

Evaluate and Improve Weight Based Pattern Detection Algorithm for Text Mining

Gagandeep Kaur¹, Hardeep Singh²

^{1,2}Department of Computer Science and Engineering, Lovely Professional University, Punjab, India

Abstract: *Semantic text mining is an abstraction of an acknowledge based on the meaning. Semantic terms are explained, phrases or words. The searching terms concern their weight is computed corresponding to their synonyms, and the term which has maximum weight is at the top. The determined technique will make use of neural technique for clustering the document present to their meaning. If various words which have similar meaning are present in document then it will cluster it in similar cluster. Increment the cluster quality neural network approach, with semantic based analyzer is used popularized.*

Keywords: data mining, text mining, text clustering, pattern taxonomy model

1. Introduction

Data mining can be expressed as an abstraction of knowledge from the large data set. This knowledge can be used for the various fields. Data mining is fact finding for information. Number of databases are studied in data mining. There are following steps in KDD: (a) Data cleaning (b) Data selection (c) Data integration (d) Data transformation (e) Data mining (f) Pattern evaluation (g) Knowledge presentation [1].

Advantages of data mining: (a) it is used in banks and marketing (b) data mining is used in retail and finance.

1.1 Text Mining

Text mining used to abstract some useful knowledge from the pages which are gather to gather. Information retrieval, abstraction, clustering, categorization are field in text mining. [1]. Techniques of Text Mining: There are two types of techniques:

- (a) **Natural processing language:** The natural processing language in which communication between human and computer. The machine learning is based on the natural processing algorithm.
- (b) **Extraction of information:** his technique is used to extract the important information from the text document [7, 20].

1.2 Text Clustering

The technique of grouping large same documents in the form of groups is called clustering. K-mean algorithm is used in text clustering (centroid-based). The advantage of clustering is that documents are classified according to their topic and subtopic, which gives the best results when searching, is done. It means the search easy [1]. Text clustering contains four main parts: (a) Text preprocessing (b) Word relativity computation (c) Word clustering (d) Text classification.

Types of Text Clustering:

- (a) Partitioning method
- (b) Hierarchical method
- (c) Grid-based method

- (d) Density-based method
- (e) Spectral method [1].

1.3 Latent Dirichlet Allocation

LDA is the generative model in natural processing language. This model is used why some part of data is same. It confess the consideration sets to be defined by intimately group. Topic mode is the example of LDA. Graphic form was presented in the first. Each document is designed as mixed of topic in LDA. The LDA model is use to give topic to each word. They are use to find the phi, error, perplexity. Alpha and Beta these two terms are used in LDA. The default value of alpha is 50/k and default value of beta is 0.1. The k value is depending upon the input dataset. The default value of LDA is 100. The information is measure by perplexity. [5].

1.4 Multinomial Naive Bayes Algorithm

Text classification and clustering in which used the algorithm. The main advantage of this algorithm, it is simplicity and computational. This algorithm is not able because performance is poor. Discriminative multinomial naive bayes will watch of likelihood as well as objectives [1, 24].

2. Literature Review

In 2012, Neelamadhab Padhy, Dr. Pragnyaban Mishra and Rasmita Panigrahi [2] in this paper discussed the number of applications used in data mining. Data mining is used for various numbers of fields. In this paper target the various approaches for searching the new areas. Data mining used in:

- universities and schools
- hospitals
- games
- domain specific.

In conclusion domain specific produce they are more correct and useful information. In 2006, Shady Shehata and Fakhri Karray [3] in this paper lized on new model that is concept-based which analysis both sentence and documents. The document analysis was based by previous model. This model

has two parts: firstly, based on analysis of term and secondly, is based on measure of similarity. F measure and entropy can be measure in cluster quality. Minimum value of entropy is beneficial for the cluster quality. More the F-measure improve the quality of results.

In 2009, Vishal Gupta and Gurpreet S. Lehal [4] text mining have developed into a great search field. The implicitly and explicitly approach to abstract in knowledge discovery in text. Number of applications is used in:

- Feature extraction
- Text base navigation
- Clustering
- Summarization
- Topic tracking
- Search and retrieval.

