

Article

A Novel Framework to Detect Irrelevant Software Requirements Based on MultiPhiLDA as the Topic Model

Daniel Siahaan ^{†*}  and Brian Risqi Paradisiaca Darnoto

Informatics Department, Institut Teknologi Sepuluh Nopember, Surabaya 60111, Indonesia

* Correspondence: daniel@if.its.ac.id; Tel.: +62-31-5939214

† Current address: Gedung Informatika, Jl. Teknik Kimia, Kampus ITS Sukolilo, Surabaya 60111, Indonesia.

Abstract: Noise in requirements has been known to be a defect in software requirements specifications (SRS). Detecting defects at an early stage is crucial in the process of software development. Noise can be in the form of irrelevant requirements that are included within a SRS. A previous study had attempted to detect noise in SRS, in which noise was considered as an outlier. However, the resulting method only demonstrated a moderate reliability due to the overshadowing of unique actor words by unique action words in the topic-word distribution. In this study, we propose a framework to identify irrelevant requirements based on the MultiPhiLDA method. The proposed framework distinguishes the topic-word distribution of actor words and action words as two separate topic-word distributions with two multinomial probability functions. Weights are used to maintain a proportional contribution of actor and action words. We also explore the use of two outlier detection methods, namely Percentile-based Outlier Detection (PBOD) and Angle-based Outlier Detection (ABOD), to distinguish irrelevant requirements from relevant requirements. The experimental results show that the proposed framework was able to exhibit better performance than previous methods. Furthermore, the use of the combination of ABOD as the outlier detection method and topic coherence as the estimation approach to determine the optimal number of topics and iterations in the proposed framework outperformed the other combinations and obtained sensitivity, specificity, F1-score, and G-mean values of 0.59, 0.65, 0.62, and 0.62, respectively.

Keywords: angle-based outlier detection; percentile-based outlier detection; multiphilda, noise; irrelevant software requirements



Citation: Siahaan, D.; Darnoto, B.R.P. A Novel Framework to Detect Irrelevant Software Requirements Based on MultiPhiLDA as the Topic Model. *Preprints* 2022, 1, 0. <https://doi.org/>

Publisher's Note: MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Copyright: © 2022 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

1. Introduction

The requirements specification process is one of the key stages in a software development project that determines its success. Researchers indicate that 40-60% of software development project failures originate in requirements specification [1]. Requirements dictate how the product is designed and implemented in the following stages [2]. Overlooking software requirements during the requirements specification process can cause future threats and failures in the operational phase [3]. Furthermore, the cost of detecting and correcting defects increases exponentially as the software progresses along the software development life cycle (SDLC) [4]. Therefore, it is essential and critical to carry out an effective requirements specification.

In an object-oriented approach, requirements engineers deliver a SRS, which is a document that serves as a guideline for the subsequent processes of software development [5]. Therefore, the SRS should have a set of quality attributes in order to maintain the quality of the end product. Nevertheless, requirements engineers often fail to comply with the quality attributes and produce ambiguity, contradictions, forward references, over-specifications, and noise in SRS [6,7]. Meyer [7] argues that the use of natural language to specify software requirements is the cause of these problems. This study focuses on noise as it is a mistake commonly made by requirements engineers. According to Meyer, noise refers to a condition where a SRS contains obscurity [8].

There are two types of noise in SRS, namely irrelevant requirements and non-requirement statements. Irrelevant requirements in SRS are requirements stated in the document that discusses topics beyond the scope of the software project. On the other hand, non-requirement

statements in SRS are elements in the document that are of relevance to the domain of the soft-ware being developed, but are not requirements. Manek and Siahaan [9] proposed a method to identify noise in SRS using spectral clustering. However, it was indicated in the study that the method was unable to effectively identify noise due to two reasons. Firstly, the method does not take into consideration the difference between irrelevant requirements and non-requirements statements. Secondly, the method disregards the different parts of requirements, namely actors and actions, during the identification process.

