# Cross Breed Clustering Algorithm for High Dimensional Data

Y.Vijay Bhaskar Reddy, L.S.S Reddy,.S.S.N.Reddy

*Abstract: Clustering plays a major role in machine learning and also in data mining. Deep learning is fast growing domain in present world. Improving the quality of the clustering results by adopting the deep learning algorithms. Many clustering algorithm process various datasets to get the better results. But for the high dimensional data clustering is still an issue to process and get the quality clustering results with the existing clustering algorithms. In this paper, the cross breed clustering algorithm for high dimensional data is utilized. Various datasets are used to get the results.*

*Keywords: clustering, data mining, deep learning.*

## I. INTRODUCTION

Clustering could be a standout amongst the foremost vital systems for dissecting info in AN unattended means, it's a large scope of utilization together with laptop vision [1, 2, 3], traditional language getting ready [4, 5, 6] and bioinformatics [7, 8]. Within the previous decades, innumerable calculations are projected to modify clustering problems [9, 10]. Be that because it might, no algorithm that matches multiple issues. This method depends upon the data to modify and also the explicit assignment. Generally, there square measure 2 arrangements of methodologies, the element based mostly grouping calculations and also the likeness based clustering calculations. The overwhelming majority of them arrange to find the inborn info structure from the primary component area or the essential topological space.

Among these algorithms, K-implies [11] and mathematician Mixture Models (GMM) [12] square measure 2 renowned component primarily based ways. K-implies makes arduous clustering that relegates every example to its nearest cluster focus. GMM settle for that info square measure created from some autonomous mathematician dispersions and makes an attempt to see these appropriations from the data. Therefore, it makes delicate assignments. Be that because it might, the 2 of them do clustering within the initial element area. Ghastly grouping [13] could be a delegate calculation of similitude primarily based clustering or topological space clustering ways.

The bulk of these methodologies begin with structure a feeling lattice and endeavour the primary info to a straight topological space. At last, clustering is completed within the topological space.

Data clustering is a necessary issue in various regions, as an example, AI, style acknowledgment, PC vision, info pressure. the target of grouping is to order comparative info into one cluster passionate about some closeness measures (e.g., geometer separation). Despite the very fact that a massive range of data clustering techniques is projected [14]–[15], ancient grouping ways for the foremost half have frightful showing on high-dimensional info, due to the wastefulness of likeness estimates utilized in these techniques. Besides, these ways by vast} expertise the sick effects of high machine unpredictability on huge scale datasets. Consequently, spatiality decrease and highlight amendment techniques are wide targeted to stipulate crude info into another element area, wherever they made info square measure easier to be isolated by existing classifiers. As a rule, existing info amendment ways incorporate straight amendment like Principal component analysis (PCA) and non-direct amendment, as an example, kernel techniques, and spectral techniques.

In any case, a deeply unpredictable dormant structure of data is heretofore testing the viability of existing clustering ways. Inferable from the advancement of deep adapting, deep neural systems (DNNs) will be utilized to amendment the data into all a lot of clustering benevolent portrayals due to its inborn property of exceptionally non-direct change. For the straightforwardness of depiction, we have a tendency to use the cross breed clustering algorithm with deep learning as deep clustering during this paper.

## II. LITERATURE SURVEY

Clustering may be a broadly speaking examined zone, and up to currently, several clustering methods are created. Here, we tend to audit in the main on the clustering methods that utilize decilitre systems, and quickly feature the points of interest.

From the many of the clustering methods, K-means and GMM square measure broadly speaking utilized in various applications. Nevertheless, they need 2 disadvantages: one is that they for the foremost half add the primary element space; the opposite is that they can not modify monumental and high-dimensional datasets are present.

Spectral clustering and its variations square measure broadly speaking thought among mathematical space clustering methods.

Builds up a sent system to tackle scanty mathematical space grouping by means that of ADMM. In any case, it worries simply direct subspaces. to deal with this issue, another methodology [16] was planned to consolidate nonlinearity into mathematical space clustering.