In this papers in which discussed the text mining applications: text mining is used in telecommunication, health care center, research center and banks.

In 2003, David M. Blei, Andrew Y. Ng and Michael I. Jordan [5] LDA is the generative model in natural processing language. This model is used why some part of data is same. The aim of this paper is to find the related topics from the given documents, then find topics of conversion and at last focus on the relationship between topics. LDA is applied they are use to find the phi, error, perplexity. Alpha and Beta is under in the hyper-meter. Alpha, Beta these two terms are used in LDA. The default value of alpha is $50/k$ and default value of beta is 0.1. The k value is depending upon the input dataset. The default value of LDA is 100. The information is measure by perplexity. The entropy is used to improve the cluster quality. The perplexity and error is calculated and take the random words then its cluster after that calculate the accuracy.

In 2009, Dr. Kanak Saxena and D.S Rajpoot [6] In this paper in which discuss the classical work of machine learning and statistics. In which occur the new problem because of perpendicular size of data. Classification rule learning is the main issue in data mining in which rules are finding that are given the partition into predefined class. Implementation of the proposed sequential patterns, improving time and space complexities of algorithms. These are the expected outcome of the proposed work.

In 2012, Shaidah Jusoh and Hejab M. Alfawareh [7] in this paper in which discussed the technique and challenging issues in text mining are used. The two main techniques for text mining:

- Natural Processing language
- Extraction of Information

The text mining used in hospitals, government organization and they are used in business is taking the right decision. There are many favorable problems faced by text mining on the one hand natural language complexity. On the other hands words can have many meanings but these meanings can be explained in different ways, this give arise certainty.

In 2012, K. Mythili and K. Yosodha [8] In this paper designed the number of mining approach for new pattern. We are used the number of patterns to make the techniques. The pattern deploying and pattern evolving are two mainly process used by this. Text mining will focus implementation for bioinformatics and it is include applying the discovered patterns for different time's series analysis.

In 2012, Ning Zhong, Yuefeng Li and Sheng- Tang Wu [9] in this paper in which discussed the number of data mining application. However we will effectively use and update discover patterns are still an open research problem. Closed sequential pattern, frequent and closed pattern are describe in this paper. Polysemy and synonymy are suffering problems in term based approach. The pattern evolving and pattern deploying are used in proposed technique.

In 2012, J. Sathya Priya and S. Priyadarshini [10] the text categorization in which use the feature selection method they are improve the efficiency and accuracy. They are deleting the duplicate information in data warehouse. In which proposed the new clustering algorithm that is text clustering with feature selection. There are 4 main modules in text clustering with feature selection:

- Text documents
- Document analysis
- Document clustering
- Feature selection

In 2012, Vandana Korde and C Namrata Mahender [11] In which discuss the text classification it means text mining in which use the text classification. In which introduce the text classification and compare the existing classifiers. The classification is providing the right class to new text document. The process of text classification is defined :

- Collection of documents
- Pre-processing
- Indexing
- Feature selection
- Classification
- Performance evaluation.

In 2012, Pratiksha Y. Pawar and S.H Gawande [12] the text classification is allowing a category through predefined ones to all text documents, consequently. This paper presents the number of approaches are discussed. There are some new approaches are proposed in the system:

- Neural text categorizer acronyms: the huge dimensionality and sparse distribution there are two problems are solved in neural text categorizer
- Improving the efficiency by self organization map's: in which approach improve the performance and reduce the dimensionality
- Soft supervised learning: the graph based SSL is new algorithm is proposed.

The multi- class problem in the text classification. The future works in which improving the performance and improving the accuracy in text classification.

In 2012, Anoop Jain, Aruna Bajpai, Manish Kumar Rohila [13] in this paper in which discussed the number of clustering methods. Methods of clustering:

- Hierarchical technique
- Optimization technique
- Density technique
- Clumping technique.

The new clustering method is implemented in this paper.

