

A NOVEL METHOD FOR FUZZY BAG-OF-WORDS BASED ON WORD CLUSTERS

¹G MOUNICA, ²N NAGA BHAVANI, ³P.A.K ROHIT, ⁴Dr.D.SUNEETHA

^{1,2,3}Student, NRI Institute of Technology, Pothavarappadu (V), Via Nunna, Agiripalli (M), PIN-521 212

⁴Professor, department of CSE, NRI Institute of Technology, Pothavarappadu (V), Via Nunna, Agiripalli (M), PIN-521 212.

ABSTRACT:

One key issue in text mining and language process (NLP) is the way to effectively represent documents mistreatment numerical vectors. One classical model is that the Bag-of-Words (BoW). In a very BoW-based vector illustration of a document, every component denotes the normalized variety of prevalence of a basis term within the document. To count the amount of prevalence of a basis term, BoW conducts actual word matching, which might be considered a tough mapping from words to the premise term. BoW illustration suffers from its intrinsic extreme sparseness, high spatiality, and inability to capture high-level linguistics meanings behind text knowledge. To deal with the on top of problems, we have a tendency to propose a brand new document illustration methodology named Fuzzy Bag-of-Words (FBoW) during this project. Fuzzy Bag-of-Words model uses basis words for representation and clusters are formed as cluster-item pairs. Document representation is done by using this clusters. Since word semantic matching instead of exact word string matching is used, the FBoW could encode more semantics into the numerical representation. In addition, we propose to use word clusters instead of individual words as basis terms and develop Fuzzy Bag-of-WordClusters (FBoWC) models.

INTRODUCTION:

As the net grows, an outsized range of text data is currently out there within the kind of machine readable electronic documents. The method of Automatic data Retrieval so got abundant importance in recent years because of the exponential growth of the amount of documents in digital kind. Text classification or Text categorization is that the method to reason the digital documents into individual classes that describe the contents of the documents. Every document will belong to

1 or a lot of clusters supported by their contents. The Preprocessing step is followed by a Text classification rule in a very Text Categorization method. Feature extraction and have choice area unit the 2 main steps of Preprocessing. In Feature extraction, tokenization, stop word removal and stemming area unit carried. In Feature choice, the term coefficient strategies area unit administrated to search out the foremost relevant data from a collection of documents. For text classification, there area unit completely

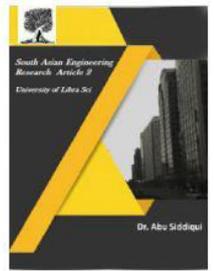


2581-4575

International Journal For Recent Developments in Science & Technology



A Peer Reviewed Research Journal



different algorithms. K nearest neighbor is one amongst the foremost economical and straightforward classification algorithms used for text classification. used. In this paper, we propose fuzzy BoW models to learn more dense and robust document representations encoding more semantics. To overcome the limitations of the original BoW model, we propose to replace the original hard mapping by a fuzzy mapping, and develop the FuzzyBoW (FBoW) model.

RELATED WORK:

The planned model is fuzzy bag of words model which may be terribly helpful in classification by reducing the additional effort. exploitation fuzzy bag of words the classified words ar keep within the bag which is able to be utilized in classification. There ar several measures to implement fuzzy bag of words model. this could be free from meagerness, high spatial property, and inability of capturing the linguistics meanings of the text. to create this linguistics matching of words ar replaced the word to word actual matching by linguistics matching as this can be additional outstanding. The fuzzy bag of words may inscribe additional linguistics into the numerical illustration. during this we tend to mentioned regarding K-Nearest neighbor as this has it's own importance in forming the clusters, thus by exploitation KNN with fuzzy bag of words will offer the most effective results with within the limit. TF-idf is employed to come up

with the candidate keys and people ar employed by knn with fuzzy bag of words to categorise the document consequently to accumulate highest classification accuracies.

