A chunking-based text-tokenization framework incorporating domain knowledge

Anonymous ACL submission

Abstract

1

One of the crucial activities in Natural 2 Language Processing (NLP) is to tokenize 3 text to extract features so that data mining 4 models can be applied. Many widely-used 5 tokenization algorithms take the approach 6 of using words as tokens. This approach 7 suffers from the following limitations: (a) 8 Using words as features leads to high ۵ dimensionality of the data file generated 10 from text, (b) These algorithms use a one-11 size fits all approach on the text and extract 12 tokens uniformly without consideration of 13 the prior knowledge available in the 14 domain. Here a novel method is proposed 15 which extracts features by tokenizing text 16 as chunks. Domain specific knowledge is 17 used to generate syntactic rules which are 18 then used to split text documents into Finite 19 State Rule Based Chunks. The chunks are 20 the tokens on which data mining models are 21 then applied to generate insights. The 22 effectiveness of chunk-based tokenization 23 is demonstrated by extracting chunk-based 24 tokens as well as word-based tokens from a 25 document corpus, and performing 26 clustering in both cases. A comparison of 27 the clustering of corpus with chunk tokens 28 vis-à-vis word tokens shows marked 29 improvement in clustering performance in 30 the former. 31

32 Introduction

³³ With the advent of machine learning and natural ³⁴ language processing (NLP) tools, huge volumes of 35 text data can be mined and patterns, sentiments, 76 ³⁶ latent topics identified with reduced human effort. ⁷⁷ ³⁷ One of the major problems in NLP is generating 78 methodology and presents a case study showcasing ³⁸ features from text. Some of the popular methods ⁷⁹ the application of this model. Section 2 covers the ³⁹ involve embedding either the individual words in ⁸⁰ gaps and challenges in application of NLP models 40 the text (also called tokens) or a combination of 81 and techniques. Section 3 describes the concept of 41 two, three, or n-words (also called bigrams, 82 chunking and the broad methodology. A

42 trigrams, n-grams, etc). The limitations in these ⁴³ approaches are threefold: a) Embedding individual 44 words as tokens leads to generation of high-45 dimensional vectors for each document which 46 leads to creation of sparse matrices and overfitting 47 risk; b) While bi-grams, tri-grams or n-grams help 48 in reduction of dimensions of the embedding 49 vectors, every feature thus generated consists of 50 two, three or n-words respectively as a one-size-51 fits-all solution, without considering that different 52 semantic meanings might require different 53 customizations; c) The involvement of prior 54 domain knowledge in the traditional token 55 embedding or n-gram embeddings is minimal. 56

In this paper, a text featurization framework is 57 58 proposed that uses user's prior domain knowledge 59 to extract specific customized phrases from the text 60 as features, also called as chunks. Using domain 61 specific attributes, various syntactic rules are 62 developed which are then used for splitting text 63 documents into chunks. This method is called 64 Finite State Rule Based Chunking, where the 65 syntactic rules generated are used to identify 66 chunks from text documents using tag patterns 67 (sequence of part of speech tags). Once the chunks 68 are identified, clustering is applied to group the 69 chunks based on their semantic similarity, thus ⁷⁰ reducing the chunks to a standardized list of cluster 71 labels. Using the cluster labels as features, the 72 documents are then embedded. This method helps 73 avoid the high dimensionality and sparsity of the 74 other embedding frameworks and provides better 75 performance.

