

# Nomenclature and Contemporary Affirmation of the Unsupervised Learning in Text and Document Mining By

Annaluri Sreenivasa Rao<sup>1</sup> and Prof. S. Ramakrishna<sup>2</sup>

<sup>1</sup> MRCE, Hyderabad, Telangana State, India

Received: 15 December 2014 Accepted: 31 December 2014 Published: 15 January 2015

6

## Abstract

Document clustering is primarily a method applied for an uncomplicated, document search, analysis and review of content or is a process of automatic classification of documents of similar type categorized to relevant clusters, in a clustering hierarchy. In this paper a review of the related work in the field of document clustering from the simple techniques of word and phrase to the present complex techniques of statistical analysis, machine learning etc are illustrated with their implications for future research work.

14

**Index terms**— Document clustering is primarily a method applied for an uncomplicated, document search, analysis and review of content or is a process of automatic c

## 1 Introduction

Document clustering [1], [2], [3], [4] techniques find relevance in a wide range of tasks from a simple search with a few terms to vast information retrieval processes. The early document clustering techniques used were developed for typically enhancing information retrieval systems [5], were designed to find documents according to the query type, however could not perform the task of creating a query, generate a synopsis of the documents, or provide an interface to the search results. The progress of internet, digital libraries, news sources and companywide intranets has made available huge volumes of text documents. The tremendous increase in the already quantum size of web data and the classification of the web documents into relevant and moderate number of clusters has led to the development of large number of web clustering engines and high performing clustering algorithms.

The process of document clustering involves four stages which are, i) Data collection, crawling to accumulate the documents, indexing the set of documents in a structured fashion, filtering of data with techniques of tokenization, stop words removal and stemming, lemming etc. ii) preprocessing where the data is represented in suitable form, vector etc. and measurable factors applied to determine the similarity, iii) Document clustering where a clustering technique and an efficient clustering algorithm are identified for clustering based on preset criteria and iv) Post processing involving applications of business and scientific requirements adaptation of the document clustering technique.

The applications of document clustering are of diverse nature such as, i) Creation of document taxonomies ii) IR process of search, accessing and collection [6],

Similar documents identification, review and classification of results [7], automatic topic extraction [8], content summarization iii) Recommendation System, iv) Search Optimization, etc. For instance the processes are used enormously in the data classification process such as Google Web Directory, Social media data classification etc.

The clustering techniques though being studied since several years, still face many of the same challenges. These challenges [9,10] of document clustering are mostly of, i) Huge volume of data, ii) The high dimensionality of the feature space, iii) A feasible clustering method in terms of constraints such as cluster quality and performance and iv) Representing the results in an effective browsing interface. The current challenges associated with text clustering are the requirement of dynamic clustering techniques to incrementally update clusters as new data is added [11,12]. For instance the social media has to generate user specific content [13] instantly and this requires real time data clustering methodologies.

## 4 CONTEMPORARY AFFIRMATION OF THE

---

45 The remainder of this paper is organized as follows. In Section 2 we discuss the "Taxonomy" of document  
46 clustering, in Section 3 the "Contemporary literature work of clustering techniques" are evaluated and Section 4  
47 gives the "Conclusion" of the paper.

### 48 2 II. Taxonomy

49 The clustering functionality can be expressed as a function comprising of a document set mapped to a D set  
50 of clusters. Based on specified constraints the minimum and maximum of the function defines the clustering  
51 difficulty and algorithms applied over the similarity criteria determine the clustering quality.

52 The preprocessing step of clustering for finding the document similarity is determined with methods based  
53 on the following strategies, (i) phrase or pair-wise methodology, (ii) tree form data depiction, (iii) component  
54 dependent data depiction, (iv) semantic relation dependent documents depiction, (v) concept and feature vector  
55 dependent depiction.

56 The clustering methods of are generally of two types, 1) Word patterns and phrases based 2) Feature based.

57 The clustering methods algorithms are mostly of two types 1) hierarchical methods and 2) partitioning methods  
58 (non hierarchical) [14,15,16]. The hierarchical algorithms for clustering represent data sets as a cluster tree and  
59 are of two types 1-1) agglomerative [17] 1 -2) divisive hierarchical clustering methods. Partitional clustering  
60 algorithms [17] are of two types, 2-1) iterative 2 -2) single pass methods. K means and its variants etc. are the  
61 popular partitioning methods. The hierarchical clustering algorithms are considered efficient than the remaining  
62 algorithms [18] however due to their inherent complexity they are not applicable to huge document sets.

63 The techniques for determining inter-cluster similarity in classification ??19 20] ex. single link and for  
64 enhancing the value of the clusters where the cluster size differs or fluctuates by a huge factor [17], especially in  
65 case of high performing clustering algorithms have been studied widely in recent years.