Natural Language Processing (NLP) has played an important role in human-computer interaction and artificial intelligence. It is often implemented in the domain of language structure analysis, technology-based machine translation, and language recognition [10]. Several unsupervised methods have been proposed to identify and model topics within textual documents. One of the most popular methods for topic modeling is Latent Dirichlet Allocation (LDA) [11,12]. The basic idea of the LDA method is that each document is represented as a random mixture of hidden topics, where each topic is characterized based on the distribution of words that lie within the topic itself. The LDA model generates a distribution of topics for the documents that can be immediately used as features for classification. The combination of topic and word frequency information may result in an enhanced set of features and in turn increase the performance of text classifiers [13]. Schröder [14] developed an LDA-based document classification method for accurate topic-based tagging of scientific articles. The proposed method is called CascadeLDA. It was indicated in the study that CascadeLDA required significantly fewer number of iterations compared to another extension of LDA, namely Labeled LDA (L-LDA), but relies heavily on the label structure and general characteristics of the dataset. In the studies conducted by Suri and Roy [15] and Chen *et al.* [12], the LDA method was compared to Negative Matrix Factorization (NMF) for data detection with a large text flow. It was shown that the semantic results produced by LDA were more meaningful than the NMF method.

Irrelevant requirements can be viewed as outliers since they discuss topics that are not within the scope of a system, in which the word components of irrelevant topics deviate from the word components of relevant topics. Several studies have proposed methods to detect outliers. Kriegel *et al.* [16] proposed the ABOD method for the detection of outliers. ABOD works on high-dimensional data. This method evaluates the degree of outliers based on the variance of the angle (VOA) formed by the target object and all other object pairs. The smaller the angular variance of the object, the more likely the object is an outlier. A recent study combined the EBOD method with several different methods and proposed an ensemble-based outlier detection (EBOD) method to detect outliers in noisy datasets [17]. In addition to ABOD, another method for outlier detection is the PBOD method. Sharma *et al.* [18] developed a PBOD method to reduce noise in large data sets and explored the use of machine learning methods, namely Gradient Boosting Regression (GBR) and Random Forest (RF), to predict short-term waiting times at the US-Mexico border. The results encourage to incorporate more advanced prediction algorithms and methods into the data set. Another research went further and combined percentile and LDA to predict Days on Market, which is a measure of liquidity in the real estate industry [19].

This study is a continuation of our previous work [20]. The aim of this study is to identify noise in the form of irrelevant requirements within SRS documents. This study proposes a new noise detection framework which relies on a modification of LDA, called MultiPhiLDA. The proposed framework distinguishes actors and actions as separate word distributions with two multinomial probability functions. Therefore, the process of clustering topics becomes more efficient due to the omission of irrelevant requirements within the SRS documents.

2. Related Works

Schröder [14] developed an LDA-based document classification model to aid in accurately tagging academic publications. The proposed LDA model is named CascadeLDA. CascadeLDA utilizes hierarchical labels and adopts different Hierarchically Supervised

LDA (HSLDA) mechanisms. The prediction problem with the extensions discussed earlier is that the topic-word distribution (θ_k) does not seem to differ between labels. CascadeLDA aims to differentiate features between the same nodes by enlarging one label (non-leaf) at a time.

Manek and Siahaan [9] proposed a method to detect noise in SRS using spectral clustering. The method was tested on 648 requirements that were manually extracted from 14 SRS documents. However, based on their experiment, the method showed moderate reliability. They suggested that this is due to the suboptimal pre-processing text process, in which they did not preserve the actors and actions in the requirement statements as a distinctive element of a requirement. Hence, actors or actions with noise can cause a requirement to become irrelevant.

