



RESEARCH ARTICLE

Clustering Sentence-Level Text Using a Fuzzy Back- Propagation Clustering Algorithm

¹M.Padmavathi, ²T.V.N.Sudheer

^{1,2}CSE&QISCET

¹padmavathimudigonda@gmail.com, ²tvn.sudheer99@gmail.com

Abstract— In comparison with hard clustering methods, in which a pattern belongs to a unique cluster, clustering algorithms with fuzziness allow patterns with differing degrees of membership to belong to all clusters. This is important in domains such as sentence clustering, as a sentence may belong to more than a topic present within a document or set of documents. Since most sentence similarity measures do not represent sentences in a common metric space, traditional fuzzy clustering approaches are generally not applicable to sentence clustering. This paper presents a back propagation fuzzy clustering algorithm. The algorithm uses a graph representation of the data, and operates in an Back Propagation framework in which the graph centrality of an object in the graph is interpreted as a likelihood. Results of applying the algorithm to sentence clustering tasks demonstrate that the algorithm is suitable of identifying more clusters of related sentences, and that it is therefore of potential use in a variety of text mining tasks.

Keywords— Sentence Clustering, Fuzzy clusters, Back Propagation, Page ranks, Membership values

I. INTRODUCTION

In last two years information technology developed the way for world full data. Potentially these data are not that much useful. It makes it order to useful one, we need information or knowledge underlying the data to extract in large amount. Inside the huge amount of data, Data mining is a process of extracting the valuable information. In the data discovery and data analysis can help clustering technique. Information Retrieval (IR) Process is mainly useful in clustering the sentences. The sentence level and document level has various clustering text in many differences. Document clustering divide the documents into various parts and cluster those parts depend on theme. It doesn't give that much of importance to the similar sentence in each document. In the multi document summarization there may be content overlap or bad coverage of theme. One clustering in each data element is assign in the hard clustering method. The most important unsupervised learning framework is declared as a group of data item in clusters, similar and dissimilar between them to the object belonging to other clusters. In variety of text mining applications used sentence clustering. The output of cluster was specified by the user which should be related to the query. Similar distance between the sentences is calculated by using some distance function such as Euclidean distance. In sentence clustering the recently used methods are represent sentence in the document matrix and performing clustering algorithm. The work is described in the fuzzy relationships which are used to increase the breadth and scope of problems can be applied successfully in sentence clustering. Such as document, clustering text at the sentence level accept the specific challenges not present when clustering large segment of text. The examiner some existing approaches to fuzzy clustering was highlighting some important differences between clustering at theses two levels. The data element may be belonging to more than one cluster with various degrees of membership in the Fuzzy C-Means (FCM) algorithm. Fuzzy set theory and robust statistic connections are establish for analyzing a various popular robust clustering method. The rough based FCM algorithm use arbitrary dissimilarity of data. All kind of dataset containing outlier and deal with all kind of relational data can handle fuzzy relational algorithm smoothly. The parameter of fuzzification degree greatly affect on the performance

of FCM. A suitable kernel function having a key to the success of configuration for kernel method. A single kernel that is choosing from predefined group is sufficient to represent the data. The multiple kernels are combine from the set of basis kernel have adoption for refining the results of single kernel learning. A hierarchical organization is an organizational structure is subordinate to single other entity and it represent in the form of hierarchy. In hierarchy structure consist of singular or group of power at the top level. The members of hierarchical structures are communicate with their instant superior and with their instant subordinate and it can reduce overhead communication for limiting information flow. A modern computational technology which is a method of examining and calculating and estimating the claims about human language itself is known as Natural Language Processing (NLP). Applying NLP to the data mining and text mining previously unknown information can be discovered. The text mining refers to the process of extracting high quality information from text. Document clustering (or Text clustering) is the process of automatically organizing the documents, extraction of topics and for fast information retrieval or filtering.