The remainder of this paper covers the

⁸³ demonstration of model application is given in a ¹³⁴ probabilistic topic modelling. While topic models ⁸⁴ case study in section 4 where safety incident ¹³⁵ are useful for identifying the prevalent themes in a ⁸⁵ investigation reports generated in a steel plant are ¹³⁶ document saving time and effort, owing to their ⁸⁶ analyzed to identify chain of events model of ¹³⁷ generalization they face many limitations: (a) They 87 accident causation. TFIDF vector space model is 138 do ⁸⁸ used for converting the chunks into feature vectors ¹³⁹ relationships (b) They are difficult to be expressed ⁸⁹ and perform clustering. The findings and ¹⁴⁰ in a meaningful way. Hence some variations of ⁹⁰ comparison of the chunking results are also ¹⁴¹ topic models have been proposed e.g. supervised 91 presented. Section 5 provides a comparative 142 topic models for classification and regression 92 analysis of clustering performance of chunking vis- 143 respectively ⁹³ à-vis original documents. Finally, sections 6 and 7 ¹⁴⁴ heterogeneity and bias (Rodrigues et. al., 2017), ⁹⁴ provide the conclusions, limitations and future ¹⁴⁵ discriminative relational topic models (Chen et. al., 95 scope of work.

96 2 Literature Review

115

97 Text mining refers to knowledge discovery by ³⁸ identifying high-relevance patterns from text ¹³⁰ differences and similarities. 99 (Feldman and Dagan, 1995). Allahyari et. al. 100 (2017) classified various text mining approaches as 153 101 Information Retrieval (IR), text summarization, 154 supervised and unsupervised models is in text 102 Natural Language Processing (NLP), sentiment 155 representation. Currently the most widely used 103 analysis, Extraction (IE). Zhang et. al. (2015) classified text 157 (Stavrianou et. al., 2007) where a string of text is 105 mining algorithms into four classes: text 158 described by a vector of text features and its content 106 categorization, trend modelling, text clustering and 159 contains a function of the feature attributes (e.g. 107 association rule mining. Text mining uses both 160 frequencies with which these features appear in the supervised and unsupervised models. Some of the 161 corpus or corpora, relevance of the features etc.) nearest neighbor classifiers, decision trees, rule-111 based classifiers and probabilistic classifiers 164 extensions of the VSM. Some representations 112 (Sebastini, 2002). Some popular categorization 165 focus on phrases instead of single words (Blake and 113 techniques in supervised models include KNN, 166 Pratt, 2001; Caropreso et. al, 2001; Mitra et. al., 114 SVM, Naïve Bayes etc (Kadhim, 2019).

While supervised algorithms offer advantages 116 in terms of obtaining precise, user-friendly insights 118 and labels from historical data (in form of labelled 119 training data), in big datasets, there is requirement of a large annotated corpus and it becomes 121 expensive and time consuming to do annotations. 122 Secondly, when there is some modification in the 123 domain, it requires re-annotation of the training 124 data all over again. Hence unsupervised algorithms 125 become potentially important. Use of probabilistic 126 techniques for text mining include topic modelling 127 using probabilistic Latent Semantic Analysis 128 (PLSA) (Hofmann, 1999) and Latent Dirichlet 129 Allocation (LDA) (Blei, 2012). The objective of 130 topic modelling is to retrieve suitable latent topics ¹³¹ from document collections. For example, Chang et ¹³⁴ address some of the problems faced during text 132 al (2009) used quantitative methods for evaluating 133 semantic meaning in topics inferred from

not provide insights on hierarchical accounting for annotators' 146 2014) which is a generative process that combines 147 modeling of the document link structure with 148 document contents, and differential topic models 149 (Chen et al, 2014]) that use hierarchical Bayesian 150 nonparametric techniques to model both topic

Another major challenge faced by both opinion mining and Information 156 representation is the vector space model (VSM) common models used in supervised models include 162 (Salton et. al., 1975). Various kinds of text are 167 1997) and some give importance to semantics of 168 words or relations between them (Cimiano et. al., 169 2005; Kehagias et. al., 2003; Rajman and 170 Besancon, 1999) while some take advantage of the 171 hierarchial structure of the text (Antonellis and 172 Gallopoulos, 2006). These models either use term 173 by sentence matrix (Antonellis and Gallopoulos, 174 2006), association rules (Blake and Pratt, 2001), 175 combination of bag of words and concept hierarchy 176 (Bloehdorn et. al., 2006), n-grams (Caropreso et. 177 al., 2001), concept hierarchy (Cimiano et. al., 178 2005), supervised term weighing method (TF-RF) 179 (Lan et. al., 2008), Latent Semantic Indexing where 180 the semantic structure of the documents is used to 181 improve detection of relevant documents on the 182 basis of query terms or sense-based vectors 183 (Kehagias et. al., 2003). Thus, while these models 185 representation, they lack in some other areas e.g. (i) 186 Lack of incorporation of prior knowledge during 237 approach based chunking. In statistical approach to 187 text representation (ii) fixed n-gram models (iii) 238 chunking, a training dataset is annotated with POS 188 word hierarchy is unaddressed in some of the 239 tags and chunk class tags and the chunker model is 189 models.