66 The widely used document clustering methods are Spectral Clustering, LSI dependent cluster development  
67 and NMF technique based clustering. The Spectral clustering methods [21] are LPI, LSI etc. Latent semantic  
68 indexing (LSI) [22] a feature extraction approach [23] tries to optimize the documents space compared to the  
69 given document and is a widely used linear document indexing method [24]. LSI is inapplicable for processes  
70 with a high range of documents [24] and similarly spectral clustering when used in a large dimensional space the  
71 dimensionality reduction is very costly which limits its usability.

72 The word patterns and phrases based approaches are the traditional strategies where the clustering is dependent  
73 on the documents features such as words, phrases and sequences [25,26]. These methods are of four types, 1-  
74 1) Clustering with Frequent Word Patterns 1-2) Application of Word Clusters in Document Clusters 1-3) Co-  
75 clustering Words and Documents, Co-clustering with graph partitioning and Information-Theoretic Co -clustering  
76 1-4) Clustering based on Frequent Phrases. The technique VSM is used in almost all the document clustering  
77 methods used nowadays [27]. The vector space model is a data model for representing the terms related to the  
78 words in a document as a feature vector.

79 The features based clustering approaches are of two types 2-1) Feature Extraction 2-2) Feature Selection.

80 The Feature Extraction approaches are based on the algorithm of two types i) linear and ii) nonlinear  
81 techniques. The models of linear type algorithms are unsupervised PCA, OCA, MMC etc. The examples of non  
82 linear algorithms are LLE, Laplacian Eigenmaps, and ISOMAP etc. The linear methods show better operational  
83 performance in contrast to nonlinear approaches, however underperform in the clustering of huge and complicated  
84 data of the internet. The feature extraction technique finds applications in the fields of IR based on human  
85 language learning ability, comparing reviewed and submitted papers, of various languages or networks and filter  
86 of data. Feature selection algorithms are of two types, 2-2-1) Feature Ranking that is metric based and 2-2-2)  
87 Subset Selection from the possible features. The feature selection algorithms are of two categories, i) supervised  
88 and ii) unsupervised. The supervised feature selection algorithms are the most researched as well as used and  
89 they are IG, CHI, and MI. The unsupervised methods that are most popular are, i) DF-based selection dependent  
90 on term strength and ranking dependent on entropy or term contribution, ii) LSI-based method and iii) NMF  
91 based method. These techniques of unsupervised approach such as, decision trees, statistics, NLP and ML are  
92 being used in BI or analytics, in neural networks for developing AI or bio neural networks, for developing systems  
93 of AI that are rule based for intelligent content development, database development, information retrieval and  
94 automatic grouping of web documents with Enterprise Search engines or open source software's in web mining  
95 or text mining.

96 The strategies of feature selection used mostly are i) wrapper, ii) filter and iii) embedded methods [28] however  
97 a study [29] has shown, the methods of supervised feature selection dependent on algorithms using the filter metric  
98 IG, are most efficient over others techniques.

### 99 3 III.

## 100 4 Contemporary Affirmation of the

101 Recent Literature

102 An approach of bisecting k-means algorithm proposed by Steinbach, M, Karypis, G, & Kumar, V [14] breaks  
103 up a large cluster into small clusters repetitively to generate k numbers of clusters of huge similarity for filtering  
104 the clusters and collecting similar texts based on the method.

---

105 A technique called CCA [30] widely used in the emerging technologies of ML etc applies correlation for  
106 measuring the similar features in a document. However, CCA has its own limitations in clustering.

## 107 5 C

108 An approach of spectral clustering based on graph partitioning strategy called LPI [31] proposed however fails  
109 in feature selection and comprises of the existing problems of distance based clustering documents.

110 An approach for document clustering called Frequent Term based Clustering or HFTC [32] is a topic of  
111 extensive research. However it is not scalable for huge data or of documents.

112 A technique known as Hierarchical Document Clustering using Frequent itemsets (FIHC) approach proposed  
113 by Fung, B., Wang, K., Ester, M, is discussed in [33]. The strategy of FIHC though performs better than HFTC  
114 underperforms in clustering efficiency when compared to existing approaches such as UPGMA and Bisecting  
115 K-means.

116 The TDC algorithm technique based on closed frequent itemsets for clustering is proposed by Yu, H., Searsmit, D., Li, X., Han, J [34]. The algorithm performs better compared to HFTC and FIHC however the use of closed  
117 itemsets makes it avoidable.

118 A strategy of Hierarchical Clustering using Closed Interesting Itemsets, referred to as HCCI proposed by Malik, H.H., Kender, J.R [35], is the best clustering method available. However the technique may cause information  
119 loss.