In 2012, Divya Nasa [14] the web mining and text mining is the application of the data mining. They are abstract the important information the text pages. The web mining is essentially determining the knowledge from the web data. Web mining can be separated into classifications web connection mining, web structure mining and web utilization mining.

In 2013, B. Drakshayani and E. V. Prasad [15] in this paper meaning idioms based is cluster. Processing idioms, POS tagging, the documents of pre-processing, the terms of the semantic based, the terms of calculation semantic weight, semantic grammar, similarities of document, applying clustering algorithm these are steps used under in model consisting. The clustering is used to find the beneficial possible result by hierarchical algorithm. The clustering is used by chameleon algorithm. F measure and entropy is used in measuring efficiency.

In 2013, Charushila Kadu, Praveen Bhanodia and Pritesh Jain [16] Text mining used to abstract some useful knowledge from the pages which are gather to gather. Reduced the dimensional dataset in proposed work of hybrid approach. Reduce the dimension with the feature based analysis. The term frequency is describing the system. The clustering is which used the weight of documents. The pattern discovery result will be improved in this model. There are number of phases are described in proposed system. Phases are:

- data preprocessing
- semantic based analysis
- similarity based analysis
- Pattern evolving and pattern mining.

They are describing the useful pattern to develop the new approach using for future work.

In 2013, Dipti S. Charjan and Mukesh A. Pund [17] in this paper we target on large data collection is discovering the patterns by deploying efficient algorithm. Extracting and non trivial is used to refer in text mining. In which implement the number of approaches:

- Pattern Taxonomy model
- Inner pattern
- Stemmer.

Better mean exploration will be investigated by us in the long pattern. In the conclusion patterns of repetitions have been focused by us.

In 2013, Anisha Radhakrishnan [18] in this paper there are

many mining application used for pattern developing. Polysemy and synonymy are main problem facing in term based approach. Recently research gives us pattern evolving and pattern deploying is useful patterns.

In Salton G. and McGill M.J. [19] collection of objects is considered by us for information retrieval. One or more properties are characterizing each object associated. In literature to used measure vector similarities is widely by dice and jaccard. A large mass of data is basing the most of the information retrieval. Such data are difficult by the manipulation. Classification is grouping of similar items into common classes. The classification methods are mainly used for two purposes:

- To classify the set of index keywords
- To classify the documents into subject classes.

The clustering process improves the search process.

In 2013, Inje Bhushan. V and Ujwalapatil [20] text mining is used for the research areas. In this papers in which discussed the information extraction, pattern taxonomy model. Effective discovery approach has been designed to overcome the problem; the problem is low frequency and misinterpretation. This paper is solving to cover all challenging in data mining. Pattern method is used to identify the pattern. The successful techniques are discussed in this research paper. We will make strong application for the solving problems.

In 201X, Fanghuai Hu, Zhiqing Shao and Tong Ruan [21] Synonym and self supervised learning are discussed in this paper. The synonym means word in the same language. The synonym is used in which NLP. There are two types of tool are used in NLP:

- Lightweight
- Stand one.

Performance is perfectly excellent and low time complexity by using lightweight tool.

In 2013, Nihal M. Abdel Hamid, M.B Abdel Halim, M. Waleed Fakhr [22] the document clustering is main part of data mining. High dimensionality, scalability and accuracy is provide in the clustering. The BA is population based algorithm. Local and global search both perform in Bees algorithm. There are three methods are used in document clustering:

- Graph-based method
- Hierarchical method
- Partitioning method.

In which discuss the main two optimization technique: first one is Evolutionary algorithm (EA) and second one is Swarm-based optimization algorithm. The optimization issue is solving with the EA technique. The optimal solution is finding the SOA. In future work we are doing main concentrate to Bees algorithm of hybridizing and comparison to other hybridizing model.

In 2013, R. Jensi and Dr. G. Wiselin Jiji [23] the document

text clustering is used in the research center. They store the large amount of information. The document clustering is developing the number of techniques. Localized and globally search in the text document cluster. In which introduce the soft computing technique. The soft computing is represented in tree structure.

