



Semi-Supervised Clustering for High Dimensional Data Clustering

Tejashree V Patil¹, Govind S. Pole²

Department of Computer Engineering, MES College of Engineering, Pune, India

ABSTRACT: Cluster formation has three types as supervised clustering, unsupervised clustering and semi supervised. This paper reviews traditional and state-of-the-art methods of clustering. Clustering algorithms are based on active learning, with ensemble clustering-means algorithm, data streams with flock, fuzzyclustering for shape annotations, Incremental semi supervised clustering, Weakly supervised clustering, with minimum labeled data, self organizing based on neural networks. Incremental semi-supervised clustering ensemble framework (ISSCE) which makes utilization of the advantage of the arbitrary subspace method, the limitation spread approach, the proposed incremental ensemble member choice process, and the normalized cut algorithm to perform high dimensional information clustering. The incremental ensemble member choice process is recently intended to sensibly evacuate excess gathering individuals in light of a recently proposed neighborhood cost work and a worldwide cost work, and the standardized slice calculation is received to serve as the accord work for giving more steady, hearty, and precise results.

KEYWORDS: Cluster Ensemble, Semi-Supervised Clustering, Random Subspace, Cancer Gene Expression Profile, Clustering Analysis.

I. INTRODUCTION

The bunch troupe methodologies are more points of interest and more consideration because of its valuable applications in the regions of example acknowledgment, data mining, bioinformatics, and more one. At the point when contrasted and customary single grouping calculations, bunch gathering methodologies can coordinate various grouping arrangements got from various information sources into a bound together arrangement, and give a more hearty, steady and precise last result. In any case, conventional cluster ensemble approaches have a few statutes of impediments: First they don't consider how to make utilization of earlier information given by specialists, which are spoken to by Pair savvy limitations. Match shrewd requirements are regularly characterized as the must-connect limitations and they can't interface imperatives. The must-interface limitation implies that two component vectors ought to be doled out to a similar group, while they can't connect requirements implies that two element vectors can't be appointed to a similar cluster. First most of the cluster ensemble methods cannot procure acceptable results on high dimensional datasets. Third not all the ensemble members add to the last result. So as to address the 1 and 2 restrictions, we first propose the random subspace based semi-supervised clustering ensemble framework (RSSCE), joins the irregular subspace method, the imperative proliferation approach [12], and the normalized cut algorithm [13] into the cluster ensemble framework to perform high dimensional information grouping. At that point, the incremental semi-supervised clustering ensemble framework (ISSCE) is intended to expel the copy ensemble members. At the point when contrasted and customary with traditional semi-supervised clustering algorithm, ISSCE is elements by the incremental ensemble member selection (IEMS) handle in view of an as of late proposed worldwide target work and a nearby target work, which decision ensemble individuals dynamically. The nearby target capacity is ascertained in view of an as of late planned closeness work which chooses how comparative two arrangements of properties are in the subspaces. Besides, the computational cost and the space utilization of ISSCE are dissected hypothetically. Labeled data can classify easily, but unlabeled data classification is very challenging task.

In incremental data clustering data is updated so at every time new clusters need to form for better result. It is very difficult in semisupervise to form a cluster when unnamed data is coming. All things considered, we take various nonparametric tests to think about number of semi supervised clustering ensemble approaches more than a few datasets. The test outcomes demonstrate the change of ISSCE over customary semi-supervised clustering ensemble approaches



International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Website: www.ijirccce.com

Vol. 4, Issue 12, December 2016

or traditional cluster ensemble methods on six true datasets from UCI machine learning repository [14] and 12 true datasets of tumor quality expression profiles. While there are few kinds of cluster ensemble techniques, little of them consider how to handle high dimensional information clustering, and how to make utilization of earlier learning. High dimensional datasets have too huge number of ascribes in respect to the quantity of tests, which will prompt to the over fitting issue. The greater part of the ordinary cluster ensemble methods do not consider how to handle the over fitting issue, and cannot acquire agreeable results when taking care of high dimensional information. Our strategy embraces the arbitrary subspace procedure to produce the new datasets in a low dimensional space. Incremental semisupervised clustering would be gives better results because it works on mixed type datasets. Different methods are as follows

II. LITERATURE SURVEY

semi supervised clustering approaches are:

1. Active Learning to improve semi supervised clustering

Semi supervised clustering is major tasks grouping the data objects into meaningful clusters that the same of objects within clusters is maximized and the similarity of objects to minimized clusters.

