

EMPLOYABILITY OF NEURAL NETWORK ALGORITHM IN MEASURING SIMILARITY IN TEXT DOCUMENT CLUSTERING

Jahnvi Gupta

Student, University Institute of Engineering and Technology, Panjab University, Chandigarh, India

ABSTRACT

Broad utilization of the World Wide Web for data search utilizing famous web indexes has turned numerous specialists to focuses on text mining issues. NLP required effective techniques to take the real needs of the client during Machine Learning. Using genetic algorithm and similarity measures for text mining during document clustering yields critical outcomes for WordSim353 informational collections. Tests show that using Echo State Neural Network and Radial Basis Function to the informative preparation index gives the better classification of text archives dependent on the stored-weights to stay away from the recovery of unnecessary documents.

I. INTRODUCTION

The idea of similarity measure in web indexes is valuable and is used in numerous applications. To recover the ideal web reports dependent on the inquiry, the similitude measure assumes a significant part. The idea of the archives is unique; it tends to be either organized or semi-organized, which causes trouble in taking care of. Along these lines, utilize the concept of clustering to develop the recovery approach further. It is a deliberate method of managing the reports given to the clients dependent on the inquiry. It tends to be finished by creating dependent on substance and connections. Digital psychological warfare examination, subject point, email steering, and language considering are applications where the vast majority of the searches are achieved. Pre-processing and archive representation investigations are done to work on the accuracy and characterization. In pre-processing, to discover the ideas from many words, highlights are achieved using Verb-Argument Structures. In some examination regions, a pack of terms is found from the content reports. This enormous arrangement of words should be decreased utilizing highlight clustering techniques. The resultant is additionally analyzed for the archive nearest, and reports are grouped if they are comparative.

Numerous fuzzy nearest based models and predictions have been presented with the basic idea of its participation limits, fuzzy association, fuzzy c-implies, creation rules . The demonstration of a data framework relies upon the calculations, and the appraisal of the presentation dwells in the actions applied that could be completed adequately. The client will consistently like accurate and wonderful site pages to the question given to the framework. Thus our key goal is to nourish a decent assumed match of the investigation. The specific game is feasible to decide if the worth of a given term coordinates with the worth determined in the inquiry, which requires an action to fulfill the data need.

The rapid development of the World Wide Web has forced difficulties to group the reports over the web and work on their effectiveness. Web indexes are confronting trouble getting sorted out the

relevant messages among enormous list items that got back to a basic question. We require a powerful strategy to tackle the issue by clustering the comparable archives, which helps the client recognize the applicable information effectively.

We overviewed related works and referenced them in Section 2, trailed by preparation the means engaged with different periods of our proposed technique in Section 3. It is followed by a mathematical examination of the outcomes acquired during the preparation cycle for Word-Sim353 informational indexes. The human evaluations of the cosine similarity measures are obtained utilizing Kardi assets. Section 4 presents the outcomes acquired from the Weka tool on utilization of the information-mining procedures for the results obtained in the preparation cycle. At last, in Section 5, we finish up with the outline, featuring the downsides and bearings for future work.

Clustering a document is significant for the fast retrieval of important reports on the web for a given pair of words. The time taken for recovering the records ought to be in under seconds so the client can check whether the document applies to the query. Programming and calculation are being utilized via web indexes, which can fulfil the clients' questions and return pertinent reports accurately on schedule. Consequently, more productive new archive grouping calculations are needed than ordinary bunching calculations.

Ling Zhuang Honghua Dai 2004 presented the underlying focuses fork-implies calculation as irregular focuses. This bunching approach is hard for understanding as it is unstructured and noisy.

Benjamin Fung et al., 2003 has presented a technique for clustering the reports. A tree is built dependent on the subjects, and likeness is produced among clusters. Sharma et al., 2009 has presented this methodology for huge word sets.

