

# New Substructure Based on Regularized Clustering Functions

<sup>1</sup> B. Ramakrishna, M.Tech, Asst.Professor, & <sup>2</sup> S. Krishna Reddy, M.Tech, Asst.Professor,  
Dept.of.CSE, Kallam Haranadhareddy Institute of Technology, Guntur, Andhra Pradesh, India.  
<sup>1</sup> [ramakrishna.badiguntla@gmail.com](mailto:ramakrishna.badiguntla@gmail.com), <sup>2</sup> [kreshnareddy@gmail.com](mailto:kreshnareddy@gmail.com)

**Abstract:** Feature selection is elimination irrelevant and redundant dimensions by analyzing the entire dataset. Subspace clustering algorithms localize find relevant dimensions clusters backend many possibly overlapping subspaces. In latest years a lot of works is using deep neural networks is learn a clustering-friendly representation retorts in a significant deployments clustering results. We define a clustering objective function using relative entropy minimization regularized by a prior for the frequency of cluster results. An alternating strategy is then derived to modify objective and updating parameters and estimating cluster assignments. Many outlier detection methods is proposed till date. This proposed system is broadly categorized as distribution based clustering-based, density-based and model based approaches.. We introduce a nonparametric density estimation system is modify anisotropic kernels. It is shown that this method provides accurate configurationally entropies for dimensions is improving on the quasiharmonic approximation. We compare the two main methods to subspace clustering using empirical scalability and accuracy tests and discuss some potential applications where subspace clustering could be particularly useful.

**Index Terms:** projected clustering, high dimensional data, deep learning, data representation, network architecture, Outlier detection,.

## 1. INTRODUCTION

Entropies are key quantities in physics, chemistry and biology is energy optimize govern the direction of all chemical processes including reaction equilibrium entropy changes are the underlying driving forces of legend binding, security folding other phenomena driven by hydrophobic forces [1]. In very high dimensions it is common for all of the objects in a dataset to be nearly equidistant from each other, completely masking the clusters. Feature selection application is employed some new successfully to improve cluster quality [2]. In present deep embedding clustering is proposed and followed by other novel methods making deep clustering become a popular research field [3]. The clustering problem is extensively studied in many applications the performance of standard clustering algorithms is adversely affected when dealing with high-dimensional data, and their time complexity dramatically increases when working with large-scale datasets [4]. Outlier detection has new applications in numerous fields. Some of them are: detecting fraudulent applications for credit cards detecting deceptive usage of credit cards and mobile phones to detect fraudulent applications potentially problematical customers [5]. In traditional clustering each dimension is equally weighted when computing the distance between points. Most of these algorithms perform well in clustering low-dimensional data sets in higher dimensional feature spaces, their performance and efficiency deteriorate to a greater

extent due to the high dimensionality [6]. Direct method is additional advantages such methods are independent from finding suitable perturbation pathways between the states of interest and/or analytical [7].



(a) Raw Data (b) NonJoint DEPICT (c) Joint DEPICT

Figure 1: Visualization discriminative capability of embedding subspaces

## 2. RELATED WORKS

These principals are very successful and uncovering latent structure in datasets. However the preserve relative distances between objects they are less effective when there are large numbers of irrelevant attributes that hide the clusters noise [8]. Also the new features are combinations of the originals and may be very difficult to interpret the new features in the context of the domain [9]. A fully-connected network (FCN) consists is multiple layers of neurons every neuron is connected to every neuron in the previous layer, and each connection has its own weight. The FCN is also known as multi-layer perception (MLP) [10]. . Our DEPICT algorithm is

discriminative clustering model is auxiliary reconstruction task to alleviate this model training of our discriminative clustering algorithm [11]. The normal class is taught algorithm learns to recognize abnormality. It model gradually as new data arrives, tuning the model to improve the fit as each new epitome becomes available [12]. This algorithm combines density and grid-based clustering and uses an APRIORI style technique to find cluster subspaces. This may be because the unit densities vary in different subspace cardinalities such as the identification of the dense units in all subspaces is used absolute unit density threshold [13]. The classical treatment of those stiff degrees of freedom yields unphysical entropies by allowing unlimited sharpness of the threads [14].

### 3. SYSTEM MODELS

The bottom-up search method is advantage of the downward closure property of density to reduce the search space using an APRIORI methods. Algorithms first create a histogram for each dimension and selecting those bins with densities of threshold [15]. The algorithm starts with an arbitrary dense unit and greedily grows a maximal region in each dimension until the union of all the regions covers the entire cluster [16]. The regularization is encourages balanced cluster assignments and removed allocating clusters to outlier samples. Furthermore the reconstruction loss of auto encoder is also employed to prevent corrupted feature representation [17].