II. LITERATURE REVIEW

The system A. Skabar *et al* [1] used general graph centrality measure by using page rank algorithm and review of the Gaussian mixture model approach. Page Rank can be used within an Expectation-Maximization framework to construct a complete relational fuzzy clustering algorithm. The name FRECCA was given to the Page Rank centrality which can be viewed as a special case of eigenvector. The important part of their paper was novel fuzzy relational clustering algorithm.

This algorithm motivated by the mixture model approach and it also gather all the data as combination of component. To determine the model parameter they can use the Expectation- Maximization framework by applying page rank algorithm to each cluster. This framework was interpreting the page rank score of an object. The relationship between object was express in term of pair wise similarities can be applied in any domain as per result of fuzzy relational clustering algorithm.

K. Sathish kumar *et al* [2] there was a sentence level clustering algorithm used for text data as per the survey represent. The measuring sentence similarity special treatment is necessary. It was describe topic or themes which defined as the clusters in highly related sentence. It was also avoid redundancy and cover more diverse information co-clustering. In both intrinsic clustering evolution and extrinsic summarization evolution shows clear advantage in clustering algorithm. Text mining operation was used to identify outlier document in micro-level contradiction analysis techniques.

D. Wang, *et al* [3] there was proposed a new multi-document summarization framework based on sentence-level similar analysis and non-negative matrix factorization. By using semantic analysis it construct similarity matrix. For group sentence into cluster they were used similar matrix factorization. There had been shown to be equally normalized spectral clustering. It was given benefit from sentence-level similar understanding and the clustering over similar matrix by using SNMF algorithm. Multiple document summarizations have two type of summarization; extractive summarization and abstractive summarization. The term inverse sentence frequency, sentence or term place and number of keywords generally ranks the sentences in the documents. According to their scores calculated by a set of preen features in extractive summarization. Their paper was proposed a new framework based on sentence level semantic analysis (SLSS) and symmetric non-negative matrix factorization (SNMF). The relationship between sentences in a semantic manner was better capture by SLSS and SSNF can divide the similarity matrix to obtain meaningful cluster of sentence.

In system J. Saranya [4] event detection was treated as a sentence level text classification problem. There was a given comparison in between the performance of discriminative and generative approaches: namely, a Support Vector Machine (SVM) classifier versus a Language Modeling (LM) approach. The term derived from worse net was uses handcrafted lists of „trigger“ by rule-based method for investigating. The effective feature selection and proper choice of algorithm for the task at hand are requiring for good clustering of text. The handling document clustering was depending upon the different distance measures, a number of method have been proposed to handle document clustering. The Euclidean distance was typical and widely used distance measure. Euclidean distance used k means method which minimizes sum of the squared Euclidean distance between data points and their similar cluster center. It was advantageous to finding the low dimensional low dimensional presenting the documents to reduce calculation complexity.

Kamal Sarkar [5] proposed cluster which represent the sentence in multi-document text summarization depend on the factors such as clustering the sentences, cluster ordering. The uni-gram Matching-based similarity measure after a preprocessing in a similar sentence. During preprocessing stemming was not applied on input and properties such as length, sentence position, and cue phrase are not incorporated to make system effective and portable in domain and language.

S.V.Wazarkar [6] proposed Rough set clustering whose exact border line cannot be defined due to incomplete information gives another way of representing datasets. Rough sets have been conventional used and can be equally useful in clustering for classification of a sets. The crisp boundary line did not necessary in data mining.

Amit Pimpalkar[7] this system collects the number of reviews from various online websites. The given text sentences at document level was checked by all the detail of that particular product. It clusters the contents of the documents +ve, -ve or neutral. The output for any product reviews Rule based method approach was used for proper filter. Sentiment of the product was used for selecting directly and it can also accept the smiley"s of the product. To select the best product between the two it compares two products.

D. McLean[8] proposed the semantic and word order information presents method for measuring the similarity between sentences or very short text. The lexical knowledge base and corpus has given by Semantic similarity . Word order similarity measures the number of different words as well as word pairs in different order. This method was inefficient and requires human input and was not adaptable to all application domains.