Report list chart based record clustering is advanced by Momin et al., 2006. Record clustering techniques depend on a solitary term assessment of the informational report index. Archive Index Graph (DIG) permits records to be encoded utilizing phrases. It centres around further developing expression based closeness measures. Additionally, a Document Index Graph-based Clustering (DIGBC) calculation is likewise being proposed. It steadily structures clusters dependent on the group report similarity measure, and the archives can be allocated to more than one group.

Muflikhah et al., 2009 presented space and cosine similitude estimation. The researcher utilized Latent Semantic Index (LSI) approach with Principle Component Analysis (PCA) or Singular Vector Decomposition (SVD). The strategy diminishes the framework measurement by recognizing the example in the record collection, which alludes to synchronous terms. Each method utilizes recurrence of the term as weight in Vector Space Model (VSM) with the assistance of fluffy c-means procedure for clustering. Shyu et al., 2004 introduced Web report grouping on partiality-based closeness measure. The idea is enhanced to be utilized in web archive clustering by building up the fondness-based closeness measure methodology and uses access examples to discover the similitudes through a probabilistic model. Different analyses are dissected dependent on the continuous dataset. The trial results represented that the introduced likeness measure beats Euclidean distance and cosine coefficient under other report grouping methods.

Eldest et al., 2009, given a novel similitude measure for archive grouping dependent on point phrases. The most recent pattern utilizes the expression to be a more educational component and considering the issue that drove in expanding the exhibition of record grouping. This paper used a strategy for breaking down the similarity proportion of VSM by considering the theme expressions of the record and applying them to the Buckshot procedure. These techniques increase the upsides of metrics to develop grouping viability further.

Cobos et al., 2010 utilized k-means, term sets, Bayesian data for web record grouping. Haojun et al., 2008 created advanced calculations for report clustering using group covering rates to develop productivity further. Subgroups (say two) are combined if they have a high covering rate. In order to number the parameter, this paper utilized the Gaussian combination as a model in the Expectation-Maximization technique.

II. TOOLS AND TECHNIQUES

Should perform Pre-handling assignments before information mining techniques are utilized on the informational index. These incorporate the information in filtering, gathering, recognizing proof, way fulfilment, and arranging. It is done to work on the nature of the outcomes. This segment talks about different periods of the proposed technique, as displayed in Figure 1.