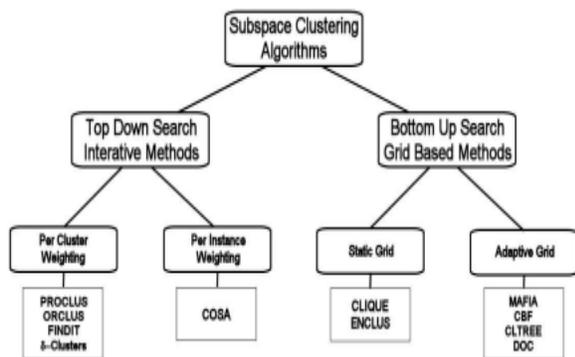


Figure 2: Hierarchy of Subspace Clustering.

### 4. PROPOSED SYSTEM

New genetic algorithm, is particular class of evolutionary algorithms is recognized to be well suited to multi-objective modifications problems. In our work we employ multi objective subspace clustering algorithm for clustering data sets based on subspace approach [18]. They provide AntiHub method using reverse nearest neighbor counts for outlier detection. This method can efficiently find

outliers in high dimensional dat. DEPICT consists of a soft-max layer stacked on top of a multilayer convolutional autoen coder[19]. The clustering loss can be k-means loss, agglomerative clustering loss, locality-preserving loss and so on. For deep clustering methods based on AE network, the network loss is essential. But some other work designs a specific clustering loss to guide the optimization of networks, in which case the network loss can be removed [20].

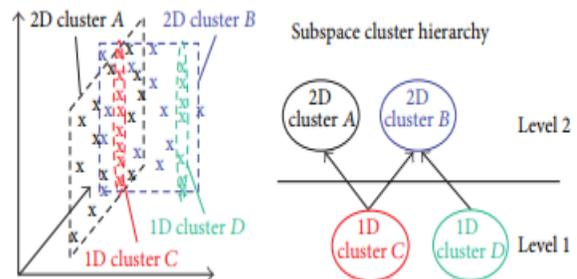


Figure 3: An example of subspace clusters.

### 5. OUTLIER DETECTION TECHNIQUES

1) **Unsupervised:** It is the process in which no information about the dataset class distribution is available beforehand. This approach is widely used now a day.

2) **Supervised:** The dataset consists of class objects is classified as normal or abnormal. But the limitation of FMN method is that, user has to tune the parameters to get good recognition accuracy. The recognition accuracy at the cost of recall time is increased in the above stated method.

3) **Semi-supervised:** This method is use pre-classified data but only learns data which is marked normal. The normal class is taught but the algorithm learns to recognize abnormality. It can learn the model gradually as new data arrives, tuning the model to improve the fit as each new epitome becomes available [21]

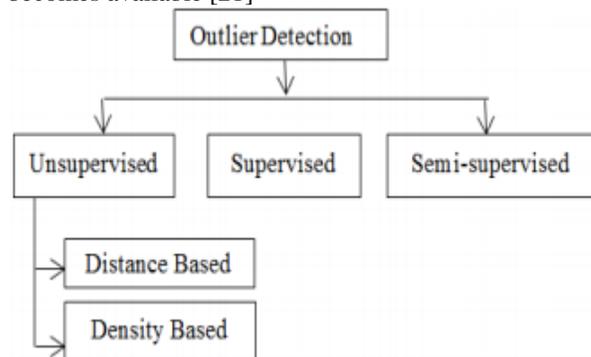


Fig -4: Modes of operation of outlier detection techniques

### 6. MULTI OBJECTIVE SUBSPACE CLUSTERING

The goal of preprocessing step is to identify all dimensions in a data set which exhibit some cluster structure by discovering dense regions and their location in each dimension is identified dimensions represent potential candidates for relevant dimensions of the subspace clusters [22] [23]. Individual's genetic material being passed to the next generation. Here we adopt the tournament selection method because its time complexity is low [24].

1. begin
2. for  $i=1$  to
3.  $P_{pop}$  do  $P \square$  best fitted item among  $N$  tour elements randomly
4. selected from  $P$ ;
5. return  $P \square$
6. end begin for  $i=1$  to  $P_{pop}$  do  $P \square$  best fitted item among  $N$  tour elements randomly selected from  $P$ ; return  $P \square$  en

### Deep Embedded Regularized Clustering (DEPICT):

DEPICT is a sophisticated method consisting of multiple striking tricks. It consists of a soft ax layer stacked on top of a multi layer convolution auto encoder. It minimizes a relative entropy loss function with a regularization term for clustering [25]. For another, the optimization procedures of these methods involve discrete reconfigurations of the objective, which require updating the clustering parameters and network parameters alternatively. We divide CDNN-based deep clustering algorithms into three categories according to the ways of network initialization, i.e., unsupervised pre-trained, supervised pre-trained and randomly initialized [26].