III. PROPOSED WORK

We proposed a work, in the form of similarity relationships between pairs of objects are available when the data to be clustered. We analyze advantage of the capability and stability of clusters. The new back propagation fuzzy relational clustering algorithm which does not require any limitation on relational matrix is depending on the given fuzzy C-means (FCM) algorithm. This algorithm is applied in the form of text files for the clustering of the text data which is represent in the given document for the output as cluster which are grouped from text data. The similarity measure is finding out by using page rank algorithm in the algorithm. Our objective is to mine the relevant information from the closely related clusters that we have created using the back propagation fuzzy clustering relational algorithm. We will give a document as an input to search a sentence and then process further. To mine the data from a document they were taken as an input by using mining algorithm such as k-means algorithm. Find keywords, from clusters and rank them by implementing back propagation algorithm by using page rank algorithm. The back propagation fuzzy clustering is used for partitioning of the data items into collection of clusters.

The similarity measure is finding out by using page rank algorithm in the algorithm. Our objective is to mine the relevant information from the closely related clusters that we have created using the hierarchical fuzzy clustering relational algorithm. We will give a document as an input to search a sentence and then process further. To mine the data from a document they were taken as an input by using mining algorithm such as k-means algorithm. Find keywords, from clusters and rank them by implementing the algorithm using page rank algorithm. The back propagation fuzzy clustering is used for partitioning of the data items into collection of clusters. Using the Algorithm we will reduce the complexity of the system as compared to ordinary Fuzzy relational algorithm and avoid the overlapping. The page rank and Gaussian mixture model approach are used. Page rank is used for graph centrality measure and used to determine the importance of particular node within graph. The numerical score assign to every node in this algorithm it is known as page rank score. The Expectation-Maximization algorithm is used to optimize the parameter value and to formulate the cluster in Page Rank algorithm.

Expectation-Maximization is a framework which is general purposed method for learning knowledge. The maximum possibility of its parameter it is an unsupervised method is used for finding the parameters of the probability distribution that has to finding the maximum likelihood parameter of the model it is used iterative method. The E-step include the calculation of cluster membership probabilities and calculated from E-step are estimated with parameters in M-step. Producing clusters with sentences are each of them relates to some content is used fuzzy relational clustering approach. The connectivity of the association among the data element indicate output of clustering. Many existing technique have difficulties in handling extreme outlier to overcome this drawback we used hierarchical fuzzy relational clustering algorithm. The Fig. shows that we have a text document and a sentence to be searched as an input. Clustering text at the document level is well established in the Information Retrieval (IR) literature, where documents are typically represented as data points in a high dimensional vector space in which each dimension corresponds to a unique keyword, leading to a rectangular representation in which rows represent documents and columns represent attributes of those documents. In hierarchical fuzzy relational clustering algorithm we use Text mining algorithm for extracting keywords from the document. After the extracting use a clustering algorithm on the basis of extracted set of keywords and form basic clusters. In the cluster for ranking the keyword ranking algorithm is use.