120 An approach based on PSSM histogram by Gad and Kamel [36] combines the text semantic with the process  
121 of incremental clustering and measures the similarity of the documents for adjusting the insertion order of the  
122 documents in the cluster for quality.

123 An improved incremental clustering technique for an efficient clustering algorithm proposed by Gavin and Yue  
124 [37] improves categorization of web data incrementally. The method based on cluster specific multiple information  
125 anew document is assigned to a cluster.

126 An approach for improving text clustering mining by Shehata, S, Fakhri, K, & Mohamed S, S. [38] outperforms  
127 the existing techniques such as HAC, k-NN etc.

128 A progressive clustering algorithm by Liu, Y, Ouyang, Y, Sheng, H, & Xiong, Z. ( ??008) [39] based on Cluster  
129 Average Similarity Area determines the cluster coherence and progressively assigns the new data items to the  
130 clusters.

131 A technique for enhancing the clustering functionality based on the partial disambiguation of words by means  
132 of their PoS [40] is recommended by the developers as the approach finds the inefficiency of considering synonyms  
133 and hypermy my for selecting the right sense of the word disambiguated solely by PoS tags.

134 The CFWS technique proposed by Y. LI, and S.M. Chung, enhances the capability to process the document,  
135 considering the word sequences apart from the words [41].

136 The technique of non linear representation of the data by J.B. Tenenbaum, V. de Silva, and J.C. Langford  
137 [42] keeps specific local data simultaneously based on the optimization factors however is associated with high  
138 complexity.

139 A study of the approaches for reducing the complexity of feature extraction based on a new technique called  
140 approximation algorithm [43], [44], [45] is found to be good.

141 A software for automatically retrieving information from websites by Zamir O Etzioni [46] is designed for  
142 websites comprising of vast amount of data

143 The approach of integrating clustering and feature selection for text clustering based on the semantic relation  
144 of the text documents with ontology was proposed by Thangamani.M and P.Thangaraj in [47]. The approach  
145 minimizes dimensionality and improves feature selection.

146 The clustering technique, for finding the clustering quality based on WordNet [48] phrasal noun and semantic  
147 relationships [49] shows better performance with hyperny my based strategy compared to other noun phrases.

148 A system for determining the ontology related semantic relations of the term or word and associated weight  
149 measure is given by Prof. K. Raja, C. Prakash Narayanan [9]. However the technique has dimensionality and  
150 other problems.

151 A description of the task of Ontology based automatic categorizing of web documents [50] and the scope  
152 of Ontology in improving the current machine learning and IR approaches is given by Andreas Hotho. The  
153 integration of ontology's for combining various information types of multiple resources by Young-Woo et al. in  
154 the paper [51].

155 The process of using domain specific ontology's for enhancing performance of text classification where text  
156 learning and IR are used to generate ontology's with minimum user interaction is given in [52,53].

157 The methods utilizing Wikipedia ontology for improving primarily the document depiction and cluster quality  
158 by Gabrilovich and Markovitch [54] and a further extension provided a structure based on the Wikipedia guidelines  
159 and groups [55,56]. The Wikipedia ontology is most relevant as it is applicable to a large cross section of domains  
160 and also restructured on a regular basis.

161 A technique for feature selection in text clustering based on supervised feature selection on the intermediary  
162 clustering outcomes by Xu, J. Xu, B [57] generates a efficient subset for classification. The suggested techniques  
163 performance is efficient compared to manual process.

164 A technique of feature selection dependent on the ACO algorithm by M. Janaki Meena,K.R.

## 6 Year 2015

167 Global Journal of Computer Science and Technology Volume XV Issue II Version I ( ) C Chandran, J. Mary  
168 Binda," [58] is a unique method. Comparative tests of the approach with existing chisquare and CHIR techniques  
169 shows the proposed approach achieves better performance in FS.

171 An entropy based FS approach i.e. a filter solution [59] tested with various data types that reduces  
172 dimensionality and is efficient in finding the subset of major features.

173 A feature co-selection method called MFCC (multi type feature co-selection), proposed by Shen huang, Zheng  
174 Chen, Yong Yu, and Wei-Ying main [60] shows enhanced clusters performance of web documents based on the  
175 outcomes of intermediate clustering.

176 A method to remodel the matrix of data similarity as a bi-stochastic matrix prior to executing algorithms by  
177 F. Wang, P. Li, and A. C. K Aonig showed better clustering performance [61].

178 The techniques of document clustering that are term based for clustering in dynamic environments, is given in  
179 [11] by Wang, X, Tang, J, & Liu, H, synonyms and hypernyms by Bharathi and Vengatesan [62], Synonyms and  
180 Hyponyms, Nadig, R, Ramanand, J, & Bhattacharyya, P in [12]. These approaches are however not applicable  
181 to technically similar documents.

