

1 Feature Extraction and Duplicate Detection for Text Mining: A 2 Survey

3 Ramya R S¹ and Venugopal K R²

4 ¹ University Visvesvaraya College of Engineering, UVCE

5 *Received: 12 December 2015 Accepted: 31 December 2015 Published: 15 January 2016*

6

7 **Abstract**

8 Text mining, also known as Intelligent Text Analysis is an important research area. It is very
9 difficult to focus on the most appropriate information due to the high dimensionality of data.
10 Feature Extraction is one of the important techniques in data reduction to discover the most
11 important features. Processing massive amount of data stored in a unstructured form is a
12 challenging task. Several pre-processing methods and algorithms are needed to extract useful
13 features from huge amount of data. The survey covers different text summarization, classi-
14 fication, clustering methods to discover useful features and also discovering query facets which
15 are multiple groups of words or phrases that explain and summarize the content covered by a
16 query thereby reducing time taken by the user. Dealing with collection of text documents, it is
17 also very important to filter out duplicate data. Once duplicates are deleted, it is
18 recommended to replace the removed duplicates. Hence we also review the literature on
19 duplicate detection and data fusion (remove and replace duplicates). The survey provides
20 existing text mining techniques to extract relevant features, detect duplicates and to replace
21 the duplicate data to get fine grained knowledge to the user.

22

23 **Index terms**— text feature extraction, text mining, query search, text classification.

24 **1 I. Introduction**

25 Society is increasingly becoming more digitized and as a result organisations are producing and storing vast amount
26 of data. Managing and gaining insights from the produced data is a challenge and key to competitive advantage.
27 Web based social applications like people connecting websites results in huge amount of unstructured text data.
28 These huge data contains a lot of useful information. People hardly bother about the correctness of grammar
29 while forming a sentence that may lead to lexical syntactical and semantic ambiguities. The ability of finding
30 patterns from unstructured form of text data is a difficult task.

31 Data mining aims to discover previously unknown interrelations among apparently unrelated attributes of
32 data sets by applying methods from several areas including machine learning, database systems, and statistics.
33 Many researches have emphasized on different branches of data mining such as opinion mining, web mining, text
34 mining. Text mining is one of the most important strategy involved in the phenomenon of knowledge discovery. It
35 is a technique of selecting previously unknown, hidden, understandable, interesting knowledge or patterns which
36 are not structured. The prime objective of text mining is to diminish the effort made by the users to obtain
37 appropriate information from the collection of text sources [1].

38 Thus, our focus is on methods that extract useful patterns from texts in order to categorize or structure text
39 collections. Generally, around 80 percent of company's information is saved in text documents. Hence text mining
40 has a higher economic value than data mining. Current research in the area of text mining tackles problems of
41 text representation, classification, clustering, information extraction or the search for and modelling of hidden
42 patterns. Selection of characteristics, influence of domain knowledge and domain-specific procedures play an
43 important role.

3 C) TRAFFIC BASED EVENT IN TEXT MINING

44 The text documents contain large scale terms, patterns and duplicate lists. Queries submitted by the user on
45 web search are usually listed on top retrieved documents. Finding the best query facet and how to effectively use
46 large scale patterns remains a hard problem in text mining. However, the traditional feature selection methods
47 are not effective for selecting text features for solving relevance issue. These issues suggests that we need an
48 efficient and effective methods to mine fine grained knowledge from the huge amount of text documents and
49 helps the user to get information quickly about a user query without browsing tens of pages.

50 The paper provides a review of an innovative techniques for extracting and classifying terms and patterns.
51 A user query is usually presented in list styles and repeated many times among top retrieved documents. To
52 aggregate frequent lists within the top search results, various navigational techniques have been presented to
53 mine query facets.

54 The Organisation of the paper is as follows: Section 1 introduces a detailed overview of text mining frameworks,
55 application and benefits of text mining. Sections 2 and 3 reviews feature selection, feature extraction and
56 techniques of pattern extraction. Section 4 discusses various text classification and clustering algorithms in text
57 mining. Sections 5 and 6 introduce a detailed overview of discovering facets and fine grained knowledge. Section
58 7 reviews the duplicate detection in text documents. Section 8 contains the conclusions.

59 1 Year 2016

60 2 () C a) Text Mining Models

61 Text mining tasks consists of three steps: text preprocessing, text mining operations, text post processing. Text
62 preprocessing includes data selection, Many approaches [2] have been concerned of obtaining structured datasets
63 called intermediate forms, on which techniques of data mining [3] When documents contains terms with same
64 frequency. Two terms can be meaningful while the other term may be irrelevant. Inorder to discover the semantic
65 of text, the mining model is introduced. Figure 2 represents a new mining model based on concepts. The model
66 is proposed to analyse terms in a sentence from documents. The model contains group of concept analysis, they
67 are sentence based concept analysis, document based concept analysis and corpus based similarity measure [4].
68 Similarity measure concept based analysis calculates the similarity between documents. The model effectively
69 and efficiently finds Benefits of text mining are better collection development to resolve user needs, information
70 retrieval, to resolve usability and system performance, data base evaluation, hypothesis development. Information
71 professionals(IP) [8] are always in forefront for emerging technologies. Inorder to make their product and service
72 better and more efficient, usually libraries and information use these IP. The trained information professionals
73 manage both technical and semantic infrastructures which is very important in text mining. IP also manages
74 content selection and formulation of search techniques and algorithms.

75 Akilan et al., [9] pesented the challenges and future directions in text mining. It is mandatory to function
76 semantic analysis to capture objects relationship in the documents. Semantic analysis is computationally
77 expensive and operates on few words per second as text mining consists of significant language component.
78 An effective text refining method has to be developed to process multilingual text document. Trained knowledge
79 specialists are neccessary to deal with products and application of current text mining tools. Automated mining
80 operations is required which can be used by technical users. Domain Knowledge plays an important role in both
81 at text refining stage and knowledge distillation and hence helps in improving the efficiency of text mining.

82 Sanchez et al., [10] presented Text knowledge mining (TKM) based deductive inference that is usually targeted
83 on the feasible subset of texts which usually search for contradictions. The procedure obtains new knowledge
84 making a union of intermediate forms of texts from accurate knowledge expressed in the text.

85 Dai et al., [11] introduced competitive intelligence analysis methods FFA (Five Faces Frame work) and SWOT
86 with text mining technologies. The knowledge is extracted from the raw data while performing transforming
87 process that enables the business enterprises to take decisions more reliably and easily. Mining Environment for
88 Decisions (MinEDec) system is not evaluated in real business environments.

89 Hu et al., [12] presented a interesting task of automatically generating presentation slides for academic papers.
90 Using a support vector regression method, importance scores of sentences in the academic papers is provided.
91 Another method called Integer Linear Programming is used to generate well structured slides. The method
92 provides the researchers to prepare draft slides which helps in final slides used for presentation. The approach
93 does not focus on tables, graphs and figures in the academic papers.

94 3 c) Traffic based Event in Text Mining

95 Andrea et al., [13] [14] have proposed a realtime monitoring system for traffic event detection that fetches tweets,
96 classifies and then notifies the users about traffic events. Tweets are fetched using some text mining techniques.
97 It provides the class labels to each tweet that are related to a traffic event. If the event is due to an external
98 cause, such as football match, procession and manifestation, the system also discriminate the traffic event. Final
99 result shows it is capable of detecting traffic event but traffic condition notifications in real-time is not captured.

100 An efficient and scalable system from a set of microblogs/ tweets has been proposed to detect Events from
101 Tweets (ET) [15] by considering their textual and temporal components. The main goal of proposed ET system
102 is the efficient use of content similarity and appearance similarity among keywords and to cluster the related
103 keywords. Hierarchical clustering technique is used to determine the events, which is based on common co-

104 occurring features of keywords [16]. ET is evaluated on two different datasets from two different domains. The
105 results show that it is possible to detect events of relevance efficiently. The use of semantic knowledge base like
106 Yago is not incorporated.

107 Schulz et al., [17] proposed a machine learning algorithm which includes text classification and increasing the
108 semantics of the microblog. It identifies the small scale incidents with high accuracy. It also precisely localizes
109 microblogs in space and time which enables it to detect incidents in real time. The algorithm will not only give
110 us information about the incident and in addition give us valuable information on previous unknown information
111 about the incidents. It does not consider NLP techniques and large data.

112 ITS (Intelligent Transportation Systems) [18] recognizes the traffic panels and dig in information contained
113 on them. Firstly, it applies white and blue color segmentation and then at some point of interest it derives
114 descriptors. These images that can now be considered as sack of words and classified using Na?"ve Bayes or SVM
115 (state vector method). The kind of categorization where the images are classified based on visual appearance is
116 new for traffic panel detection and it does not recognize multiframe integration.

117 Text may be loosely organized without complete information in the documents and may also contain omitted
118 information. The text has to be scanned attentively to determine the problems. If it is not scanned and scrutinised
119 properly then it leads to poor accuracy on unstructured data and hence preprocessing is necessary.

120 Preprocessing guarantees successful implementation of text analysis, but may spend substantial processing
121 time. Text processing can be done in two basic methods. a) Feature Selection b) Feature Extraction.

122 Research in numerous fields like machine learning, data mining, computer vision, statistics and linked fields
123 has led to diversity of feature selection approaches in supervised and unsupervised surroundings.

124 Feature Selection (FS) has an important role in data mining in categorization of text. The centralized idea of
125 feature selection is the reduction of the dimension of the feature set by determining the features appropriately
126 which enhances the efficiency and the performance. FS is a search process and categorized into forward search
127 and backward search.

128 Mehdi et al., [19] [20] executed a innovative feature selection algorithm based on Ant Colony Optimization
129 (ACO).

130 Without any prior knowledge of features, a minimal feature subset is determined by applying ACO [21]. The
131 approach uses simple nearest neighbor classifier to show the effectiveness of ACO algorithm by reducing the
132 computational cost and it outperforms information gain and chi methods. Complex classifiers and different kinds
133 of datasets are not incorporated. Combining feature selection algorithm with other population-based feature
134 selection algorithms are not considered.

135 Gasca et al., [22] proposed feature selection method based on Multilayer Perceptron (MLP). Under certain
136 objective functions the approach determines and also corrects proper set of irrelevant set of attributes. It further
137 computes the relative contributions for individual attribute in reference to the units that are to be output. For
138 each output unit, contribution are sorted in the descending order. An objective function called prominence is
139 computed for each attribute. Selecting the features from large document faces problem in unsupervised learning
140 because of unnamed class labels.

141 Sivagaminathan et al., [23] [24] proposed a fixed size subset, an hybrid approach to solve feature subset
142 selection problem in neural network pattern classifier. It considers both the individual performance and subset
143 performance. Features are selected using the pheromone trail and value of heuristic by state transition rules.
144 After selecting the feature, the global updating rule takes place to increment the features, which ultimately gives
145 better classification performance without increase in the overall computational cost. Ogura et al., [28] proposed
146 an approach to reduce a feature dimension space which calculates the probability distribution for each term
147 that deviates from poissos. These deviations from poissos are non significant for the documents that does not
148 belong to category. Three measures are employed as a benchmark and by using two classifiers SVM and K-NN
149 gives better performance than other conventional classifiers. Gini index proved to be better than chisquare, IG in
150 terms of macro, micro average values of F1. These measures do not utilize the number of times the term occurs
151 in a document. The computational complexity could not be suppressed for other typical measures such as
152 information gain and CHI.