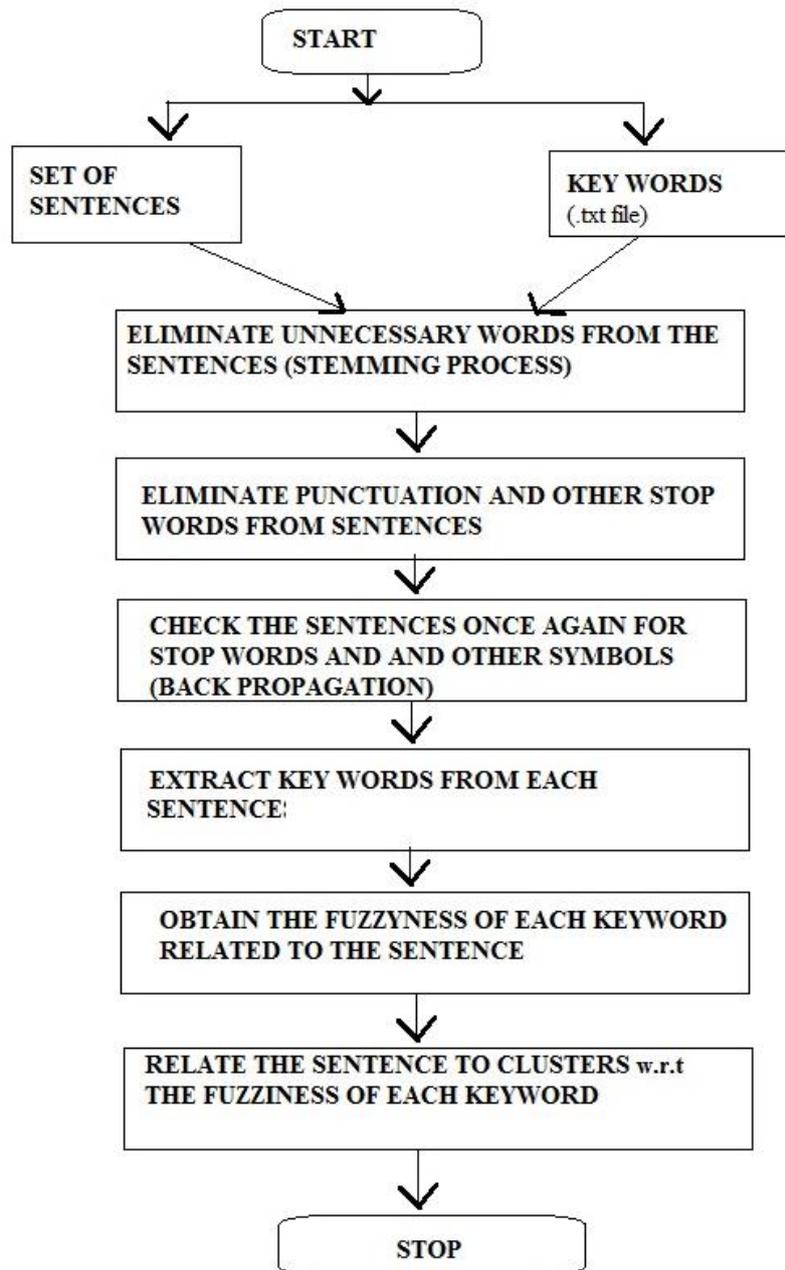


Fig., Data Flow diagram for back propagation fuzzy sentence clustering algorithm

Expectation-Maximization is a framework which is general purposed method for learning knowledge. The maximum possibility of its parameter it is an unsupervised method is used for finding the parameters of the probability distribution that has to finding the maximum likelihood parameter of the model it is used iterative method. The E-step include the calculation of cluster membership probabilities and calculated from E-step are estimated with parameters in M-step. Producing clusters with sentences are each of them relates to some content is used fuzzy relational clustering approach. The connectivity of the association among the data element indicate output of clustering. Many existing technique have difficulties in handling extreme outlier to overcome this drawback we used hierarchical fuzzy relational clustering algorithm. The Fig. shows that we have a text document and a sentence to be searched as an input. Clustering text at the document level is well established in the Information Retrieval (IR) literature, where documents are typically represented as data points in a high dimensional vector space in which each dimension corresponds to a unique keyword, leading to a rectangular representation in which rows represent documents and columns represent attributes of those documents. In hierarchical fuzzy relational clustering algorithm we use Text mining algorithm for extracting keywords from the document. After the extracting use a clustering algorithm on the basis of extracted set of keywords and form basic clusters. In the cluster for ranking the keyword ranking algorithm is use. The top ranking sentences can be extracted for the summary by using ranking sentences according to their centrality. Now apply Hierarchical Fuzzy Clustering Relation Algorithm with the ranked keywords and the input sentence and show the output for the input query.