Fahmi and Siahaan [8] compared the performance of several supervised methods in classifying non-requirement statements in SRS. They used the dataset that was previously used by Manek and Siahaan [9]. The classification methods that were analyzed were Naïve Bayes (NB), Support Vector Machine (SVM), Decision Tree, Random Forest (RF), and k-Nearest Neighbour (KNN). Their experiment showed that SVM produced a better classification model compared to the other methods to identify non-requirement statements in SRS.

Li *et al.* [21] developed a new model called the Common Semantics Topic Model (CSTM) to solve noise problems found in short texts and analyzed the performance of the model on three real-world datasets. They assumed that each short text is a mixture of the topic of the selected function and all common topics. According to Li, words that have high-frequency domains is a type of noise. CSTM defines a new type of topic to capture standard semantics and noise words.

Liu *et al.* [22] proposed Collaboratively Modeling and Embedding – Dirichlet Multinomial Mixtures (CME-DMM) to apply topic modeling on short texts that were collected from a social media platform. Liu identified major challenges in modeling latent topics in short texts, including the inadequacy of word co-occurrence instances, the need to capture local contextual dependencies in noisy short texts due to their limited length, and maintaining consistency between the latent topic models based on word co-occurrences and their representation based on local contextual dependencies. The proposed method overcomes these challenges by considering both the global and contextual word co-occurrence relationship and also ensuring a high consistency between the topic and word embeddings in the latent semantic space.

Vo and Ock [23] classified short text documents using topic modeling. They used the LDA method before performing classification. They applied three machine learning methods, namely SVM, NB, and KNN for the classification of short text documents three universal data sets and compared the performances of the three methods. The results using SVM are better in terms of accuracy compared to the other two methods.

Albalawi *et al.* [24] conducted a comparative analysis on the application of topic modeling methods on short text data. With the use of use of the LDA and Non-negative Matrix Factorization (NMF) methods, higher quality and coherence topics were obtained compared to the other methods. However, the LDA method was more flexible and provides more meaningful and logical extracted topics, especially with fewer topics. Both LDA and NMF showed similar performance, but LDA was more consistent. As a result, the LDA method outperformed other topic modeling methods with most number of features. The experimental results showed that LDA was able to obtain higher F-score values compared to the other methods. Furthermore, the LDA method was able to define the best and most precise meaning than the other topic modeling methods.

3. Noise Detection Framework

This section details the noise detection framework which relies on MultiPhiLDA to detect noise in SRS. Fig. 1 shows that the input to the framework is a list of requirements which may be irrelevant or relevant requirements. In order to differentiate between ir-

relevant and relevant requirements, this framework has six main processes. This section elaborates these processes in more detail.

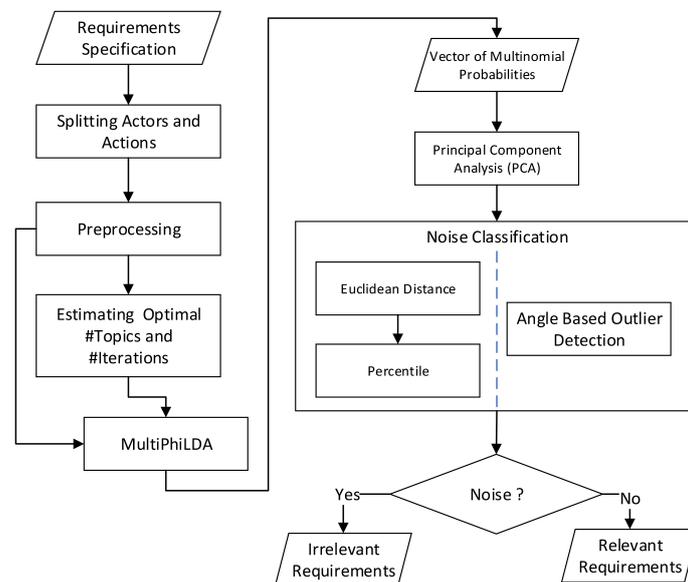


Figure 1. The noise detection framework.