182 A document clustering approach [63] dependent on phrases and the STC technique by O. Zamir, O. Etzioni,  
183 O. Madanim, and R.M. Karp builds the clusters on the common documents suffixes. The method though efficient  
184 in cluster quality however is associated with high amount of term redundancy.

185 A study of the TF-IDF method of clustering [64], term frequency dependent algorithms [65] and a review of  
186 clustering algorithms [66] showed that majority of clustering approaches are TF-IDF based, however associated  
187 with several problems.

188 The NMF (Nonnegative Matrix Factorization) technique in text classification [67], improved clustering  
189 performance compared to the existing approaches [68], relationship study of NMF techniques with earlier  
190 clustering techniques [69], [70] [71]. A review of established techniques of NMF such as multiplicative updates  
191 [72], projected gradients [73] though efficient however are associated with the problems of memory for huge  
192 datasets streamed and not disk based [74]. To overcome these problems, approaches such as random projections  
193 [61,75] and sketch/sampling algorithms [76] have been proposed. An NMF based technique by Li and Zhu in 2011  
194 [77] for research specific documents minimizes high dimensionality, finds relevant topics for clustering and shows  
195 performance efficiency in classification comparatively. A study of the online algorithm based on Nonnegative  
196 Matrix Factorization [78], a NMF based method that uses features based on weights and similar cluster property  
197 by Sun Park, Dong Un An, Choi Im Cheon [79] performs comparatively more efficiently than the remaining NMF  
198 based strategies.

199 IV.

## 7 Conclusion

200 In this paper we analyzed several techniques developed for clustering documents with their applications and  
201 relevance in terms of today's requirements. The task of developing perfect strategies for classification of varied  
202 forms and types of documents for a near optimal solution or finding accurate ways of assessing the quality of  
203 the performed clustering though is impossible and is increasing in its complex nature, the field today deals  
204 with extraordinary tasks like granular taxonomies generation, sentiment analysis and document summarization  
205 for generating reliable and relevant insights applicable to several fields. In conclusion we can say document  
206 clustering is going to be widely studied and will find relevance in a number of newer areas. <sup>1</sup> <sup>2</sup>

---

<sup>1</sup>© 2015 Global Journals Inc. (US)

