BBIJTM

# A Multi-Hierarchical Clustering on Topic Modeling using Latent Dirichlet Allocation

**Anusha Chhabra[1], Monika Arora[2]**
[1]*Research Scholar, Delhi Technological University,*
*Delhi, India*
[2]*Department of Information Technology, Bhagwan Parshuram Institute of Technology, Delhi, India,*

[1]anusha.chhabra@gmail.com, [2]monikaarora@bpitindia.com

*Abstract- With the economic crisis, it is expected that credit and down payment growth will average later on. Credit score growth will be led by spending on the infrastructure while retail credit will display an average growth. Margin demands due to lag effect of quantity cuts between attention quantity on deposits and advances, reduced treasury gains and core fee earnings and improving in conditions for NPAs is likely to put pressure on the main point here of the financial organizations. In the light of above aspects the present document tends to analyze sectoral distribution of credit in Native India Financial industry. The document studies the styles in credit growth with a perspective to project upcoming courses of growth in bank credit.*

*Key words: Unsupervised learning, K-Means Clustering, MALLET, WEKA*

## I. INTRODUCTION

Topic modeling is the technique utilized in the area of text analysis. It is used for finding the abstract topics occurring in an assortment of documents, which provides us with strategies to arrange, comprehend, and sum up enormous assortments of text based data. It is unique in relation to rule-based text analysis approaches that utilize regular expressions or word reference based keyword-finding methods. It additionally helps in finding hidden topical patterns available across the assortment of data. Topic models have numerous applications in natural processing dialects. Many papers have been published using topic modeling approaches in different subjects like online social media stages, software engineering, Cognitive science and many more. Latent Dirichlet Allocation (LDA) [1], a generative probabilistic topic model is an amazing asset of topic modeling that permits a set of perceptions to be clarified by unobserved groups. The instinct behind LDA is that data archives exhibit numerous topics. Each record shows the topics to various extent; each word in each archive is extracted from one of the topics, where the chosen topic is taken from the per-report dissemination over topics [1], [2].

The unsupervised learning via probabilistic topic model [3] has been successfully developed for document categorization[3] [1], image analysis [4], text segmentation [5], speech recognition [6], information retrieval [7], document summarization [8], [9], and many other applications. Using topic models, latent semantic topics are learned from a bag of words to capture the salient aspects embedded in data collection. If the number of clusters is considered as a number of topics and the probabilities as the extent of cluster participation, then LDA is used as a way of soft-clustering. Cluster analysis is an errand of grouping a set of objects so that objects of the same group are more similar in any sense to one another than to those in other groups. In LDA, the order of archives doesn't make any difference except hyper-parameters which rely upon whether Dirichlet distributions are assumed to be symmetric or asymmetric. For the symmetric distribution considering on alpha value, a high alpha-value is considered means that each corpus is probably going to contain a combination of maximum number of topics, not any single one explicitly [10].

On the other hand, low alpha value puts less such constraints on corpus with the meaning where it is more likely that a corpus contains the combinations of a few or even one topic. Whereas If considering beta values, a higher beta-value represents that each and every topic is likely to have a combination of maximum number of words, not any word explicitly, while a low beta value represents that a topic contains a combination of just a few of the words [1].

Table I. Parameters of LDA and their relevance with respect to documents and topics

| Hyper-parameters | Relevance | > values | < values |
|---|---|---|---|
| Alpha | Document-topic density | Documents composed of large number of topics | Documents contain fewer number of topics |
| Beta | Topic-word density | Topics composed of a large amount of words in the document | Topics composed of few words |

With the reference of parameters of LDA as in Table I, LDA model searches for rehashing term designs in the whole Document-Term framework while K-Means grouping relies a lot upon the determination of starting point of convergence [11]. So, in contrast with LDA, K-Means belongs to one cluster for each entity. So, to perform Clustering on the document containing the topics probabilities present in the corpus we are using WEKA. The algorithms can either be applied directly to a dataset or called from your own Java code. WEKA contains devices for information pre-handling, classification, regression, grouping, affiliation rules, and representation. Presently, Applying K-means clustering on output generated from LDA for example, the text record which is showing the breakdown, by rate, of every topic inside every unique text document we imported or we can say the points that form our archives. The next section begins with a review of current studies related to LDA. In Section2, discusses the Literature review, Section 3 discusses the proposed methodology. Section 4 elucidates the experimental setups and Results and Analysis has been presented in Section5. In section 6 we conclude the paper along with suggestions for future research opportunities.