153 4 Global Journal of Computer Science and Technology

154 Volume XVI Issue V Version I4 Year 2016 () C

155 Feature selection is measured based on words term and document frequency. Azam et al., [29] observes these
156 frequencies for measuring FS. The metrics of Discriminative Power Measure (DPM) and GINI index (GINI) are
157 incorporated and the term frequency based metric is useful for small feature set. The most important features
158 returned by DPM and GINI tend to discover most of the available information at a faster rate, i.e. against lower
159 number of features. The DPM and GINI are comparatively slower in covering document frequency information.

160 Yan et al., [30] presented a graph embedded framework for dimensionality reduction. The framework is also
161 used as a tool and unifies many feature extraction methods. Feature is selected based on spectral graph theory
162 and proposed framework unifys both supervised and unsupervised feature selection.

163 Zhao et al., [31] developed a framework for preserving feature selection similarity to handle redundant feature.
164 A combined optimization formulation of sparse multiple output regression formulation is used for selecting

6 2) TEXT

165 similarity preserving features. The framework do not address existing kernel, metric learning methods and
166 semi-supervise feature selection methods.

167 5 1) Feature Selection based Graph Reconstruction:A

168 Major task in efficient data mining is Feature selection. Feature selection has a significant challenge in small
169 labeled-sample problem. If data is unlabeled then it is large. If the label of data is extremely tiny, then
170 supervised feature selection algorithms fail for want of sufficient information. Zhao et al., [32] introduced graph
171 regularized data construction to overcome the problems in feature selection. The approach achieves higher
172 clustering performance in both unsupervised and supervised feature selection.

173 Linked social media crops enormous amount of unlabeled data. In the prevailing system, selecting features for
174 unlabeled data is a difficult task due to the lack of label information. Tang et al., [33] proposed an unsupervised
175 feature selection framework, LUFS(Linked Unsupervised Feature Selection), for related social media data to
176 surpass the problem. The design builds a pseudo-class labels through social dimension extraction and spectral
177 analysis. LUFS efficiently exploits association information but does not exploit link information. Computer
178 vision and pattern recognition problems are the two main problems which have inherent manifold structure. A
179 laplacian regularizer is included to smoothen the clustering process along with the scale factor. In text mining
180 applications, several existing systems incorporate a NLP-based techniques which parse the text and promote the
181 usage patterns that is used for mining and examination of the parse trees that are trivial and complex.

182 Mousavi et al., [34] have formulated a weighted graph depiction of text, called Text Graphs that further
183 captures grammar which serve as semantic dealings between words that are in textual terms. The text based
184 graphs incorporates such a framework called SemScape that creates parse trees for each sentence and uses two
185 step pattern based procedure for facilitation of extraction from parse trees candidate terms and their parsable
186 grammar.

187 Due to the absence of label information, it is hard to select the discriminative features in unsupervised learning.
188 In the prevailing system, unsupervised feature selection algorithms frequently select the features that preserve the
189 best data dissemination. Yang et al., [35] proposed a new approach that is L2, 1 -norm regularized Unsupervised
190 Discriminative Feature Selection (UDFS). The algorithm chooses the most discriminative feature subset from the
191 entire feature set in batch mode. UDFS outclasses the existing unsupervised feature selection algorithms and
192 selects discriminative features for data representation. The performance is sensitive to the number of selected
193 features and is data dependent.

194 Cai et al., [36] presented a novel algorithm, called Graph regularized Nonnegative Matrix Factorization
195 (GNMF) [37], which explicitly considers the local invariance. In GNMF, the geometrical information of the
196 data space is pre-arranged by building a nearest neighbor graph and gathering parts-based representation space
197 in which two data points are adequately close to each other, if they are connected in the graph. GNMF models
198 the data space as a sub manifold rooted in the ambient space and achieves more discriminating power than the
199 ordinary NMF approach.

200 Fan et al., [38] suggested a principled vibrational framework for unsupervised feature selection using the non
201 Gaussian data which is subjective to several applications that range from several diversified domains to disciplines.
202 The vibrational frameworks provides a deterministic alternative for Bayesian approximation by the maximization
203 of a lower bound on the marginal probability which has an advantage of computational efficiency.

204 6 2) Text

205 summarization and Dataset: Several approaches have been developed till date for automatic summarization by
206 identifying important topic from single document or clustered documents. Gupta et al., [39] describes a topic
207 representation approach that captures the topic and frequency driven approach using word probability which
208 gives reasonable performance and conceptual simplicity.

209 Negi et al., [40] developed a system that summarizes the information from a clump of documents. The
210 proposed system constructs the information from the given text. It achieves high accuracy but cannot calculate
211 the relevance of the document.

212 Debole et al., [41] initially explains the three phases in the life cycle of TC system like document indexing,
213 classifier learning and classifier evaluation. All researches takes Reuters 21578 documents for TC experiments.
214 Several researches have used Modapte split for testing. The three subsets used for the experiments are a set of
215 ten categories with more number of positive training examples.

216 Xie et al., [42] proposed an approach to the acquisition of the semantic features within phrases from a
217 single document that extracts document keyphrases. Keyphrase extraction method always performs better than
218 TFIDF and KEA. Keyphrase extraction is a basic research in text mining and natural language processing. The
219 method is developed on the concept of semantic relatedness where degrees between phrases are calculated by
220 the cooccurrences between phrases in a given document and the same is presented as a relatedness graph. The
221 approach is not domain specific and generalizes well on journal articles and is tested on news web pages.

222 To obtain any online information is an easy task. We log on to the world wide web and give simple keywords.
223 However, it is not easy for the user to read the entire information provided. Hence text summarization is needed.

224 **7 b) Feature Extraction**

225 Zhong et al., [44] has presented an effective pattern discovery technique which includes the process of pattern
226 deploying and pattern evolving as shown in Table 2, to improve the effectiveness of using and updating discovered
227 patterns for finding relevant and interesting information. The proposed model outperforms other pure data
228 mining-based methods, the concept based models and term-based state-of the-art models, such as BM25 and
229 SVM.

230 Li et al., [47] proposed two algorithms namely Fclustering and Wfeature to discover both positive and negative
231 patterns in the text documents. The algorithm Fclustering classifies the terms into three categories general,
232 positive, negative automatically without using parameters manually. After classifying the terms using Fclustering,
233 Wfeature is executed to calculate the weights of the term. Wfeature is effective because the selected terms size is
234 less than the average size of the documents. The proposed model is evaluated on RCV, Trec topics and Reuters
235 21578 dataset as shown in Table 2, the model performs much better than the term based method and pattern
236 based method. The use of

237 **8 Collection of Text Documents**

238 **9 E Extract Useful Features**

239 **10 Feature Weight Specificity Data Fusion**

240 **11 Relevant features with Duplicate Free**

241 Duplicate Detection irrelevance feedback strategy is highly efficient for improving the overall performance of
242 relevance feature discovery model.

243 Xu et al., [26] experimented on microblog dimensionality reduction-A deep learning approach. The approach
244 aims at extracting useful information from large amount of textual data produced by microblogging services.
245 The approach involves mapping of natural language texts into proper numerical representations which is a
246 challenging issue. Two types of approaches namely modifying training data and modifying training objective of
247 deep networks are presented to use microblog specific information. Meta-information contained in tweets like
248 embedded hyperlinks is not explored.

249 Nguyen et al., [49] worked on review selection using Micro-reviews. The approach consists of two steps namely
250 matching review sentences with micro reviews and selecting a few reviews which cover many reviews. A heuristic
251 algorithm performs computationally fast and provides informative reviews.

252 **12 III. Pattern Extraction**

253 Patterns which are close to their super patterns that appears in the same paragragh are termed closed relation
254 and needs to be eliminated. The shorter pattern is not considered since it is meaningless while the longer
255 pattern is more meaningful and hence these are significant patterns in the pattern taxonomy. Abonem et al.,
256 [53] presented text mining framework that discovers knowledge by preprocessing the data. Usually text in the
257 documents contains words, special characters and structural information and hence special characters is replaced
258 by symbols. It mainly focuses on refining the uninterested patterns and thus fitering decreases the time and
259 size of search space needed for the discovery phase. It is more efficient when large collection of documents are
260 considered. Postprocessing involves pruning, organizing and ordering of the results. The rule of each document
261 is to find a set of characteristics phrases and keywords i.e., length, tightness and mutual confidence. The ranking
262 of the rules within a document is measured by calculating a weight for each rule.

263 Mining entire set of frequent subsequence for every long pattern generates uncontrollable number of frequent
264 subsequence which are expensive in space and time. Yan et al., [54] proposed a solution for mining only frequent
265 closed subsequence through an algorithm Clospan-Closed Sequential Pattern Mining. Clospan efficiently mines
266 frequent closed sequences in large data sets with low minimum support but does not take advantage of search
267 space pruning property.

268 Gomariz et al., [55] presented CSpan algorithm for mining closed sequential patterns which mines closed
269 sequential patterns early by using pruning method called occurrence checking. CSpan outperforms clospan and
270 claspalgorithm.

271 **13 Global Journal of Computer Science and Technology**

272 Volume XVI Issue V Version I

273 **14 b) Mining Sequential Pattern**

274 To delimit the search and to increase the subsequence fragments Han et al., [57] proposed Freespan Frequent
275 Pattern Projected sequential pattern Mining. Freespan fuses the mining of frequent sequence with that of frequent
276 patterns and adopts projected sequence databases. Freespan runs quicker than the Apriori based GSP algorithm.
277 Freespan is highly scalable and processing efficient in mining complete set of patterns. Freespan causes page

278 thrashing as it requires extra memory. With extensive applications in data mining, mining sequential pattern
279 encounters problems with a usage of very large database.

280 Pei et al., [58] proposed a sequential pattern mining method called Prefix Span(Prefix Projected sequential
281 pattern mining). The complete set of patterns is extracted by reducing the generation of candidate subsequence.
282 Further prefix projection largely reduces projected database size and greatly improves efficiency as shown in
283 Table 3. Making use of RE(Regular Expression) [59] as a flexile constraint SPIRIT algorithm was proposed
284 by Garofalakis et al., [60] for mining patterns that are sequential. A family of four algorithms is executed for
285 forwarding a stronger relaxation of RE. Candidate sequence containing elements are pruned that do not appear
286 in RE than its predecessor in the pattern mining loop.

287 The degree to which RE constraints are enforced to prune the search space of patterns are the main distinctive
288 factor. The results on the real life data shows RE's adaptability as a user level tool for focussing on interesting
289 patterns.

290 Jian et al., developed a new framework called Pattern Growth ??PG]. PG is based on prefix monotonic
291 property. Every monotonic and anti monotonic regular expression constraints are preprocessed and are pushed
292 into a PG-based mining algorithm. PG adopts and also handles regular expression constraints which is diffi cult
293 to explore using Apriori based method like SPIRIT. The candidate generation and test framework adopted by PG
294 is less expensive and efficient in pushing many constraints than SPIRIT method. During Prefix growth various
295 irrelevant sequence can be excluded in the huge dataset. Accordingly, projected database quickly shrinks. While
296 PG outperforms SPIRIT, interesting constraints specific to complex structure mining is not be explored.

297 To filter the discovered patterns, Li et al., ??43] [61] proposed an effective pattern discovery technique that
298 deploys and evolves patterns to refine the discovered patterns. Using these discovered patterns, the relevant
299 information can be determined inorder to improve the effectiveness. All frequent short patterns and long patterns
300 are not useful and also long patterns with high specificity suffers from the low problem frequency. The problem
301 of low frequency and misinterpretation for text mining can be solved by employing pattern deploying strategies.

