Artificial Intelligence formulated this projection for compatibility purposes from the original article published at Global Journals. However, this technology is currently in beta. *Therefore, kindly ignore odd layouts, missed formulae, text, tables, or figures.*

1 2	Intuitionistic Partition based Conceptual Granulation Topic-Term Modeling
3	D. Malathi ¹ and S. Valarmathy ²
4	1 Bannari Amman Institute of Technology, Anna University, India
5	Received: 15 December 2013 Accepted: 4 January 2014 Published: 15 January 2014

7 Abstract

⁸ Document Analysis represented in vector space model is often used in information retrieval,

⁹ topic analysis, and automatic classification. However, it hardly deals with fuzzy information

¹⁰ and decision-making problems. To account this, Intuitionistic partition based cosine similarity

¹¹ measure between topic/terms and correlation between document/topic are proposed for

¹² evaluation. Conceptual granulation is emphasized in the decision matrix expressed

¹³ conventionally as tf-idf. A local clustering of topic-terms and document-topics results in

¹⁴ comparing dependent terms with membership degree using cosine similarity measure and

¹⁵ correlation. A preprocessing of documents with intuitionistic fuzzy sets results in efficient

¹⁶ classification of large corpus. But it depends on the datasets chosen. The proposed method

¹⁷ effectively works well with large sized categorized corpus.

18

19 Index terms— document analysis, intuitionistic fuzzy, topic modeling.

20 1 Introduction

ocument model in the information retrieval has three main components, namely Text Preprocessor, Topic Extractor and Corpus categoryzation. These components are integrated to deploy knowledge extraction in information system. In spite of this, the growing data and its knowledge recognition complications have considerably encouraging the extensions of machine learning algorithms.

²⁵ 2 a) Document Model

The text document Modeling is observed as latent topics model. Various prominent approaches in machine learning are used to study the model. Document model is a mixture of topics [4]. Topics are inferred by the collection of correlated words. But unsupervised learning perspective is the pulse of bubbling out the topics. By modeling, varieties of mining range can be established with various subjects. The models try to observe the likely documents and tend to focus on topics. But document models are discriminant because of random words due to

³¹ linguistic factors such as synonym, hyponym, Polysemy, etc.

32 **3 b)** Text Pre-processor

The functionalities essential for machine learning of document are document pre-processing and corpus 33 34 representation. Stop words removal, word stemming, filtering to exclude certain words, are done within each document. This process is called preprocessing of documents. Obtained vocabulary is put up in the word-35 document matrix which is generally called as bag-of-words model. The document representations may be in 36 binary (0, for nonoccurrence and 1 for occurrence of each term in a document), term frequency (tij -number of 37 occurrence of ith word in jth document) and term frequency inverse document frequency (probable occurrence 38 of tij' -distribution of ith word in jth document). Obtained data in this stage is huge in dimension, and lot of 39 techniques [15] have been proposed for dimension reduction. 40

9 PROPOSED MODEL -INTUITIONISTIC PARTITION BASED CONCEPT GRANULATION (IPCG)

41 **4 c)** Topic Extractor

42 A topic model is a probabilistic model that can be considered as a mixture of topics, represented by probability

distributions of words in a document. The latent variables or topics are the inferring components of this model.
The main objective is to learn from documents the distribution of the underlying topics in a given corpus. Topic

 $_{45}$ model is Text corpora representation by a co-occurrence matrix of words and documents. The probabilistic latent

⁴⁶ semantic analysis (PLSA) model [10] uses probability of words with given topics and probability of topics in a ⁴⁷ document, to build a topic model. The Latent Dirichlet Allocation (LDA) model [1], is another probabilistic

approach which ties the parameters of all documents through hierarchical generative model.

⁴⁹ 5 d) Corpus Categorization

Text Categorization is a classical application of Text Mining [19], and is used in email filters, social tagging and automatic labeling of documents in business libraries. Text mining applications in research and business intelligence include, latent semantic analysis techniques in bioinformatics automatic investigation of jurisdictions plagiarism detection in universities and publishing houses, cross-language information retrieval, spam filters learning, help desk inquiries, measuring customer preferences by analyzing qualitative interviews, automatic grading, fraud detection or parsing social network for ideas of new products [9].

56 6 II.