## II.    LITERATURE REVIEW

This section summarizes a few scholarly works proposed about extracting the topics using Latent Dirichlet Allocation (LDA) and various types of clustering used to cluster similar data points. Related work can also be found in the above mentioned applications in multi-domain data. Authors in a research paper [2], proposed an inference algorithm for LDA results that leads to identify documents and the relationships between the documents. [12], applied initial clustering center for K-Means as preprocessing which further makes the latent clustering center more focus on a certain topic.[10], proposed multi-grain clustering topic model which integrates topic modeling and document clustering together. [13], tried to combine LDA with document-specific word distributions for organizing large documents archives. [14], evaluated LDA with the perspective of classifying and identifying the noteworthy topics further applied to filtered collection of Twitter. [1], proposed a generative statistical model, which projects a document into topic space, and each topic contains multiple words. [15], focuses on calculating the temporal weights to reveal the importance of all the topics extracted from the probabilistic approach. [16], introduced an improved K-Means algorithm solving the limitation of traditional K-Means algorithm in terms of dependency on selecting the initial focal points by combining the largest minimum distance algorithm and the traditional K-Means algorithm. Authors in [17] conducted a comprehensive survey on already existing works. Predominantly, they concentrated on event feature learning that has ability on translating social media data into computer friendly numerical form. In [18], Different methods of document clustering and topic modeling on social media content offers a technique of categorizing, interpreting and understanding the large volume of user generated content. Authors have implemented hybrid models to cluster large document sets. Based on contextual keywords, they optimized the clustering similarity index to find the essential key document clusters. Their experimental results show that the clustering based document classification models have better performance [19].

### III.    METHODOLOGY

This section depicts the proposed methodology for finding the probabilities of topic distribution after implementing via Single-Link Type, Complete-Link Type and Average-Link Type Hierarchical clustering. The Process diagram for implementation is in Fig. 1.

A. **Single Linkage Type:** In this approach, we calculate the similarity for the closest pair. The limitation for this method is that close pair groups merge faster than the optimal, regardless of whether there is dissimilarity for those groups.

B. **Complete Linkage Type:** This approach is based on the phenomenon for calculating similarity of distant pairs along with the limitation of having outliers which causes less merging than optimal one.

C. **Average Linkage, or group linkage Type:** In this, grouping of objects are considered for calculating the similarity, rather than individual objects.
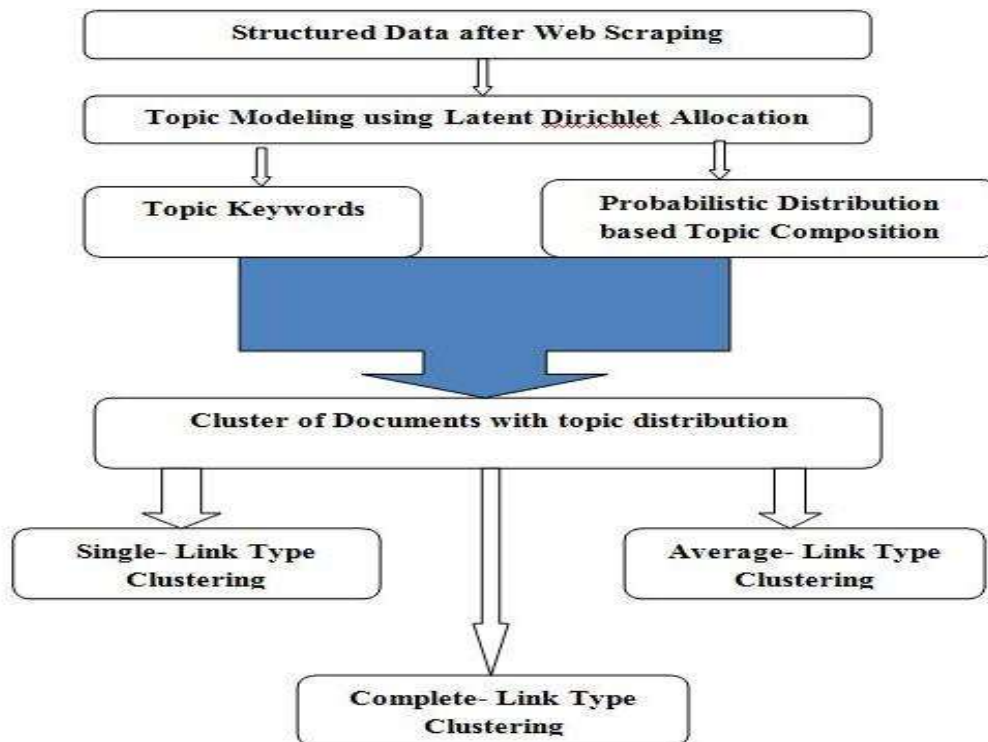