The soft computing is divided into three parts:

- Fuzzy logic
- Neural networks
- Evolutionary computation.

In future research work will make more improve the document cluster quality as compare to present research work.

In 2013, Kavitha Murugesan and Neeraj RK [24] Text mining is the part of the data mining. It is the technique of the data mining. Number of methods under in text mining. For searching the new patterns here we use the naive bayes algorithm. The prescriptive technique make the outcome organize in particularized order. Some methods are discussed:

- Pattern taxonomy model
- Classification
- Naive bayes
- Inner pattern evaluation.

In 2013, Nikunj Kansara and Shailendra Mishra [25] Fuzzy clustering and steps of preprocessing are introduce in this paper. The fuzzy clustering is mechanisms which classify the items. The fuzzy clustering is following some task:

- Resource finding
- Information selection and preprocessing
- Generalization
- Analysis.

The complexity is reducing for processing time. Easy to implement the proposed algorithm.

In 2013, Tatiane M. Nogueira, Heloisa A Camargo and Solange O. Rezende [26] The problem is solving in text mining and information retrieval by using the clustering. Text mining used to abstract some useful knowledge from the pages which are gather to gather. In fuzzy DDE method is present in this paper. There are main two problems are occurring in this method:

- First one is how to consider the fault in document cluster and how to represent query (confusion)
- Second one is how to abstract cluster description from this kindly of information.

The centroid based is not good as compare to the proposed approach. The performance is good in the proposed system.

In 2013, Aastha Joshi and Rajneet Kaur [27] Clustering is a mechanism to bring corresponding data into cluster. The unsupervised learning is part of the clustering. There are five types of techniques are defined:

- K- mean method
- Hierarchical method
- DBSCAN method

- OPTICS method
- STING method.

In 2014, V. Aswini and S. K. Lavanya [28] text mining is used for finding the effective knowledge. They are used to searching the correct information. in this paper in which discuss the pattern taxonomy model. There are two stages in pattern taxonomy model. Firstly how to abstract the pattern in text page and secondly how to improve the effectiveness. Closed and frequent pattern are used in pattern taxonomy model. They are used to abstract and update the discovered pattern. Pattern approach is used to searching new pattern. There are main issues under this pattern based approach:

- Low frequency
- Misinterpretation.

The effective pattern discovery is defined in this paper. These techniques:

- Pattern evolving
- Pattern deploying.

They are used to improve the effectiveness. These are used for searching the interesting knowledge. In proposed approach, the pattern taxonomy model is used to abstract the pattern. We can search the problem in this paper. The pattern discovery for text mining is selected base paper. We will be working in calculate the weight in the shortest time.

In 2014, Sudesh Kumar and Nancy [29] Data mining can be expressed as an abstraction of knowledge from the large data set. This knowledge can be used for the various fields. The clusters are generating in less time by using K- mean algorithm. The information searching is very easy with the help of data mining. In which introduce the concept of preprocessing and min-max normalization. The clustering concept in which used the z- score normalization. C#.net is used to implement the proposed algorithm but the WEKA tool is used in the existing algorithm. Solve the calculation issue and improved the time by using C#.net.

3. Research Methodology

Semantic text mining is an abstraction of an acknowledge based on the meaning. Semantic terms are explained, phrases or words. The dataset taken the random words, these words will cluster the words in alphabetical order. In which take the value in form of mat file. It is converted t the ASCII format then process is continue. LDA is applied they are use to find the phi, error, perplexity. Alpha, Beta these two terms are used in LDA. The default value of alpha is 50/k and default value of beta is 0.1. The k value is depending upon the input dataset. The default value of LDA is 100. The information is measure by perplexity. The entropy is used to improve the cluster quality. The perplexity and error is calculated and take the random words then its cluster after that calculate the accuracy. The technique of cluster quality is not good because they are not in sorted alphabetical order. The result is come in unsorted words. It display in randomly. In which accuracy is less and consume the more time. The Gaussian method is used in which increase the accuracy and reduce the

time as compare to existing technique. The words are display in sorted order.

the code and conclusion will be shown in the command window.