302 Rather than using individual words, some researches used phrases to discover relevant patterns from documents
303 collection. Hence there is a small improvement in the effectiveness of text mining because phrases based methods
304 have consistency of assignment and document frequency for terms to be low. Inje et al., [62] used a pattern based
305 taxonomy(is-a) relation to represent document rather than using single word. The computation cost is reduced
306 by pruning unwanted patterns and hence improves the effectiveness of system.

307 Bayardo et al., [63] evaluated Max miner algorithm inorder to mine maximal frequent itemsets from large
308 databases. Max-Miner reduces the space of itemsets considered through superset-frequency based pruning.
309 There is a performance improvements over Apriori-like algorithms when frequent itemsets are long and more
310 modest though still substantial improvements when frequent itemsets are short. Completeness at low supports
311 on complex datasets is not achieved.

312 Jan et al., [64] [65] proposed propositionalization and classification that employs long first order frequent
313 patterns for text mining. The Framework solves three text mining tasks such as information extraction,
314 morphological disambiguation and context sensitive text correction. Propositionalization approach outperforms
315 CBA by using frequent patterns as features. The performance of CBA classifiers greatly depends on number of
316 class association rules and threshold values given by the user. The proposed framework shows that the distributed
317 computation can improve performance of both method since large sample of data and a larger number of features
318 are extracted.

319 Seno et al., [66] proposed an algorithm SLP miner that finds all sequential patterns. It performs effectively
320 satisfying length decreasing support constraint and increases in average length of the sequence. It is expensive
321 as pruning is not considered in this work.

322 Nizar et al., [67] demonstrates a taxonomy of sequential pattern mining techniques. Reducing the search
323 space can be done by strongly minimizing the support count. Domain knowledge, distributed sequence are not
324 considered in the mining process.

325 15 c) Mining Frequent Sequences

326 To extract sequential patterns, various algorithms have been executed by making continuously repeated scans of
327 database and making use of hash structure.

328 16 Global Journal of Computer Science and Technology

329 Volume XVI Issue V Version I 8 Year 2016 () Zaki et al., [68] presented a new novel algorithm SPADE for
330 discovering sequential patterns at a high speed. SPADE decomposes the parent class into small subclasses.
331 These sub problems are executed without depending on other subproblems in main memory by lattice approach.
332 The lattice approach needs only one scan when having some pre-processed data. They also process depth first
333 search and breadth first search for frequent sequence enumeration within each sublattice. By using these search
334 strategies SPADE minimizes the computational costs and I/O costs by reducing number of database scans. It
335 provides pruning strategies to identify the interesting patterns and prune out irrelevant patterns.

336 BFS outperforms DFS by having more information available for pruning while constructing a set of three

337 sequence, two sequence. BFS require more main memory than DFS. BFS checks the track of idlists for all the
338 classes, while DFS needs to preserve intermediate id lists for two consecutive classes along a specific path.

339 Han et al., [69] proposed a FP(frequent pattern tree) structure where the complete set of frequent patterns can
340 be extracted by pattern fragment growth. Three techniques are used to achieve mining efficiency compression of
341 the database, (i) FP tree avoids expensive repeatedly scanning database (ii) FP tree prevents generation of large
342 number of candidates sets and uses divide and conquer method which breaks the mining task into a set of tasks
343 that lowers search space. FP growth method [70] is efficient and also scalable for extracting both long and short
344 frequent patterns and it is faster than Apriori algorithm.

345 Zhang et al., [71] executed CFP Constrained Frequent Pattern algorithm to improve the efficiency of association
346 rule mining. The algorithm is incorporated in an interrelation analysis model for celestial spectra data. The
347 module extracts correlation among the celestial spectra data characteristics. The model do not support for
348 different application domain.

349 17 d) Mining Frequent itemsets using Map Reduce

350 Database Management System have evolved over the last four decades and now functionally rich. Operating and
351 managing very large amount of business data is a challenging task. MapReduce [72] [73] is a framework that
352 process and manages a very large datasets in a distributed clusters efficiently and achieves parallelism.

353 Xun et al., [74] [75] executed Fidoop algorithm using mapreduce model. Fidoop algorithm uses frequent
354 itemset with different lengths to improve workload balance metric across clusters. Fidoop handles very high
355 dimensional data efficiently but do not work on heterogeneous clusters for mining frequent itemsets.

356 Wang et al., [76] proposed (FIMMR) Frequent Itemset Mapreduce Framework algorithm. The
357 algorithm initially extracts lower frequent itemset, applies pruning technique and later mines global frequent
358 itemset. The speedup of algorithm is satisfactory under low minimum support threshold.

359 Ramakrishnudu et al., [77] finds infrequent itemset from huge data using mapreduce framework. The efficiency
360 of framework increases as the size of the data is increased. The framework produces few intermediate items during
361 the process.

362 Ozkural et al., [78] extracts frequent item set by partitioning the graph by a vertex separator. The separator
363 mines the item distribution independently. Parallel frequent itemset algorithm replicates the items that co-relate
364 with the separator. The algorithm minimizes redundancy and load balancing is achieved. Relationship among a
365 very large number of items for real world database is not incorporated.

366 18 e) Relevance Feedback Documents

367 Xu et al., [79] presented a Expectation Maximization(EM) algorithm for relevance feedback in overlaps in feedback
368 documents. Based on dirichlet compound multinomial(DCM) distribution, EM includes a background collection
369 model reduction, by the methodology of deterministic annealing and query based regularization.

370 Several Queries which do not contain any relevance feedback needs improvisations by combining pseudo
371 relevance feedback and relevance feedback using a hybrid feed-back paradigm. Instead of using static
372 regularization, the authors adjust the regularization parameter based on the percentage of relevant feedback
373 documents [80]. Further, the design formulates the space for a much newer document progressively. The weighted
374 relevance is computed for an experimental design which further exploits the top retrieved documents by adjusting
375 the selection scheme. The relevance score algorithms need to be validated on several TREC datasets.

376 Cao et al., [81] re-examined the assumption of most frequent terms in the false feedback documents that are
377 useful and prove that it does not hold in reality. Distinguishing good and bad expansion terms cannot be done in
378 the feedback documents. The difference of term distribution between feedback documents and whole document
379 collection is exploited through the mixture model indicates that good and bad expansion terms may have similar
380 distributions that fails to distinguish. Experiments are conducted to see that each query can keep only the good
381 expansion terms. The new query model integrates the good terms, while classification of term is done to improve
382 the effectiveness of retrieval. In a final query model, the classification score is used to enhance the weight of good
383 terms. Selecting expansion terms are significantly better than traditional expansion terms by evaluating on three
384 TREC datasets. Selection of terms has to be done carefully.

385 Pak et al., [82] proposed a automatic query expansion algorithm which incorporates a incremental blind
386 approach to choose feedback documents from the top retrieved lists and further finds the terms by aggregating
387 the scores from each feedback document. The algorithm performs significantly better on large documents.

388 Algarni et al., [83] proposed the adaptive relevance feature discovery(ARFD). Using a sliding window over
389 positive and negative feedback, that ARFD updates the systems knowledge. The system provides a training
390 documents where specific features are discovered. Various methods have been used to merge and revise the
391 weight of the feature in a vector space. Documents are selected based on two categories. The first category is
392 that user provide the interested topic information and the second category is that the user changed the interest
393 topic.

394 **19 IV. Text Classification and Clustering**

395 Text categorization [84] is a significant issue in text mining. In general, the documents contains large texts and
396 hence it is necessary to classify them into specific classes. Text categorization can be broadly classified into
397 supervised and unsupervised classification. Classifying documents manually is very costly and time consuming
398 task. Hence it is necessary to construct automatic text classifiers using pre-classified sample documents whose
399 time efficiency and accuracy is much better than manual text classification.

400 Computer programs often treat the document as a sack of words. The main characteristics of text
401 categorization is feature space having high dimensionality. Even for moderate sized text documents, the feature
402 space consists of hundreds and thousands of terms.

403 Sebastiani et al., [85] reviews the standard approaches that comes under machine learning paradigm for text
404 categorization. The approach also describes the problem faced while document representation constructing
405 classifiers and evaluation of constructed classifier. The experimental study shows comparisons among different
406 classifiers on different versions of reuter dataset. Text categorization is a good benchmark for clarifying whether
407 a given learning technique can scale up to substantial sizes.

408 Irfan et al., [86] reviews different pre-processing techniques in text mining to extract various textual patterns
409 from the social networking sites. To explore the unstructured text available from social web, the basic approaches
410 of text mining like classification and clustering are provided.

411 Wu et al., [87] presents a technique consisting of three preprocessing stages to recognize the text region of
412 huge size and contrast data. A Segmentation algorithm cannot identify the changes that happen both in color
413 and illumination of character in a document image. The technique follows extracting the grayscale image such
414 as from the book cover, magazine RGB plane associated with weighted valve. A multilevel thresholding process
415 is done on each grayscale image independently to identify text region. A recursive filter is executed to interpret
416 which connects components is textual components. An approach to determine score is considered to findout the
417 probabilistic text region of resultant images. If the text region has maximum score, then it is classified as textual
418 component.

419 **20 V. Discovering Facets for Queries from Search Result**

420 Facets means a phrase or a word. A query facet is a set of items which summarize an important aspect of a
421 query. Dou et al., [88] [89] [90] explores solution of searching the set of facets for a user query. A system called
422 Query Discovery (QD) miner is proposed to mine facets automatically. Experiments are conducted for 100's of
423 queries and results shows the effectiveness of the system as shown in the table 5. It provides interesting knowledge
424 about a query and however improves searching for the users in different ways. The problem of generating query
425 suggestions based on query facets is not considered that might help users find a better query more easily.

426 Multifacted search is an important paradigm for extracting and mining applications that provides users to
427 analyze, navigate through multidimensional data.

428 Facetted search [91] can also be applied on spoken web search problem to index the metadata associated with
429 audio content that provides audio search solution for rural people. The query interface ensures that a user is
430 able to narrow the search results fastly. The approach focuses on indexing system and not generating precision
431 -recall results on a labeled set of data.

432 Kong et al., [96] incorporated the feedback of users on the query facets into document ranking for evaluating
433 boolean filtering feedback models that are widely used in conventional faceted search which automatically
434 generates the facets for a user given query instead of generating for a complete corpus. The boolean filtering
435 model is less effective than soft ranking models.

436 Bron et al., [97] proposed a novel framework by adding type filtering based on category information available
437 in wikipedia. Combining a language modelling approach with heuristic based on wikipedia's external links,
438 framework achieves high recall scores by finding homepages of top ranked entities. The model returns entities
439 that have not been judged.

440 Navarro et al., [98] develops an automatic facet generation framework for an efficient document retrieval. To
441 extract the facets a new approach is developed

442 **21 Global Journal of Computer Science and Technology**

443 Volume XVI Issue V Version I 10 Year 2016 () which is both domain independent and unsupervised. The
444 approach generates multifaceted topic effectively. The subtopics in the text collection is not investigated.

445 Liu et al., [99] presented the study of exploring topical lead lag across corpora. Selecting which text corpus
446 leads and which lags in a topic is a big challenge. Text pioneer, a visual analytic tool is introduced. The tool
447 investigates lead lag across corpora from global to local level. Multiple perspectives of results are conveyed by
448 two visualizations like global lead lag as hybrid tree, local lead lag as twisted ladder. Text pioneer donot analyze
449 topics within each corpus and across corpora.