Fig. 1. Process Flow Diagram

The Process flow diagram as in Fig.1 starts with the structured data generated after web scraping. The extracted HTML code from web scraping, we stored the data in the database. The generated dataset is then normalized and LDA is applied to extract the topic keywords and probabilistic distribution-based topic composition. The documents obtained after implementing LDA, passed through soft clustering techniques which results in clustering of documents with topic distribution classified against the types of hierarchical clustering as Single-Link type, Complete-Link type and Average- Link type clustering. The experiments have been carried out on WEKA. In this paper, the dataset has been generated by using the BeautifulSoup4 Python tool of web scraping; the python script has been employed to generate the dataset. The generated dataset in the proposed method has been stored in two text files that comprises top keywords for each topic in one file and weight distribution of each topic in another file.

# IV.   RESULTS AND ANALYSIS

The output using the commands in MALLET is stored in 2 files. 1st file- keys_n1.txt is representing what all the top keywords are for each topic as in Fig.2 and the 2nd file- composition_n1.txt indicates the breakdown, by weight distribution, of every topic extracted from original textual file as in Fig.3.



Fig. 2. Top Keywords



Fig.3. Topic Composition
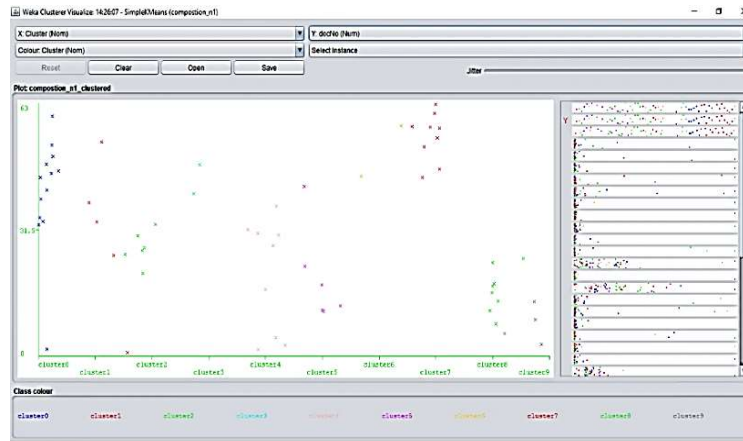
Fig. 4. K-Means Clustering



Fig. 5. Cluster of Documents with topic distribution present in the document corpus

We can also visualize the clusters that are formed. The clusters formed after applying K- means clustering on the topics composing the documents as in Fig. 5.
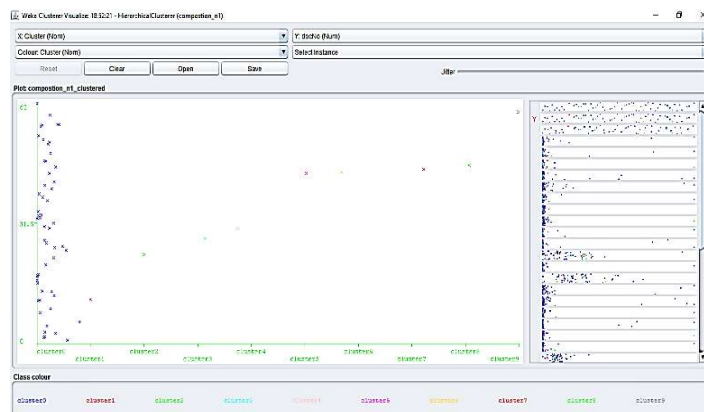