3.1. Actor-Action Splitting

In the first step, the framework begins by identifying the actor and action parts within each requirement statement with the use of syntactic dependency parsing. In this study, we implemented a python code that loads the spacy library. Listing 1 shows how the spacy library was used for syntactic dependency parsing between words within a sentence. In the example, the sentence “All staff can view the details of any student.” is parsed. A word with ‘nsubj’ tag would be classified as actor, while the rest would be classified as actions.

Listing 1: Extracting actors and actions from requirements statement.

```

1 #requirements: a textual dataset (contains requirement statements
2 )
3 #[actors , actions]: lists of actors or actions extracted from
4 requirements
5
6 for requirement in data_skpl :
7     postag = spaCy(requirement)
8     for words in postag:
9         if words.type == 'nsubj' :
10            actor.append(words)
11         else :
12            action.append(words)

```

Fig. 2 shows the result of the parsing process on the example sentence. The word ‘staff’ has a ‘nsubj’ dependency, which means that the word ‘staff’ is the subject who is doing the action and along with the determiner ‘all’ is identified as an actor part. The word ‘view’ has a ‘ROOT’ dependency, which means that the word ‘view’ and its direct object (the accusative object of ‘view’) is identified as an action part. After conducting syntactic dependency parsing on each sentence, the action and actor parts are further treated separately.

```
All tag : det
staff tag : nsubj
can tag : aux
view tag : ROOT
the tag : det
details tag : dobj
of tag : prep
any tag : det
student tag : pobj
. tag : punct
```

Figure 2. Results of syntactic dependency parsing using the SpaCy Library.

3.2. Pre-processing

In the second step, each action and actor parts are then pre-processed. The pre-processing involves case folding, punctuation removal, tokenization, lemmatization, and stopwords removal. In the case folding process, all the words are normalized to lowercase. In the punctuation removal process, all numbers, symbols, and punctuations are removed. In the tokenization process, all the words are split individually into tokens. For example, the actor part of the example sentence “all staff” is split into two tokens ‘all’ and ‘staff’, while the action part of the example sentence “can view the details of any student” is split into six tokens ‘can’, ‘view’, ‘the’, ‘details’, ‘of’, ‘any’, and ‘student’. After tokenization, lemmatization aims to reduce the words to their base form with the use of a vocabulary and morphological analysis. The last step in the pre-processing step is stopwords removal, in which non-substantial elements within a text, such as articles, conjunction, preposition, modals, and pronouns are removed.

3.3. Optimal Number of Topics and Iterations Estimation

Before determining the optimal number of topics and iterations, unique words from every actor and action parts need to be extracted separately and a set of actor words and a set of action words are obtained. Each set of unique words are stored as a separate dictionary of unique words. This process employs the gensim library to create the dictionaries. The dictionaries of unique words are later used to determine the optimal number of topics and the number of iterations in the application of MultiPhiLDA.

This study compares three estimation approaches to determine the optimal number of topics and iterations, namely perplexity, topic-coherence, and human judgment. Perplexity and topic coherence are commonly used to measure the performance of topic modeling methods [24,25]. Perplexity is considered to be a good approximation for a case where test data is similar to training data. While perplexity only considers intrinsic measurements, topic coherence also considers extrinsic measurements. These two measurements are supposed to make a distinction between semantically expoundable topics and factitious topics of statistical inference. In this study, we also compared the two estimation approaches against human judgement. Software engineers are involved in determining the optimal number of topics and iterations required to cluster the requirements based on its topics.

3.4. MultiPhiLDA

MultiPhiLDA is a modification of the LDA method. The LDA method has one phi, where phi is the distribution of words in a topic. A previous research had proposed that requirements should be divided into two parts that separately represent actors and action, as these are two critical parts of a requirement. The number of unique actor words is significantly lower than the number of unique action words. If the topic-word distribution is based on unique words extracted from requirements and the topic-word distributions of action words and actor words are not defined separately, then the unique action words will overshadow the unique actor words in the topic-word distribution. Therefore, we need to create separate topic-word distributions for action words and actor words and subsequently amalgamate a new distribution from these two distributions by assigning weights to the two distributions to overcome the aforementioned problem within the new topic-word distribution. Fig. 3 shows the plate notation of the MultiPhiLDA method. It can be seen in

Fig. 3 that ϕ is two because MultiPhiLDA processes two topic-word distributions, one for actor words and the other for action words.