190

211

191 192 we propose certain advantages: (i) Use of prior 243 themes in the document and capture the relevant 193 domain knowledge in application of unsupervised 244 phrases from the text. Hence unlike statistical 194 models will generate insights that will be relatable 245 chunking that requires annotated data, finite 195 to the domain user and allow for better 246 element chunking requires POS tags of tokens. ¹⁹⁶ interpretation (ii) rule-based chunks ¹⁹⁷ flexibility of size. Instead of rigidly defining n in n- ²⁴⁸ grammar structure which becomes the framework 198 grams, rule-based chunks can be unigrams, 249 for chunk phrase extraction. E.g. a noun phrase 199 bigrams trigrams and so on (iii) many text 250 chunk (NP chunk) will consist of a noun token, 200 documents cannot be classified into a single class 251 adjective token etc. whereas a verb token will rather contain subsections that can be classified 252 consist of verb token, adverb token and noun token. 201 into different classes. Chunking based 253 203 classification enables splitting a document into 254 ²⁰⁴ various subsections and separately classifying the ²⁵⁵ combined with unsupervised models, an ensemble 205 subsections into different classes (iv) rule-based 256 framework of Finite-State Rule Based chunker 206 chunks can incorporate word associations and 257 combined with clustering is proposed. ²⁰⁷ represent the hierarchy of insights. (v) this prevents $_{258}$ 3.2 208 high-dimensionality of the embedding vectors for 209 each document, this reducing the risk 210 overfitting.

212 213 proposes application of various NLP models on 263 phrases from text using Finite-State Rule Based 214 them for pattern recognition and knowledge ²⁶⁴ Chunking. From the rules, tag patterns are 215 discovery. In the case study presented here, in 265 identified and the POS tags of words are then used 216 analysis of incident investigation reports, to 266 to identify the chunk phrases of words that conform ²¹⁷ identify the events that form a chain that leads to an ²⁶⁷ to the tag patterns. NLTK Python library is used for 218 accident, the reports are split into various chunks ²⁶⁸ this purpose. Hence the chunks encapsulate prior ²¹⁹ where each chunk is representative of an event. ²⁶⁹ domain knowledge. These chunks then become the 220 Then K-means clustering is applied on the chunks ²⁷⁰ new tokens/features for text mining models to 221 to identify the major attributes. A comparison is 222 made between clustering of original text 272 223 documents vs clustering of text chunks, and it is 273 ²²⁵ whole documents, provides better performance.

Methodology 226 3

227 3.1 Theory of chunking

229 overlapping segments from a document. These 281 building block of the semantic structure of the 230 segments are the basic non-recursive phrases and 282 original document. A collection of such chunks can 231 named entities that correspond to major parts-of- 283 be considered as a new corpus. The words in the 232 speech. Hence a chunk is a set of tokens that form 284 chunk corpus are then TF-IDF vectorized to form a 233 a syntactically correlated phrase. 234

235 236 State Rule Based Chunking and statistical 288 then analyzed to identify primary attributes for

240 trained on it. On the contrary, in Finite-State Rule 241 Based chunking, a set of hand-crafted syntactic By choosing Finite-State Rule Based Chunking, 242 rules is used to identify the prior knowledge based have 247 These hand-crafted rules define the chunk

In this study, since domain knowledge is being

Methodology

of ²⁵⁹ The first step of the methodology involves ²⁵⁹ understanding the problem statement. This enables 261 rule crafting as per the objectives and domain Once the text is split into chunks, the model ²⁶² expertise and then use the rule for extraction of 271 apply on.