Fig. 6. Single-link type Hierarchical clustering

The results of the Single link, complete link and Average-link of the Hierarchical clustering are as in Fig 6, Fig.7 and Fig. 8.
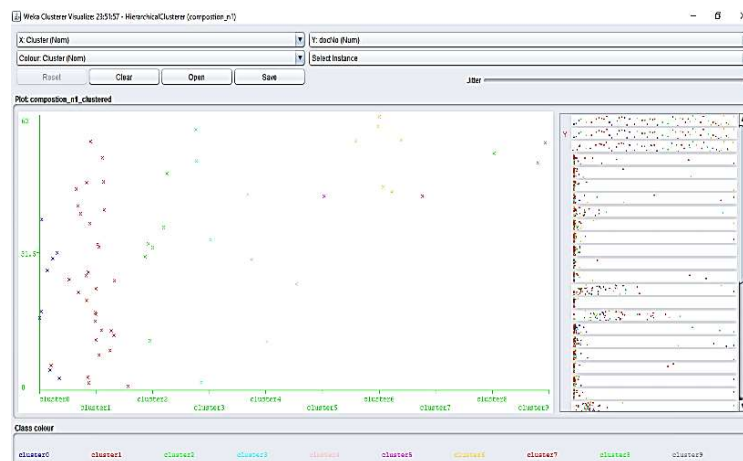


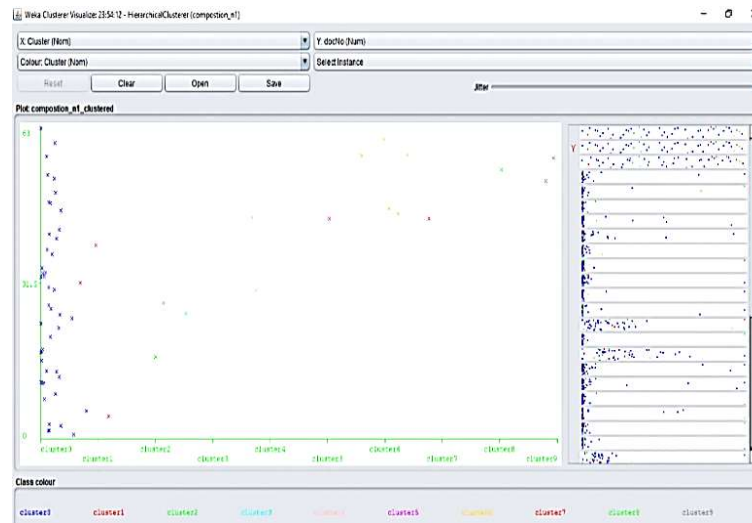Fig. 7. Complete-link type Hierarchical clustering

Fig. 8. Average-link type Hierarchical clustering

## V. CONCLUSION

From the experiments, it has been concluded that the clusters are made on the basis of the distribution of topics over which the documents are close to each other in terms of their meanings and the probability value is same for Complete-Link Type and Average-Link Type clustering for each topic in each document. Less-fold cross validation might be the reason for the same values in the probability values. In contrast to this, Single- Link type has very low probability for topics. As Topic modeling is an emerging and current field, we can increase the cross validation on dataset to have more accuracy and we can also approach different methods to improve the performance of the models for managing large archives of information. The futuristic directions for further research can Fuzzy logics and the Dynamic data-set. We can perform a hybrid approach by combining LDA and K-Means clustering on a dynamic dataset instead of static. In the Fuzzy Logic approach the simple LDA topics can be considered in a crisp set. Hence, we can fuzzify the output in results; and further the recovered topics can be in the form of a fuzzy set.