The goal is to limit the knowledge is done with data re-creation and add sparsity to improve the deep learning features. To utilize the features of deep learning, for some process the initial training is required to get the accurate results.
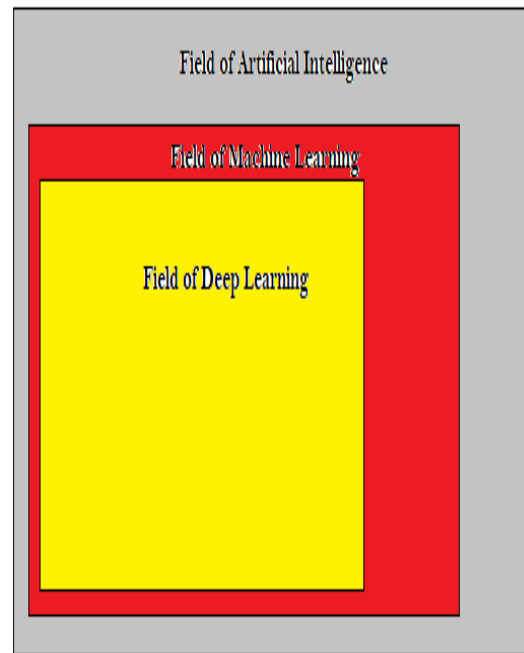
These methodologies have different component learning and clustering. To direct begin to complete grouping in deep systems, [18] proposes a model to at a similar time get acquainted with the deep portrayals and therefore the cluster focuses. It makes the exhausting task to every example and squarely will clustering on the shrouded highlights of the deep auto encoder. Associate in Nursing current endeavour is that the Deep Embedding cluster (DEC) strategy [19], that accomplishes detail of-the-workmanship results on various datasets, nonetheless it would bombs once firmly connected cluster exist.

In [20], various clustering methods are implemented experimentally such as multivariate Gaussian mixture (MGM), hierarchical clustering (HC), spectral and nearest neighbour (SNN) methods. Four closeness measures were utilized within the investigations: Pearson and Spearman relationship constant, trigonometric function alikeness and therefore the geometrician separation. The calculations were assessed with regards to 35 gene expression data info from either Affymetrix or cDNA chip stages, utilizing the balanced and record for execution assessment. The variable mathematician mix strategy gave the simplest execution in **recuperating** the real range of clusteres of the datasets. The k-implies technique showed comparative execution. during this equivalent investigation, the progressive strategy prompted affected execution, whereas the unearthly technique gave the impression to be particularly touchy to the closeness live utilized.

One issue with most element primarily based clustering methods is that they cannot scale well to high-dimensional info thanks to the scourge of spatial property. In high-dimensional info examination, it's progressively wise to consider some as reduced and delegate includes instead of the complete component area. As of late, Deep learning (DL) has been created and with an implausible accomplishment in various territories, as an example, image order and discourse acknowledgment [19]. decilitre means that to require in an incredible portrayal from the crude info through abnormal state non-straight mapping. As of late, the way to utilize a deep portrayal to boost grouping execution turns into a hot analysis subject.

## III. DEEP LEARNING

The field of AI is expansive and has been around for quite a while. Deep learning is a subset of the field of AI, which is a subfield of AI. The features that separate deep taking in systems all in all from "accepted" feed-forward multilayer systems are as per the following:



**Figure: 1 Roots of Deep Learning**

When I state "more neurons", I imply that the neuron check has ascended throughout the years to express progressively complex models. Layers additionally have developed from each layer being completely associated in multilayer systems to privately associated patches of neurons between layers in Convolutional Neural Networks and intermittent associations with a similar neuron in Recurrent Neural Networks (notwithstanding the associations from the past layer).

Deep learning that can be characterized as neural systems with an enormous number of parameters and layers in one of four key system designs

### K-Means Clustering

K-means is one of the least complex unsupervised learning calculations that take care of the notable clustering issue.

$$J(V) = \sum_{i=1}^{c} \sum_{j=1}^{c_i} \left( \|x_i - v_j\| \right)^2$$

Where,

‘$\|x_i - v_j\|$’ is the Euclidean distance between $x_i$ and $v_j$.
‘$c_i$’ is the number of data points in $i^{th}$ cluster.
‘$c$’ is the number of cluster centers.