EXISTING SYSTEM:

In topic models as well as probabilistic latent linguistics analysis and latent Dirichlet allocation, chance distributions ar introduced to explain words and also the generation method of every word during a document. the idea behind the subject model is that word selection during a document are influenced by the subject of the document probabilistically. However, in these models, the derived latent dimension lacks linguistics interpretation. for instance, LSA regards a latent dimension as a linear combination of all original terms within the vocabulary, that is counter-intuitive as a result of solely alittle a part of the vocabulary has relevancy to a precise topic. Besides, these 2 approaches each utilize the word incidence of documents to perform spatial property reduction. However, the incidence statistics might not be ready to capture actuality linguistics data underlying a document.

PROPOSED SYSTEM:

Totally different from BoW model and BoW-enhanced models like LSA and topic models that use actual word matching and arduous mapping, our planned FBoW and FBoWC models adopt linguistics matching and fuzzy mapping to project the words occurred in

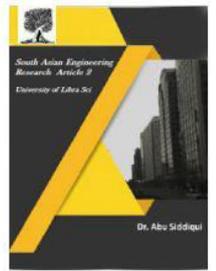


2581-4575

International Journal For Recent Developments in Science & Technology



A Peer Reviewed Research Journal



documents to the fundamental terms. In our planned fuzzy BoW models, linguistics similarity between words are based on clusters. It's believed that the captured similarity data is additional correct and general than that extracted from word incidence statistics underlying a document in previous BoW-based approaches. Besides, our planned fuzzy BoW models also can be utilized in conjunction with the LSA methodology to cut back the spatial property of the FBoW illustration. The text knowledge or data is taken as cluster-item pair where the semantics of the text is taken into account. Semantically similar words are keep together in a cluster that is served as the basic word to search documents. Based on the cluster name all the similar words are retrieved and presented. Every item is a hyperlink that generates some URLs that leads to a web page.

ARCHITECTURE:

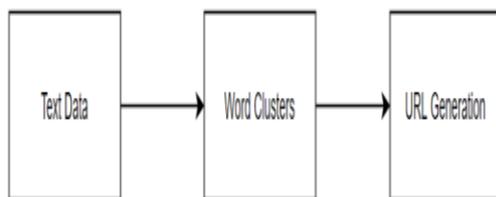


Fig. Process diagram

ALGORITHM:

Here we have a tendency to area unit mistreatment Bag of words and Fuzzy Bag of Words for document classification, algorithmic program for Bag of Words: it's employed in order to perform the term

frequency of a document.

Bag of Words:

Bag of words ignores descriptive linguistics and arrangement of words. Here we have a tendency to begin with 2 documents that is thought as corpus. a listing is formed supported the distinctive words within the corpus. Here we are going to use count Vectors to form vectors from the corpus. It counts the term frequency supported the documents. so Bag of Words is enforced.

Fuzzy Bag of Words:

FBOW is additionally called vector house model. Here a sentence is pictured as a multi set of words with none priority of descriptive linguistics. it's conjointly used for pc vision. every part here represents variety supported the frequency of the term, bag of words is exactly matched however the linguistics that means behind the information is captured because of extreme sparseness and high dimension. FBOW is employed for primarily document illustration, image classification.

RESULTS:

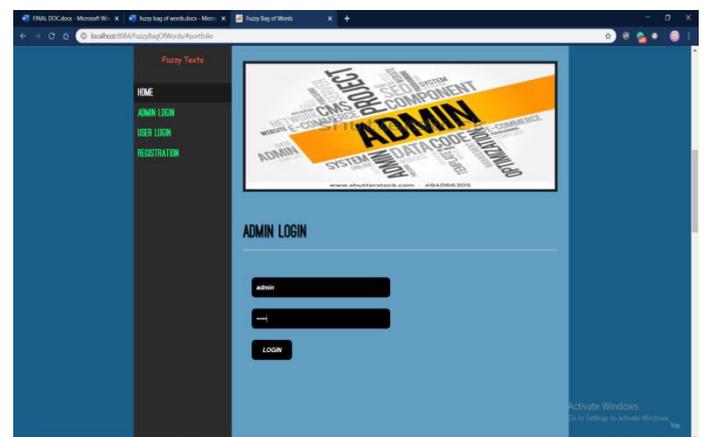


Fig. Admin Login Page



2581-4575

International Journal For Recent Developments in Science & Technology



A Peer Reviewed Research Journal

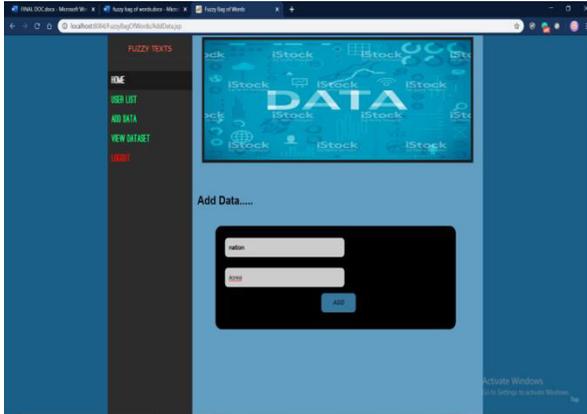
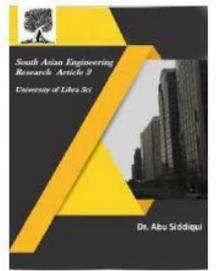