450 Jiang et al., [100] presented Cross Lingual Query Log Topic Model (CL-QLTM) to investigate query logs to
451 derive the latent topics of web search data. The model incorporates different languages by collecting cooccurrence
452 relations and cross lingual dictionaries from query log. CL-QLTM is effective and superior in discovering latent
453 topics. The model is not applied on statistical machine translation.

454 Cafarella et al., [101] exploited the interesting knowledge from webpages which consists of higher relevance to
455 user when compared to traditional approach. The system records co-occurrences of schema elements and helps
456 user in navigating, creating synonyms for schema matching use.

457 Wordnet Domains text document. The queries given by the user is free text queries. Mapping keywords to
458 different attributes and their values of a given entity is a challenging task. Castanet is simple and effective that
459 achieves higher quality results than other automated category creation algorithms. WordNet is not exhaustive
460 and few other mechanism is needed to improve coverage for unknown terms.

461 Pound et al., [102] proposed a solution that exploits user faceted search behaviour and structured data to find
462 facet utility. The approach captures values and conditional values that provides attributes and values according
463 to user preferences. Experi

464 Space Efficient Framework Robust Multiple patterns are not handled ment results show that the approach
465 is scalable and also outperforms popular commercial systems. Altingovde et al., [103] demonstrate static index
466 pruning technique by incorporating query views like document and term centric. The technique improves the
467 quality of top ranked result. When the web pages changes frequently the original index is not updated.

468 Koutris et al., [104] proposed a framework for pricing the data based on queries. The polynomial time algorithm
469 is executed for a conjunctive queries of large class and the result shows that the data complexity instance based
470 determinicy is CO NP complete. The framework do not explore interaction between pricing and privacy.

471 Liu et al., [106] developed a tool that automatically differentiate structured data from search results. A feature
472 type based approach is introduced which identifies a valid features and evaluates the quality of features using
473 exact and heuristics computation methods. The method achieves local optimality avoids dependency on random
474 initialization. Result differentiation (whether the selected features is interest to users are not) is not incorporated.

475 Liu et al., [107] proposed matrix representation to discover collection of documents based on user interest. The
476 multidimensional visualization is presented to overcome the difficulty for users to compare across different facet
477 values. The approach further enables visual ordering based on facet values to support cross facet compa risons of
478 items and also support users in exploring tasks. The intradocument details are unavailable and visual scalability
479 is not incorporated. [105] proposed two methodology for extracting user tasks when they search for relevant data
480 from search engine. The method identifies user query logs and further aggregate same kind of users tasks based
481 on supervised and unsupervised approaches. The method is effective in detecting similar latent needs from a
482 query log. Users task by task search behaviour is not represented in the model.

483 22 C

484 Colini et al., ??111] [112] proposes multiple keyword method that provides search auctions with budgets and
485 bidders. Bidders is bounded by multiple slots per keyword. Bidders which have cumulative valuations are
486 click through rates and budgets that confine the overall study of multiple keyword method. Multiple keywords
487 mechanism is compatible, optimal and rational with expectation. In combinatorial setting, each bidder is having
488 a direct involvement in a subset of keywords. Deterministic mechanisms with temper marginal valuations are
489 incompatible.

490 Wu et al., [113] introduced the concept of safe zones. It studies the moving of top K keyword query. The safe
491 zones saves the time and communication cost. The approach computes safe zone in order to optimize server side
492 computation. It is also used to establish the client server communication. Spatial keyword is not processed and
493 also the safe zone do not compute future path of moving the query.

494 Lu et al., [114] proposed reverse spatial keyword K nearest neighbour to find the query of object which is
495 similar to one of the neighbour. The query search is based on spatial location and also text associated with it.
496 The algorithm is used to prune unnecessary objects and also computes the lists. The method do not considers
497 textual description of two different objects.

498 Cao et al., [115] demonstrates the concept of weighing a query. The spatial keywords match considers both the
499 location and the text. The method focuses more on finding queries to group of objects by grouping spatial objects.
500 Top K spatial keyword and weighing of query improves the performance and efficiency. The computational time
501 is reduced but partial coverage of queries is not considered.

502 23 VI. Fine Grained Knowledge

503 Guan et al., [116] suggested "tcpdump" method to capture the web surfing activities from users. Web surfing
504 activities reflects persons fine grained knowledge by recognizing the semantic structures. Further by using
505 Dirichlet process infinite Guassian mixture model is adopted. D-iHMM process is employed for mining the
506 fine grained aspect in each part by session clustering. Discovering fine grained knowledge Feature Extraction and
507 Duplicate Detection for Text Mining: A Survey Hon et al., [95] developed space efficient frame works for top k
508 string retrieval problem that considers two metrics for relevance features which includes frequency and proximity.
509 The threshold based approach on these metrics are also been used. Compact space and sufficient space indexes
510 are derived that results index space and query time with significant robustness. The framework is robust but do
511 not index an the cache oblivious model and also the index takes twice the size of the text. Multiple patterns are
512 not handled.

513 Zhang et al., [94] proposed (SPP) Space Partition and Probing to keep track of object position and relevance
514 to the query and also to find the vector space. Quality is achieved by using MMR which is one of the important
515 diversification algorithm. The method identifies the next top K object very quickly. SPP helps in reducing object
516 axis and also increases the performance. Fixed bounded region is not considered. Zhang et al., [93] proposed
517 inverted linear quadtree index structure to accomplish both spatial and keyword based techniques to effectively
518 decreases the search space. Spatial keyword queries having two disputes: top k spatial keyword search(TOPK-
519 SK) and batch top k spatial keyword search(BTOPK-SK), in which top-sk fetch the closest k objects which
520 contains all keywords in the query. BTOPK-SK contains set of top k queries. Existing techniques in IL-quadtree
521 presents firstly Keyword first index, which is to extract the related inverted indexes. Partition based method is
522 proposed to further enhance the filtering capabilities of the signature of linear quadtree.

523 Efstathiades et al., [92] presents Link of Interest (LOI) to improve the quality of users queries. K Relevant
524 Nearest Neighbor(K-RNN) queries is based on query processing method is proposed to analyse LOI information
525 to retrieve relevant location based point of interest as shown in Table 3. The method captures the relevance
526 aspect of data. Relevance score is not computed.

527 Catallo et al., [108] proposed probabilistic k-Sky band to process subset of sliding window objects, that are
528 most recent data objects. The algorithm out performs for parameter of large values of K parameter both in
529 memory consumption and time reduction. Adaptive top K processing is not incorporated in the approach.

530 Bast et al., [109] presented pre-processing techniques to achieve interactive query times on large text collections.
531 Two similarities measures are considered which includes, firstly, query terms match -similar terms in collection.
532 Secondly, Query terms match -terms with similar prefix in collection which display the results quickly and are
533 more efficient and scalable.

534 Termehchy et al., [110] introduced the XML structure for searching the keyword effectively. Traditional
535 keyword search techniques does not support effectively. In order to overcome these problems for data-centric,
536 XML put forth the Coherency Ranking(CR), which is a database design self sustained ranking method for XML
537 keyword queries that is based on prolonging concepts of data dependencies and mutual information. With the
538 concepts of CR, that analyze the prior approaches to XML keyword search. Approximate coherency ranking
539 and current potent algorithm process queries and rank their responds using CR. CR shows better precision and
540 recall, provides better ranking than prior approaches.

541 reflected from people's interaction made knowledge sharing in collaborative environment much easier. Although
542 privacy is major issue.

543 Wang et al., [117] analysed user's searching behaviors and considered inter-query dependencies. A semi-
544 supervised clustering model is proposed based on the SVM framework. The model enables a more comprehensive
545 understanding of user's searchbehaviors via query search logs and facilitates the development search-engine
546 support for long-term tasks. The performance of the model is superior in identifying cross-session search. User
547 modeling and long-term task based personalization is not considered.

548 Kotov et al., [118] proposed a method for creating a semi automatically labeled data set that can be used for
549 identifying user's query searches from earlier sessions on the same task and to predict whether a user returns to
550 the same task during his later session. Using logistic regression and MART classifiers the method can effectively
551 model and crosssession of user's information needs. The model is not incorporated in commercial search engines.

552 24 VII. Duplicate Detection and Data Fusion

553 Duplicate detection is the methodology of identification of multiple semantic representation of the existing and
554 similar real world entities. The present day detection methods need to execute larger datasets in the least amount
555 of time and hence to maintain the overall quality of datasets is tougher.

556 Papenbrock et al., [119] proposed a strategic approach namely the progressive duplicate detection methods as
557 shown in Table 4 which finds the duplicates efficiently and reduces the overall processing time by reporting most
558 of the results as shown in table ?? than the existing classical approaches.

559 Bano et al., [120] executed innovative windows algorithm that adapts window for duplicates and also which
560 are not duplicates and unnecessary comparisons is avoided.

561 The duplicate records are a vital problem and a concern in knowledge management [124]. To Extract duplicate
562 data items an entity resolution mechanism is employed for the procedure of cleanup. The overall evaluation reveals
563 that the clustering algorithms perform extraordinarily well with accuracy and f-measure being high.

564 Whang et al., [125] investigates the enhancement of focusing on several matching records. Three types of
565 hints that are compatible with different ER algorithms:(i) an ordered list of records, (ii) a sorted list of record
566 pairs, (iii) a hierarchy of record partitions. The underlying disadvantage of the process is that it is useful only
567 for database contents.

568 13 Year 2016

569 25 () C

570 Duplicate records do not share a strategic key but they build duplicate matching making it a tedious task. Errors
571 are induced because the results of transcription errors, incomplete information and lack of normal formats.

572 Abraham et al., [126] [127] provides survey on different techniques used for detecting duplicates in both XML
573 and relational data. It uses elimination rule to detect duplicates in database.

574 Elmagarmid et al., [128] present intensive analysis of the literature on duplicate record for detection and covers
575 various similarity metrics, which will detect some duplicate records in exceedingly available information. The
576 strengths of the survey analysis in statistics and machine learning aims to develop a lot of refined matching
577 techniques that deem probabilistic models.

578 Deduplication is an important issue in the era of huge database [129]. Various indexing techniques have been
579 developed to reduce the number of record pairs to be compared in the matching process. The total candidates
580 generated by these techniques have high efficiency with scalability and have been evaluated using various data
581 sets.

582 The training data in the form of true matches and true non matches is often unavailable in various real-world
583 applications. It is commonly up to domain and linkage experts for decision of the blocking keys. Papadakis et
584 al., [122] presented a blocking methods for clean-clean ER over Highly Heterogeneous Information Spaces (HHIS)
585 through an innovative framework which comprises of two orthogonal layers. The effective layer incorporates
586 methods for construction of several blockings with small probability of hits; the efficiency layer comprises of a
587 rich variety of techniques which restricts the required number of pairwise matches.

588 Papadakis et al., [123] focuses to boost the overall blocking efficiency of the quadratic task on Entity Resolution
589 among large, noisy, and heterogeneous information areas.