The chunks are then converted into feature 224 observed that clustering of chunks, instead of 274 vectors. As the chunk phrases extracted will 275 contain differences owing to presence of different 276 words or morphological variants of the same word, 277 it is necessary to group semantically similar chunks 278 and for that clustering is utilized. For that purpose, 279 we postulate that each chunk/phrase can be 228 Chunking is the process of extracting non- 280 considered as separate document representing a 285 global TF-IDF matrix. K-Means clustering is ²⁸⁶ applied on the TF-IDF vectors and the chunks are A chunker is developed in two ways: Finite- 287 classified into clusters. The cluster word-clouds are

each cluster and give an appropriate cluster label. 339 ²⁹⁰ The label is then applied to annotate all the chunks ³⁴⁰ unsafe event is as follows: 291 belonging to that cluster. Once the chunks are 341 annotated with the parent cluster labels, the chunks 342 <adjective><noun><verb><noun/adverb> 292 are re-assembled to form the original document and ³⁴³ ²⁹⁴ the chunks are then substituted with the respective ³⁴⁴ ²⁹⁵ cluster labels. These cluster labels then become the ³⁴⁵ representing the unsafe event can be extracted. ²⁹⁶ new tokens of the document. Then considering that ³⁴⁶ 297 the document is made of these new cluster label 347 4.2 tokens, the tokens are TF-IDF vectorized and K-298 299 Means clustering is carried out to cluster the original documents. 300

4 301 events at a steel plant 302

305 happened at a steel plant from 2015-2018. 636 356 1923 unique tokens were obtained. Hence the 306 incident descriptions were collected. These 357 vectorization leads to a 1923 x 1 vector generated 307 described the root cause behind accidents and what 358 for each chunk. The global corpus embedding 309 310 events/activities. 311

312

Finite-State Rule-Based Chunking 313 **4.1**

314 After pre-processing of the text data, the POS tags 315 of words are noted. Identifying the POS tag of each ³¹⁶ word is crucial to do chunking which will use these ³¹⁷ POS tags to identify the tag patterns of the various 318 chunks.

319

An accident is caused by escalation of a 320 321 sequence of events that lead a Hazardous Element to cause harm. These activities an actor and an 322 323 action. The action is represented by a verb phrase and the actor is represented by a noun phrase. So 324 there are two parts to this: A noun phrase 326 (represented by adjective-noun combination) that 327 cause the verb phrase (verb-adverb combination) to 328 occur.

329

For Example, if the following incident was 330 331 recorded "Toxic gas leaked without being 332 detected" the verb-adverb combination comprising 333 the Verb Phrase "leaked without being detected" 334 represents the action that can escalate to an 335 accident. The actor is represented by "Toxic gas" 336 which is an adjective-noun combination forming 337 the Noun Phrase (NP). The Verb Phrase was caused 338 by the Noun Phrase.

Hence the tag pattern to be developed for the

Using this tag pattern all such chunks/phrases

K-Means clustering using TF-IDF vector

349 In this method, each chunk is considered as a 350 separate document and corpus from all these **Case Study: Identification of chain of**³⁵¹ chunks are formed. These chunk corpus is then TF-352 IDF vectorized to obtain the TF-IDF matrix of all ³⁵³ the unique words in the corpus. In the steel plant ³⁰³ To demonstrate chunk effectiveness, we have ³⁵⁴ accident report corpus, 1708 chunks were considered a corpus of accident reports that 355 identified. Listing the 1708 chunks into a corpus, were the sequence of unsafe events/activities that ³⁵⁹ matrix is thus of size 1708 x 1923. Using elbow escalated to the accident. The objective is to cluster 360 method, the optimal number of clusters obtained is the documents by using the properties of the unsafe 361 28, and K-means clustering is performed. As the ³⁶² clusters are formed in a 1923-dimensional space, to ³⁶³ plot them on a 2D space, t-SNE algorithm is used. ³⁶⁴ The clusters are shown in Fig 1. Word-clouds are 365 obtained for the clusters to identify the keywords 366 that identify the major attributes of the clusters. The 367 word-clouds are shown in Fig 2.