Fig. Add Data Page

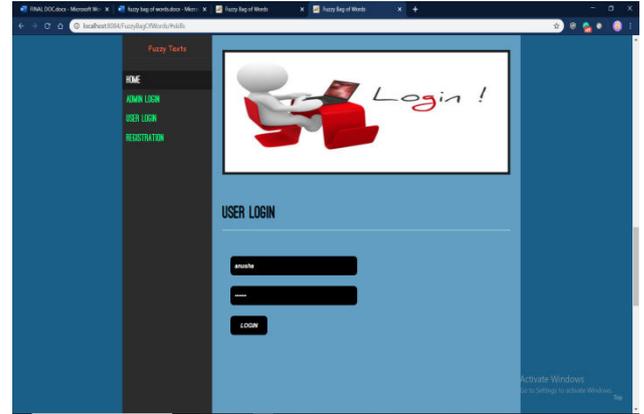


Fig. User Login Page

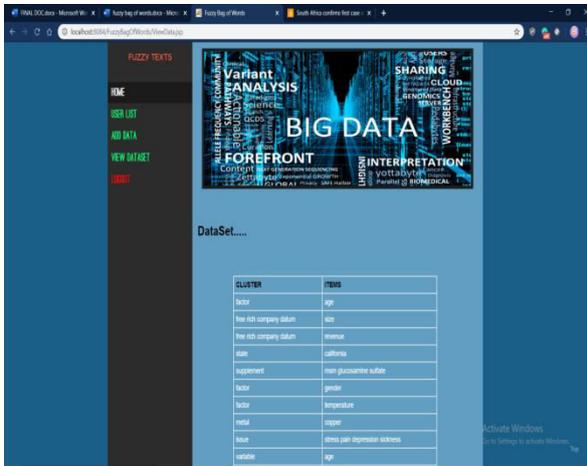


Fig. View Dataset Page

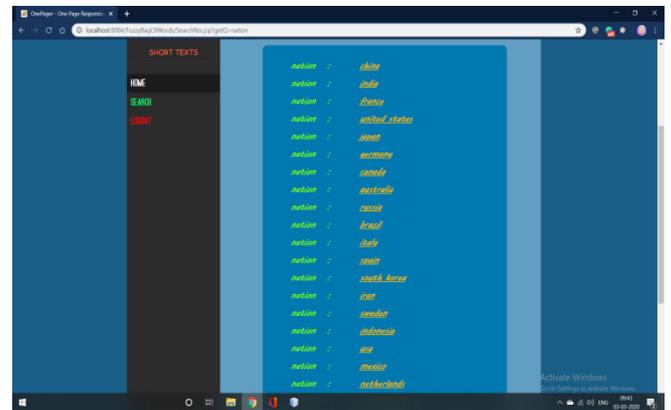


Fig. Search Results Page

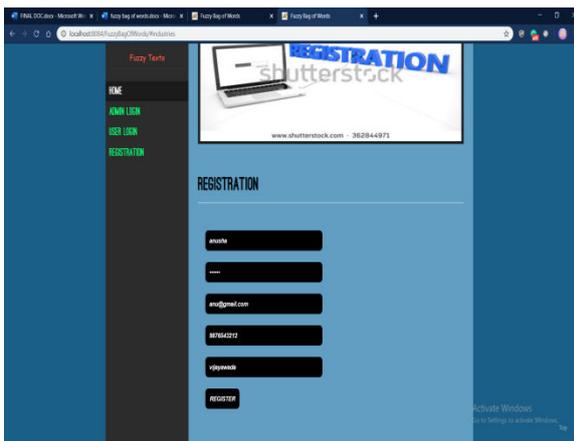


Fig. Registration Page

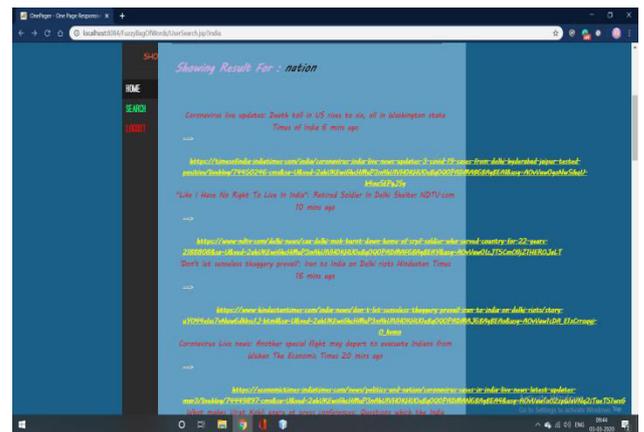


Fig. Search Results Page.

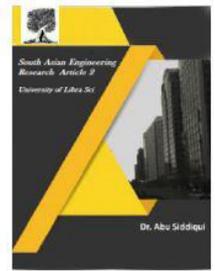


Fig. Search Results Page

CONCLUSION AND FUTURE ENHANCEMENT:

In this paper, we have proposed Fuzzy Bag-of-Words models including FBoW and FBoWC to address issues of sparsity and lack of high-level semantics of BoW representation. Word embeddings are utilized to measure semantic similarity among words and construct fuzzy membership functions of basis terms in BoW space over words in the task-specific corpus. Since word2vec embeddings can be trained over billions of words, word embeddings adopted in our methods are able to capture high-quality and meaningful semantic information that are not contained by the task-specific corpus alone. To determine basis terms in BoW space, FBoWC utilizes word clusters, while FBoW directly regards high-term-frequencies words as original BoW does. The adoption of word clusters in FBoWC can reduce feature redundancy and improve feature discrimination. Three different measures have been designed to evaluate similarity between clusters and words, and three corresponding variants of FBoWC models as FBoWC_{mean}, FBoWC_{max} and FBoWC_{min} have been

developed. The performance of our approaches has been experimentally verified through seven multi-class document categorization tasks. As a next step work, document structure or word order information will be considered in document representation learning. In addition, the effects of multi-sense word embeddings and different term weighting schemes will be explored in future.

REFERENCES:

- [1] Y. Bengio, A. Courville, and P. Vincent, "Representation learning: A review and new perspectives," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 35, no. 8, pp. 1798–1828, 2013.
- [2] K. Nigam, A. K. McCallum, S. Thrun, and T. Mitchell, "Text classification from labeled and unlabeled documents using em," *Machine learning*, vol. 39, no. 2-3, pp. 103–134, 2000.
- [3] R. Zhao and K. Mao, "Supervised adaptive-transfer pls for cross domain text classification," in *Data Mining Workshop (ICDMW), 2014 IEEE International Conference on. IEEE*, 2014, pp. 259–266.
- [4] B. Pang and L. Lee, "Opinion mining and sentiment analysis," *Foundations and trends in information retrieval*, vol. 2, no. 1-2, pp. 1–135, 2008.
- [5] M. Steinbach, G. Karypis, V. Kumar et al., "A comparison of document clustering techniques," in *KDD workshop on text mining*, vol. 400, no. 1. Boston, 2000, pp. 525–526.
- [6] R. Zhao and K. Mao, "Cyberbullying detection based on semanticehanced marginalized denoising auto-encoder," *IEEE*