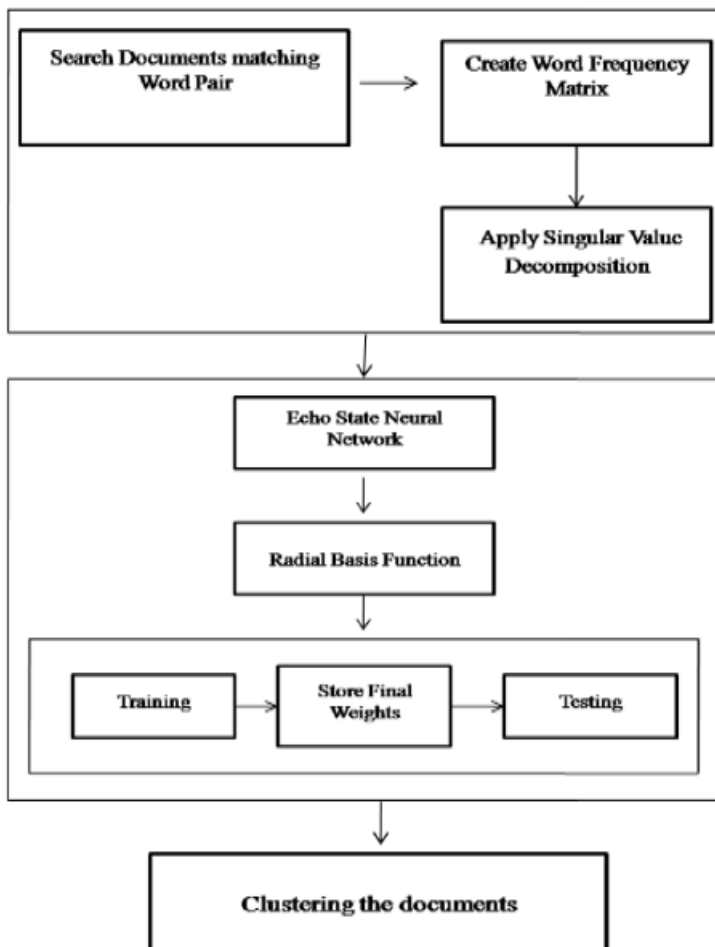


Figure1. System Architecture

Document are records that contain ASCII and non-ASCII characters. Test classifications are postulation, paper, diary, Invoice, quote, RFP, Proposal, Contract, Packing slip, Manifest, Report nitty-gritty and rundown, Spreadsheet, Waybill, Bill of Lading, Financial articulation, Nondisclosure agreement, Mutual nondisclosure arrangement, summons, testament, permit periodical, white paper, application structures, client guide, brief, model, script. These reports can be made using designs. Can modernize the pieces by using various editors, beginning from fundamental editors who employ just ASCII characters to complex editors who insert illustrations. Unique characters address the images in a report.

Looking at comparable records includes handling archives where disposal of data ought to be tried not to deal with word pair search. The libraries should initially be preprocessed. When the documents are appropriately preprocessed, the proportions of similitude will be awesome, and the nature of the search will be most excellent. Can change each library over into a vector D. Every vector contains multiple measurements. Each component of a vector is a one of a kind copy of a report. When the substance of two words is practically comparative, then, at that point, the mathematical upsides of the vector are almost the same. Each element of the vector will address a specific term. The term is a solitary word or expression.

2.1 Preprocessing

Stage 1. Pair of words is given in the query region.

Stage 2. Each record is preprocessed and changed over into vectors of mathematical qualities. If the vectors of mathematical rates are now accessible in the looking through organizer comparing to the accessible records, then, at that point, preprocessing of the archives and changing over into vectors need not be done.

Stage 3. The packs of the description of words (certain single words) are made. During this cycle, unimportant information is precluded for computation purposes.

Stage 4. Ordering the bag of words is finished. In this cycle, tokens are made by portioning strings by blank area and accentuations called Tokenization. Every token stems from its root structure by changing over a thing into its particular design and eliminating linguistic words like articles, conjunctions, pronouns called Stop Word Elimination.

Stage 5. Changing over vectors into a grid is framed. In this grid presence of nonattendance of words compared to an archive is shown utilizing 1 or 0. Then again, the recurrence of words is addressed in portions in the grid by standardization.

2.2 Computing ANN calculation

During the preparation cycle, different organization boundaries are streamlined, and because of which the organization goes through a learning stage. The sorts of calculations for preparing the ANN

geography are Radial Basis Function (RBF) and Echo State Neural Network (ESNN). The grids created above are utilized to register Echo State Neural Network calculation (ESNN).

Reverberation State Neural Network can best fit in multi-dimensional space in the wake of guaranteeing the best match to the preparation information. This one of a kind element empowers it to deal with issues identified with order and Clustering. The spline work applied to the informational preparation collection is:

where x addresses hubs of the organization, which is given as a contribution to the neural organization calculation. Reverberation State Network gives better bunching capacities to a recently made informational index with a total group determination. Steadily, a gauged framework is produced utilizing reverberation state property after the preparation interaction, and the entire preparing set is gone through the organization.

Outspread premise work utilizes remarkable capacity as its enactment esteems in the hidden layer of the ANN geography. In this calculation, preparing designs with input highlights addressing a record is used. Focus designs are made from preparation designs. The summation of the distance between each preparation example and every one of the middle examples are found. The added values are ignored an actuation capacity to get RBF yield in the hidden layer. An inclination worth of „1“ is utilized in the secret layer yield. Also, for all the leftover preparing designs, the RBF yields are acquired as the yields in the personal layer. All the RBF yields from all the preparation designs are handled with the relegated target esteems to get many conclusive loads.