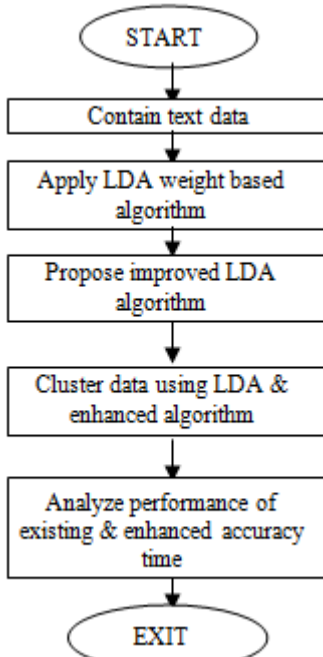


Figure 1: Proposed Technique Work Flow

4. Experimental Results

We will use the MATLAB tool to implement the flowchart. MATLAB tool is used for the area of research. It is provide the benefit for mathematical equation. It means mathematics equation is easily solved in this tool. MATLAB is the simplest programming language in which solves the linear equation. Mathematical programs are written in this tool. It has number of tool boxes that are the beneficial for optimization. When we open the matlab then show the graphical user interface.

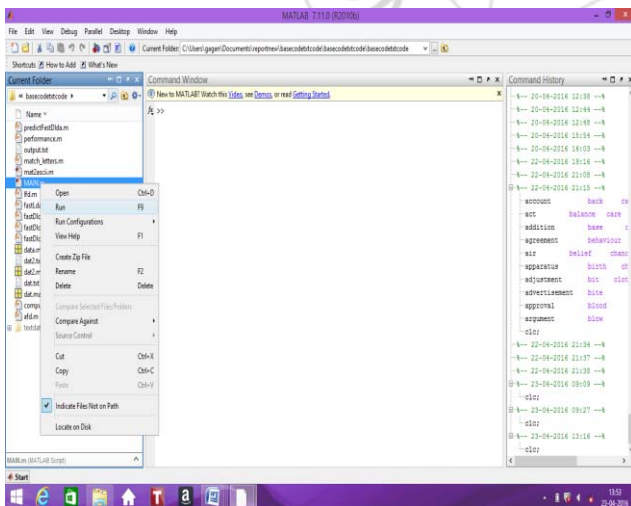


Figure 2: Loading Code

The code is run in MATLAB, firstly right click on the folder which contain the code, then open the new pop up window in which various option to run the code and display the result. The code is executed in other alternative way. Double click on the code file a new editor window will be open. The toolbar is top to the editor window there will be option to run

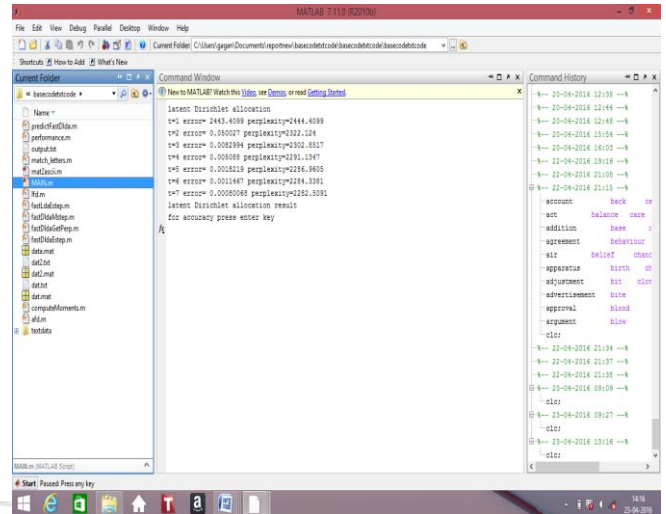


Figure 3: LDA Function Applied

The base code of text clustering is implemented with the LDA function which will cluster the words in dataset. Perplexity and error is calculated in which LDA function. The iterations are display one by one but they consume more time to display the result.