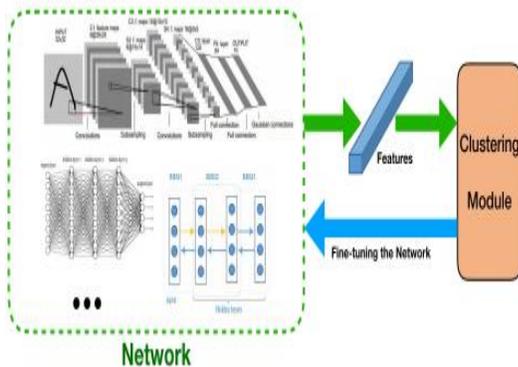


FIGURE 5. Architecture of CDNN-based deep clustering algorithms

### 7. EMPIRICAL COMPARISON

We representative top-down and bottom-up algorithms is well each of the algorithms were able ever cluster. The implementation of each algorithm was provided by the respective authors.

**Data:** To facilitate the comparison of the two algorithms, we chose to use synthetic datasets so that we have control over of the characteristics of the data is comparing the output of the algorithms to the known input clusters.

**Clustering Metrics:** We have used 2 of the most popular evaluation criteria widely used for clustering algorithms normalized mutual information (NMI). The best mapping between cluster assignments and true labels is computed using the Hungarian algorithm [16] to measure accuracy.

**Scalability with Average Cluster Dimensionality:** The dependency of the execution time on the average cluster dimensionality where the latter increases from 10 to 50 in a 100-dimensional data space. In each case, the dataset has 50000 points distributed over 5 clusters with a fixed 10 percent level of outliers.

The experiments show that our algorithm is able to detect clusters and their relevant dimensions accurately in various situations. The results of PROCLUS are less accurate than those given by MOSCL.

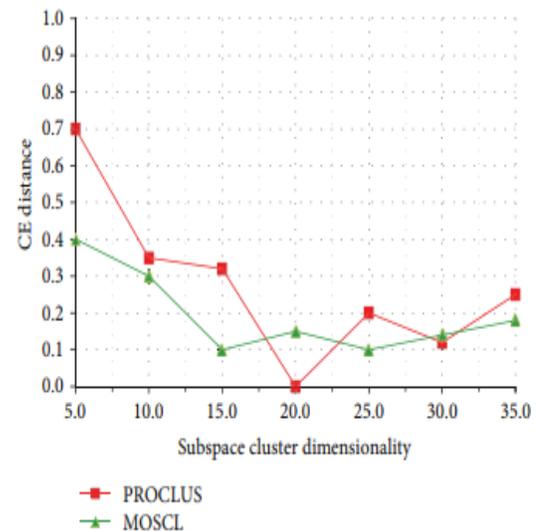


Figure 6: Distances between MOSCL output and the true clustering

## 8. CONCLUSIONS AND FUTURE OPPORTUNITIES

We have introduced a kernel-based density estimation method and showed that it provides accurate results in even the high-dimensional and quite complex configurationally space density generated by the dynamics of biological macromolecules. We have presented a robust MOSCL algorithm for the challenging problem of high-dimensional clustering and illustrated the suitability of our algorithm in tests and comparisons with previous work. Furthermore, a joint learning framework was introduced to train all network layers simultaneously and avoid layer-wise pertaining. Based on the existing literature and review, we argue that the following perspectives of deep clustering are worth being implemented further.

## REFERENCES

[1] A. K. Jain, M. N. Murty, P. J. Flynn, "Data clustering: A review", *ACM Comput. Surv.*, vol. 31, no. 2, pp. 264-323, 2015

[2] C. C. Aggarwal. Re-designing distance functions and distance-based applications for high dimensional data. *ACM SIGMOD Record*, 30(1):13–18, 2001.

[3] C. C. Aggarwal. Towards meaningful high-dimensional nearest neighbor search by human-computer interaction. In *Data Engineering, 2002. Proceedings. 18th International Conference on*, pages 593–604, 2002.

[4] C. C. Aggarwal, A. Hinneburg, and D. A. Keim. On the surprising behavior of distance metrics in high dimensional space. In *Database Theory, Proceedings of 8th International Conference on*, pages 420–434, 2001.

[5] J. A. Hartigan and M. A. Wong, "Algorithm AS 136: A k-means clustering algorithm," *J. Roy. Stat. Soc. C, Appl. Stat.*, vol. 28, no. 1, pp. 100–108, 1979.