The calculations are created dependent on various weight refreshing principles. In each weight vital standard, mistakes are determined in the forward cycle of the ANN, and weight updation is finished during the opposite or intermittent interaction. In the preparation cycle of the ANN calculations, the associations among the hubs between layers are addressed by grids. In a large portion of the calculations, the frameworks are instated with arbitrary numbers. Toward the finish of the preparation cycle, the networks contain the last weight. During the testing of the ANN or testing the recovering of the records relating to word pair, last loads are utilized for handling with the vector comparing to word pair. The previous loads are gotten from the actual example with no introduction of the weight networks in another strategy.

2.3 Clustering the archives

In light of the put away last loads, records are bunched utilizing Expectation-Maximization. This is finished using a weka tool. It delivers a most disproportionate number of fifteen groups to perform possible recovery of reports.

III. OUTCOMES AND DISCUSSION

Figure 2 shows the word sets with human assessments. These words are discovered to be semantically comparative in the larger part of the records. This closeness is used for building the similarity lattice.

By and large, it contains 1.0s along with the crooked. Even though two words are discovered to be indistinguishable, they may compare in their significance. In any case, their worth remaining parts as before (for example, 1.0) is opposing and can be considered for additional exploration in semantic similarity.

Figure 2. Word pairs and their similarity scores

Relation: term_sim_measure			
No.	1: Word 1 Nominal	2: Word 2 Nominal	3: Human (mean) Numeric
330	problem	challenge	6.75
331	size	prominence	5.31
332	country	citizen	7.31
333	planet	people	5.75
334	develop...	issue	3.97
335	experience	music	3.63
336	music	project	3.63
337	glass	metal	5.56
338	aluminum	metal	7.83
339	chance	credibility	3.88
340	exhibit	memorabilia	5.31
341	concert	virtuoso	6.81
342	rock	jazz	7.59
343	museum	theater	7.19
344	observat...	architect...	4.38
345	space	world	6.53
346	preserva...	world	6.19
347	admission	ticket	7.69
348	shower	thunders...	6.31
349	shower	flood	6.03
350	weather	forecast	8.34
351	disaster	area	6.25
352	governor	office	6.34
353	architect...	century	3.78

Table 1 below contains test word sets explored different avenues regarding by us from the WorSim353 informational index. The relationship coefficient and the other error rate is determined using the weka device. Various examinations tracked down that human scoring has reliably high connections. In both RG and WordSim353, the confidence levels are huge and critical in their distinction.

Table 1. Error Rates with Correlation Coefficient

Lowest Value = -0.34971644612476366 Highest Value = 0.6502835538752364	
Inverted Covariance Matrix * Target-value Vector: Lowest Value = -0.2588713208324032 Highest Value = 0.27371229405845615	
Time taken to build model: 0.66 seconds Scheme: weka.classifiers.functions.LinearRegression	
Correlation coefficient	0.9632
Mean absolute error	0.6996
Root mean squared error	0.87
Relative absolute error	39.0574 %
Root relative squared error	40.0484 %
Coverage of cases (0.95 level)	100 %
Total Number of Instances	353

Table 2 Applying RBF for training of Data

Scheme: weka.classifiers.functions.RBFNetwork -B 2 -S 1 -R 1.0E-8 -M -1 -W 0.1	
Relation:	Data- weka.filters.supervised.attribute.AddClassification-Wweka.classifiers.rules.ZeroR
Instances: 353	Attributes: 4
Test mode :10-fold cross-validation	
Radial basis function network (Linear regression applied to K-means clusters as basic functions):	
Linear Regression Model SM = -2.1219 * pCluster_0_0 + 2.1212 * pCluster_0_1 + 5.5318	
Time taken to build model: 0.08 seconds	

Straight relapse order calculation yields a normal worth of 0.58 around with an outright mistake pace of 40 %. This proposed approach is appropriate for any unique situation, as question preparation doesn't need an arrangement plan. Table 2 sums up the error rates for 353 ordered examples for the given word-sim353 information.