> 369 These attributes will be used to annotate the 28 370 clusters with appropriate cluster labels which will 371 automatically annotate all the chunks present in 372 that cluster and derive the information present in 373 various chunks. The cluster labels are given in 374 Table 1.



Figure 1: Distribution of clusters

368

due area gerator Line	olax metal- slide gate adle sicture closing energency	ladles	gas leakage observed	starteohot Talling metal
Cluster0	Cluster1	Cluster2	Cluster3	Cluster4
seal injury Water turn injury		cast house tap hole closed cast	fireS1te	file file caught calle extinguished
Cluster5	Cluster6	Cluster7	Cluster8	Cluster9
pressure flow mig lance 835 mean value Solie Solie Solie Solie Solie Solie	car sus Slag retain service of slag pot	pressure furnace	started burner	servicetaken Line pipe Lances walve
Clus-	Clus-	Clus-	Clus-	Clus-
ter10	ter11	ter12	ter13	ter14
observed conner area stoke-observed	reatment on error weather the released of the	tank startedeid	energency turret with position time the lade	metal hot metal
Clus-	Clus-	Clus-	Clus-	Clus-
ter15	ter16	ter17	ter18	ter19
CO gas ppm co	cooling _{due} SVSTEM	Vessel	steam fine	happened joint Incident
Clus-	Clus-	Clus-	Clus-	Clus-
ter20	ter21	ter22	ter23	ter24
	battery Sound explosion	pressure walvefound 211 fan		
Clus-	Clus-	Clus-		
ter25	ter20	ter2/		

Figure 2: Cluster word-clouds

375 376 word cloud keywords

Clusters	Major attributes
Cluster0	Equipment malfunction
Cluster1	Slide gate operations issue
Cluster2	Crane, ladle and car and related
	operations related issues
Cluster3	Metal, oil and gas leakage
Cluster4	Falling of liquid metal
Cluster5	Hazardous water related
	incidents
Cluster6	Gas leakage, exposure, ignition
	related incidents
Cluster7	Cast and tap hole operations
	issues
Cluster8	Related to personnel activities
	and workers' exposure to various
	safety issues
Cluster9	Fire related issues

Cluster10	Miscellaneous activities
	(cluster can be ignored in analysis)
Cluster11	Slag and slag equipment related
	issues
Cluster12	Furnace related issues
Cluster13	Equipment operation issues
Cluster14	Pipe, line and pipe equipment
	related problem
Cluster15	Flame and smoke related issues
Cluster16	Administering of first aid
Cluster17	Tanks and chemicals related
	issues
Cluster18	Emergency situations created
	w.r.t equipment, personnel,
	hazardous substances
Cluster19	Hot metal/molten metal related
	issues
Cluster20	Carbon monoxide gas
	leakage/exposure
Cluster21	Miscellaneous activities of
	equipment (cluster can be ignored)
Cluster22	Vessel related issues
Cluster23	Valve related issues
Cluster24	Various kinds of incidents
	happening
Cluster25	Issues due to high/low pressure
Cluster26	Battery and explosion related
	issues
Cluster27	Issues due to air equipment, air
	pressure etc

377

383 384

385

A sample case is demonstrated here where the 378 Table 1: Annotation of the clusters based on ³⁷⁹ cluster labels for the chunks are used to define the 380 new tokens. The original document and chunks ³⁸¹ identified are given in Table 2. The chunk trees are ³⁸² developed as demonstrated in Fig 3.