[6] S. Wold, K. Esbensen, and P. Geladi, "Principal component analysis," *Chemometrics Intell. Lab. Syst.*, vol. 2, nos. 1–3, pp. 37–52, 1987.

[7] T. Hofmann, B. Schölkopf, and A. J. Smola, "Kernel methods in machine learning," *Ann. Stat.*, vol. 36, no. 3, pp. 1171–1220, 2008.

[8] A. Y. Ng, M. I. Jordan, and Y. Weiss, "On spectral clustering: Analysis and an algorithm," in *Proc. Adv. Neural Inf. Process. Syst.*, 2002, pp. 849–856.

[9] J. Schmidhuber, "Deep learning in neural networks: An overview," *Neural Netw.*, vol. 61, pp. 85–117, Jan. 2015.

[10] J. R. Hershey, Z. Chen, J. Le Roux, and S. Watanabe, "Deep clustering: Discriminative embeddings for segmentation and separation," in *Proc. IEEE Int. Conf. Acoust., Speech Signal Process. (ICASSP)*, Mar. 2016, pp. 31–35.

[11] K. Kailing, H.-P. Kriegel, and P. Kroger. Density-connected subspace clustering for high-dimensional data. In *Proceedings of the 2004 SIAM International Conference on Data Mining*, pages 246–256. SIAM, 2004.

[12] E. Keogh, S. Chu, D. Hart, and M. Pazzani. An online algorithm for segmenting time series. In *Data Mining, 2001. ICDM 2001, Proceedings IEEE International Conference on*, pages 289–296. IEEE, 2001.

[13] D. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

[14] D. P. Kingma, S. Mohamed, D. J. Rezende, and M. Welling. Semi-supervised learning with deep generative models. In *Advances in Neural Information Processing Systems (NIPS)*, pages 3581–3589, 2014.

[15] A. Krause, P. Perona, and R. G. Gomes. Discriminative clustering by regularized information maximization. In *Advances in neural information processing systems (NIPS)*, pages 775–783, 2010.

[16] H. W. Kuhn. The hungarian method for the assignment problem. *Naval research logistics quarterly*, 2(1-2):83–97, 1955.

[17] K. Zhang, M. Hutter, and H. Jin, "A new local distancebased outlier detection approach for scattered realworld data," in *Proc 13th Pacific-Asia Conf. Knowl. Discovery Data Mining*, 2009, pp. 813–822.

[18] H.-P. Kriegel, P. Kroger, E. Schubert, and A. Zimek, "LoOP: Local outlier probabilities," in *Proc 18th ACM Conf. Inform. Knowl. Manage.*, 2009, pp. 1649–1652.

[19] Liangwei Zhang, Jing Lin, Ramin Karim, "An anglebased subspace anomaly detection approach to highdimensional data: With an application to industrial fault detection." in *Reliability Engineering and System Safety* 142 (2015) 482-497 Elsevier.

[20] Jayanta K. Dutta, Bonny Banerjee, Member, IEEE, and Chandan K. Reddy, Senior Member, IEEE, "RODS: Rarity based Outlier Detection in a Sparse Coding Framework" in *IEEE Transactions on Knowledge and Data*

Engineering, Volume: PP Issue: 99 September 2015.

[21] Jabez J, Dr.B.Muthukumar, "Intrusion Detection System (IDS): Anomaly Detection using Outlier Detection Approach" in International Conference on Intelligent Computing, Communication & Convergence (ICCC-2015) Elsevier.

[22] W. Wang, J. Yang, and R. Muntz, "STING: a statistical information grid approach to spatial data mining," in Proceedings of the 23rd International Conference on Very Large Data Bases (VLDB '97), vol. 97, pp. 186–195, 1997.

[23] C. C. Aggarwal, C. Procopiuc, J. L. Wolf, and J. S. Park, "Fast algorithms for projected clustering," in Proceedings of the 1999 ACM SIGMOD International Conference on Management of Data, vol. 28, no. 2, pp. 61–72, 1999.

[24] A. Patrikainen and M. Meila, "Comparing subspace clusterings," IEEE Transactions on Knowledge and Data Engineering, vol. 18, no. 7, pp. 902–916, 2006.

[25] R. Agrawal, J. Gehrke, D. Gunopulos, and P. Raghavan, "Automatic subspace clustering of high dimensional data for data mining applications," in Proceedings of the ACM-SIGMOD Conference on the Management of Data, vol. 27, no. 2, pp. 94–105, 1998.

[26] K. Kailing, H. P. Kriegel, and P. Kroger, "Density-connected subspace clustering for high-dimensional data," in Proceedings 4th SIAM International Conference on Data Mining (SDM '04), p. 4, 2004