Table 2 shows the preparation results procured using RBF utilizing the Weka device for five cycles. Further lists show that the determined scores would permit us a correlation for the scientific classification based sparkles.

Figure 3 shows the archives containing the necessary word sets are grouped utilizing the Expectation-Maximization technique. It uses 353 examples for testing ascribes. It delivers a limit of fourteen quantities of groups to perform compelling recovery of records.

Scheme: weka.clusters.EM-I100-N-M1.0E-6-S100															
Relation: Aishu_dataweka.filters.supervised.attribute.AddClassification-															
Wweka.classifiers.rules.ZeroR															
Instances: 353															
Attributes: 4															
Word 1															
Word 2															
HR															
PM															
Test mode: user supplied test set: 353 instances															
==== Model and evaluation on training set====															
EM= Number of clusters selected by cross validation: 15															
Cluster															
Attribute	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
	0.1	0.12	0.06	0.01	0.08	0.02	0.06	0.09	0.07	0.06	0.02	0.08	0.04	0.15	0.02

Fig 3. Using Expectation-Maximization for Clustering of documents

IV. CONCLUSION

This paper accentuates the relevance of fake neural groups in bunching reports. Should consider the significant three stages indicated in this paper for grouping the information is utilizing neural organization calculations. Trials led on 353-word sets show that the proposed strategy outflanks. The yields of ANN can also be deciphered whether the reports recovered or bunched are pertinent to the words. Future work incorporates the utilization of other ANN calculations that can execute for comparative grouping records. The constraint is that it's anything but a qualification among essential and auxiliary classifications.

REFERENCES

- [1] F. Beil, M. Ester, and X. Xu, "Frequent term-based text clustering", In Proceedings of 8th International Conference on Knowledge Discovery and Data Mining, 2002.
- [2] A. Budanitsky, and G. Hirst, "Evaluating wordnet-based measures of semantic distance", Comput. Linguistics, vol. 32, no. 1, pp. 13–47, 2006.
- [3] C.M. Benjamin Fung, Wang Ke, and Ester Martin, "Hierarchical document clustering using frequent item sets", In Proceedings SIAM International Conference on Data Mining, pp. 59-70, 2003.

- [4] C. Cobos, J. Andrade, W. Constain., M. Mendoza., and E. Leon, "Web document clustering based on global-best harmony search, k-means, frequent term sets and bayesian information criterion", IEEE Congress on Evolutionary Computation, pp. 1-8, 2010.
- [5] A.E. Eldesoky, M. Saleh, and N.A. Sakr, "Novel similarity measure for document clustering based on topic phrases", International Conference on Networking and Media Convergence, pp. 92-96, 2009.
- [6] Haojun Sun, Zhihui Liu, and Lingjun Kong, "A document clustering method based on hierarchical algorithm with model clustering", 22nd International Conference on Advanced Information Networking and Applications, pp. 1229-1233, 2008.
- [7] Han-Saem Park, Si-Ho Yoo, and Sung-Bae Cho, "Evolutionary Fuzzy Clustering Algorithm with Knowledge-Based Evaluation and Applications for Gene Expression Profiling", Journal of Computational and Theoretical Nanoscience, vol. 11, no. 4, pp. 524-53, 2005.
- [8] S. Karthick, S.M. Shalinie, A. Eswarimeena, P.Madhumitha, T.N. Abhinaya, "Effect of multi-word features on the hierarchical clustering of web documents", Recent Trends in Information Technology(ICRTIT) International Conference, pp. 1 – 6, 2014.
- [9] Ling Zhuang, and Honghua Dai, "A maximal frequent item set approach for web document clustering", In Proceedings of the IEEE Fourth International Conference on Computer and Information Technology, 2004.
- [10] B.F. Momin, P.J. Kulkarni, and A. Chaudhari, "Web document clustering using document index graph", In Proceedings IEEE International Conference on Advanced Computing and Communications, 2006.
- [11] L. Muflikhah, and B. Baharudin, "Document clustering using concept space and cosine similarity measurement", International Conference on Computer Technology and Development, vol. 1, pp. 58-62, 2009.
- [12] N. Narayanan, J. E. Judith and J. JayaKumari, "Enhanced distributed document clustering algorithm using different similarity measures", Information & Communication Technologies (ICT), IEEE Conference, pp. 545-550, 2013.
- [13] H.A. Nguyen, and H. Al-Mubaid, "New ontology-based semantic similarity measure for the biomedical domain", In Proc. IEEE GrC, pp. 623–628, 2006.
- [14] Peipei Li, Haixun Wang, K.Q. Zhu Zhongyuan Wang, Xuegang Hu and Xindong Wu, "A large probabilistic semantic network based approach to compute term similarity", IEEE Transaction on Knowledge and Data Engineering, vol. 27, pp. 2604-2617, 2015.
- [15] J. Prasannakumar, and P. Govindarajulu, "Duplicate and near duplicate documents detection", A Review European Journal of Scientific Research ISSN 1450-216X vol. 32, no. 4, pp. 514-527, 2009.
- [16] G.S. Reddy, T.V. Rajinikanth, A.A. Rao, "A frequent term based text clustering approach using novel similarity measure", Advanced Computer Conference (IACC), IEEE International, pp. 495-499, 2014.
- [17] Ruxixu and Donald Wunsch, "A survey of clustering algorithms", IEEE Transactions on Neural Networks, vol. 16, no. 3, pp. 645-678, 2005.
- [18] S. Satwardhan, Incorporating dictionary and corpus information into a context vector measure of semantic relatedness, Master's thesis, Univ. Minnesota, Minneapolis, 2003.