<sup>2</sup>© 2015 Global Journals Inc. (US) 1



Figure 1: Year 2015 Global

## **7 CONCLUSION**

---

- 
- 208 [Chizi] , Barak Chizi . Tel-Aviv University, Israel.
- 209 [Rokach] , Lior Rokach . Ben-Gurion University, Israel.
- 210 [Lsa @ Cu and Boulder ()] , Lsa @ Cu , Boulder . <http://lsa.colorado.edu/> 2010.
- 211 [ Information Science and Applications ICISA ()] , *Information Science and Applications ICISA* 2010 6 February 212 2010.
- 213 [Yang and Pedersen ()] 'A Comparative Study on Feature Selection in Text Categorization'. Y Yang , J O 214 Pedersen . *Proc. 14th Int'l Conf. Machine Learning*, (14th Int'l Conf. Machine Learning) 1997. p. .
- 215 [Tenenbaum et al. ()] 'A Global Geometric Framework for Nonlinear Dimensionality Reduction'. J B Tenenbaum 216 , V Silva , J C Langford . *Science* 2009. 290 p. .
- 217 [Xu et al. ()] 'A new feature selection method for text clustering'. J Xu , B Xu , W Zhang , Z Cui , W Zhang . 218 *wuhan university journal of natural sciences* 2007. 12 p. .
- 219 [pp. John Wang (ed.) ()] *a survey of feature selection techniques*, 10.4018/978-1-60566-010-3.ch289. pp. John 220 Wang (ed.) 2009. 13 p. 9781605660103. Oded Maimon (Tel-Aviv University, Israel ; Montclair State University, 221 USA
- 222 [Carpinetto et al. ()] 'A survey of web clustering engines'. C Carpineto , S Osi'nski , G Romano , D Weiss . *ACM 223 Comput. Surv* 2009. 41 (3) p. .
- 224 [Prathima and Supreethi ()] 'A survey paper on concept based text clustering'. Y Prathima , K P Supreethi . 225 *International Journal of Research in IT & Management* 2011. 1 (3) p. .
- 226 [Lee and Seung ()] 'Algorithms for nonnegative matrix factorization'. D D Lee , H S Seung . *Advances in Neural 227 Information Processing System (NIPS)*, 2000. p. .
- 228 [Shehata et al. ()] 'An efficient concept-based mining model for enhancing text clustering'. S Shehata , K Fakhri 229 , S Mohamed , S . *IEEE Transactions On Knowledge And Data Engineering* 2010. 22 (10) p. .
- 230 [Iu et al. ()] 'An Incremental Algorithm for Clustering Search Results'. Y Iu , Y Ouyang , H Sheng , Z Xiong . 231 *IEEE International Conference on Signal Image Technology and Internet Based Systems*, 2008. p. .
- 232 [Nadig et al. ()] 'Automatic evaluation of Word Net synonyms and hypermy my India'. R Nadig , J Ramanand , 233 P Bhattacharyya . *Proceedings of ICON-2008, 6th International Conference on Natural Language Processing*, 234 (ICON-2008, 6th International Conference on Natural Language Processing) 2008.
- 235 [Weng et al. (2003)] 'Candid Covariance-Free Incremental Principal Component Analysis'. J Weng , Y Zhang , 236 W.-S Hwang . *IEEE Trans. Pattern Analysis and Machine Intelligence* Aug.2003. 25 (8) p. .
- 237 [Hardoon et al. ()] 'Canonical Correlation Analysis: An Overview with Application to Learning Methods'. D R 238 Hardoon , S R Szedmak , J R Shawetaylor . *J. Neural Computation* 2004. 16 (12) p. .
- 239 [Liu and Croft ()] 'Cluster-based retrieval using language models'. X Liu , W B Croft . *Proceedings of 240 the 27th annual International ACM SIGIR Conference on Research and Development in Information 241 Retrieval (SIGIR)*, (the 27th annual International ACM SIGIR Conference on Research and Development in 242 Information Retrieval (SIGIR)) 2004. p. .
- 243 [Prof et al.] 'Clustering Technique with Feature Selection for Text Documents'. . K Prof , C Prakash Raja , 244 Narayanan . *Proceedings of the Int. Conf. on*, (the Int. Conf. on)
- 245 [Gabrilovich and Markovitch ()] 'Computing Semantic Relatedness Using Wikipedia-based Explicit Semantic 246 Analysis'. E Gabrilovich , S Markovitch . *Proc. of The 20th Intl. Joint Conf.on Artificial Intelligence*, (of 247 The 20th Intl. Joint Conf.on Artificial Intelligence) 2007.
- 248 [Dhillon and Modha ()] 'Concept decompositions for large sparse text data using clustering'. I S Dhillon , D S 249 Modha . *Machine Learning*, 2001. 42 p. .
- 250 [Ding et al. ()] 'Convex and seminonnegative matrix factorizations'. C Ding , T Li , M I Jordan . *IEEE 251 Transactions on Pattern Analysis and Machine Intelligence* 2010.
- 252 [Khalilian and Mustapha ()] 'Data Stream Clustering: Challenges and Issues'. M Khalilian , & N Mustapha . 253 *Proceedings of the International Multiconference of Engineers and Computer Scientists IMECS 2010*, (the 254 International Multiconference of Engineers and Computer Scientists IMECS 2010Hong Kong) 2010. p. .
- 255 [Cao et al. ()] 'Detect and track latent factors with online nonnegative matrix factorization'. B Cao , D Shen 256 , J Sun , X Wang , Q Yang , Z Chen . *Proc. International Joint Conference on Artificial Intelligence*, 257 (International Joint Conference on Artificial Intelligence) 2007. p. .
- 258 [Silva et al. ()] 'Document clustering and cluster topic extraction in multilingual corpora'. J Silva , J Mexia , A 259 Coelho , G Lopes . *Proceedings of the 1st IEEE International Conference on Data Mining (ICDM)*, (the 1st 260 IEEE International Conference on Data Mining (ICDM)) 2001. p. .
- 261 [Xu and Gong (2004)] 'Document Clustering by Concept Factorization'. W Xu , Y Gong . *Proc. Int'l Conf. 262 Research and Development in Information Retrieval*, (Int'l Conf. Research and Development in Information 263 Retrieval) July 2004. p. .

