



OUTLIER DISCERNMENT USING K-MEANS

Madhav Bokare

Research Scholar

Priyadarshini Institute of Engineering & Technology,
Nagpur.

V. M Thakare

Professor

Sant Gadge Baba Amravati University,
Amravati.

ABSTRACT

In this paper we propose a clustering based method to capture outliers. We apply K-means clustering algorithm to divide the data set into clusters. We present a unified approach for simultaneously clustering and discovering outliers in data. Our approach is formalized as a generalization of the k-means problem. In this paper shows the outlier of sample data set by using K-means algorithm.

Keywords— Apriori, Bary, CCT.

I. INTRODUCTION

Outlier detection is a fundamental issue in data mining; specifically it has been used to detect and remove anomalous objects from data. Outliers, also called contaminant observations, are data points that deviate from other data points that they seem to be generated because of any of the faulty condition in the experimental setup. When observations that are taken in the experimental setup are subjected to analysis there is a possibility of the two conditions based on the outlier. Either (i) outliers negatively influence the results of analysis, or (ii) the search for outliers is the main task of data analysis. In data mining outlier detection is also regarded as the detection of anomaly. In many applications, a set of training values is required to define „normality“. Security applications are examples, in which a typical behaviour by people or technical systems has to be detected. Outlier tests are used with the statistical model for generating the observations and presume some knowledge of the number of assumed outliers. Many of them can only cope with a single outlier. Outlier detection is used in various domains in data mining. This has resulted in a huge and highly diverse literature of outlier detection techniques. A lot of these techniques have been developed in order to solve problems based on some of the particular features, while others have been

developed in a more generic fashion. It has been argued by many researchers whether clustering algorithms are an appropriate choice for outlier detection. For example, in (Zhang and Wang, 2006) [2], the authors reported that clustering algorithms should not be considered as outlier detection methods. This might be true for some of the clustering algorithms, such as the k-means clustering algorithm (Mac Queen, 1967) [3]. This is because the cluster means produced by the k-means algorithm is (Laan, 2003). But here we propose an algorithm that uses clustering efficiency of the k-means algorithm and uses a hybridized outlier detection technique of finding outlier through clustering.

In this paper we will propose a generalization of the k-means problem with the aim of simultaneously clustering data and discovering outliers. A naive approach to apply the k-means algorithm and list as outliers the top points that are the furthest away from their nearest cluster centres. However, there is a subtle point that needs to be noted: the k-means algorithm itself is extremely sensitive to outliers, and such outliers may have a disproportionate impact on the final cluster configuration. This can result in many false negatives: i.e., data points that should be declared outliers are masked by the clustering and also false positives: data points that are incorrectly labelled as outliers.

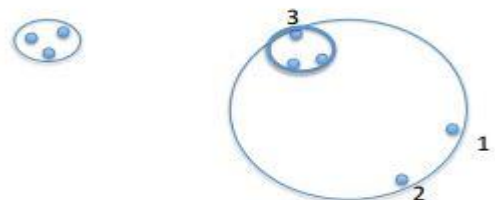


Figure 1: The k-means algorithm is extremely sensitive to outliers. By removing two points (1) and (2), we can obtain much tighter clusters (the

In the figure 3 shows the some options such as Select Excel File, Copy to Array, CCT, K-means. Here we are using sample data set. In left panel the selected data highlighted. Then click on Copy to Array option then CCT and then K-means. Finally output will display as graphical manner.

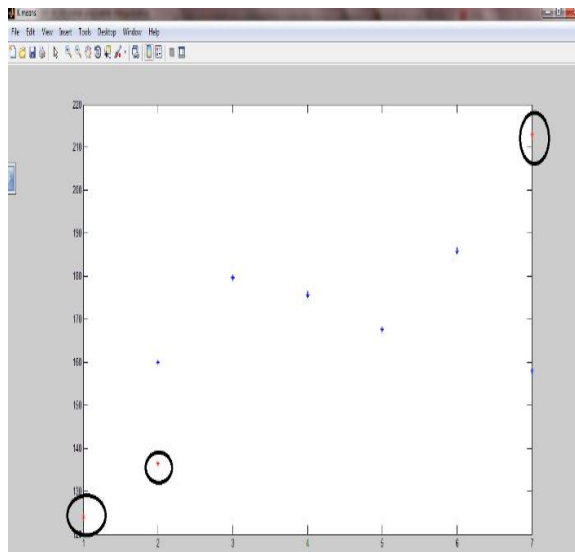


Figure 4Outlier using K-means

IV. CONCLUSIONS

In this paper we have shown the outlier of sampled data set by applying the k-means algorithm..

ACKNOWLEDGMENT

We wish to thank Dr. S. B. Thorat, Director,ITM,and Nanded.

REFERENCES

- [1] “An Effective Clustering-Based Approach for Outlier Detection” by Moh’d Belal Al- Zoubi European Journal of Scientific Research Vol.28 No.2 (2009).
- [2] “Detecting outlying subspaces for high-dimensional data: the new Task, Algorithms, and Performance, Knowledge and Information Systems”, Zhang, J. and H. Wang, 2006.10 (3): 333-355.
- [3] “Some methods for classification and analysis of multivariate observations” by MacQueen, J.,1967. Proc. 5th Berkeley Symp. Math. Stat. and Prob, pp. 281 -97.s.
- [4] Web-link:
<https://sites.google.com/site/dataclusteringalgorithms/k-means-clustering-algorithm>