Keyword Extraction:

The technique which is used for document renewal, Web page renewal, document clustering and review of data set, text mining and other is nothing but Keyword extraction. They can simply select which document to read and learn the relationship between documents by using extracting suitable keywords. A famous algorithm indexing is used for extracting keywords that arrive randomly in a document, but that don't arrive randomly in the remainder of the corpus. The text mining is "keyword extraction" in the form of context.

Ranking Algorithm:

This algorithm performs better than fuzzy clustering and gives the description of the application of algorithm to data set. The description of the use of Page Rank and use the Gaussian mixture model approaches are given in proposed algorithm. Graph centrality measure used in Page Rank. For determining a specific node within graph is by using Page Rank algorithm. The measure of centrality uses significance of node. This algorithm gives each every node from 0 to 1 numerical score in graph and it is also known as Page Rank Score. It's gives similar value between sentence and represent node on a graph and edges are weighted. The Expectation- Maximization algorithm to raising the parameter values and to produce the clusters is used in Page Rank. Along with the Page Rank algorithm the graph representation of data objects is used. It is a framework for learning knowledge from the insufficient data which is a common purpose method. The document is indicated by a node in the directed graph and the objects with weights represented the object equality in each sentence.

Fuzzy Relational Clustering

Unlike Gaussian mixture models, which use a likelihood function parameterized by the means and covariances of the mixture components, the proposed algorithm uses the PageRank score of an object within a cluster as a measure of its centrality to that cluster. These PageRank values are then treated as likelihoods. Since there is no parameterized likelihood function as such, the only parameters that need to be determined are the cluster membership values and mixing coefficients. The algorithm uses Expectation Maximization to optimize these parameters. We assume in the following that the similarities between objects are stored in a similarity matrix $S = \{s_{ij}\}$, where s_{ij} is the similarity between objects i and j . Initialization. We assume here that cluster membership values are initialized randomly, and normalized such that cluster membership for an object sums to unity over all clusters. Mixing coefficients are initialized such that priors for all clusters are equal.

Expectation step. The E-step calculates the PageRank value for each object in each cluster. PageRank values for each cluster are calculated, with the affinity matrix weights w_{ij} obtained by scaling the similarities by their cluster membership values; i.e.,

$$W_{ij}^m = s_{ij} \times p_i^m \times p_j^m$$

where W_{ij}^m is the weight between objects i and j in cluster m , s_{ij} is the similarity between objects i and j , and p_i^m and p_j^m are the respective membership values of objects i and j to cluster m . The intuition behind this scaling is that an object's entitlement to contribute to the centrality score of some other object depends not only on its similarity to that other object, but also on its degree of membership to the cluster. Likewise, an object's entitlement to receive a contribution depends on its membership to the cluster. Once PageRank scores have been determined, these are treated as likelihoods and used to calculate cluster membership values.

Maximization step. Since there is no parameterized likelihood function, the maximization step involves only the single step of updating the mixing coefficients based on membership values calculated in the Expectation Step.

Algorithm: To construct the cluster using back propagation fuzzy clustering algorithm.

Input: Raw Clusters for clustering and a set of sentences.

Output: Clusters to which a sentence belongs to.

Process:

1. Partitioning of sentences.
2. Removing stop words.
3. Applying back propagation.
4. Applying fuzziness to the sentences.
5. Calculating the membership values.
6. for $i = 1$ to N
7. for $m = 1$ to C
8. $p_i^m = \text{rnd}$ // random number on $[0, 1]$
9. end for
10. for $m = 1$ to C
11. $p_i^m = p_i^m / \sum_{j=1}^C p_j^m$ // normalize
12. end for
13. end for
14. for $m = 1$ to C
15. $\Pi_m = 1/C$ // equal priors
16. end for
17. repeat until convergence
18. Remove duplicate clusters.
19. Show the result.

IV. CONCLUSION

Clustering on the basis of meaningful sentences will be done and the relationship similarity values will be shown. The clustering techniques are depending upon the type of input data set and similarity measure the performance. The feature selection and its increasing good clustering of text are based on effectiveness of the algorithm.

