDOF<sup>11</sup> which is actually pruning out the data items which are similar and only the outliers are detected and taken into consideration, we have tried to extend this work by using ODF function in order to make pruning more effective by enhancing the inner radius, which will be quantified according to their deviations.

In the past, many outlier detection methodologies have been proposed<sup>20,21,22</sup>. The outlier methodologies are broadly categorised into distribution based, distance based and density based methods. Statistical based methods usually follow the underlying assumptions of the data and this type of approach aims to find the outliers which deviate from such distributions. Most of the statistical distribution models are univariate, as the multivariate distributions lack the robustness. The solutions of the statistical models suffer from noise present in the data as the assumptions or the prior knowledge of the data distribution is not easily determined for the practical problems. In case of the distance based methods the distance between each point of interest and its neighbours

are calculated. This distance is compared with the predetermined threshold and if the data points lie beyond the threshold, then those points are termed as outliers, otherwise they are considered as the normal points<sup>13</sup>. In case of the multiclustered structured data, the data distribution is quite complex as no prior knowledge of the data distribution is needed. In such cases, improper neighbours are determined which enhances the false positive rate.

To alleviate this problem, the density based methods are proposed. LOF is one of the important density based methods to measure the outlierness of each data instance as LOF not only determines the outliers rather highlight the degree of outlierness, which provide the suspicious ranking scores for all the samples. Although it identifies the outliers in the local data structure via density estimation and also awares the user of the outliers that are sheltered under the global data structures, yet the estimation of the local density for each instance is computationally expensive, usually when the data size is very large. Apart from the above work, there are other outlier detection approaches that are recently proposed<sup>6,9,10</sup>. Among them, ABOD i.e., angle based outlier detection is the one which is very crucial and unique as it calculates the variation of the angles between each target instance and the remaining data points and it's usually observed that the outliers always produce a smaller angle variance than the normal one. Then further it's extension named as fast ABOD was proposed to generate the approximation of the original ABOD solution. K-means algorithm is an important clustering algorithm to cluster  $n$ -objects based on the attributes into  $k$ -partitions where  $k < n$ . It is quite similar to the expectation-maximisation algorithm for the different Gaussian mixture in a manner that they both attempt to find the centres of the natural clusters in the data. Moreover, it assumes that the object attribute forms a vector space. A major drawback of the K-means algorithm is its constant attempt to find out the local are used systematically. Projected clustering is a usual data mining task for unmanned grouping object<sup>17</sup>. In paper<sup>18</sup> they presented a detailed probabilistic approach to  $k$ -nearest neighbour classification. The feasibility of incorporating the hubness data for Bayesian class prediction is examined<sup>15</sup>. In<sup>16</sup> it is shown that the nearest neighbour methods can be further enhanced by taking a point in the hubness. In HDD, it's very difficult to handle the ordinary machine learning algorithms, which particularly characterize the problem of curse of dimensionality. In<sup>19</sup> they provided better assured proposed labels by exposing the fuzzy nearest neighbour classification and they also enlarged the already existing crisp hubness based approach into a fuzzy counterpart. Moreover, hybrid fuzzy functions are available and tested in detail in<sup>19</sup>.

### **3. CHALLENGES IN HIGHER DIMENSIONAL DATA**

#### **A. Dimensionality Reduction**

In higher dimensional space, unsupervised methods detect every point equally a good outlier because the distance becomes indiscernible as the dimensionality increases. All the data points become equidistant from each other showing that all of them carry useful and important information. Due to the sparsity of the data, the outliers are hidden in the lower dimensional subspaces and are usually do not prominently highlighted while dealing with independent components.

#### **B. Prediction of Relevant Attribute**

In case of the higher dimensional clustering, the relevant attributes contain the projected clusters and the irrelevant ones contain outliers or noise<sup>20</sup>. Cluster structure is the region with the higher density of the points than its surroundings. Such dense regions represent the 1-dimensional projection of the cluster. So, by detecting the dense regions in each dimension, it becomes easier to differentiate between the dimensions of the relevant and irrelevant clusters. The dense regions can be distinguished from the sparse one using the predefined density threshold, but in such a case, this value of density threshold may affect the accuracy and the goodness of the cluster.