## 7 CONCLUSION

---

- 264 [Li and Zhu ()] 'Document clustering in research literature based on NMF and testor theory'. F Li , Q Zhu .  
265 *Journal of Software* 2011. 6 (1) p. .
- 266 [Park et al. ()] 'Document Clustering Method Using Weighted Semantic Features and Cluster Similarity'. Sun  
267 Park , Dong Un An , Choi Im Cheon . *Third IEEE International Conference on Digital Game and Intelligent*  
268 *Toy Enhanced Learning*, 2010. 2010. p. . (digitel)
- 269 [Shahnaz et al. ()] 'Document clustering using nonnegative matrix factorization'. F Shahnaz , M W Berry , V P  
270 Pauca , R J Plemmons . *Information Processing and Management* 2006. 42 (2) p. .
- 271 [Wang and Li ()] 'efficient non-negative matrix factorization with random projections'. F Wang , P Li .  
272 *Proceedings of the 10th SIAM International Conference on Data Mining (SDM)*, (the 10th SIAM International  
273 Conference on Data Mining (SDM)) 2010. p. .
- 274 [Hung and Xiaotie (2008)] 'Efficient Phrase-Based Document Similarity for Clustering'. C Hung , D Xiaotie .  
275 *IEEE Transaction on Knowledge and Data Engineering* September. 2008. 20 p. .
- 276 [Zhong ()] 'Efficient streaming text clustering'. S Zhong . *Neural Networks* 2005. 18 (5-6) p. .
- 277 [Gavin and Yue ()] 'Enhancing an incremental clustering algorithm for Web page collections'. S Gavin , X Yue .  
278 *ACM International Joint Conferences on Web Intelligence and Intelligent Agent Technologies*, 2009. p. .
- 279 [Hu et al. ()] 'Enhancing Text Clustering by Leveraging Wikipedia Semantics'. J Hu , L Fang , Y Cao . *Proc. of*  
280 *31st Annual Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval*, (of 31st Annual  
281 Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval) 2008.
- 282 [Everitt et al. ()] Brian S Everitt , Sabine Landau , Morven Leese . *Cluster Analysis*, 2001. Oxford University  
283 Press. (fourth edition)
- 284 [Zheng ()] 'Exploiting noun phrases and semantic relationships for text document clustering'. Kang Zheng , Kim  
285 . *Information Science* 2009. 179 p. .
- 286 [Hu et al. ()] 'Exploiting Wikipedia as External Knowledge for Document Clustering'. X Hu , X Zhang , C Lu .  
287 *Proc. of Knowledge Discovery and Data Mining*, (of Knowledge Discovery and Data Mining) 2009.
- 288 [Dash et al.] 'Feature Selection for Clustering -A Filter Solution'. Manoranjan Dash , Kiseok Choi , Peter  
289 Scheuermann , Huan Liu . ICDM'02)0-7695-1754-4/02 © 2002 IEEE. *Proceedings of the 2002 IEEE*  
290 *International Conference on Data Mining*, (the 2002 IEEE International Conference on Data Mining)
- 291 [Seo et al. (2004)] *Feature Selections for Extracting Semantically Rich Word for Ontology Learning*, Young-Woo  
292 Seo , Anupriya Ankolekar , Katia Sycara . CMU-RI-TR-04-18. March 2004.
- 293 [Beil et al. ()] 'Frequent Term-based Text Clustering'. F Beil , M Ester , X Xu . *Proc. of Intl. Conf. on Knowledge*  
294 *Discovery and Data Mining*, (of Intl. Conf. on Knowledge Discovery and Data Mining) 2002.
- 295 [Das et al. ()] 'Google news personalization: Scalable online collaborative filtering'. A Das , M Datar , A Garg  
296 , S Rajaram . *Proceedings of the 16th International Conference on World Wide Web (WWW)*, (the 16th  
297 International Conference on World Wide Web (WWW)) 2007. p. .
- 298 [Zamir and Etzioni ()] 'Grouper: A Dynamic Clustering Interface to Web Search Results'. O Zamir , O Etzioni  
299 . *Computer Networks* 1999. 31 p. .
- 300 [Wang et al. ()] 'H Document clustering via matrix representation'. X Wang , J Tang , Liu . *11th IEEE*  
301 *International Conference on DataMiningICDM2011*, 2011. p. .
- 302 [Fung et al. ()] *Hierarchical Document Clustering Using Frequent Itemsets*, B C M Fung , K Wan , M Ester .  
303 2003. p. 3.
- 304 [Fung et al. ()] 'Hierarchical document clustering using frequent Itemsets'. B C M Fung , K Wang , M Ester .  
305 *Proceedings of SIAM International Conference on Data Mining*, (SIAM International Conference on Data  
306 Mining) 2003.
- 307 [Malik and Kender ()] 'High Quality, Efficient Hierarchical Document Clustering Using Closed Interesting  
308 Itemsets'. H H Malik , J R Kender . *Proc. of IEEE Intl. Conf. on Data Mining*, (of IEEE Intl. Conf. on  
309 Data Mining) 2006.
- 310 [Yan et al. ()] 'IMMC: Incremental Maximum, Marginal Criterion'. J Yan , B S Zhang , Z Yan , W Chen , Q  
311 Fan , W Y Yang , Q Ma , Cheng . *Proc. 10th ACM SIGKDD*, (10th ACM SIGKDD) 2004. p. .
- 312 [Bharathi and Vengatesan ()] 'Improving information retrieval using document clusters and semantic synonym  
313 extraction'. G Bharathi , D Vengatesan . *Journal of Theoretical and Applied Information Technology* 2012.  
314 36 (2) p. .
- 315 [Gad and Kamel ()] 'Incremental clustering algorithm based on phrase-semantic similarity histogram'. W K Gad  
316 , M S Kamel . *Proceedings of the Ninth International Conference on Machine Learning and Cybernetics*, (the  
317 Ninth International Conference on Machine Learning and Cybernetics) 2010. 11 p. .
- 318 [Deerwester et al. ()] 'Indexing by Latent Semantic Analysis'. S C Deerwester , S T Dumais , T K Landauer , G  
319 W Furnas , R A Harshman . *J. Am.Soc. Information Science* 1990. 41 (6) p. .