## REFERENCES

[1] D. M. Blei, A. Y. Ng, and M. I. Jordan, "Latent Dirichlet allocation," J. Mach. Learn. Res., vol. 3, no. 4–5, pp. 993–1022, 2003, doi: 10.1016/b978-0-12-411519-4.00006-9.

[2] M. Griffiths, T.L.; Steyvers, "'Finding scientific topics,'" in Proceedings of the National Academy of Sciences, 2004, pp. 5228--35.

[3] D. M. Blei, "Probabilistic Topic Models," pp. 77–84.

[4] D. Blei, L. Carin, and D. Dunson, "and applications to document and image analysis ] [ A focus on graphical model design," no. November, pp. 55–65, 2010.

[5] J. Chien, S. Member, and C. Chueh, "Topic-Based Hierarchical Segmentation," vol. 20, no. 1, pp. 55–66, 2012.

[6] J. Chien, S. Member, C. Chueh, and S. Member, "Dirichlet Class Language Models for Speech Recognition," vol. 19, no. 3, pp. 482–495, 2011.

[7] J. Chien and M. Wu, "Adaptive Bayesian Latent Semantic Analysis," vol. 16, no. 1, pp. 198–207, 2008.

[8] Y. Chang and J. Chien, "LATENT DIRICHLET LEARNING FOR DOCUMENT SUMMARIZATION," pp. 1689–1692, 2009.

[9]  N. Chiao, "HIERARCHICAL THEME AND TOPIC MODEL FOR SUMMARIZATION Jen-Tzung Chien and Ying-Lan Chang Department of Electrical and Computer Engineering," 2013.

[10] E. P. Xie, P.; Xing, "'Integrating Document Clustering and Topic Modeling,'" 2013.

[11] W. H. Li Y., "A clustering Method Based on K- Means Algorithm," in 2012 International Conference on Solid State Devices and Materials, Science Direct: Physics Procedia 25, pp. 1104–1109.

[12] Z. Guan, P.; Wang, M.; Chen, B.; Fu, "'K-means Document Clustering Based on Latent Dirichlet Allocation,'" Forty-Fifth Annu. Meet. West. Decis. Sci. Inst., 2016.

[13] M. Chemudugunta C.; Smyth P.; Steyvers, "'Modeling General and Specific Aspects of Documents with a Probabilistic Topic Model,'" Adv. neural Inf. Process. Syst., 2006.

[14] O. D.A., "'Using Latent Dirichlet Allocation for Topic Modeling in Twitter,'" in In: Proceedings of the 2015 IEEE 9th International Conference on Semantic Computing, CA, USA, 2015, pp. 978-1-4799-7935–6.

[15] L. J. Z. D. Chen H., Zhang G., "A Fuzzy Approach for measuring development of topics in patents using Latent Dirichlet Allocation," in In: 2015 IEEE International Conference on Fuzzy Systems, 2015, p. ISBN: 978-1-4673-7428-6., doi: 10.1007/s10489-014-0595-0.

[16] Z. Yang, Q. Li, Z. Lu, Y. Ma, Z. Gong, and H. Pan, "Semi-supervised Multimodal Clustering Algorithm Integrating Label Signals for Social Event Detection," Proc. - 2015 IEEE Int. Conf. Multimed. Big Data, BigMM 2015, pp. 32–39, 2015, doi: 10.1109/BigMM.2015.26.

[17] H. Zhou, H. Yin, H. Zheng, and Y. Li, "A survey on multi-modal social event detection," Knowledge-Based Syst., vol. 195, p. 105695, 2020, doi: 10.1016/j.knosys.2020.105695.

[18] S. A. Curiskis, B. Drake, T. R. Osborn, and P. J. Kennedy, "An evaluation of document clustering and topic modelling in two online social networks: Twitter and Reddit," Inf. Process. Manag., vol. 57, no. 2, 2020, doi: 10.1016/j.ipm.2019.04.002.

[19] S. A. Devi and S. S. Kumar, "A hybrid document features extraction with clustering based classification framework on large document sets," Int. J. Adv. Comput. Sci. Appl., vol. 11, no. 7, pp. 364–374, 2020, doi: 10.14569/IJACSA.2020.0110748.