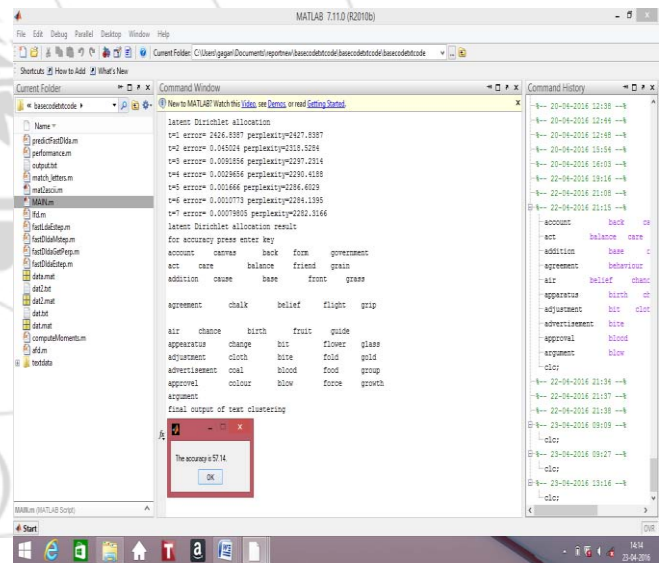


Figure 4: Shows Accuracy and Outcome Cluster

In this figure all iterations are done. After that we have press the enter and display the result. The accuracy is shown which is 57.14 and display the cluster of words. The words are not display in serial order they display in unsorted order.

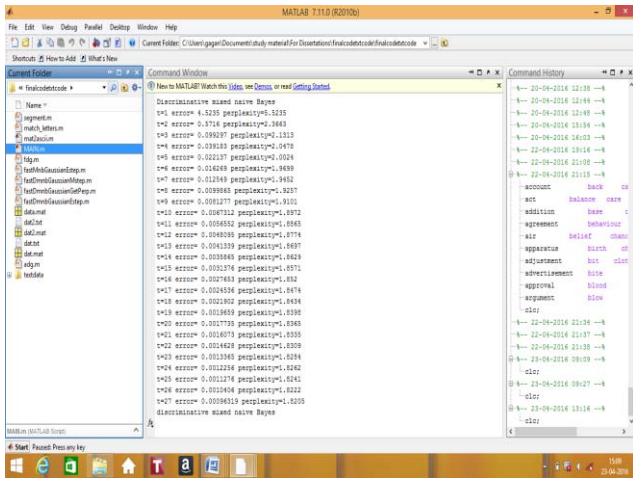


Figure 5: Results of DMNB Algorithm

DMNB method is start in which calculate the error, perplexity and accuracy of result. All the iterations are done then DMNB algorithm is complete. After that press the enter and display the result. In this method in which calculate the accuracy and cluster the words.

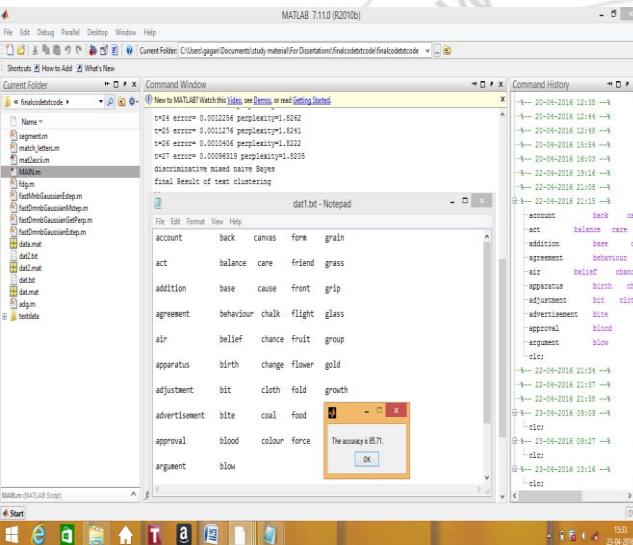


Figure 6: Clustering A-Z series and Increase Accuracy

After pressing the enter the accuracy is display. Accuracy will be 85.71 and clustering the words in the serial order. The words are display in the sorted order.