Active learning algorithm is classifying data of similar wide research, also in application, the target domain with active learning algorithm, to simplify the point label complexity. Important application of active learning in the NLP (Natural Language Processing), focus how to obtain high quality training dataset.

The two phase approaches are phase incrementally selection and expands the neighbor data nodes these two approaches increases efficiency in classification. Active learning is facilitate where the aim is to cluster group of objects by actively querying the distances between many pairs of points.

Active learning use to minimize the query to obtain a cluster.

So for this the link based algorithm is used

The linking is provided in between categorized dataset and Numerical Dataset from these both clusters final output is combine to get final result as following diagram

2. Clustering ensemble Semi supervised clustering:

Clustering ensemble recent and advanced in unsupervised learning. To combine the clustering multiple data partitions improve the accuracy of clustering. Many semi-supervised algorithms were proposed in various methodologies, some based on EM with generative mixture models, self-training, co-training, Ideally we should use a method whose EM with generative mixture models good choice if the classes produce well clustered data; features split into two sets; graph-based methods can be used with similar features and same class. But there is no way for of semi-supervised algorithm. Improve clustering accuracy for the results, supervision provided: either by using semi-supervised algorithms in the clustering ensemble and a feedback used in the function stage.

supervisor the can tune the clustering process the clustering that fits the type of the input data clustering Ensemble is to integrate clustering partitions obtained using various methods.

Clustering ensemble algorithms are usually divided into two ways. At the first different partitions of same dataset are using independent runs of different clustering algorithms. Another consensus function is used to find partition Combining supervision with clustering ensemble to give higher level of accuracy. The supervision at ensemble generation step can aid and bias different clustering to produce better and high quality base partitions.

The consensus function can also take benefit from user feedback about the base partitions to produce higher quality target partition. This approach gives user flexibility of choosing multiple types of supervision and feedback in both steps. This type of weighting scheme can be applied to other consensus functions as well.

III. CONSTRAINT PARTITIONING K-MEANS ALGORITHM

Data clustering high dimension dataset using Constraint-Partitioning K-Means (COP-KMEANS) clustering algorithm which not fit cluster high dimensional data sets in effectiveness and efficiency, because of intrinsic sparse of high



International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Website: www.ijirccce.com

Vol. 4, Issue 12, December 2016

dimensional input and resulted in producing indefinite and inaccurate clusters. So two steps for clustering high dimension dataset. First we perform dimensionality reduction on the high dimension dataset using Principal Component Analysis (PCA) as preprocessing step to data clustering. we integrate the COP-KMEANS clustering algorithm to dimension reduced to produce good and correct clusters. The experimental results very effective in producing accurate and precise clusters.

Clustering with grouping objects which are similar to each other and dissimilar to the other clusters Cluster is used to assemble that appear to fall naturally simultaneously

IV. FUZZY CLUSTERING FOR SHAPE ANNOTATIONS:

A fuzzy clustering algorithm is used group shapes into clusters. Each cluster is represented by a prototype that is manually labeled and used to unlabeled shapes that cluster. To capture the evolution of the image set over time, the previously discovered prototypes are added as pre-labeled objects to the current shape set and semi-supervised clustering is used.

Each selected object, its to a class is derived according to some similarity measures. To classify an object and to unnamed data, firstly the object has to be numerically described. Image is classify by considering its shape,color,and texture. When new image is come then previous history of classification is to be consider. Same shapes get added in single cluster. It is different a little but difficult to form cluster than text dataset. In clustering at testing phase to unlabel data if star shape image is coming then it can be classify in flower label cluster ;like this similar type shape to be consider in this type of clustering.