57 7 Literature Support

The theory of fuzzy set is Consider as a degree of membership assigned to each element, where the degree of non-membership is just automatically equal to D its complement. However, human interpretation often does not express the corresponding degree of nonmembership as the complement to 1. So, Atanassov [1][2] [3] introduced the concept of intuitionistic fuzzy set that is meant to reflect the fact that the degree of nonmembership is not always equal to 1 minus degree of membership, but there may be some hesitation degree.

Intuitionistic fuzzy set is a generalized constructive logic applied in fuzzy set. It is defined on a X of objects, with each object x is described by the degrees of membership and non-membership to a certain property, () () $\{ \} X x x x X A A ?, , ? \mu (1) () () 1 0 ? + ? x x A A ? \mu X x ? ? (2)$

Therefore the degree of non determinacy of the object x with respect to the intuitionistic fuzzy set A is imposed as, () () () x x x A A A ? μ ? + = X x ? (3)

The model is well suited to represent a classification problem with high dimension. The confusion matrix of high dimension can be probably reduced to concept matrix of low dimension. The similarity measures [14] and distance measures [21] [20] between two intuitionistic fuzzy sets can be applied in pattern recognition.

71 In this paper, a Partition based approach [16] inspired by Hierarchical segmentation [8] and topic based segmentation [6] are extended using Intuitionistic fuzzy set approach [23] for local centralization of conceptual 72 words. The intuitionistic fuzzy set theory is applied in conceptual term/topic detection. A cosine similarity 73 and correlation are taken into for defining membership degree and the non-membership degree respectively. The 74 results using this measure found better with respect to the dataset chosen. In literature a intuitionistic fuzzy 75 representation of images for clustering [18] [12] by utilizing a novel similarity metric are defined. But a minimal 76 support is extended for text classification. So, a local centralization of conceptual terms using Intuitionistic 77 logical clustering has been applied in the work. 78

79 **8 III.**

9 Proposed Model -Intuitionistic Partition based Concept ⁸¹ Granulation (IPCG)

Intuitionistic logic is a natural deduction system [13],that have introduction rules μ and elimination rules ? for the logical connectives and quantifiers. The { }) (), (, ij i ij i ij w w w A ? μ = where 1 0 < < ij w (4)

The similarity between words and on a topic is calculated by the cosine measure. Each document vector is normalized with the weight and length of terms in k partitions. Then the optimal term ij w [16] should The intuitionistic angular or cosine similarity [22] measure between the m terms in a partitioned set is given as follows:()?? = = = = m i i B m i i A m i i B i A x x x x B A C 1 2 1 2 1)()()()()()() (, $\mu \mu \mu \mu \mu (6)$

The intuitionistic correlation [7] of rows all fuzzy numbers are included from the samples of tf-idf (Partition Model). The crisp set is modified intuitionistically with the sample mean and variance of membership function as:()()()()()()??? = = = ??????????? = n i B i B n i A i A n i B i B A i A I x x x x B A CR 1 1 1)()()()()()(, $\mu \mu \mu \mu \mu \mu \mu \mu \mu (7)$

The effectiveness of the intuitionistic classification of corpus is approximately studied and analyzed using the following entropy [22] specifically used for Intuitionist Fuzzy Set 'A'.()()()()()()()()) () () () A i A i A i A i A i A x x v x x v x n E 1, max, min 1 ? μ ? μ (8)

95 IV.

Datasets a) Newspaper Article collection 10 96

The newspaper articles under different topics are collected. The categories are marked. The training and testing 97 documents are randomly chosen. The growing social media made essential to include newspaper article collection 98 to include in this work. News are generally categorized by topic area ("politics," "business," etc.) written in 99 clear, correct, "objective," and somewhat schematized language [5]. This would pave way to extend the research 100 towards social networking and marketing. The collection includes about 780 documents with 25 categories. All 101 new social relevant topics ("mobile", "opinion", etc.) are included for categorizing. 102

b) Reuters-21578 Data Set 11 103

The Reuters-21578 Data Set collection provides a classification task with challenging properties. There are 104 multiple categories, the categories are overlapping and non exhaustive, and there are relationships among the 105 categories. There are interesting possibilities for the use of domain knowledge. There are many possible feature 106 sets that can be extracted from the text, and most plausible feature/example matrices are large and sparse [11]. 107

c) Movie Review Dataset 12108

The Movie Review Dataset, Polarity dataset v0.9 with 900 positive and 900 negative reviews is used. Using movie 109 reviews as data, the problem of classifying documents using standard machine learning techniques definitively 110 111 outperform human-produced baselines processed reviews [17]. The training cases are chosen randomly from each 112 class about 100 documents. Which means about 500 cases are considered for training. V.