590 The problem of merging many large databases is often encountered in KDD. It is usually referred to as the
591 Merge/Purge problem and is difficult to resolve in scale and accuracy. The Record linkage [130] is a wellknown
592 data integration strategy that uses sets for merging, matching and elimination of duplicate records in large and
593 heterogeneous databases. The suffix grouping methodology facilitates the causal ordering used by the indexes

594 26 C

595 for merging blocks with least marginal extra cost resulting in high accuracy. An efficient grouping similar suffixes
596 is carried out with incorporation of a sliding window technique. The method is helpful in various health records for
597 understanding patient's details but is not very efficient as it concentrates only on blocking and not on windowing
598 technique. Additionally the methodology with duplicates that are detected using the state of the art expandable
599 paradigm is approximate [131]. It is quite helpful in creating clusters records. Bronselaer et al., [132] focused on
600 Information aggregation approach which combine information and rules available from independent sources into
601 summarization. Information aggregation is investigated in the context of inferencing objects from several entity
602 relations. The complex objects are composed of merge functions for atomic and subatomic objects in a way that
603 the composite function inherits the properties of the merge functions.

604 Sorted Neighborhood Method (SNM) proposed by Draisbach et al., [133] partitions data set and comparison
605 are performed on the jurisdiction of each partition. Further, the advances in a window over the data is done by
606 comparison of the records that appears within the range of same window. Duplicate Count Strategy (DCS) which
607 is a variation of SNM is proposed by regulating the window size. DCS++ is proposed which is much better than
608 the original SNM in terms of efficiency but the disadvantage is that the window size is fixed and is expensive for
609 selection and operation. Some duplicates might be missed when large window are used.

610 The tuples in the relational structure of the database give an overview of the similar real world entities such
611 tuples are described as duplicates. Deleting these duplicates and in turn facilitating their replacement with
612 several other tuples represents the joint informational structure of the duplicate tuples up to a maximum level.
613 The incorporated delete and then replacement mode of operation is termed as fusion. The removal of the original
614 duplicate tuples can deviate from the referential integrity.

615 Bronselaer et al., [121] describes a technique to maintain the referential integrity. The fusion Propogation
616 algorithm is based on first and second order fusion derivatives to resolve conflicts and clashes. Traditional
617 referential integrity strategies like DELETE cascading, are highly sophisticated. Execution time and recursively
618 calling the propagation algorithm increases when the length of chain linked relations increases.

619 Bleiholder et al., proposes the SQL Fuse by inducing the schema and semantics. The existential approach is
620 towards the architecture, query languages, and query execution. The final step of actually aggregating data from
621 multiple heterogeneous sources into a consistent and homogeneous dataset and is often inconsiderable.

622 Naumann et al., [134] observes that amount of noisy data are in abundance from several data sources. Without
623 any suitable techniques for integrating and fusing noisy data with deviations, the quality of data associated with
624 an integrated system remains extremely low. It is necessary for allowing tentative and declarative integration of
625 noisy and scattered data by incorporating schema matching, duplicate detection and fusion. Subjected to SQL-
626 like query against a series of tables instance, oriented schema matching covers the cognitive bridge of the varied
627 tables by alignment of various corresponding attributes. Further, a duplicate detection technique is used for
628 multiple representations of several matching entities. Finally, the paradigm of data fusion for resolving a conflict
629 in turn merges around Bleiholder et al., [135] explains a conceptual understanding of classification of different
630 operators over data fusion. Numerous techniques are based on standard and advanced operators of algebraic
631 relations and SQL. The concept of Co-clustering is explained from several techniques for tapping the rich and
632 associated meta tag information of various multimedia web documents that includes annotations, descriptions
633 and associations. Varied Coclustering mechanisms are proposed for linked data that are obtained from multiple

27 VIII. CONCLUSIONS

634 sources which do not matter the representational problem of precise texts but rather increase their performance
635 up to the most minimally empirical measurement of the multi-modal features.

636 The two channel Heterogeneous Fusion ART (HF-ART) yields several multiple channels divergently. The
637 GHF-ART [136] is designed to effectively represent multimedia content that incorporates Meta data to handle
638 precise and noisy texts. It is not trained directly using the text features but can be identified as a key tag by
639 training it with the probabilistic distribution of the tag based occurrences. The approach also incorporates a
640 highly and the most adaptive methodology for active and efficient fusion of multimodal.

641 27 VIII. Conclusions

642 The paper presents different techniques and framework to extract relevant features from huge amount of
643 unstructured text documents. The paper also reviews a survey on various text classification, clustering,
644 summerization methods.

645 To guarantee the quality of extracted relevant features in a collection of text documents is a great challenge.
646 Many text mining techniques have been proposed till date. However how effectively the discovered features is
647 interesting and useful to the user is an open issue.

648 Our future work is to efficiently utilize relevant documents from non relevant documents. Effective filtering
649 model is required to automatically generate facets. The security and time to extract the useful features that is
650 duplicate free and fine grained knowledge helps the user to reduce time in searching various web pages needs to
be addressed. ¹

Clustering: The process of grouping similar kind of information is called clustering that results in finding interesting knowledge. The new discovered knowledge can be used by an industry for further development and helps in competing with their competitors.

Question Answering: For separating and combining terms we use standard text searching techniques that use boolean operators. Sophisticated search in text mining executes the searching process in sentence or phrase level and verbal connection identification between various search terms, which is not possible in traditional search. The result obtained by sophisticated search can be used for providing specific information that can be influenced by an organization.

[Note: Concept linkage: The results obtained from sophisticated search are linked together to produce a new hypothesis. The linking of concepts is called concept linkage. Hence, new domain of knowledge can be generated by making use of concept linkage application.]

Figure 1:

651

¹Feature Extraction and Duplicate Detection for Text Mining: A Survey

1

Sl.no	Authors	Feature Selection (FS)	selection algorithms.		
			Algorithm	Advantages	Disadvantages
1. [25]	Zhao et al.,(2016)	Unsupervised Gradient	Preserve similarity and dis-clustering performance is criminant information, high achieved		Supervised FS is not considered
2. [26]	Xu et al.,(2016)	Deep Learning	Performs better than traditional dimensional reduction method		Meta data information of tweet is not considered
3. [27]	Wang et al.,(2015)	Supervised Global and Unsupervised	Features are more compact and redundancy and discriminant,Superior performance without minimization parameter		-

Figure 2: Table 1 : shows comparison of feature

3) PCA and Random Projection RP: Principal identifiers that are useful for retrieving the important Component Analysis (PCA) is a simple technique used to explore and visualize the data easily. PCA extracts useful information from complicated data sets using non parametric method. It determines a lower dimension space by statistical method. Based on eigen value decomposition of the covariance matrix transformation matrix of PCA is calculated and thereby computation cost is more and it is also not suitable for very high dimensional data. The strength of PCA is that there are no parameters to

Year 2016

6
)
(C

fine tune and also no co-efficient is required to adjust.

Fradkin et al., [51] [52] reported a number of experiments by evaluating random projection in supervised learning. Different datasets were tested to compare random projection and PCA using several machine learning methods. The results show PCA outperforms RP in efficiency for supervised learning. The results also shows that RP's are well suited to use with nearest neighbour and with SVM classifier and are less satisfactory with decision trees.

2) Feature Extraction for Classification: Khadhim et al., [50] [21] developed two weighting methods TF-IDF means clustering algorithm is used for feature extraction for classification.

[Note: and TF-IDF (Term Frequency/Inverse Document Frequency) global to reduce dimensionality of datasets because it is very difficult to process the original features i.e, thousands of features. Fuzzy c Global Journal of Computer Science and Technology Volume XVI Issue V Version I Feature Extraction and Duplicate Detection for Text Mining: A Survey 1) Feature Mining for Text Mining: Li et al.,[43] designed a new technique to discover patterns i.e., positive and negative in text document. Both relevant and irrelevant document contains useful features. Inorder to remove the noise, negative documents in the training set is used to improve the effectiveness of Pattern Taxonomy Model PTM. Two algorithms HLF mining and N revision was introduced.]

Figure 3:

2

Figure 4: Table 2 :

3

Figure 5: Table 3 :

4

Sl.no	Authors	Algorithm	Window selection	Advantages	Disadvantages
1.	Papenbrock et (2015) [119]	PSNM al.,	Adaptive	Efficient with limited execution time	Delivers results moderately
2.	Bano al., (2015) [120]	Innovative et Window	Adaptive	Unnecessary Comparison is avoided	Do not support on Datasets
3.	Bronselaer et (2015) [121]	Fusion al., Propogation	-	Conflicts in relationship attributes are resolved	Multiple More Expensive
4.	Papadakis et al., (2013) [122]	Attribute Clus- tering	-	Effective on real datasets	low quality workloads, Parallelizing is not adopted
5.	Papadakis et al., (2011) [123]	-	Adaptive	Time Complexity is reduced	Process is very slow

[Note: © 2016 Global Journals Inc. (US) 1]

Figure 6: Table 4 :

652 .1 Global Journal of Computer Science and Technology

653 Volume XVI Issue V Version I

654 [Li et al.] , R Li , K H Lei , R Khadiwala , K. C.-C Chang . *Tedas: A Twitter-based Event Detection*

655 [Yang et al. ()] '1-norm Regularized Discriminative Feature Selection for Unsupervised Learning'. Y Yang , H
656 T Shen , Z Ma , Z Huang , X Zhou . *Proceedings of the International Joint Conference on Artificial*
657 *Intelligence(IJCAI)*, (the International Joint Conference on Artificial Intelligence(IJCAI)) 2011. 2 p. .

658 [Papadakis et al. ()] 'A Blocking Framework for Entity Resolution in'. G Papadakis , E Ioannou , T Palpanas ,
659 C Nieder'ee , W Nejdl . *Highly Heterogeneous C Transactions on Database Systems (TODS)*, 2014. 2013. 39
660 p. . (Information Spaces)

661 [Yang and Pedersen ()] 'A Comparative Study on Feature Selection in Text Categorization'. Y Yang , J O
662 Pedersen . *ICML* 1997. 97 p. .

663 [Bhat et al. ()] 'A Data Mining Approach for Data Generation and Analysis for Digital Forensic Application'.
664 V H Bhat , P G Rao , R Abhilash , P D Shenoy , K R Venugopal , L Patnaik . *International Journal of*
665 *Engineering and Technology* 2010. 2 (3) p. .

666 [Agrawal and Batra ()] 'A Detailed Study on Text Mining Techniques'. R Agrawal , M Batra . *International*
667 *Journal of Soft Computing and Engineering (IJSCE) ISSN* 2013. 2 (6) p. .

668 [Song et al. ()] 'A Fast Clusteringbased Feature Subset Selection Algorithm for High-Dimensional Data'. Q Song
669 , J Ni , G Wang . *IEEE Transactions on Knowledge and Data Engineering* 2013. 25 (1) p. .

670 [Sivagaminathan and Ramakrishnan ()] 'A Hybrid Approach for Feature Subset Selection using Neural Networks
671 and Ant Colony Optimization'. R K Sivagaminathan , S Ramakrishnan . *Expert systems with Applications*
672 2007. 33 (1) p. .

673 [Bhat et al. ()] 'A Novel Data Generation Approach for Digital Forensic Application in Data Mining'. V H Bhat
674 , P G Rao , R Abhilash , P D Shenoy , K R Venugopal , L Patnaik . *Proceedings of Second International*
675 *Conference on Machine Learning and Computing (ICMLC)*, (Second International Conference on Machine
676 Learning and Computing (ICMLC)) 2010. p. .

677 [Zhang et al. ()] 'A Review on Text Mining'. Y Zhang , M Chen , L Liu . *Proceedings of 6th IEEE International*
678 *Conference on Software Engineering and Service Science (ICSESS)*, (6th IEEE International Conference on
679 Software Engineering and Service Science (ICSESS)) 2015. p. .