Table 2:	Example	of chunk	tree	formation

Incident Description				ACM	[chunks	
L1	level indicator		['level indicator			
shown	5	%,	all	dept1	shown	pump',
pumps	trip	oped	and	diesel	'diesel	pump
pump started			started']			

After identifying the chunks in the document 386 387 the chunks are substituted with respective cluster ³⁸⁸ label to obtain the new standardized text document, 389 as shown in Table 3:

390



Figure 3: Chunk tree for accident causing mechanism

Table 3: Clusters of chunks using K-means 391 392 clustering

Chunk	Cluster
level indicator shown pump	0
diesel pump started	13

Hence the corpus of original text will look as 394 430 395 shown in Table 4: 431

396

393

432 Table 4: Documents with cluster labels as 397 433 398 tokens

Description	434
C6 C8 C0 C26 C24	435
C0 C13	437
C25 C23	438
C6 C26 C0 C23 C23 C23 C23 C0	439
C26 C11 C9	440
C19	441
	442

399

So now each document in the corpus is 400 444 401 considered to consist of cluster labels only as the 402 new words. These labels are then tokenized and 403 embedded using TF-IDF to generate embedding 447 404 matrices of size 636 x 28, where 636 is the number 405 of documents, 28 is the number of cluster labels ⁴⁰⁶ representing the tokens and for each document, the 407 IF-IDF gives the term frequency-inverse document ⁴⁰⁸ frequency of each cluster label. On this embedding 409 matrix, K-means clustering is carried out. This 410 clustering is based on chunk as features which 411 incorporate information about unsafe events, hence 412 documents showing similarity in unsafe events will 448 413 get clustered together.

Results and Discussion 414 5

415 First, to analyze if splitting a text document into 416 chunks provides any advantages, clustering of the 452 Traditional NLP models classify a text document 417 corpus obtained by listing of all chunks is 453 into various classes. However, at times many text 418 compared with clustering of the corpus of original 454 documents have a structure where various portions 419 documents. Two cluster validity indices- Davies- 455 of the document belong to various classes and it is

used to evaluate the clustering models. The results are showcased in Table 5: 422

Table 5. DBI and SI	values	faluctor	ing models
Table J. DDI allu SI	values 0	i ciusici	mg moucis

		TFIDF	based
	TFIDF based	clustering	of
	clustering of	original	
	chunks	documents	
DBI	4.596	5.405	
SI	0.031	0.017	

The following observations can be made:

- A lower DBI score indicates better a) clustering. Hence it can be observed that chunking + clustering is giving better clusters compared to clustering original documents
- If the SI value is closer to 1, the b) clustering is better. Higher value of SI for clustering of chunks suggests that chunking and clustering is better than clustering of original documents

Then clustering performance of original and documents based on the 636 x 28 embedding 40 matrix is compared with clustering of TF-IDF ¹⁴¹ embedding of the original documents. The SI and ¹⁴² DB index of both the clustering is shown in Table 443 6:

Table 6: Clustering performance of chunk-446 documents vs Original Documents

	Clustering of documents using the cluster-label embedding matrix	TFIDF clustering original documents	based of
DBI	2.7885	5.405	
SI	0.1208	0.017	

Hence it can be seen that clustering based on 449 450 chunks is giving better results.

451 6 Conclusion

420 Bouldin index (DBI) and Silhouette Index (SI)- are 456 not feasible to classify an entire text document into 457 a class. While identifying chain of events,

458 identification of these subclasses is important. 511 represent help in customized information discovery 459 Hence this framework has been proposed which 512 from the document. 460 enables classification of various subsections into 513 461 various classes and enables development of chain 514 462 of events. Demonstration of the framework is 515 corpus embedding using TF-IDF generated 1708 x 463 conducted in form of a case study where the study 516 1923 embedding matrix, while using chunking + 464 was conducted to find out how to do text mining of 517 clustering and then doing TF-IDF embedding 465 safety text data generated by an organization and 518 generates 1708 x 28 embedding matrix. This 466 generate chain of events/accident paths with 519 reduces the risk of high dimensionality and 467 minimal human involvement. In the case study, 520 overfitting. 468 unlike other NLP models which generate numerical ⁴⁶⁹ values of tokens based on a global corpus, here the ⁵²¹ 7 470 rules obtained from hazard theory have been used 471 to guide the formation of initial dataset which is 522 Despite the advantages, the model suffers from a 472 then analyzed. On changing domains, the rules also 523 few drawbacks. Removal of stop words and 473 change, hence this model is useful to obtain domain 524 connecting words also helps in chunks becoming 474 specific insights.