5. Conclusion

The proposed work will link text mining with natural language processing and minimum the gap between them. The new proposed hypothesis of semantic model will improve the cluster quality. The cluster which is formed after applying proposed technique will be in alphabetical order either in A-Z series. The introduced method will improve the efficiency and accuracy of the existing algorithm. The proposed technique is better than existing technique.

The existing work can be extended by linking it to web documents. This technique works with the words only. You can improve this work by implementing on document which

contains phrases or sentences. The semantic algorithm will improve the efficiency. The improvement can be done by working with the concept of hypernyms or working on both synonyms and hypernyms. Further to increase the cluster quality we can use hybrid clustering.

6. Acknowledgment

I am using this opportunity to express my gratitude to everyone who supported me in the research work .First I offer my sincerest gratitude to my supervisor, Hardeep Singh, who has supported me throughout my Dissertation. Without him this Dissertation would not have been completed or written. I am thankful for his aspiring guidance, invaluable constructive criticism and friendly advice during the research work. I am sincerely grateful to him for sharing their truthful and illuminating views on a number of issues related to the research. Finally, I thanks to my parents for supporting me throughout all my studies at university.

References

Books

- [1] Data Mining: concepts and Techniques, Second Edition Jiawei Han, Jian Pei and Micheline Kamber.

References

- [2] Neelamadhab Padhy, Dr. Pragnyaban Mishra and Rasmita Panigrahi, "The Survey of Data Mining Applications And Feature Scope", International Journal Of Computer Science Engineering And Information Technology (IJCSIEIT), Vol.2, No.3, June 2012.
- [3] Shehata Shady and Fakhri Karray, "Enhancing Text Clustering using Concept-based Mining Model", Proceedings of the Sixth International Conference on Data Mining (ICDM'06) 0-7695-2701-9/06,2006.
- [4] Vishal Gupta and Gurpreet S. Lehal, "A Survey of Text Mining Techniques and Applications", Journal of Emerging Technologies in Web Intelligence, Vol.1, No.1, August 2009.
- [5] David M. Blei, Andrew Y. Ng and Michel I. Jordan, "Latent Dirichlet Allocation", Journal of Machine Learning Research 3,993-1022,2003.
- [6] Dr. Kanak Saxena and D.S Rajpoot, "A Way to Understand Various Patterns of Data Techniques for Selected Domains", International Journal of Computer Science and Information Security, Vol.6, No.1, 2009.
- [7] Shaidah Jusoh and Hejab M. Alfawareh, "Techniques, Applications and Challenging Issue in Text Mining", International Journal of Computer Science Issues, Vol.9, Issue 6, No.2, November 2012.
- [8] K. Mythili and K. Yosodha, Research Scholar, "A Pattern Taxonomy Model with New Pattern Discovery Model for Text Mining", International Journal of Science and Applied Information Technology, Vol. 1, No.3, August 2012.
- [9] Ning Zhong, Yuefeng Li and Sheng- Tang Wu, "Effective Pattern Discovery for Text Mining", IEEE Transactions on Knowledge and Data Engineering, vol.24, No.1, January 2012.