In our work, we have used two measures LDOF<sup>12</sup> and PLDOF<sup>11</sup> in which former indicates how much a point is deviating from its neighbours and is a probable candidate of outliers and later prunes out the data items of clusters and focuses only on the outlier candidates only. The factor ldof is calculated as:

Ldof of  $m$ : The local distance based outlier factor of  $m$  is defined as:

$$\text{Ldof}(m) = \frac{d_m}{D_m} \quad (1)$$

$d_m$  is the KNN of  $m$ . If  $N_m$  is the set of  $k$ -nearest neighbours of object  $m$ . Let  $\text{dist}(m, n) \geq 0$  be the distance measure between objects  $m$  and  $n$ . The  $k$  nearest neighbour distance of the object  $m$  is

$$d_m = \frac{1}{k} \sum_{q \in N_m} \text{dist}(m, n) \quad (2)$$

$D_m$  is the kNN inner distance of  $m$ . Inner distance  $D_m$  is defined as:

$$D_m = \frac{1}{k(k-1)} \sum_{qq' \in N_m, q \neq q'} \text{dist}(q, q') \quad (3)$$

Based upon this ldof, PLDOF i.e., pruning based local outlier detection measure was proposed<sup>11</sup>. The main idea underlying the pruning based algorithm is to first cluster the complete data set into clusters and then the points which are not the outliers are pruned out.

#### 4. PROPOSED METHODOLOGY

In this section we will describe our proposed method which is a further improvement over ldof and pldof. The main idea of the effective pruning based outlier detection is to effectively prune the data items which are basically the part of the clusters and only consider the outliers which are deviating from the nearest neighbours. In this algorithm, the pruning is increased effectively by increasing the threshold distance due to which only the genuine outlier candidate are taken into account. Moreover, the KNN inner distance is also decreased due to which the percentage of pruned data items is enhanced. We briefly describe the steps that we need to perform by our pruning based algorithm

1. Generation of Clusters: Initially, K-Means algorithm is used to cluster the entire data set into  $c$  clusters and then the radius of each cluster is calculated. If any of the cluster contains less no. of points than the required no. of outliers, then the radius pruning will not be done for that cluster.
2. Pruning the points with the increased Threshold: Firstly, calculate the distance of each point from the centroid of the cluster and half the distance of the radius from the centroid so as to enhance the threshold value and then effectively prune out the elements which are the part of the clusters and the outliers are taken into consideration.
3. Quantify the Outliers: All the outliers detected are then quantified using ODF function which is based upon ldof in such a way that a numerical value is associated with it in such a way that a numerical value is associated with each outlier and hence the top  $n$  outliers are taken into consideration. Hence, the outlier points are computed in an effective way using ODF function.

The complexity of the K-Means algorithm is  $c \cdot it \cdot N$  where  $c$  is the no. of clusters to be formed,  $it$  is the no of iterations and  $N$  is the no. of the data points. Total computations in our method are  $c \cdot it \cdot N + c \cdot n_p + (x + N)^2$  where  $n_p$  represents average no. of points in each cluster and  $x$  indicates the fraction of the data points after pruning, which depends upon the threshold value as in our case. As the outliers are very less in our case because most

of the data points which are the cluster candidates are pruned out due to which the value of  $x$  is very small and hence the complexity is very much reduced even from  $O(N^2)$ .

## 5. DISCUSSION

### A. Outlier Metric-Outlier Defining Function

All the arbitrary points are divided into three discrete clusters  $c_i, c_j, c_k$ . Then the distance of the centroids  $c_1, c_2, c_3$  as shown in Figure 1 of the clusters  $c_i, c_j, c_k$  respectively is calculated from the origin using the ordinary distance formula of coordinate geometry. Then we define a metric known as ODF i.e., Outlier Defining Factor which defines and quantifies the outlieriness of the points. The higher value of the ODF indicates and tells that how much a point is deviating from its neighbours and probably it can be an outlier and provides its ranking among all the points.

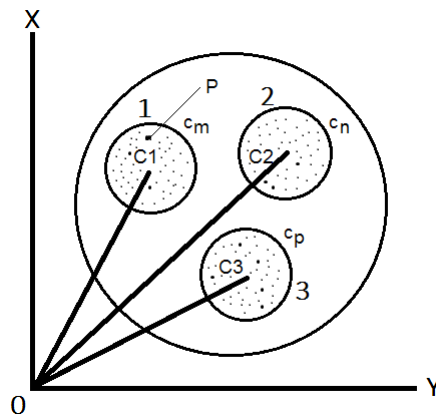


Figure 1: Shows a dataset with three clusters with their centroids and inner radius

We have considered the unsupervised way of learning here because we want that the outliers should be quantified well enough even though we do not have the training classes with us. To consider that the point  $p$  is an outlier or not, consider the first cluster, its distance  $\text{dist}(c_1, c_m)$  where  $c_m, c_n, c_p$  are the points on the cluster circumference or the farthest point on the cluster.  $c_i, c_j, c_k$  respectively.