475 476 477 specific lexical rule-based chunking 478 unsupervised NLP tools, this model aims to bridge 529 written in a haphazard manner, then this model will 479 the drawbacks of both supervised and unsupervised 530 not give optimal results. Thirdly, using K-Means 480 approaches:

481

503

One major drawback of supervised modelling is 533 482 483 annotation of huge corpus of training documents. 534 484 By splitting a document into chunks and 535 future research directions are proposed: 485 performing clustering on them, all similar chunks, 536 486 representing a particular concept, get annotated at 537 487 once just by annotating the clusters. Secondly, text 538 488 documents can have multiple annotations for 539 489 various sub-sections which might not get captured 540 490 during manual annotation. Splitting the documents 541 491 into chunks and annotating using clustering allows 542 492 for a text document to have multiple annotations 543 for different subsections. 544 493 494 545

Unsupervised models are domain independent, 495 496 hence combining these models with chunking 546 References enables the user to incorporate user's prior 547 Ronen Feldman, and Ido Dagan. 1995. Knowledge knowledge of domain features into generation of 548 498 499 insights. This helps in domain relevant 549 500 interpretation of insights and thus resolve the 550 difficulty of interpretation of insights obtained 551 Mehdi Allahyari, Seyedamin Pouriyeh, Mehdi Assefi, 501 552 from unsupervised models. 502 553

554 While unsupervised models are useful to model 504 505 latent topics and themes, one of the key issues is 556 506 difficulty of modelling hierarchy and relationship 557 Yu Zhang, Mengdong Chen, and Lianzhong Liu. 2015. 507 between themes. Chunking of a document 558 508 organizes topics and themes in a hierarchial 559 $_{\rm 509}$ manner, which will then be discovered using NLP $^{\rm 560}$ 510 models. Chunks and the relationships they 561

Finally, in the case study, the original document

Limitations and Future Work

525 less user friendly to understand. Secondly, the 526 model makes an assumption that the descriptions By adopting an ensemble approach of domain 527 are written in a sequence, and the sequence is used with 528 to obtain the chain of events. If the descriptions are 531 clustering does not always capture the semantic 532 nature, even in chunks.

Keeping in mind the limitations, the following

- a) Future work will involve how to remove these drawbacks while doing chunking analysis.
- Future work will include how to better b) cluster chunks capturing the semantic nature.
- Application of other feature vector c) generation models on chunking to improve insights generated.

- Textual Databases Discovery in (KDT). In KDD (Vol. 95, pp. 112-117).
- Saied Safaei, Elizabeth D. Trippe, Juan B. Gutierrez, and Krys Kochut. 2017. A brief survey of text mining: Classification, clustering and extraction techniques. arXiv preprint arXiv:1707.02919.
- A review on text mining. In 2015 6th IEEE International Conference on Software Engineering and Service Science (ICSESS), pp. 681-685. IEEE.