### Dis-Advantages of K-Means Clustering

1) The learning algorithm requires apriori specification of the number of cluster centers.
2) Unable to handle noisy data and outliers.
3) Algorithm fails for non-linear data set.

### Cross Breed Clustering Algorithm (CBCA)

In this paper, after analyzing the issues in k-means clustering algorithm the batch normalization is adopted to improve the quality of clustering.

Batch normalization is the deep learning technique which is used to tuning of weight initialization and learning parameters.

Whatever the introduction of weights, be it arbitrary or experimentally picked, they are far from the educated weights. Think about a small scale clump, during beginning ages, there will be numerous anomalies as far as required component initiations.

The deep neural system without anyone else's input is not well presented, for example, a little bother in the underlying layers prompts an enormous change in the later layers.

During back-propagation, these wonders cause a diversion to slopes, which means the angles need to repay the exceptions, before learning the loads to create required yields. This prompts the necessity of extra ages to merge.

Batch normalization regularizes this slope from diversion to exceptions and streams towards the shared objective (by normalizing them) inside the scope of the scaled small batch.

**Learning rate issue:** Generally, learning rates are kept low, to such an extent that solitary a little part of slopes revises the loads, the reason is that the angles for anomaly actuations not to influence learned enactments. By cluster standardization, these anomaly enactments are diminished and thus higher learning rates can be utilized to quicken the learning procedure.

**Input:** Values of $x$ over a mini-batch: $\mathcal{B} = \{x_{1...m}\}$;
Parameters to be learned: $\gamma, \beta$

**Output:** $\{y_i = \mathrm{BN}_{\gamma,\beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^{m} x_i \qquad \text{// mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^{m} (x_i - \mu_{\mathcal{B}})^2 \qquad \text{// mini-batch variance}$$

$$\widehat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \qquad \text{// normalize}$$

$$y_i \leftarrow \gamma \widehat{x}_i + \beta \equiv \mathrm{BN}_{\gamma,\beta}(x_i) \qquad \text{// scale and shift}$$

## IV. EVOLUTION RESULTS

We compare the performance of k-means and Cross Breed Clustering Algorithm with various datasets. For the evolution and results used one hand written digits dataset called as USPS. Fig:2 show the dataset.
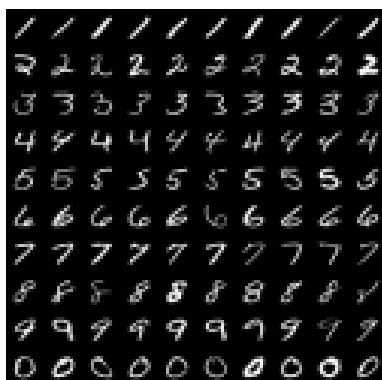


**Figure: 2 USPS dataset**

USPS: this is the handwritten digits (0-9) dataset and this contains 1100 samples in each class. 11000 16 *16 images are present in dataset.

**Table 1: Dataset Information**

| Dataset | #Classes | #Dims | #Samples |
|---------|----------|-------|----------|
| USPS | 10 | 256 | 11,000 |

The performance is calculated on the bases of accuracy is the one parameter.

$$ACC = \max_m \frac{\sum_{i=1}^{N} \mathbf{1}(r_i = m(c_i))}{N}$$

**Table2: Performance comparison on the one real dataset.**

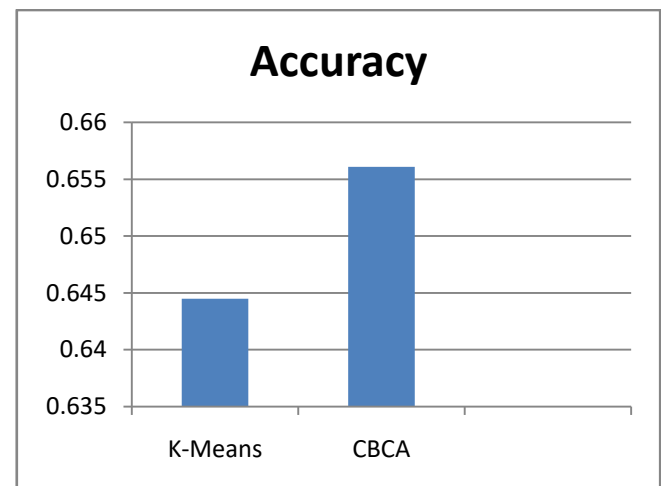| Algorithm Name | Accuracy |
|----------------|----------|
| K-Means | 0.6445 |
| CBCA | 0.6561 |