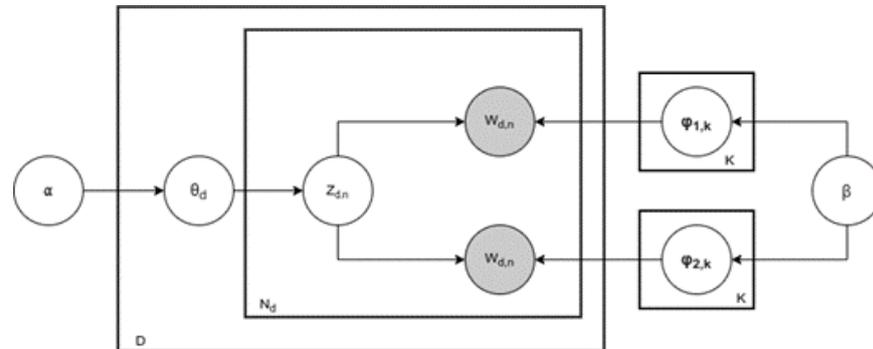


Figure 3. MultiPhiLDA plate notation.

Listing 2: MultiPhiLDA pseudocode.

```

1 dataset = data[actors , actions]
2 for all actors
3   randomly assign a topic to the actor
4 end for
5 for all actions
6   randomly assign a topic to the action
7 end for
8 for iteration
9   for all actors in dataset
10    for all actors
11      theta calculation
12      phi calculation
13      phi calculation with weight =  $\frac{\phi * actor\_weight}{\sum(\phi * actor\_weight)}$ 
14      topic determination using phi and theta
15    end for
16  end for
17  for all actions in dataset
18    for all actions
19      theta calculation
20      phi calculation
21      phi calculation with action weight =  $\frac{\phi * action\_weight}{\sum(\phi * action\_weight)}$ 
22      topic determination using phi and theta
23    end for
24  end for
25 end for

```

The proposed framework is expected to find irrelevant requirements in SRS. Requirements that have been processed and has topics assigned to it is further clustered with respect to topics. Clustering can be carried out using the K-means algorithm or other methods. The result of clustering is visualized to identify relevant and irrelevant requirements.

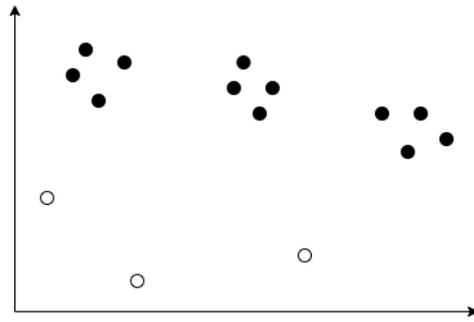


Figure 4. Visualization of MultiPhiLDA.

Fig. 4 illustrates the results of clustering by topic using MultiPhiLDA. The black dots represent relevant topics, while the white dots represent irrelevant topics. In an SRS, relevant requirements discuss relevant topics that are within the scope of the system, while irrelevant requirements discuss irrelevant topics that are outside the scope of the system. For example, a hospital information system may have relevant topics that are within the scope of the system such as patients, pharmacies, outpatients, poly, diagnosis, and disease. The visualization shown in Fig. 4 can be used to indicate which topics are considered relevant or irrelevant.