562 Fabrizio Sebastiani. 2002. Machine learning in 622 automated text categorization. ACM computing 623 563 surveys (CSUR) 34, no. 1: 1-47. 564 624 565 625 566 Ammar Ismael Kadhim. 2019. Survey on supervised 626 Mandar Mitra, Chris Buckley, Amit Singhal, and machine learning techniques for automatic text 627 567 classification. Artificial intelligence review 52, no. 628 568 1:273-292. 569 629 570 571 Thomas Hofmann, 1999. Probabilistic latent semantic 631 indexing. In Proceedings of the 22nd annual 632 572 international ACM SIGIR conference on Research 633 573 and development in information retrieval. 574 634 575 635 2012. 576 David M. Blei. Probabilistic topic 636 models. Communications of the ACM, 55(4), pp.77- 637 577 84 578 638 579 639 580 Jonathan Chang, Sean Gerrish, Chong Wang, Jordan 640 Boyd-Graber, and David Blei. 2009. Reading tea 641 581 leaves: How humans interpret 582 models. Advances in neural information processing 643 583 systems, 22. 584 644 645 585 586 Filipe Rodrigues, Mariana Lourenco, Bernardete 646 Ribeiro, and Francisco C. Pereira. 2017. Learning 647 587 supervised topic models for classification and 648 Ioannis Antonellis, and Efstratios Gallopoulos. 2006. 588 regression from crowds. IEEE transactions on 649 589 pattern analysis and machine intelligence 39, no. 650 590 12:2409-2422. 591 651 592 593 Ning Chen, Jun Zhu, Fei Xia, and Bo Zhang. 2014. 653 594 Discriminative relational topic models. IEEE 654 transactions on pattern analysis and machine 655 595 intelligence 37, no. 5: 973-986. 596 656 657 597 Changyou Chen, Wray Buntine, Nan Ding, Lexing 658 598 Xie, and Lan Du. 2014. Differential topic 659 599 models. IEEE transactions on pattern analysis and 660 600 machine intelligence 37, no. 2: 230-242. 601 602 662 Anna Stavrianou, Periklis Andritsos, and Nicolas 663 603 Nicoloyannis. 2007. Overview and semantic issues 664 604 of text mining. ACM Sigmod Record 36, no. 3: 23- 665 605 34. 666 606 607 608 Gerard Salton, Anita Wong, and Chung-Shu Yang. 668 1975. A vector space model for automatic 669 609 indexing. Communications of the ACM 18, no. 11: 670 610 613-620. 611 671 612 613 Catherine Blake, and Wanda Pratt. 2001. Better rules, 672 fewer features: a semantic approach to selecting 614 features from text. In Proceedings 2001 IEEE 673 615 International Conference on Data Mining, pp. 59-616 66. IEEE. 617 618 619 Maria Fernanda Caropreso, Stan Matwin, and Fabrizio Sebastiani. 2001. A learner-independent evaluation 620 of the usefulness of statistical phrases for automated 621

text categorization. Text databases and document management: Theory and practice 5478, no. 4: 78-102.

- Claire Cardie. 1997. An analysis of statistical and syntactic phrases. In RIAO, vol. 97, pp. 200-214.
- 630 Philipp Cimiano, Andreas Hotho, and Steffen Staab. 2005. Learning concept hierarchies from text corpora using formal concept analysis. Journal of artificial intelligence research 24: 305-339.
 - Athanasios Kehagias, Vassilios Petridis, Vassilis G. Kaburlasos, and Pavlina Fragkou. 2003. A comparison of word-and sense-based text categorization classification using several algorithms. Journal of Intelligent Information Systems 21: 227-247.
- topic 642 Martin Rajman, and Romaric Besançon. 1999. Stochastic distributional models for textual information retrieval. In Proc. of 9th International Symposium on Applied Stochastic Models and Data Analysis (ASMDA-99), pp. 80-85.
 - Exploring term-document matrices from matrix models in text mining. arXiv preprint cs/0602076.
 - 652 Stephan Bloehdorn, Philipp Cimiano, and Andreas Hotho. 2006. Learning ontologies to improve text clustering and classification. In From Data and Information Analysis to Knowledge Engineering: Proceedings of the 29th Annual Conference of the Gesellschaft für Klassifikation eV University of Magdeburg, March 9-11, 2005, pp. 334-341. Berlin, Heidelberg: Springer Berlin Heidelberg.
 - 661 Man Lan, Chew Lim Tan, Jian Su, and Yue Lu. 2008. Supervised and traditional term weighting methods for automatic text categorization. IEEE transactions on pattern analysis and machine intelligence 31, no. 4:721-735.
 - 667 Scott Deerwester, Susan T. Dumais, George W. Furnas, Thomas K. Landauer, and Richard Harshman. 1990. Indexing by latent semantic analysis. Journal of the American society for information science 41, no. 6: 391-407.