- 320 [Thangamani and Thangaraj] 'integrated clustering and feature selection scheme fo textdocuments'. P M  
321 Thangamani , Thangaraj . 10.3844/jcssp.2010.536.54. DOL:10.3 844/jcssp.2010.536.541. <http://www.thescipub.com/abstract/10.3844/jcssp.2010.536.54> *J.Comput.Sci* 6 p. 536.
- 323 [Meena et al.] 'integrating swarm intelligence and statistical data forfeature selection in text categorization'. M  
324 Meena , K R Chandran , J Mary Brinda . ©2010 *International Journal of Computer Applications* 1 (11) p. .
- 325 [Kumar and Srinathan ()] N Kumar , K Srinathan . *A New Approach for Clustering Variable Length Docu-*  
326 *ments(Proceedings of the Advanced computing Conference*, 2009. IEEE. p. .
- 327 [Dumais ()] 'Latent Semantic Indexing (LSI) and TREC-2'. S T Dumais . *Proc.Second Text Retrieval Conf.*  
328 *(TREC)*, (.Second Text Retrieval Conf. (TREC)) 1993. p. .
- 329 [Wang et al. ()] 'Learning a bistochastic data similarity matrix'. F Wang , P Li , A C Käonig . *Proceedings of*  
330 *the 10th IEEE International Conference on Data Mining (ICDM)*, (the 10th IEEE International Conference  
331 on Data Mining (ICDM)) 2010.
- 332 [Lee and Seung ()] 'Learning the parts of objects with nonnegative matrix factorization'. D D Lee , H S Seung .  
333 *Nature* 1999. 401 p. .
- 334 [Cai et al. (2005)] 'Locality Preserving Indexing'. D Cai , X He , J Han . *Document Clustering Using Knowledge*  
335 *and Data Eng* Dec. 2005. 17 (12) p. . (IEEE Trans)
- 336 [Sebastiani (2002)] 'Machine Learning in Automated Text Categorization'. Fabrizio Sebastiani . *ACM Computing*  
337 *Surveys* March 2002. 34 (1) .
- 338 [Shen Huang and Ma (2006)] 'multitype features coselection for web document clustering'. Wei-Ying Shen Huang  
339 , Ma . 1041-4347/06/\$20.00. *ieee transactions on knowledge and data engineering* april 2006. 2006. 18 (4) .  
340 (ieee published by the ieee computer society)
- 341 [Ng et al. ()] 'On Spectral Clustering: Analysis and an Algorithm'. A Y Ng , M Jordan , Y Weiss . *Advances in*  
342 *Neural Information Processing Systems* 2001. MIT Press. 14 p. .
- 343 [Hiraoka and Hamahira ()] 'On Successive Learning Type Algorithm for Linear Discriminant Analysis'. K  
344 Hiraoka , M Hamahira . *IEIC Technical Report* 1999. 99 p. . (in Japanese)
- 345 [Ding et al. ()] 'On the equivalence of nonnegative matrix factorization and spectral clustering'. C Ding , X He  
346 , H D Simon . *Proceedings of the 5th SIAM Int'l Conf. Data Mining (SDM)*, (the 5th SIAM Int'l Conf. Data  
347 Mining (SDM)) 2005. p. .
- 348 [Li et al. ()] 'One sketch for all: Theory and application of conditional random sampling'. P Li , K W Church ,  
349 T Hastie . *Advances in Neural Information Processing System (NIPS)*, 2008. p. .
- 350 [Hotho et al. (2001)] 'Ontologybased text clustering'. A Hotho , S Staab , A Maedche . *Proceedings of the IJCAI-*  
351 *2001 Workshop Text Learning: Beyond Supervision*, (the IJCAI-2001 Workshop Text Learning: Beyond  
352 SupervisionSeattle,USA) August 2001.
- 353 [Lin] 'Projected gradient methods for nonnegative matrix factorization'. C J Lin . *Neural Computation* 19 (10)  
354 p. .
- 355 [Martinez-Morais et al. ()] 'Providing QoS with the Deficit Table Scheduler'. R Martinez-Morais , F J Alfaro-  
356 Cortes , & J L Sanchez . *IEEE Transactions on Parallel and Distributed Systems* 2010. 21 (3) p. .
- 357 [Kotsiantis and Pintelas ()] 'Recent Advances in Clustering: A Brief Survey'. S Kotsiantis , P Pintelas . *WSEAS*  
358 *Trans. Information Science and Applications* 2004. 1 (1) p. .
- 359 [Gaussier and Goutte ()] 'Relation between plsa and nmf and implications'. E Gaussier , C Goutte . *Proceedings*  
360 *of the 28th Annual International ACM SIGIR Conference on Research and Development in Information*  
361 *Retrieval (SIGIR)*, (the 28th Annual International ACM SIGIR Conference on Research and Development in  
362 Information Retrieval (SIGIR)) 2005. p. .
- 363 [Siersdorfer and Sizov (2004)] 'Restrictive Clustering and Metaclustering for Self-Organizing Document Collec-  
364 tions'. S Siersdorfer , S Sizov . *Proc. Int'l Conf. Research and Development in Information Retrieval*, (Int'l  
365 Conf. Research and Development in Information Retrieval) July 2004. p. .
- 366 [Yu et al. ()] 'Scalable Construction of Topic Directory with Nonparametric Closed Termset Mining'. H Yu , D  
367 Searsmith , X Li , J Han . *Proc. of Fourth IEEE Intl. Conf.on Data Mining*, (of Fourth IEEE Intl. Conf.on  
368 Data Mining) 2004.
- 369 [Steinbach et al. ()] M Steinbach , G Karypis , V Kumar . *A comparison of document clustering techniques. KDD*  
370 *Workshop on Text Mining*, 2000.
- 371 [Xu Rui ()] 'Survey of Clustering Algorithms'. Xu Rui . *IEEE Transactions on Neural Networks* 2005. 16 (3) p. .
- 372 [Berkhin ()] *Survey of clustering data mining techniques*, P Berkhin . [http://www.accrue.com/products/rp\\_cluster\\_review.pdf](http://www.accrue.com/products/rp_cluster_review.pdf) 2004.
- 374 [Salton and Buckley ()] 'Term-weighting approaches in automatic text retrieval'. G Salton , C Buckley .  
375 *Information Processing & Management* 1998. 24 (5) p. .