- [10] J. Sathya Priya and S. Priyadarshini, "Clustering Technique in Data Mining for Text Documents", International Journal of Computer Science and Information Technologies, Vol.3(1),2943-2947, 2012.
- [11] Vandana Korde and C Namrata Mahender, "Text Classification and Classifiers: A survey", International Journal of Artificial Intelligence and Applications, Vol.3, No.2, March 2012.
- [12] Pratiksha Y. Pawar and S.H Gawande, "A Comparative Study on different Types of Approach to Text Categorization", International Journal of Machine Learning and Computing, Vol.2, No.4, August 2012.
- [13] Anoop Jain, Aruna Bajpai, Manish Kumar Rohila, "Efficient Clustering Technique for Information Retrieval in Data Mining", International Journal of Emerging Technology and Advanced Engineering, Vol.2, Issue 6, 2250-2459, June 2012.
- [14] Divya Nasa, "Text Mining Techniques- A Survey", International Journal of Advanced Research Computer and Software Engineering, Vol.2, Issue 4, April 2012.
- [15] B. Drakshayani and E.V.Prasad, "Semantic Based Model for Text Document Clustering with Idioms", International Journal of Data Engineering (IJDE), Vol.4, Issue 1, 2013.
- [16] Charushila Kadu, Praveen Bhanodia and Pritesh Jain, "Hybrid Approach to Improve Pattern Discovery in Text Mining", International Journal of Advanced Research in Computer and Communication Engineering, Vol.2, Issue 6, June 2013.
- [17] Dipti S. Charjan and Mukesh A. Pund, "Pattern Discovery for Text Mining using Pattern Taxonomy", International Journal of Engineering Trends and Technology (IJETT), Vol.4, Issue 10, October 2013.
- [18] Anisha Radhakrishnan, "Efficient Updating of Discovered Patterns for Text Mining: A Survey", International Journal of Computer Science and Network Security, Vol.13, No.10, October 2013.
- [19] Salton G. and McGill M. J., Introduction to Modern Information Retrieval.
- [20] Inje Bhushan.V and Ujwalapatil, "A Comparative Study on Different Types of Effective in Text Mining: A Survey", International Journal of Computer Engineering and Technology (IJCET), March-April 2013.
- [21] Fanghuai Hu, Zhiqing Shao and Tong Ruan, "Self-Supervised Synonym Extraction from Web", Journal of Information Science and Engineering, XX, XXX-XXX, 201X.
- [22] Nihal M. Abdel Hamid, M.B Abdel Halim, M. Waleed Fakhr, "Bees Algorithm- Based Document Clustering", 6th, International Conference on Information Technology, 2013.
- [23] R. Jensi and Dr. G. Wiselin Jiji, "A Survey on Optimization Approaches to Text Document Clustering", International Journal on Computational Sciences and Applications, Vol.3, No.6, December 2013.
- [24] Kavitha Murugesan and Neeraj RK, "Discovering Patterns to Produce Effective Output Through Text Mining Using Naive Bayes Algorithm", International Journal of Innovative Technology and Exploring Engineering, Vol.2, Issue 6, 2278-3075, May 2013.
- [25] Nikunj Kansara and Shailendra Mishra, "An Improved Fuzzy Clustering Technique for User's Browsing Behaviors", International Journal of Emerging Trends and Technology in Computer Science, Vol.2, Issue 2, March-April 2013.
- [26] Tatiane M. Nogueira, Heloisa A Camargo and Solange O. Rezende, "Fuzzy- DDE: A Fuzzy Method for the Extraction of Document Cluster Descriptors", International Journal of Computer Information Systems and Industrial Management Applications, Vol. 5, 2150-7988, 2013.
- [27] Aastha Joshi and Rajneet Kaur, "A Review: Comparative Study of Various Clustering Techniques in Data Mining", International Journal of Advanced Research in Computer Science and Software Engineering, Vol.3, Issue 3, March 2013.
- [28] V. Aswini and S. K. Lavanya, "Pattern Discovery for Text Mining", International Conference on Computation of Power, Energy, Information and Communication, 2014.
- [29] Sudesh Kumar and Nancy, "Efficient K- Mean Clustering Algorithm for Large Datasets Using Data Mining Standard Score Normalization", International Journal on Recent and Innovation Trends in Computing and Communication, Vol.2, Issue 10, 2321-8169, October 2014.

Author Profile

Gagandeep Kaur received the B. Tech degree in Computer Science Engineering from Sant Baba Bhag Singh Institute of Engineering and Technology in 2013 and M.Tech degrees in Computer Science Engineering from Lovely Professional University (Jalandhar) in 2014 to 2016, respectively.