113

13 **Results and Analysis** 114

The machine learning classification methods, such as Bayesian, Naïve Bayes, J48, Support Vector Machines, LMT 115 are strong enough to support classifications. 116

In the case of concept granulation in document classification, the feature selection is fine tuned to achieve 117 categories strictly connected to the human perception. Before imposing the features into the classifier, some 118 form of selection must be chosen. The proposed method, selects the features according to the intuitionist logic. 119 The features tf-idf matrix has been The proposed Concept Granulation Using Intuitionistic Partition Based 120 Classification Model is implemented administered in the Java based system and analyzed for its significance. The 121 intuitionistic correlation is applied to the specified datasets. In which the chosen dataset and the partitions play 122

the very important role in finding the result of the model. The tfidf-IP is favorable for Reuter dataset than for 123 Newspaper and Movies. This is represented in the Figures 2(a The perplexity is depicted in Figure 3 and Table1. 124 So the analysis can be interpreted or inferred in the following ways: 125

Intuitionistic approach is in favor of the classified documents or corpus chosen Partition plays the important 126 role in the proposed model. Out of four types of partition, k=8 plays a smoothened strong support for the 127 proposed model k=16, the highest partition yield only a very moderate result and more confusions. 128

129 k=4, the least partition model yield the smooth but less significant support for all the datasets. k=8, yield the partially smooth but supportive significant for the movie dataset. (Than other partitions) 130

The results are focused to average training datasets and micro f-measure (Table 2) to show up the IPCG 131 performs better with dimension reduction for categorization of corpus. Every datasets chosen for analysis behaves 132 to the pull and push of various stages of the proposed model. 133

Conclusions 14 134

In this paper, we have proposed a intuitionistic partition based concept granulation topic-term model for a 135 nominal tf-idf vector space model which is often used in information retrieval, topic analysis, and automatic 136 classification. The cosine distance and correlation treatment to the tf-idf reduces the dimension and improves 137 the efficiency of bag of words/terms in topics. However, it is priory treated using the intuitionistic partition for 138 fitting the model into decision-making problems. To account this, Intuitionistic partition based cosine similarity 139 measure between topic/terms and correlation between document/topic are included. The proposed fuzzy model 140 is tailored with normal combinational approach to fetch intuitionistic fuzzy crisp set. Yet, it is observed the 141 model is well behaving and promising for the categorized documents and not so bad support for the low inference 142 corpus collections like movie review. So, this make us clear that the social media documents should be specially 143 treated before introducing this model. It is felt that aggregation of social media topic-terms is needed. This is 144 taken for future work or extension of the proposed work. 145

¹© 2014 Global Journals Inc. (US) Intuitionistic Partition based Conceptual Granulation Topic-Term Modeling



Figure 1:



Figure 2: ?



Figure 3:



Figure 4: Figure 1 :



Figure 5:



Figure 6: Figure 2 :

	Training with 300	Dimension	Perplexity Correlation			
	Doc	Reduction				
	Newspaper	26%	0.231	0.582		
	Reuters	22%	0.311	0.520		
	Movie	16%	0.483	0.480		
Dataset		tf-idf			IPCG	
Classifiers Reuters News Paper Movie				Reuters	News	Movie
					Paper	
SVM	0.482	0.422	0.321	0.844	0.841	0.799
NB	0.401	0.369	0.297	0.872	0.834	0.810
J48	0.400	0.399	0.381	0.798	0.797	0.784
Bayes'	0.541	0.411	0.399	0.831	0.854	0.829
LMT	0.442	0.541	0.587	0.878	0.798	0.722

Figure 7: Table 1 :

 $\mathbf{2}$

Figure 8: Table 2 :