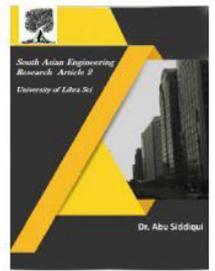


2581-4575

International Journal For Recent Developments in Science & Technology



A Peer Reviewed Research Journal



Transactions on Affective Computing, vol. PP, no. 99, pp. 1–1, 2016.

[7] K. Sparck Jones, “A statistical interpretation of term specificity and its application in retrieval,” *Journal of documentation*, vol. 28, no. 1, pp.11–21, 1972.

[8] M. Lan, C. L. Tan, J. Su, and Y. Lu, “Supervised and traditional term weighting methods for automatic text categorization,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 31, no. 4, pp. 721–735, April 2009.

[9] G. Salton and C. Buckley, “Term-weighting approaches in automatic text retrieval,” *Information processing & management*, vol. 24, no. 5, pp. 513–523, 1988.

[10] H. C. Wu, R. W. P. Luk, K. F. Wong, and K. L. Kwok, “Interpreting tf-idf term weights as making relevance decisions,” *ACM Transactions on Information Systems (TOIS)*, vol. 26, no. 3, p. 13, 2008.

[11] T. Mikolov, K. Chen, G. Corrado, and J. Dean, “Efficient estimation of word representations in vector space,” in *ICLR 2013*. ICLR, 2013.[Online]. Available: <https://arxiv.org/pdf/1301.3781.pdf>

[12] S. Dumais, G. Furnas, T. Landauer, S. Deerwester, S. Deerwester et al., “Latent semantic indexing,” in *Proceedings of the Text Retrieval Conference*, 1995.

[13] T. Hofmann, “Unsupervised learning by probabilistic latent semantic analysis,” *Machine learning*, vol. 42, no. 1-2, pp. 177–196, 2001.

[14] D. M. Blei, A. Y. Ng, and M. I. Jordan, “Latent dirichlet allocation,” *the Journal of*

machine Learning research, vol. 3, pp. 993–1022, 2003.

[15] Y. Bengio, R. Ducharme, P. Vincent, and C. Janvin, “A neural probabilistic language model,” *The Journal of Machine Learning Research*, vol. 3, pp. 1137–1155, 2003.

[16] R. Collobert, J. Weston, L. Bottou, M. Karlen, K. Kavukcuoglu, and P. Kuksa, “Natural language processing (almost) from scratch,” *The Journal of Machine Learning Research*, vol. 12, pp. 2493–2537, 2011.

[17] A. Mnih and G. E. Hinton, “A scalable hierarchical distributed language model,” in *Advances in neural information processing systems*, 2009, pp.1081–1088.

[18] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean, “Distributed representations of words and phrases and their compositionality,” in *Advances in Neural Information Processing Systems*, 2013, pp. 3111–3119.

[19] R. Socher, J. Pennington, E. H. Huang, A. Y. Ng, and C. D. Manning, “Semi-supervised recursive auto encoders for predicting sentiment distributions,” in *Proceedings of the Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 2011, pp. 151–161.

[20] R. Socher, A. Perelygin, J. Y. Wu, J. Chuang, C. D. Manning, A. Y. Ng, and C. Potts, “Recursive deep models for semantic compositionality over a sentiment tree bank,” in *Proceedings of the conference on empirical methods in natural language processing (EMNLP)*, vol. 1631. Citeseer, 2013, p. 1642.



2581-4575

International Journal For Recent Developments in Science & Technology



A Peer Reviewed Research Journal