## 7 CONCLUSION

---

- 376 [Aas and Eikvil ()] *Text Categorisation: A Survey*, K Aas , L Eikvil . 941. 1999. Oslo Norway: Norwegian  
377 Computing Center. (Technical Report) ([iteseer.ist.psu.edu/aas99text.html](http://iteseer.ist.psu.edu/aas99text.html))
- 378 [Soon and John ()] 'Text document clustering based on frequent word meaning sequences'. M C Soon , D H John  
379 , Yanjun , L . *Data & Knowledge Engineering* 2008. 64 p. .
- 380 [Li and Chung (2005)] 'Text Document Clustering Based on Frequent Word Sequences'. Y Li , S M Chung .  
381 *Proceedings of the. CIKM*, (the. CIKMBremen, Germany) 2005. 2005. October 31-November 5.
- 382 [Berendt et al. ()] 'Towards semantic web mining'. B Berendt , A Hotho , G Stumme . *Proceedings of International*  
383 *Semantic Web Conference (ISWC)*, (International Semantic Web Conference (ISWC)) 2002. p. .
- 384 [Hotho (2005)] *Using Ontologies to Improve the Text Custering and Classification Task*, Andreas Hotho . January  
385 14, 2005. Knowledge and Data Engineering Group, University of Kassel
- 386 [Van Rijsbergen ()] Van Rijsbergen . *London: Butterworth*, 1989. (Secondth ed.)
- 387 [Zamir et al.] 'Web Document Clustering, A Feasibility Demonstration'. O Zamir , K Development , Mugun-  
388 thadevi . *Proceedings of the 21st International ACM SIGIR Conference on Research*, (the 21st International  
389 ACM SIGIR Conference on Research) IJCSE.
- 390 [Miller ()] 'Wordnet: A lexical database for English'. G Miller . *CACM* 1995. 38 (11) p. .
- 391 [Sedding and Kazakov ()] *Wordnetbased text document clustering*, J Sedding , D Kazakov . 2004. p. . (3rd  
392 Workshop on Robust Methods in Analysis of Natural Language Data)