- [19] K. Selvi, and R.M. Suresh, "Context similarity measure using fuzzy formal concept analysis", In Proc. of The Second Int'l conference On Computer Science and Engineering and Information Technology CCSEIT, pp. 416-423, 2012.
- [20] K. Selvi, and R.M. Suresh, "An efficient technique to implement similarity measures in text document clustering using artificial neural network algorithm", Research Journal of Applied Sciences Engineering and Technology, vol. 8(23), pp. 2320-2328, 2014.
- [21] A. Sharma, and R. Dhir, "A wordsets based document clustering algorithm for large datasets", In Proceeding of International Conference on Methods and Models in Computer Science, 2009.
- [22] M.L. Shyu, S.C. Chen, M. Chen, and S.H. Rubin, "Affinity-based similarity measure for web document clustering", IEEE International Conference on Information Reuse and Integration, pp. 247-252, 2004.
- [23] K.M. Sim, and P.T. Wong, "Toward agency and ontology for web-based information retrieval", IEEE Trans. Syst., Man, Cybern. C, Appl. Rev., vol. 34, no. 3, pp. 257-269, 2004.
- [24] Y. Syed Mudhasir, and J. Deepika, "Near duplicate detection and elimination based on web provenance for efficient web search", In the Proceedings of International Journal on Internet and Distributed Computing Systems, vol. 1, no. 1, pp. 22-32, 2011.
- [25] Ted Pedersen, V.S. Serguei. Pakhomov, Siddharth Patwardhan, and Christopher G. Chute, "Measures of semantic similarity and relatedness in the biomedical domain", Journal of Biomedical Informatics, vol. 40, pp. 288-299, 2007.
- [26] Thanh Van Le, Trong Nghia, Hong Nam Nguyen, Tran Vu Pham, "An efficient pretopological approach for document clustering", Intelligent Networking and Collaborative Systems (INCoS), 5th International Conference , pp. 114 – 120, 2013.
- [27] <http://people.revoledu.com/kardi/tutorial/Similarity/Stringinstance.html#TextSimilarityCalculator>.
- [28] Xinjuan Peng, Lijun Cai, Bo Liao, Haowen Chen, and Wen Zhu, "Detecting the Maximum Similarity Bi-Clusters of Gene Expression Data with Evolutionary Computation", Journal of Computational and Theoretical Nanoscience, vol. 11, no. 7, pp. 1585-1591, 2014.