3.4.1. Separate Topic-Word Distribution for Actor Words and Action Words

After the identification of actor words and action words, two topic-word distributions are created, one for actor words and the other for action words. Firstly, the topics and the words are determined. Once the topics and words have been determined, randomly assign topics to words. The probability of topics for a document is calculated using Eq. 1, while the probability of words belonging to a topic is calculated using Eq. 2 as follows:

$$\theta_{d,k} = \frac{n_{(d)}^k + \alpha}{n_{(d)} + K\alpha} \quad (1)$$

$$\phi_{k,v} = \frac{n_{(k)}^v + \beta}{n_{(k)} + V\beta} \quad (2)$$

The next step is to apply the weight for the actor words and action words. Here, the weight is used to determine the topic-word distribution for both sets of actor words and action words. The weight for action and actor words are predetermined, where the sum of both weights equal to 1. Once the phi value is obtained, weighting is carried out using Eq. 3 as follows:

$$\phi = \frac{\phi_i \text{ weight}}{\Sigma(\phi \text{ weight})} \quad (3)$$

Eq. 4 is used to calculate the probability of a document belonging to a topic by taking into consideration the number of topics that the document may belong to. The likelihood of a word in the document belonging to a topic (w_n) is calculated by taking into consideration the number of topics that the word may belong to. For a specific w_n in document d , find the topic k whose probability has the maximum value and reassign the word to topic k .

$$p(z_i = k | z_i, w) = \phi_{k,v} \theta_{d,k} \quad (4)$$

Do as many iterations as previously estimated until two topic-word distributions of actor words and action words are obtained based on the predetermined number of topics.

3.4.2. Amalgamation of the Topic-Word Distributions

In this process, the two topic-word distributions of action words and actor words are merged into one topic-word distribution based on the topics that each word may belong

to. Each topic consists of a few actor words and action words. The words that represent a topic are determined through their probabilistic values for the given topic. Amalgamation is carried out by sorting the probabilistic values of action and actor words for a given topic, and only a few words with the highest probabilistic values are selected to represent the topic. After amalgamating the two distributions into a single distribution, the next process is to identify the irrelevant requirements in the SRS documents.

3.5. Identification of Irrelevant Requirements

This process discriminates the requirements statements that are irrelevant. In this study, we compared the performance of two approaches, namely PBOD and ABOD. Percentiles are values below a certain percentage of observations [26]. In this study, we calculated the percentile value using Euclidean distance. Euclidean distance calculates the average distance between each member of a cluster to their centroid. If the value is greater than the percentile value, then the requirement statements that belongs to the cluster are considered irrelevant requirements.

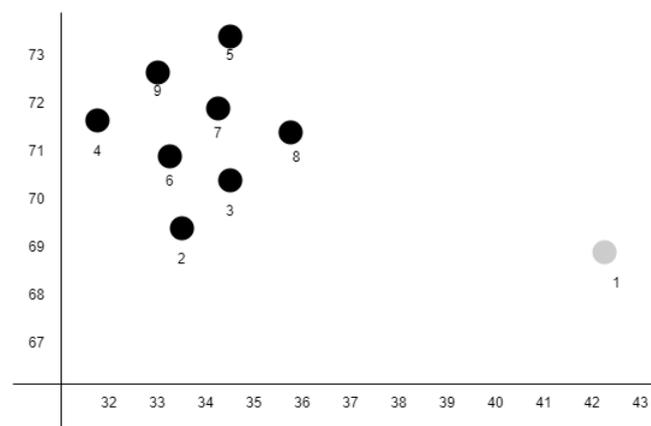


Figure 5. Point ranking in the sample dataset.

ABOD is a method to detect outliers. ABOD assigns an Angle-Based Outlier Factor (ABOF) to each point in the database and returns a list of points that had been sorted based on the ABOF of the points [16]. Fig. 5 illustrates the ranking of each point within a dataset. The top-ranking point (number 1) is the farthest outlier. The subsequent ranks are determined by the next farthest cluster boundary point. Since distance is only taken into account as a weight for the variance of angles, ABOD can succinctly detect outliers even in high-dimensional data.