**Figure: 3 Performance of Existing and proposed system**

## V. CONCLUSION

In this paper, the CBCA is the proposed algorithm utilized the USPS hand written dataset which is used to get the quality of results. The comparison is done between k-means and CBCA and improved the accuracy of the system. The deep learning algorithm plays the major role to improve the results. The accuracy is the parameter is shown in this system.

## REFERENCES

1. Joulin, A., Bach, F., Ponce, J.: Discriminative clustering for image co-segmentation. In: Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on. pp. 1943–1950. IEEE (2010).
2. Liu, H., Shao, M., Li, S., Fu, Y.: Infinite ensemble for image clustering. In: Proceedings of ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (2016).
3. Wang, R., Shan, S., Chen, X., Gao, W.: Manifold-manifold distance with application to face recognition based on image set. In: Computer Vision and Pattern Recognition, 2008. CVPR 2008. IEEE Conference on. pp. 1–8. IEEE (2008).
4. Aggarwal, C.C., Zhai, C.: A survey of text clustering algorithms. In: Mining text data, pp. 77–128. Springer (2012).
5. Beil, F., Ester, M., Xu, X.: Frequent term-based text clustering. In: Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and data mining. pp. 436–442. ACM (2002).

6. Xu, J., Wang, P., Tian, G., Xu, B., Zhao, J., Wang, F., Hao, H.: Short text clustering via convolutional neural networks. In: Proceedings of NAACL-HLT. pp. 62–69 (2015).

7. Tian, K., Shao, M., Wang, Y., Guan, J., Zhou, S.: Boosting compound-protein interaction prediction by deep learning. Methods 110, 64–72 (2016).

8. Zhang, R., Cheng, Z., Guan, J., Zhou, S.: Exploiting topic modeling to boost metagenomic reads binning. BMC Bioinformatics 16, S2 (2015)

9. Dueck, D., Frey, B.J.: Non-metric affinity propagation for unsupervised image categorization. In: Computer Vision, 2007. ICCV 2007. IEEE 11th International Conference on. pp.1–8. IEEE (2007)

10. Ng, A.Y., Jordan, M.I., Weiss, Y., et al.: On spectral clustering: Analysis and an algorithm. Advances in neural information processing systems 2, 849–856 (2002)

11. Kanungo, T., Mount, D.M., Netanyahu, N.S., Piatko, C.D., Silverman, R., Wu, A.Y.: An efficient k-means clustering algorithm: Analysis and implementation. IEEE transactions on pattern analysis and machine intelligence 24(7), 881–892 (2002).

12. Bishop, C.M.: Pattern recognition. Machine Learning 128, 1–58 (2006)

13. Liu, H., Shao, M., Li, S., Fu, Y.: Infinite ensemble for image clustering. In: Proceedings of ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (2016)

14. T. Kohonen, The self-organizing map, Neurocomputing, vol. 21,nos. 1–3, pp. 1–6, 1998.

15. J. A. Hartigan and M. A. Wong, k-means clustering algorithm, Stat. Soc. C, Appl. Stat., vol. 28, no. 1, pp. 100–108,1979.

16. S. Wold, K. Esbensen, and P. Geladi, ''Principal component analysis,'' Chemometrics Intell. Lab. Syst., vol. 2, nos. 1–3, pp. 37–52, 1987.

17. T. Hofmann, B. Schölkopf, and A. J. Smola, ''Kernel methods in machine learning,'' Ann. Stat., vol. 36, no. 3, pp. 1171–1220, 2008.

18. [18] 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.

19. [19] J. Schmidhuber, ''Deep learning in neural networks: An overview,'' Neural Netw., vol. 61, pp. 85–117, Jan. 2015.

20. [20] de Souto MC, Costa IG, de Araujo DS, Ludermir TB, Schliep A. Clustering cancer gene expression data: a comparative study. BMC bioinformatics. 2008;9(1):497. pmid:19038021.