Unlabeled shape is classify by using nearest matching type in testing phase. clustering algorithms can group unnamed data so that similar shapes are arranged into the one cluster. When new shapes entered, we have to re-process the entire then by using training dataset entered data is tested and classify in cluster. It is physical level clustering method, it also gives effective resulted not only text data but also in shaped data we can form a cluster. Different shapes are classify from using fuzzy clustering method.

V. SYSTEMATIC APPROACH

In many machine learning, there is a large incoming of unlabeled data but limited labeled data, which sometimes hard to generate cluster .Semi-supervised learning, learning is combination of both labeled and unlabeled data,. It Overcomes limitations of the Traditional cluster There is no need of prior knowledge of the datasets given by experts. Traditional cluster ensemble methods cannot obtain satisfactory results when handling high dimensional data.Remove redundant ensemble members based on a newly proposed local cost function and a global cost function, Finally, a set of tests are compare multiple semi-supervised clustering ensemble approaches over different datasets to produce the satisfactory result

VI. USING MULTIPLE CLUSTERINGS

In supervised clustering tested data is labeled so it can easily handle. But in unsupervised learning testing data is difficult to form cluster. and to create label by test this data is very difficult. Labeling is critical and take more time, very limited number of objects get label. However, designing approaches able to work efficiently with a very limited number of labeled samples is highly challenging. Semi supervised database operation deal with both labeled and unlabeled data. Pre label and post label these two types of labeling. Pre labeled data means supervised data easily classifying in testing phase. Post labeling is when unnamed data is inserted and by using testing data get label and then classification perform in proper cluster. Challenging task is to label data properly and minimize label and to form minimum clusters. Cluster form by using nearest neighbour similar objects. Two clusters have different label,and one cluster contain similar datatypes,similar object's behaviour .



International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Website: www.ijirccce.com

Vol. 4, Issue 12, December 2016

[IV] Comparisons of different approaches:

no	Methods	ADVANTAGE	DISADVANTAGES
1	Active Learning	Grouping data in meaningful clusters	Do not achieves higher clustering results
2	Clustering Ensemble	Mixed model can use EM and Self training based clustering	Feedback clustering gives incorrect resulted and complex work on this
3	K-Means Algorithm	Dimention reduce for better cluster result	Not effective an efficient for high dimensional database
4	Fuzzy Clustering for Shape Annotation	It is an important task when managing large image collections.	It is an challenging task when managing large image collection
5	A Systematic Approach	Obtain satisfactory results. An use on mixed data type	The output quality of the process is very low Time consuming to seperate mixed type data
6	Limited labeled data	Minimize cluster formation	Updatation clustering criticallyimpliment

V. CONCLUSION

From above these contents we can conclude that there are various methods we can use to form cluster in semi supervised clustering. Each method has its own some benefits and limitations. For constant dataset all methods are ok, but for updated data incremental semi supervised clustering would be more useful, because in this the data is continuously entered in system, continuously update data, and form new clusters as per their contents, and some times changes clusters as per user demands. This data is labeled or unlabeled or in shape so incremental can work on all these type of data than other methods. So incremental semi supervised clustering is can be useded method of clustering. ble approach. Which create correct cluster on given mixed type datasets.

REFERENCES

- [1] An Improved Semi-Supervised Clustering Algorithm Based on Active Learning S.Shalini1, R.Raja International Journal of Innovative Research in Computer and Communication Engineering Vol.2, Special Issue 1, March 2014
 [2] Semi-supervised Clustering Ensemble by Ashraf Mohammed Iqbal1, AbidalrahmanMoh' d2, and Zahoor Ali Khan3Voting University, Halifax, Canada
 [3] Efficient High Dimension Data Clustering using Constraint-Partitioning K-Means Algorithm Aloysius GeorgeThe International Arab Journal of Information Technology, Vol. 10, No. 5, September 2013