4. Experimental Setup

This section explains how the experiment was carried out in this study. The first section describes the dataset used in the experiment. The second section describes the three different evaluation scenarios used to evaluate the performance of the proposed framework. These three scenarios are performed as a means of providing empirical data to indicate the performance of the proposed framework.

4.1. Dataset

In this study we used the same dataset that was used by a previous research conducted by Manek and Siahaan [9]. The dataset is publicly available [27]. The dataset contains fourteen SRS from various software development projects. However, the original dataset contains typos and grammatical errors. Therefore, we proofread each document beforehand. Table 1 shows the statistical information of each SRS document within the dataset. The number of requirement statements (#Reqs) varies highly between documents (between 6 to 106). The average length of each requirement statement (the average number of words, word) also varies. The average number of verbs (verb) in a requirement statement is between 1-2 words. The SRS documents of the different projects possess different number

of actors. It can also be seen that the dataset is highly imbalanced, where the number of irrelevant requirements (#Irr) is significantly smaller (9.26%) than the number of relevant requirements.

Table 1. Recapitulation of the Requirements in Each SRS

ID	Project Name	#Reqs	# μ -word	# μ -verb	#Actor	#Irr
DA-1	Submit Job	13	6	1	1	1
DA-2	System Administrator	17	12	1	4	5
DA-3	Archive Administrator	39	12	2	4	2
DA-4	SPG Application	6	6	1	1	3
DA-5	System Administrator Staff	33	40	1	8	3
DA-6	Software Development	65	17	2	33	5
DA-7	Display System	106	13	2	17	3
DA-8	Internet Access	64	18	2	22	7
DA-9	Meeting Initiator	33	24	1	7	4
DA-10	Online System	17	15	1	7	2
DA-11	Library System	86	17	2	32	9
DA-12	IMSETY System	70	12	2	13	3
DA-13	Manage Student Information	24	22	1	9	3
DA-14	PHP Project	75	28	1	50	10

4.2. Evaluation Scenarios

In this study, we conducted three evaluation scenarios to evaluate the performance of the proposed framework. Scenario 1 aims to evaluate the topics produced by the proposed framework and estimate the optimal number of topics and iterations. Scenario 2 aims to measure the performance of the proposed framework using four performance metrics, namely specificity, sensitivity, F1 score, and G-mean. Sensitivity is the ratio of true positive predictions to the total number of actual positive data. Specificity is the ratio of true negative predictions to the total number of actual negative data. G-mean combines specificity and sensitivity and balances both concerns. F1 Score combines precision and recall values produced by the proposed framework by taking their harmonic mean. Scenario 3 aims to measure the reliability of the proposed framework by using two reliability measurements, namely average variance extracted (AVE) and composite reliability (CR) to evaluate the consistency of the proposed framework in identifying irrelevant requirements over the fourteen SRS documents. In all the evaluation scenarios, all possible combinations of estimation approaches (perplexity, topic coherence, and human judgment) and outlier detection methods (ABOD and PBOD) within the proposed framework were individually evaluated. Furthermore, in scenario 2, a comparison of performance was carried out between the proposed framework and the method proposed by Manek and Siahaan [9].

5. Result

In this section, we evaluated the performance of the proposed framework based on the three evaluation scenarios.