$$\text{dist}(c_1, c_m) = \text{Radius}/2 = R_{0.5} \quad (4)$$

Now,  $\text{dist}(c_1, p) < R_{0.5}$ , then  $p$  is an outlier, but if  $\text{dist}(c_1, p) > R_{0.5}$ , then  $p$  is not an outlier and ODF is not calculated for it. If the later holds, then the data item is pruned out, otherwise the point  $p$  is termed as an outlier and the ODF has to be calculated using Euclidean or simple distance formulae of co-ordinate geometry. Similarly, various points of cluster  $c_i$  are taken and their corresponding distances are calculated and then the mean of all the points are taken and is termed as  $D(p)$  where,

$$\text{ODF} = D(p) = 1/k \sum_{p \in c_i} \text{dist}(c_1, p) \quad (5)$$

This ODF function is based on  $\text{ldof}$  [12] and hence it provides the information that how much a point is deviating from its neighbours but with a difference that we have calculated the inner radius  $R_{0.5}$  and every point will be compared with it in order to decide whether it is an outlier or not. In case of outliers, its ODF will be calculated otherwise the data objects will be pruned out.

### B. Effective Pruning Outlier Detection Algorithm

Dataset  $X$ , iteration: no. of loops, cluster\_no, assump\_outlier

```

Step 1. Set  $X \leftarrow K\text{ Means}(c, k, S)$ 
for  $i=0$  to  $cluster\_no\_1$ 
Do
Cluster_centrei = Point  $P_i \in \text{Dataset}$ 
Done
End for
While (iteration-----)
Do
For  $i=0$  to  $\text{Dataset} - 1$ 
Do
If ( $dist\_mean_j[i] < dist\_mean_{j+1}[i] < dist\_mean_{j+2}[i]$ )
cluster_no[i] = j
cluster  $c_j \leftarrow \text{Point } P_i$ 
endif
 $\forall$  Cluster  $c_i$ 
do
cluster_centrei =  $\sum_{j=1}^n P_j \in C_i / C_i$ 
done
 $\forall$  Cluster ( $c_j$ )
Do
Radiusi  $\leftarrow$  Radius ( $c_j$ )
Do
If ( $c_i.\text{elemnt} > \text{assump\_outlier}$ )
For  $j= 0$  to  $c_i.\text{elemnt}-1$ 
Do
If ( $dist(P_i \in c_i, cluster\_centre) \leq (\text{radius} * 0.8)$ )
Prune ( $P_i$ )
Else
Move  $P_i$  to resultant_cluster
End if
Done
Else

```

$\forall$  point  $P_i \in C_i$

Do

Move  $P_i$  to resultant \_cluster

Done

End if

Done

$\forall$  Point  $P_i$ ,  $T$  is the resultant cluster

Do

ODF ( $P_i$ ) // Applying the ODF function

Find  $n$  points with higher ODF values and the desired outliers.

.According to the Implementation, the graph is plotted for the desired outliers and cluster points.

Done

## 6. EXPERIMENTAL RESULTS

In this section, we have compared the outlier detection performance of our Effective Pruning Based Outlier Detection method with the PLDOF and LDOF methods

*WDBC (Medical Diagnosis Data)*: To validate our experiment, we have used the medical data set WDBC, (diagnosis) from UCI repository which has already been further used by the nuclear feature extraction for the Breast Cancer Diagnosis. This data set contains 569 medical diagnosis records, each with 32 real valued input features. This diagnosis is Binary i.e., cancer data can be Benign or Malignant. We assume that the objects labelled as benign are normal data whereas the malignant are considered as abnormal one or an outlying data. In our experiment, initially, we use all the 360 Benign Diagnosis records as normal objects and added five malignant records in them as outliers. This process is repeated a no. of times by varying the values of the neighbourhood objects. Every time, the value of the neighbourhood size i.e.,  $k$  is varied (may be increased or decreased). The three measures that are highly affected by varying the neighbourhood size  $k$  are (i) The  $n$ -top potential outliers (ii) The detection of the precision (iii) The percentage of the data pruned. In order to verify this effect, we will repeat this experiment 10 times by adding the random no. of outliers (Malignant cancer data) every time. With the help of the various independent runs, say from 10 to 60, the average detection precision is calculated and both the top  $n$ -outliers and the percentage of the pruned data are varied.