ISSN(Online): 2320-9801
ISSN (Print): 2320-9798

International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Website: www.ijirccce.com

Vol. 4, Issue 12, December 2016

- [4] Incremental semi-supervised clustering in a data stream with a flock of agents Pierrick Bruneau, Fabien Picarougne, Marc Gelgon 978-1-4244-2959-2/09/\$25.00 c 2009 IEEE
- [5] Incremental Semi-Supervised Fuzzy Clustering for Shape Annotation Giovanna Castellano, Anna Maria Fanelli and Maria Alessandra Torsello
- [6] a systematic approach for analyzing the patient's future diseases using incremental semi supervised clustering r.anitha (assistant professor), m.r.ramya (pgscholar), international Journal on engineering technology and science ijets™ issn(p): 2349-3968, issn (o): 2349-3976 volume iii, issue xi, november- 2016
- [7] Efficient Active Learning Constraints for Improved Semi-Supervised Clustering Performance Ramkumar Eswaraprasad1, and Shanmugam Vengidusamy International Journal of Computer Science and Electronics Engineering (IJCSSEE) Volume 3, Issue 4 (2015) ISSN 2320-4028 (Online)
- [8] An Active Learning for Weakly Supervised Clustering Ms.A.Savithamani, Mr.M.Mohanraj International Journal of Science, Engineering and Technology Research (IJSETR), Volume 3, Issue 12, December 2014
- [9] Semi-supervised learning using multiple clusterings with limited labeled data Germain Forestiera, Cédric Wemmertb,* aMIPS, University of Strasbourg, France
- [10] An Online Semi-Supervised Clustering Algorithm Based on a Self-organizing Incremental Neural Network Youki Kamiya, Toshiaki Ishii, Shen Furao, and Osamu Hasegawa
- [11] Proceedings of International Joint Conference on Neural Networks, Orlando, Florida, USA, August 12-17, 2007
- [12] E. Akbari, H.M. Dahlan, R. Ibrahim, and H. Alizadeh, "Hierarchical cluster ensemble selection," Eng. Appl. Artif. Intell., vol. 39, pp. 146-156, 2015.
- [13] H. Zeng and Y.-M. Cheung, "Semi-supervised maximum margin clustering with pairwise constraints," IEEE Trans. Knowl. Data Eng., vol. 24, no. 5, pp. 926-939, May 2012.
- [14] S. Anand, S. Mittal, O. Tuzel, and P. Meer, "Semi-supervised kernel mean shift clustering," IEEE Trans. Pattern Anal. Mach. Intell., vol. 36, no. 6, pp. 1201-1215, Jun. 2014.
- [15] C.-L. Liu, W.-H. Hsiao, C.-H. Lee, and F.-S. Gou, "Semi-supervised linear discriminant clustering," IEEE Trans. Cybern., vol. 44, no. 7, pp. 989-1000, Jul. 2014.
- [16] L. Zheng and T. Li, "Semi-supervised hierarchical clustering," in Proc. IEEE 11th Int. Conf. Data Mining, 2011, pp. 982-991.
- [17] S. Xiong, J. Azimi, X. Z. Fern, "Active learning of constraints for semi-supervised clustering," IEEE Trans. Knowl. Data Eng., vol. 26, no. 1, pp. 43-54, Jan. 2014.
- [18] N. M. Arzeno and H. Vikalo, "Semi-supervised affinity propagation with soft instance-level constraints," IEEE Trans. Pattern Anal. Mach. Intell., vol. 37, no. 5, pp. 1041-1052, May 2015.
- [19] D. Wang, X. Gao, and X. Wang, "Semi-supervised nonnegative matrix factorization via constraint propagation," IEEE Trans. Cybern., 2015, Doi: 10.1109/TCYB.2015.2399533.
- [20] D. Wang, X. Gao, and X. Wang, "Semi-supervised nonnegative matrix factorization via constraint propagation," IEEE Trans. Cybern., 2015, Doi: 10.1109/TCYB.2015.2399533.