14 CONCLUSIONS

- 146 [Feinerer et al. ()], I Feinerer, K Hornik, D Meyer. Journal of Statistical Software 2008. 25 p. .
- 147 [Malathi and Valarmathy ()] 'A Comprehensive Survey on Dimension Reduction Techniques for Concept Ex-
- traction from a Large Corpus'. D Malathi , S Valarmathy . International Journal of Computing Information
 Systems 2011. 3 p. .
- [Atanassov ()] K T Atanassov . Intuitionistic Fuzzy Sets, Theory, and Applications, Series in Fuzziness and Soft
 Computing, 1999. Phisica-Verlag.
- [Xu et al. ()] 'Clustering Algorithm for Intuitionistic Fuzzy Sets'. Z Xu , J Chen , J Wu . Information Sciences
 2008. 178 p. .
- 154 [Chiang and Lin ()] Correlation of fuzzy sets, Fuzzy Sets and Systems, D A Chiang, N P Lin. 1999. 102 p. .
- [Szmidt and Kacprzy ()] Distance Between Intuitionistic Fuzzy Set, Fuzzy Sets System, E Szmidt , J Kacprzy .
 2000. 114 p. .
- [Malathi and Valarmathy ()] 'Domain Classifier using Conceptual Granulation and Equal Partition Approach'.
 D Malathi , S Valarmathy . *Indian Journal of Engineering* 2013. 7 p. . (Science and Technology)
- 159 [Szmidt and Kacprzyk ()] Entropy for Intuitionistic Fuzzy Set, Fuzzy Sets System, E Szmidt , J Kacprzyk . 2001.
 118 p. .
- [Pelekis et al.] 'Fuzzy Clustering of Intuitionistic Fuzzy Data'. N Pelekis , D K Iakovidis , E K Evangelos , I
 Kopanakis . International Journal of Business Intelligence and Data Mining 3 (1) p. .
- 163 [Iakovidis et al. ()] 'Intuitionistic Fuzzy Clustering with Applications in Computer Vision. Advanced Concepts
- for Intelligent Vision Systems'. D K Iakovidis , N Pelekis , E K Evangelos , I Kopanakis . Lecture Notes in
 Computer Science 2008. 5259 p. .
- [Atanassov and Stoeva ()] 'Intuitionistic Fuzzy Set'. K T Atanassov , S Stoeva . Polish Symposium on Interval
 and Fuzzy Mathematics 1993. p. .
- 168 [Atanassov ()] Intuitionistic Fuzzy Set, Fuzzy Sets System, K T Atanassov . 1986. p. .
- [Blei et al. ()] 'Latent Dirichlet Allocation'. D M Blei , A Y Ng , M I Jordan . Journal of Machine Learning
 Research 2003. 3 p. .
- [Ye ()] 'Multicriteria Decision-making Method Based on a Cosine Similarity Measure between Trapezoidal Fuzzy
 Numbers'. J Ye . International Journal of Engineering, Science and Technology 2011. 3 p. .
- [Li and Cheng ()] 'New Similarity Measures Of Intuitionistic Fuzzy Sets And Application To Pattern Recognition'. D Li , C Cheng . *Pattern Recognition Letter* 2002. 23 p. .
- [Pang et al. ()] B Pang , L Lee , S Vaithyanathan . Thumbs up? Sentiment Classification using Machine Learning
 Techniques, Proceedings of EMNLP, 2002. p. .
- [Hofmann ()] 'Probabilistic Latent Semantic Analysis'. T Hofmann . Proceedings of the 15th Annual Conference on Uncertainty in Artificial Intelligence (UAI-99), (the 15th Annual Conference on Uncertainty in Artificial Intelligence (UAI-99)San Francisco, CA) 1999. Morgan Kaufmann. p. .
- [Sebastiani ()] F Sebastiani . 10.1145/505282.505283. Machine Learning in Automated Text Categorization, 2002.
 34 p. .
- [Le et al. ()] 'Semantics and Aggregation of Linguistic Information, Based on Hedge Algebras'. V H Le, C H
 Nguyen, F Liu. The 3rd International Conference on Knowledge, Information, and Creativity Support
 Systems, 2013.
- [Berendt (ed.) ()] Text Mining for News and Blogs Analysis, B Berendt . C. Sammut, & G. I. Webb (ed.) 2010.
 London: Springer. p. . (Encyclopedia of Machine learning)
- [Chien and Chueh ()] 'Topic-Based Hierarchical Segmentation'. J T Chien , C H Chueh . *IEEE Transactions on* Audio, Speech and Language Processing 2012. 20 p. .
- 189 [Brants et al. ()] 'Topicbased document segmentation with Probabilistic Latent Semantic Analysis'. T Brants
- F Chen , I Tsochantaridis . the proceeding of International Conference on Information and Knowledge
 Management, 2002. p. .