REFERENCES

- [1] Andrew Skabar and Khalid Abdalgader “Clustering Sentence-Level Text Using a Novel Fuzzy Relational Clustering Algorithm”, IEEE Transactions on Knowledge and Data Engineering, Vol. 25, No. 1, January 2013.
- [2] K. Sathishkumar, M. Ramalingam, V. Azhaharasan, “A Thorough Investigation on the Sentence Level Clustering Approaches and its Issues in Various Applications”, International Journal of Applied Research and Studies (iJARS) ISSN: 2278-9480 Volume 2, Issue 7 July- 2013.
- [3] D. Wang, T. Li, S. Zhu, and C. Ding, “Multi-Document Summarization via Sentence-Level Semantic Analysis and Symmetric Matrix Factorization,” Proc. 31st Ann. Int’l ACM SIGIR Conf. Research and Development in Information Retrieval, pp. 307-314, 2008
- [4] Saranya .J, “Survey on Clustering Algorithms for Sentence Level Text”, International Journal of Computer Trends and Technology (IJCTT) – Volume 10 Number2 – Apr 2014.
- [5] Kamal Sarkar, ”Sentence Clustering-based Summarization of Multiple Text Documents”, TECHNIA –International Journal of Computing Science and Communication Technologies, Vol. 2, No. 1, year 2009.
- [6] Seema V. Wazarkar, Amrita A. Manjrekar, “Text Clustering Using HFRECCA and Rough K-Means Clustering Algorithm”, International Conference on Advances in Computer Engineering & Applications (ICACEA-2014) at IMSEC, GZB, Volume 15, Number 40, April 8, 2014.
- [7] Amit Pimpalkar, “Review of Online Product using Rule Based and Fuzzy Logic with Smiley’s”, International Journal of Computing and Technology, Volume 1, Issue 1, February 2014 [8] Y. Li, D. McLean, Z.A. Bandar, J.D. O’Shea, and K. Crockett , “Sentence Similarity Based on Semantic Nets and Corpus Statistics,” IEEE Trans. Knowledge and Data Eng., Vol. 8, No. 8, pp., year 2006
- [9] K.Sathishkumar, E.Balamurugan and D. Kavin, “Sentence Level Clustering Approaches and its Issues in Various Applications”, International Journal of Applied Research and Studies, 2278-9480 Volume 2 Issue 9, 2013.
- [10] E.H. Ruspini, “A New Approach to Clustering,” Information and Control, vol. 15, pp. 22-32, 1969. [11] T. Geweniger, D. Zuhlke, B. Hammer, and T. Villmann, “MedianFuzzy C-Means for Clustering Dissimilarity Data,” Neurocomputing, vol. 73, nos. 7-9, pp. 1109-1116, 2010.
- [12] G. Erkan and D.R. Radev, “LexRank: Graph-Based Lexical Centrality as Salience in Text Summarization,” J. Artificial Intelligence Research, vol. 22, pp. 457-479, 2004.
- [13] P. Corsini, F. Lazzarini, and F. Marcelloni, “A New Fuzzy Relational Clustering algorithm Based on the Fuzzy C - Means Algorithm,” Soft Computing, vol. 9, pp. 439-47, 2005.
- [14] R. Vasanth Kumar Mehta, B. Sankarasubramaniam, S. Rajalakshmi, “An algorithm for fuzzy-based sentence-level document clustering for micro-level contradiction analysis”, Proceeding ICACCI '12 Proceedings of the International Conference on advances in Computing, vol 10, No.2, Year 2012.
- [15] R.M. Aliguyev, “A New Sentence Similarity Measure and Sentence Based Applications”, An International Journal of Expert Systems with Applications, vol. 36, pp. 7764- 7772, 4 May 2009.
- [16] G.Thilagavath, ” Sentence-Similarity Based Document clustering Using Fuzzy Algorithm”, International Journal of Advance foundation and Research in Compute, Vol 1, Issue 3, March 2014.