680 [Christen ()] 'A Survey of Indexing Techniques for Scalable Record Linkage and Deduplication'. P Christen .
681 *IEEE Transactions on Knowledge and Data Engineering* 2012. 24 (9) p. .

682 [Gupta ()] 'A Survey of Text Summarization Extractive Techniques'. V Gupta , GS . *Journal of Emerging*
683 *Technologies in Web Intelligence* 2010. 2 (3) p. .

684 [Niharika et al. ()] 'A Survey on Text Categorization'. S Niharika , V S Latha , D Lavanya . *International Journal*
685 *of Computer Trends and Technology* 2012. 3 (1) p. .

686 [Irfan et al. ()] 'A Survey on Text Mining in Social Networks'. R Irfan , C K King , D Grages , S Ewen , S U
687 Khan , S A Madani , J Kolodziej , L Wang , D Chen , A Rayes . *The Knowledge Engineering Review* 2015.
688 30 (2) p. .

689 [Abraham et al. ()] 'A Survey on Various Methods used for Detecting Duplicates in 127'. A A Abraham , S D
690 Kanmani , ; J J Tamilselvi , C B Gifta . *International Journal of Computer Applications* 2011. 15 (4) p. .
691 (Handling Duplicate Data in Data Warehouse for Data Mining)

692 [Mabroukeh and Ezeife ()] 'A Taxonomy of Sequential Pattern Mining Algorithms'. N R Mabroukeh , C I Ezeife
693 . *ACM Computing Surveys (CSUR)* 2010. 43 (1) p. .

694 [Xu and Akella ()] 'Active Relevance Feedback for Difficult Queries'. Z Xu , R Akella . *Proceedings of the 17th*
695 *ACM Conference on Information and Knowledge Management*, (the 17th ACM Conference on Information
696 and Knowledge Management) 2008. p. .

697 [Draisbach et al. ()] 'Adaptive Windows for Duplicate Detection'. U Draisbach , F Naumann , S Szott , O
698 Wonneberg . *Proceedings of IEEE 28th International Conference on Data Engineering (ICDE)*, (IEEE 28th
699 International Conference on Data Engineering (ICDE)) 2012. p. .

700 [Debole and Sebastiani ()] 'An Analysis of the Relative Hardness of Reuters-21578 Subsets'. F Debole , F
701 Sebastiani . *Journal of the American Society for Information Science and Technology* 2005. 56 (6) p. .

702 [Wang et al. ()] 'An Efficient Algorithm of Frequent Itemsets Mining based on Mapreduce'. L Wang , L Feng , J
703 Zhang , P Liao . *Journal of Information and Computational Science* 2014. 11 (8) p. .

704 [Shehata et al. ()] 'An Efficient Concept-based Mining Model for Enhancing Text Clustering'. S Shehata , F
705 Karray , M S Kamel . *IEEE Transactions on Knowledge and Data Engineering* 2010. 22 (10) p. .

706 [Ahonen et al. ()] 'Applying Data Mining Techniques for Descriptive Phrase Extraction in Digital Document
707 Collections'. H Ahonen , O Heinonen , M Klemettinen , A I Verkamo . *Proceedings of IEEE International
708 Forum on Research and Technology Advances in Digital Libraries*, (IEEE International Forum on Research
709 and Technology Advances in Digital Libraries) 1998. p. .

710 [Bronselaer and De Tr'e ()] 'Aspects of object Merging'. A Bronselaer , G De Tr'e . *Annual Meeting of the North
711 American Fuzzy Information Processing Society (NAFIPS)*, 2010. p. .

712 [Vu et al. ()] 'Automatic Extraction of Text Regions from Document Images by Multilevel Thresholding and
713 kmeans Clustering'. H N Vu , T A Tran , I S Na , S H Kim . *Proceedings of IEEE/ACIS 14th International
714 Conference on Computer and Information Science (ICIS)*, (IEEE/ACIS 14th International Conference on
715 Computer and Information Science (ICIS)) 2015. p. .

716 [Dou et al. ()] 'Automatically Mining Facets for Queries from Their Search Results'. Z Dou , Z Jiang , S Hu ,
717 J.-R Wen , R Song . *IEEE Transactions on Knowledge and Data Engineering* 2016. 28 (2) p. .

718 [Gomariz et al. ()] 'Clasp: An Efficient Algorithm for Mining Frequent Closed Sequences'. A Gomariz , M
719 Campos , R Marin , B Goethals . *Advances in Knowledge Discovery and Data Mining*, 2013. p. .

720 [Tiwari et al. ()] 'Classification Framework of Mapreduce Scheduling Algorithms'. N Tiwari , S Sarkar , U Bellur
721 , M Indrawan . *ACM Computing Surveys (CSUR)* 2015. 47 (3) p. .

722 [Joshi et al. ()] 'Classification of Alzheimer's Disease and Parkinson's Disease by using Machine Learning and
723 Neural Network Methods'. S Joshi , D Shenoy , P Rashmi , K R Venugopal , L Patnaik . *Proceedings of Second
724 International Conference on Machine Learning and Computing (ICMLC)*, (Second International Conference
725 on Machine Learning and Computing (ICMLC)) 2010. p. .

726 [Bhat et al. ()] 'Classification of Email using Beaks: Behavior and Keyword Stemming'. V H Bhat , V R Malkani
727 , P D Shenoy , K R Venugopal , L Patnaik . *Proceedings of IEEE Region 10 Conference TENCON*, (IEEE
728 Region 10 Conference TENCON) 2011. p. .

729 [Pei et al. ()] 'Closet: An Efficient Algorithm for Mining Frequent Closed Itemsets'. J Pei , J Han , R Mao .
730 *ACM SIGMOD Workshop on Research Issues in Data Mining and Knowledge Discovery* 2000. 4 (2) p. .

731 [Yan et al. ()] 'Clospan: Mining Closed Sequential Patterns in Large Datasets'. X Yan , J Han , R Afshar . *SDM*,
732 2003. p. .

733 [Azam and Yao ()] 'Comparison of Term Frequency and Document Frequency based Feature Selection Metrics
734 in Text Categorization'. N Azam , J Yao . *Expert Systems with Applications* 2012. 39 (5) p. .

735 [Computer Science and Technology Volume XVI Issue V Version I C Analysis System Proceedings of IEEE 28th International Co
736 'Computer Science and Technology Volume XVI Issue V Version I C Analysis System'. *Proceedings of IEEE
737 28th International Conference on Data Engineering (ICDE)*, (IEEE 28th International Conference on Data
738 Engineering (ICDE)) 2012. p. .

739 [Hassanzadeh and Miller ()] 'Creating Probabilistic Databases from Duplicated Data'. O Hassanzadeh , R J
740 Miller . *The VLDB Journal-The International Journal on Very Large Data Bases*, 2009. 18 p. .

741 [Jiang et al. ()] 'Cross-Lingual Topic Discovery from Multilingual Search Engine Query Log'. D Jiang , Y Tong
742 , Y Song . *ACM Transactions on Information Systems (TOIS)* 2016. 35 (2) p. .

743 [Bleiholder and Naumann ()] 'Data Fusion'. J Bleiholder , F Naumann . *ACM Computing Surveys (CSUR)* 2009.
744 41 (1) p. .

745 [Naumann et al. ()] 'Data Fusion in Three Steps: Resolving Schema, Tuple, and Value Inconsistencies'. F
746 Naumann , A Bilke , J Bleiholder , M Weis . *International Journal of Engineering Research and Technology*
747 2014. 2006. 3 (1) p. . (Engineering and Management)

748 [Liu and Chen ()] 'Differentiating Search Results on Structured Data'. Z Liu , Y Chen . *ACM Transactions on
749 Database Systems (TODS)* 2012. 37 (1) p. .

750 [Lucchese et al. ()] 'Discovering Tasks from Search Engine Query Logs'. C Lucchese , S Orlando , R Perego , F
751 Silvestri , G Tolomei . *ACM Transactions on Information Systems (TOIS)* 2013. 31 (3) p. .

752 [Li et al. ()] 'Distributed Data Management using Mapreduce'. F Li , B C Ooi , MT , S Wu . *ACM Computing
753 Surveys (CSUR)* 2014. 46 (3) p. .

754 [Elmagarmid et al. ()] 'Duplicate Record Detection: A Survey'. A K Elmagarmid , P G Ipeirotis , V S Verykios
755 . *IEEE Transactions on Knowledge and Data Engineering* 2007. 19 (1) p. .

756 [Shenoy et al. ()] 'Dynamic Association Rule Mining using Genetic Algorithms'. P D Shenoy , K Srinivasa , L M
757 Venugopal , Patnaik . *Intelligent Data Analysis* 2005. 9 (5) p. .

758 [Zhong et al. ()] 'Effective Pattern Discovery for Text Mining'. N Zhong , Y Li , S.-T Wu . *IEEE Transactions
759 on Knowledge and Data Engineering* 2012. 24 (1) p. .

760 [Zhong et al. ()] 'Effective Pattern Discovery for Text Mining'. N Zhong , Y Li , S.-T Wu . *IEEE Transactions
761 on Knowledge and Data Engineering* 2012. 24 (1) p. .

762 [Lu et al.] *Efficient Algorithms and Cost Models for Reverse Spatial-Keyword K Nearest Neighbor Search*, Y Lu
763 , J Lu , G Cong , W Wu , C Shahabi . ACM.

764 [Papadakis and Nejdl ()] ‘Efficient Entity Resolution Methods for Heterogeneous Information Spaces’. G Papadakis , W Nejdl . *Proceedings of IEEE 27th International Conference on Data Engineering Workshops (ICDEW)*, (IEEE 27th International Conference on Data Engineering Workshops (ICDEW)) 2011. p. .

767 [Bast and Celikik ()] ‘Efficient Fuzzy Search in Large Text Collections’. H Bast , M Celikik . *ACM Transactions on Information Systems (TOIS)* 2013. 31 (2) p. .

769 [Efstathiades et al. ()] ‘Efficient Processing of Relevant Nearest-Neighbor Queries’. C Efstathiades , A Efentakis
770 , D Pfoser . *ACM Transactions on Spatial Algorithms and Systems (TSAS)* 2016. 2 (3) p. .

771 [Cao et al. ()] ‘Efficient Processing of Spatial Group Keyword Queries’. X Cao , G Cong , T Guo , C S Jensen ,
772 B C Ooi . *ACM Transactions on Database Systems (TODS)* 2015. 40 (2) p. .

773 [Bayardo ()] ‘Efficiently Mining Long Patterns from Databases’. R J BayardoJr . *ACM Sigmod Record* 1998. 27
774 (2) p. .

775 [Gasca et al. ()] ‘Eliminating Redundancy and Irrelevance using a New Mlp-based Feature Selection Method’. E
776 Gasca , J S S’anchez , R Alonso . *Pattern Recognition* 2006. 39 (2) p. .

777 [Parikh and Karlapalem ()] ‘Et: Events from Tweets’. R Parikh , K Karlapalem . *Proceedings of the 22nd
778 International Conference on World Wide Web Companion*, (the 22nd International Conference on World
779 Wide Web Companion) 2013. p. .

780 [Fradkin and Madigan ()] ‘Experiments with Random Projections for Machine Learning’. D Fradkin , D Madigan
781 . *Proceedings of the Ninth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*,
782 (the Ninth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining) 2003. p. .