In EPLDOF, the percentage of the pruned data is more as compared to the PLDOF. This is basically because of the ODF function that we have used in our EPLDOF algorithm. Due to the increase in the threshold value for the inner radius of the clusters, more data is pruned out and when the data gets much more pruned, the time and the space complexity are much more suppressed and the outliers will be more prominently highlighted. Moreover, when the precision is compared with LDOF and PLDOF, EPLDOF reaches at par at  $k = 30$  even though the 60% of the data is pruned initially. Hence the time complexity and computation time are further decreased.

## 7. CONCLUSION

In this paper, we have proposed an effective and an efficient outlier detection algorithm which is based on already existing methods LDOF and PLDOF but with the major difference that the pruning has been done very effectively

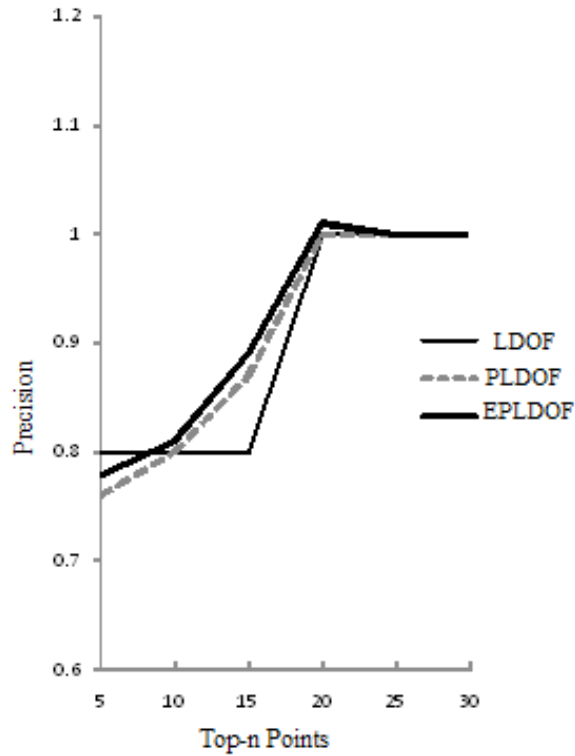


Figure 2: Graphical Comparison of LDOF, PLDOF and EPLDOF

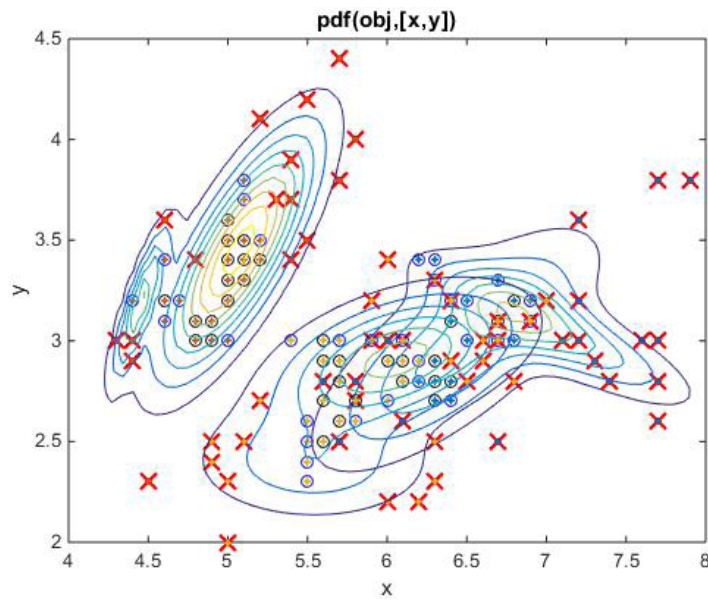


Figure 2: Clusters Using Gaussian Mixture at Max Threshold

using ODF function. Figure 2 shows the precision of the three methods against top n points and neighbourhood size respectively which gives the exact comparison of these methods. All the points which are not the probable candidates of the outliers are pruned out and the remaining points which are termed as outliers are taken into consideration and are further quantified for their outlieriness. Secondly we have also tried to use the Gaussian Mixture model to optimise our clustering algorithm as shown in Figure 3. The time and the space complexities



are drastically reduced and hence the computation cost is also suppressed than the previously existing methods and the accuracy is quite high even though our pruned data is comparatively higher.

### **Acknowledgement**

Authors acknowledge the opportunity and support provided by I.K Gujral Punjab Technical University, Jalandhar to conduct the present work.

### **REFERENCES**

- [1] B. V and L. T, Wiley John and Sons “*Outliers in Statistical Data*”. New York, 1994.
- [2] F. Angiulli, S. Basta, and C. Pizzuti “*Distance-based detection and prediction of outliers*” *IEEE Transactions on Knowledge and Data Engineering*, 18:145–160, 2006.
- [3] F. Angiulli and C. Pizzuti. Outlier mining in large high dimensional data sets. *IEEE Transactions on Knowledge and Data Engineering*, 17:203–215, 2005.
- [4] E. M. Knorr, R. T. Ng, and V. Tucakov. Distance-based outliers: algorithms and applications. *The VLDB Journal*, 8(3-4):237–253, 2000.
- [5] E. M. Knorr and R. T. Ng. Algorithms for mining distance based outliers in large datasets. In *Proc. 24th Int. Conf. Very Large Data Bases, VLDB*, pages 392–403, 1998.
- [6] S. Ramaswamy, R. Rastogi, and K. Shim. *Efficient algorithms for mining outliers from large data sets. pages 427–438, 2000.*
- [7] M. M. Breunig, H.-P. Kriegel, R. T. Ng, and J. Sander. *LoF: identifying density-based local outliers. SIGMOD Rec., 29(2):93–104, 2000.*
- [8] R. T. Ng and J. Han. *Efficient and effective clustering methods for spatial data mining. pages 144–155, 1994.*
- [9] S. Guha, R. Rastogi, and K. Shim. *CURE: An efficient clustering algorithm for large databases. SIGMOD Rec., 27(2):73–84, 1998.*
- [10] T. Zhang, R. Ramakrishnan, and M. Livny. *Birch: an efficient data clustering method for very large databases. SIGMOD Rec., 25(2):103–114, 1996. 256*
- [11] Rajendra Pamula, Jatindra Kumar Deka, Sukumar Nandi. *An Outlier Detection Method based on Clustering, Second International Conference on Emerging Applications of Information Technology, 2011.*
- [12] K. Zhang, M. Hutter, and H. Jin. *A new local distance-based outlier detection approach for scattered real-world data. In PAKDD '09: Proceedings of the 13th Pacific Asia conference on Advances in Knowledge Discovery and Data Mining pages 813–822, 2009.*
- [13] M. Breunig, H.-P. Kriegel, R.T. Ng, and J. Sander, “*LOF: Identifying Density-Based Local Outliers,*” *Proc. ACM SIGMOD Int’l Conf. Management of Data*, 2000.
- [14] H.-P. Kriegel, M. Schubert, and A. Zimek, “*Angle-Based Outlier Detection in High- Dimensional Data,*” *Proc. 14th ACM SIGKDD Int’l Conf. Knowledge Discovery and data Mining*, 2008.
- [15] F. Angiulli, S. Basta, and C. Pizzuti, “*Distance-Based Detection and Prediction of Outliers,*” *IEEE Trans. Knowledge and Data Eng.*, vol. 18, no. 2, pp. 145-160, 2010
- [16] V. Barnett and T. Lewis, *Outliers in Statistical Data*. John Wiley&Sons, 1994.
- [17] W. Jin, A.K.H. Tung, J. Han, and W. Wang, “*Ranking Outliers Using Symmetric Neighborhood Relationship,*” *Proc. Pacific-Asia Conf. Knowledge Discovery and Data Mining*, 2006.

- [18] N.L.D. Khoa and S. Chawla, "Robust Outlier Detection Using Commute Time and Eigenspace Embedding," Proc. Pacific-Asia Conf. Knowledge Discovery and Data Mining, 2010.
- [19] E.M. Knox and R.T. Ng, "Algorithms for Mining Distance-Based Outliers in Large Data Sets," Proc. Int'l Conf. Very Large Data Bases, 1998.
- [20] H.-P. Kriegel, P. Kroger, E. Schubert, and A. Zimek, "Outlier Detection in Axis-Parallel Subspaces of High Dimensional Data," Proc. Pacific-Asia Conf. Knowledge Discovery and Data Mining, 2009.
- [21] C.C. Aggarwal and P.S. Yu, "Outlier Detection for High Dimensional Data," Proc. ACM SIGMOD Int'l Conf. Management of Data, 2012
- [22] C.C. Aggarwal and P. S. Yu, —Finding generalized projected clusters in high dimensional spaces, in Proc. 26th ACM SIGMOD Int. Conf. on Management of Data, 2000, pp. 70–8
- [23] K. Kailing, H.-P. Kriegel, P. Kroger, and S. Wanka, Ranking subspaces for clustering high dimensional data, in Proc. 7th European Conf. on Principles and Practice of Knowledge Discovery in Databases (PKDD), 2003, pp. 241–252.