783 [Liu et al. ()] ‘Exploring Topical Lead-Lag Across Corpora’. S Liu , Y Chen , H Wei , J Yang , K Zhou , S M
784 Drucker . *IEEE Transactions on Knowledge and Data Engineering* 2015. 27 (1) p. .

785 [Kong and Allan ()] ‘Extending Faceted Search to the General Web’. W Kong , J Allan . *Proceedings of the 23rd
786 ACM International Conference on Information and Knowledge Management*, (the 23rd ACM International
787 Conference on Information and Knowledge Management) 2014. p. .

788 [Pound et al. ()] ‘Facet Discovery for Structured Web Search: A Query-Log Mining Approach’. J Pound , S
789 Paparizos , P Tsaparas . *Proceedings of the ACM SIGMOD International Conference on Management of
790 Data*, (the ACM SIGMOD International Conference on Management of Data) 2011. p. .

791 [Diao et al. ()] ‘Faceted Search and Browsing of Audio Content on Spoken Web’. M Diao , S Mukherjea , N
792 Rajput , K Srivastava . *Proceedings of the 19th ACM International Conference on Information and Knowledge
793 Management*, (the 19th ACM International Conference on Information and Knowledge Management) 2010.
794 p. .

795 [Kadhim et al. ()] ‘Feature Extraction for Co-occurrence-based Cosine Similarity Score of Text Documents’. A I
796 Kadhim , Y Cheah , N H Ahamed , L A Salman . *Proceedings of IEEE Student Conference on Research and
797 Development (SCoReD)*, (IEEE Student Conference on Research and Development (SCoReD)) 2014. p. .

798 [Wang et al. ()] ‘Feature Selection via Global Redundancy Minimization’. D Wang , F Nie , H Huang . *IEEE
799 Transactions on Knowledge and Data Engineering* 2015. 27 (10) p. .

800 [Ogura et al. ()] ‘Feature Selection with a Measure of Deviations from Poisson in Text Categorization’. H Ogura
801 , H Amano , M Kondo . *Expert Systems with Applications* 2009. 36 (3) p. .

802 [Xun et al. ()] ‘Fidoop: Parallel Mining of Frequent Itemsets Using Mapreduce’. Y Xun , J Zhang , X Qin . *IEEE
803 Transactions on Systems, Man, and Cybernetics: Systems* 2016. 46 (3) p. .

804 [Guan et al. ()] ‘Fine-Grained Knowledge Sharing in Collaborative Environments’. Z Guan , S Yang , H Sun , M
805 Srivatsa , X Yan . *IEEE Transactions on Knowledge and Data Engineering* 2015. 27 (8) p. .

806 [Hassanzadeh et al. ()] ‘Framework for Evaluating Clustering Algorithms in Duplicate Detection’. O Hassanzadeh
807 , F Chiang , H C Lee , R J Miller . *Proceedings of the VLDB Endowment*, (the VLDB Endowment) 2009. 2
808 p. .

809 [Han et al. ()] ‘Freespan: Frequent Pattern-Projected Sequential Pattern Mining’. J Han , J Pei , B Mortazavi-
810 Asl , Q Chen , U Dayal , M.-C Hsu . *Proceedings of the Sixth ACM SIGKDD International Conference on
811 Knowledge Discovery and Data Mining*, (the Sixth ACM SIGKDD International Conference on Knowledge
812 Discovery and Data Mining) 2000. p. .

813 [Srinivasa et al. ()] ‘Generic Feature Extraction for Classification using Fuzzy C-means Clustering’. K Srinivasa ,
814 A Singh , A Thomas , K R Venugopal , L Patnaik . *Proceedings of 3rd International Conference on Intelligent
815 Sensing and Information Processing*, (3rd International Conference on Intelligent Sensing and Information
816 Processing) 2005. p. .

817 [Muni et al. ()] ‘Genetic Programming for Simultaneous Feature Selection and Classifier Design’. D P Muni , N
818 R Pal , J Das . *IEEE Transactions on Systems, Man, and Cybernetics* 2006. 36 (1) p. . (Part B (Cybernetics))

819 [Yan et al. ()] ‘Graph Embedding and Extensions: A General Framework for Dimensionality Reduction’. S Yan
 820 , D Xu , B Zhang , H.-J Zhang , Q Yang , S Lin . *IEEE Transactions on Pattern Analysis and Machine*
 821 *Intelligence* 2007. 29 (1) p. .

822 [Zhao et al. ()] ‘Graph Regularized Feature Selection with Data Reconstruction’. Z Zhao , X He , D Cai , L
 823 Zhang , W Ng , Y Zhuang . *IEEE Transactions on Knowledge and Data Engineering* 2016. 28 (3) p. .

824 [Zhao et al. ()] ‘Graph Regularized Feature Selection with Data Reconstruction’. Z Zhao , X He , L Zhang , W
 825 Ng , Y Zhuang . *IEEE Transactions on Knowledge and Data Engineering* 2016. 28 (3) p. .

826 [Cai et al. ()] ‘Graph Regularized Nonnegative Matrix Factorization for Data Representation’. D Cai , X He , J
 827 Han , T S Huang . *IEEE Transactions on Pattern Analysis and Machine Intelligence* 2011. 33 (8) p. .

828 [Mousavi et al. ()] ‘Harvesting Domain Specific Ontologies from Text’. H Mousavi , D Kerr , M Iseli , C
 829 Zaniolo . *Proceedings of IEEE International Conference on Semantic Computing (ICSC)*, (IEEE International
 830 Conference on Semantic Computing (ICSC)) 2014. p. .

831 [Schulz et al. ()] ‘I See a Car Crash: Real-Time Detection of Small Scale Incidents in Microblogs’. A Schulz , P
 832 Ristoski , H Paulheim . *The Semantic Web: ESWC Satellite Events*, 2013. p. .

833 [Paik et al. ()] ‘Incremental Blind Feedback: An Effective Approach to Automatic Query Expansion’. J H Paik ,
 834 D Pal , S K Parui . *ACM Transactions on Asian Language Information Processing (TALIP)* 2014. 13 (3) p. .

835 [Bartoli et al. ()] ‘Inference of Regular Expressions for Text Extraction from Examples’. A Bartoli , A Lorenzo ,
 836 E Medvet , F Tarlao . *IEEE Transactions on Knowledge and Data Engineering* 2016. 28 (5) p. .

837 [Bano and Azam ()] ‘Innovative Windows for Duplicate Detection’. H Bano , F Azam . *International Journal of*
 838 *Software Engineering and Its Applications* 2015. 9 (1) p. .

839 [Zhang et al. ()] ‘Interrelation Analysis of Celestial Spectra Data using Constrained Frequent Pattern Trees’. J
 840 Zhang , X Zhao , S Zhang , S Yin , X Qin , I Senior , Member . *Knowledge-Based Systems* 2013. 41 (4) p. .

841 [Zhang et al. ()] ‘Inverted Linear Quadtree: Efficient Top k Spatial Keyword Search’. C Zhang , Y Zhang
 842 , W Zhang , X Lin . *IEEE Transactions on Knowledge and Data Engineering* 2016. 28 (7) p. .

843 [Xie et al. ()] ‘Keyphrase Extraction based on Semantic Relatedness’. F Xie , X Wu , X Hu . *Proceedings of*
 844 *9th IEEE International Conference on Cognitive Informatics (ICCI)*, (9th IEEE International Conference on
 845 Cognitive Informatics (ICCI)) 2010. p. .

846 [Wang et al. ()] ‘Learning to Extract Cross-Session Search Tasks’. H Wang , Y Song , M.-W Chang , X He ,
 847 R W White , W Chu . *Proceedings of the 22nd International Conference on World Wide Web*, (the 22nd
 848 International Conference on World Wide Web) 2013. p. .

849 [Sebastiani ()] ‘Machine Learning in Automated Text Categorization’. F Sebastiani . *ACM Computing Surveys*
 850 (*CSUR*) 2002. 34 (1) p. .

851 [Venugopal and Buyya ()] *Mastering c++*, K R Venugopal , R Buyya . 2013. Tata McGraw-Hill Education.

852 [Xu et al. ()] ‘Microblog Dimensionality Reduction-A Deep Learning Approach’. L Xu , C Jiang , Y Ren , H.-H
 853 Chen . *IEEE Transactions on Knowledge and Data Engineering* 2016. 28 (7) p. .

854 [Dai et al. ()] ‘Minedec: A Decision-Support Model that Combines Textmining Technologies with Two Compet-
 855 itive Intelligence Analysis Methods’. Y Dai , T Kakkonen , E Sutinen . *International Journal of Computer*
 856 *Information Systems and Industrial Management Applications* 2011. 3 (10) p. .

857 [Raju and Varma ()] ‘Mining Closed Sequential Patterns in Large Sequence Databases’. V P Raju , G S Varma
 858 . *International Journal of Database Management Systems* 2015. 7 (1) p. .

859 [Han et al. ()] ‘Mining Frequent Patterns without Candidate Generation: A Frequent-Pattern Tree Approach’.
 860 J Han , J Pei , Y Yin , R Mao . *Data Mining and Knowledge Discovery* 2004. 8 (1) p. .

861 [Ramakrishnudu and Subramanyam ()] ‘Mining Interesting Infrequent Itemsets from Very Large Data based on
 862 Mapreduce Framework’. T Ramakrishnudu , R Subramanyam . *International Journal of Intelligent Systems*
 863 *and Applications* 2015. 7 (7) p. .

864 [Li et al. ()] ‘Mining Positive and Negative Patterns for Relevance Feature Discovery’. Y Li , A Algarni , N Zhong
 865 . *Proceedings of the 16th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*,
 866 (the 16th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining) 2010. p. .

867 [Shi and Yang ()] ‘Mining Related Queries from Web Search Engine Query Logs using an Improved Association
 868 Rule Mining Model’. X Shi , C C Yang . *Journal of the American Society for Information Science and*
 869 *Technology* 2007. 58 (12) p. .

870 [Chen et al. ()] ‘Mining Temporal Patterns in Time Interval-Based Data’. Y.-C Chen , W.-C Peng , S.-Y. Lee .
 871 *IEEE Transactions on Knowledge and Data Engineering* 2015. 27 (12) p. .

872 [Kotov et al. ()] ‘Modeling and Analysis of Cross-Session Search Tasks’. A Kotov , P N Bennett , R W White
 873 , S T Dumais , J Teevan . *Proceedings of the 34th International ACM SIGIR Conference on Research and*
 874 *Development in Information Retrieval*, (the 34th International ACM SIGIR Conference on Research and
 875 Development in Information Retrieval) 2011. p. .

876 [Wu et al. ()] 'Moving Spatial Keyword Queries: Formulation, Methods, and Analysis'. D Wu , M L Yiu , C S
877 Jensen . *ACM Transactions on Database Systems (TODS)* 2013. 38 (1) p. .

878 [Colini-Baldeschi et al. ()] 'On Multiple Keyword Sponsored Search Auctions with Budgets'. R Colini-Baldeschi
879 , S Leonardi , M Henzinger , M Starnberger . *ACM Transactions on Economics and Computation* 2016. 4 (1)
880 p. .

881 [Zhao et al. ()] 'On Similarity Preserving Feature Selection'. Z Zhao , L Wang , H Liu , J Ye . *IEEE Transactions
882 on Knowledge and Data Engineering* 2013. 25 (3) p. .

883 [Inje and Patil ()] 'Operational Pattern Revealing Technique in Text Mining'. A Inje , U Patil . *Proceedings of
884 IEEE Students' Conference on Electrical, Electronics and Computer Science (SCEECS)*, (IEEE Students'
885 Conference on Electrical, Electronics and Computer Science (SCEECS)) 2014. p. .

886 [Ozkural et al. ()] 'Parallel Frequent Item Set Mining with Selective Item Replication'. E Ozkural , B Ucar , C
887 Aykanat . *IEEE Transactions on Parallel and Distributed Systems* 2011. 22 (10) p. .

888 [Whang et al. ()] 'Pay-asyougo Entity Resolution'. S E Whang , D Marmaros , H Garcia-Molina . *IEEE
889 Transactions on Knowledge and Data Engineering* 2013. 25 (5) p. .

890 [Hu and Wan ()] 'Ppsgen: Learning-Based Presentation Slides Generation for Academic Papers'. Y Hu , X Wan
891 . *IEEE Transactions on Knowledge and Data Engineering* 2015. 27 (4) p. .

892 [Pei et al. ()] 'Prefixspan: Mining Sequential Patterns Efficiently by Prefix-Projected Pattern Growth'. J Pei , J
893 Han , B Mortazavi-Asl , H Pinto , Q Chen , U Dayal , M.-C Hsu . *ICCN* 2001. p. .

894 [Papenbrock et al. ()] 'Progressive Duplicate Detection'. T Papenbrock , A Heise , F Naumann . *IEEE
895 Transactions on Knowledge and Data Engineering* 2015. 27 (5) p. .

896 [Bronselaer et al. ()] 'Propagation of Data Fusion'. A Bronselaer , D Van Britsom , G De Tre . *IEEE Transactions
897 on Knowledge and Data Engineering* 2015. 27 (5) p. .

898 [Sejal et al. ()] 'Qrgqr: Query Relevance Graph for Query Recommendation'. D Sejal , K Shailesh , V Tejaswi , D
899 Anvekar , K R Venugopal , S Iyengar , L Patnaik . *Proceedings of IEEE Region 10 Symposium (TENSYMP)*,
900 (IEEE Region 10 Symposium (TENSYMP)) 2015. p. .

901 [Sejal et al. ()] 'Query Click and Text Similarity Graph for Query Suggestions'. A Sejal , K Shailesh , V Tejaswi
902 , D Anvekar , K R Venugopal , S Iyengar , L Patnaik . *International Workshop on Machine Learning and
903 Data Mining in Pattern Recognition*, 2015. p. .

904 [Koutris et al. ()] 'Query-Based Data Pricing'. P Koutris , P Upadhyaya , M Balazinska , B Howe , D Suciu .
905 *Journal of the ACM (JACM)* 2015. 62 (5) p. .

906 [Bron et al. ()] 'Ranking Related Entities: Components and Analyses'. M Bron , K Balog , M De Rijke .
907 *Proceedings of the 19 th ACM International Conference on Information and Knowledge Management*, (the 19
908 th ACM International Conference on Information and Knowledge Management) 2010. p. .

909 [D'andrea et al. ()] 'Real-Time Detection of Traffic from Twitter Stream Analysis'. E D'andrea , P Ducange , B
910 Lazzerini , F Marcelloni . *IEEE Transactions on Intelligent Transportation Systems* 2015. 16 (4) p. .

911 [Li et al. ()] 'Relevance Feature Discovery for Text Mining'. Y Li , A Algarni , M Albathan , Y Shen , M A
912 Bijaksana . *IEEE Transactions on Knowledge and Data Engineering* 2015. 27 (6) p. .

913 [Nguyen et al. ()] 'Review Selection Using Micro-Reviews'. T.-S Nguyen , H W Lauw , P Tsaparas . *IEEE
914 Transactions on Knowledge and Data Engineering* 2015. 27 (4) p. .

915 [Vries et al. ()] 'Robust Record Linkage Blocking using Suffix Arrays and Bloom Filters'. T De Vries , H Ke , S
916 Chawla , P Christen . *ACM Transactions on Knowledge Discovery from Data (TKDD)* 2011. 5 (2) p. .

917 [Algarni et al. ()] 'Selected New Training Documents to Update User profile'. A Algarni , Y Li , Y Xu .
918 *Proceedings of the 19 th ACM International Conference on Information and Knowledge Management*, (the 19
919 th ACM International Conference on Information and Knowledge Management) 2010. p. .

920 [Cao et al. ()] 'Selecting Good Expansion Terms for Pseudo-Relevance Feedback'. G Cao , J.-Y Nie , J Gao
921 , S Robertson . *Proceedings of the 31st Annual International ACM SIGIR Conference on Research and
922 Development in Information Retrieval*, (the 31st Annual International ACM SIGIR Conference on Research
923 and Development in Information Retrieval) 2008. p. .

924 [Meng et al. ()] 'Semi-Supervised Heterogeneous Fusion for Multimedia Data Co-Clustering'. L Meng , A.-H Tan
925 , D Xu . *IEEE Transactions on Knowledge and Data Engineering* 2014. 26 (9) p. .

926 [Seno and Karypis ()] 'Slpminer: An Algorithm for Finding Frequent Sequential Patterns using Length-
927 Decreasing Support Constraint'. M Seno , G Karypis . *Proceedings of IEEE International Conference on
928 Data Mining*, (IEEE International Conference on Data Mining) 2002. p. .

929 [Venugopal et al. ()] *Soft Computing for Data Mining Applications*, K R Venugopal , K Srinivasa , L M Patnaik
930 . 2009. Springer.

931 [Hon et al. ()] 'Space-Efficient Frameworks for Top-k String Retrieval'. W.-K Hon , R Shah , S V Thankachan ,
932 J S Vitter . *Journal of the ACM (JACM)* 2014. 61 (2) p. .

933 [Navarro ()] 'Spaces, Trees, and Colors: The Algorithmic Landscape of Document Retrieval on Sequences'. G
934 Navarro . *ACM Computing Surveys (CSUR)* 2014. 46 (4) p. .

935 [Zaki ()] 'Spade: An Efficient Algorithm for Mining Frequent Sequences'. M J Zaki . *Machine learning* 2001. 42
936 (1-2) p. .

937 [Garofalakis et al. ()] *Spirit: Sequential Pattern Mining with Regular Expression Constraints*, M N Garofalakis ,
938 R Rastogi , K Shim . 1999. 99 p. . (VLDB)

939 [Altingovde et al. ()] 'Static Index Pruning in Web Search Engines: Combining Term and Document Popularities
940 with Query Views'. I S Altingovde , R Ozcan , O" Ulusoy . *ACM Transactions on Information Systems (TOIS)*
941 2012. 30 (1) p. .

942 [Gonzalez et al. ()] 'Text Detection and Recognition on Traffic Panels from Street-Level Imagery using Visual
943 Appearance'. A Gonzalez , L M Bergasa , J J Yebes . *IEEE Transactions on Intelligent Transportation
944 Systems* 2014. 15 (1) p. .

945 [Aghdam et al. ()] 'Text Feature Selection Using Ant Colony Optimization'. M H Aghdam , N Ghasem-Aghaee
946 , M E Basiri . *Expert Systems with Applications* 2009. 36 (3) p. .

947 [Sanchez et al. ()] 'Text Knowledge Mining:An Alternative to Text Data Mining'. D Sanchez , M J Martin-
948 Bautista , I Blanco , C Torre . *Proceedings of IEEE International Conference on Data Mining Work-
949 shops(ICDMW)*, (IEEE International Conference on Data Mining Workshops(ICDMW)) 2008. p. .

950 [Verma et al. ()] 'Text Mining and Information Professionals: Role, Issues and Challenges'. V K Verma , M
951 Ranjan , P Mishra . *Proceedings of 4th International Symposium on Emerging Trends and Technologies
952 in Libraries and Information Services (ETTLIS)*, (4th International Symposium on Emerging Trends and
953 Technologies in Libraries and Information Services (ETTLIS)) 2015. p. .

954 [Brown ()] 'Text Mining the Contributors to Rail Accidents'. D E Brown . *IEEE Transactions on Intelligent
955 Transportation Systems* 2015. 27 (5) p. .

956 [Akilan ()] 'Text mining: Challenges and Future Directions'. A Akilan . *Proceedings of Second International
957 Conference on Electronics and Communication Systems (ICECS)*, (Second International Conference on
958 Electronics and Communication Systems (ICECS)) 2015. p. .

959 [Negi et al. ()] 'Text Summarization for Information Retrieval using Pattern Recognition Techniques'. P S Negi
960 , M Rauthan , H Dhami . *International Journal of Computer Applications* 2011. 21 (10) p. .

961 [Arguello and Capra ()] 'The Effects of Aggregated Search Coherence on Search Behavior'. R Arguello , Capra .
962 *ACM Transactions on Information Systems (TOIS)* 2016. 35 (1) p. .

963 [Sakr et al. ()] 'The Family of Mapreduce and Large-Scale Data Processing Systems'. S Sakr , A Liu , A G
964 Fayoumi . *ACM Computing Surveys (CSUR)* 2013. 46 (1) p. .

965 [Pripu?zi?c et al. ()] 'Timeand Space-Efficient Sliding Window Top-k Query Processi-ng'. K Pripu?zi?c , I P
966 Zarko , K Aberer . *ACM Transactions on Database Systems (TODS)* 2015. 40 (1) p. .

967 [Catallo et al. ()] 'Top-k Diversity Queries Over Bounded Regions'. E Catallo , P Ciceri , D Fraternali , M
968 Martinenghi , Tagliasacchi . *ACM Transactions on Database Systems (TODS)* 2013. 38 (2) p. .

969 [Tang and Liu ()] 'Unsupervised Feature Selection for Linked Social Media Data'. J Tang , H Liu . *Proceedings of
970 the 18th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, (the 18th
971 ACM SIGKDD International Conference on Knowledge Discovery and Data Mining) 2012. p. .

972 [Cai et al. ()] 'Unsupervised Feature Selection for Multi-cluster Data'. D Cai , C Zhang , X He . *Proceedings of
973 the 16th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, (the 16th ACM
974 SIGKDD International Conference on Knowledge Discovery and Data Mining) 2010. p. .

975 [Fan et al. ()] 'Unsupervised Hybrid Feature Extraction Selection for High-Dimensional Non-Gaussian Data
976 Clustering with Variational Inference'. W Fan , N Bouguila , D Ziou . *IEEE Transactions on Knowledge
977 and Data Engineering* 2013. 25 (7) p. .

978 [Desai et al. ()] 'User Feedback Session with Clicked and Unclicked Documents for Related Search Recommen-
979 dation'. S Desai , V Chandrasheker , V Mathapati , K R V Rajuk , S S Iyengar , L M Patnaik . *IADIS-
980 International Journal on Computer Science and Information Systems* 2016. 11 (1) p. .

981 [Termehchy and Winslett ()] 'Using Structural Information in Xml Keyword Search Effectively'. A Termehchy ,
982 M Winslett . *ACM Transactions on Database Systems (TODS)* 2011. 36 (1) p. .

983 [Thai et al. ()] 'Visual Abstraction and Ordering in Faceted Browsing of Text Collections'. V Thai , P.-Y Rouille
984 , S Handschuh . *ACM Transactions on Intelligent Systems and Technology (TIST)* 2012. 3 (2) p. .

985 [Cafarella et al. ()] 'Webtables: Exploring the Power of Tables on the Web'. M J Cafarella , A Halevy , D Z Wang
986 , E Wu , Y Zhang . *Proceedings of the VLDB Data to find Endowment*, (the VLDB Data to find Endowment)
987 2008. 1 p. .