5.1. Evaluation Scenario 1: Estimate the Optimal Number of Topics and Iterations

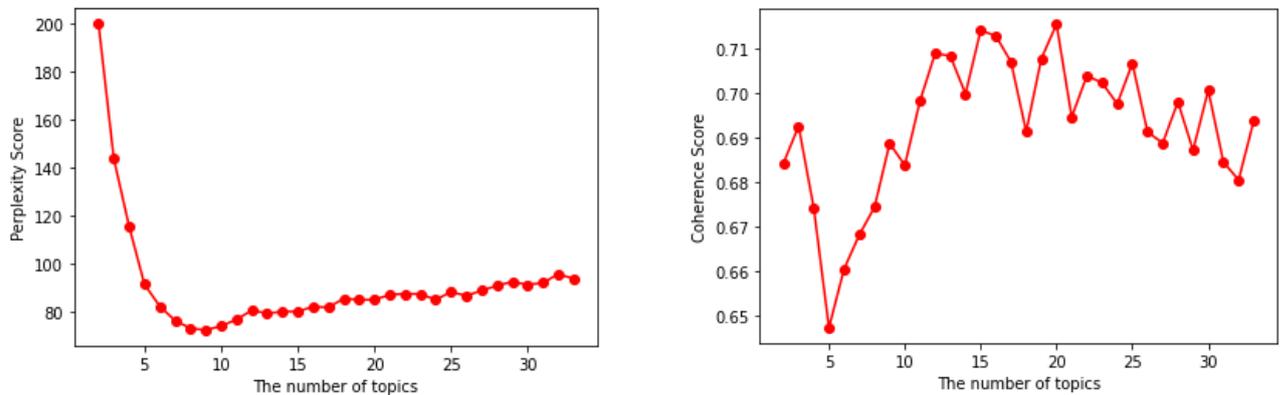


Figure 6. Visualization of topic coherence and perplexity values.

Since the MultiPhiLDA method requires the number of topics and number of iterations as its input, this scenario aims to evaluate the topics produced by the proposed framework with the use of different estimation approaches (perplexity, topic coherence, and software engineer) and estimate the optimal number of topics and iterations. In this evaluation, the minimum number of topics for each SRS is two, while the maximum number of topics is predetermined by the software engineer. Furthermore, in this evaluation, the minimum number of iterations is one, while the maximum number of iterations is 500. The evaluation was carried out on two models for each estimation approach, one that uses ABOD and another that uses PBOD as the outlier detection method. The values obtained from the two models are added and divided by two. After obtaining one principal value, the value is visualized to estimate the optimal number of topics and iterations. Fig. 6 shows the perplexity scores and the coherence scores obtained from the two models which are added and divided by two.

With the use of topic coherence, the most satisfactory model is indicated by the highest coherence score. On the other hand, with the use of perplexity, the most satisfactory model is indicated by the lowest perplexity score. The results of the optimal number of topics and iterations can be seen in Table 2.

Table 2. The Optimal Number of Topics and Number of Iterations for Each SRS Document

Noise Classifier	Perplexity		Topic Coherence		Software Engineer		Status
	#Topic	# Iteration	#Topic	#Iteration	#Topic	#Iteration	
DA-01	5	224	4	273	4	273	Not Accepted
DA-02	15	62	13	493	9	424	Not Accepted
DA-03	10	425	19	490	9	59	Not Accepted
DA-04	5	22	4	323	6	62	Not Accepted
DA-05	18	466	12	286	19	86	Not Accepted
DA-06	23	198	12	473	27	201	Not Accepted
DA-07	9	122	21	203	35	187	Accepted
DA-08	14	74	20	123	26	400	Not Accepted
DA-09	13	96	18	294	18	294	Not Accepted
DA-10	7	400	6	107	11	239	Not Accepted
DA-11	19	433	30	390	36	132	Not Accepted
DA-12	12	330	9	375	48	430	Not Accepted
DA-13	17	472	10	409	20	500	Not Accepted
DA-14	23	430	39	89	53	320	Not Accepted

After obtaining the optimal number of topics, the next step is to determine the optimal number of iterations with the use of the optimal number of topics that had been estimated by the three estimation approaches. The SRS document used to determine the optimal number of iterations must have a graph of perplexity values whose ups and downs are not too extreme and possess a gentle trend value. If the SRS document does not meet these requirements, the SRS document cannot be used. Fig. 7 is an example of a SRS document (DA-07) with a graph of perplexity values and trend values that meets these conditions. The number of requirements within the DA-07 SRS document is the largest among all the SRS documents.

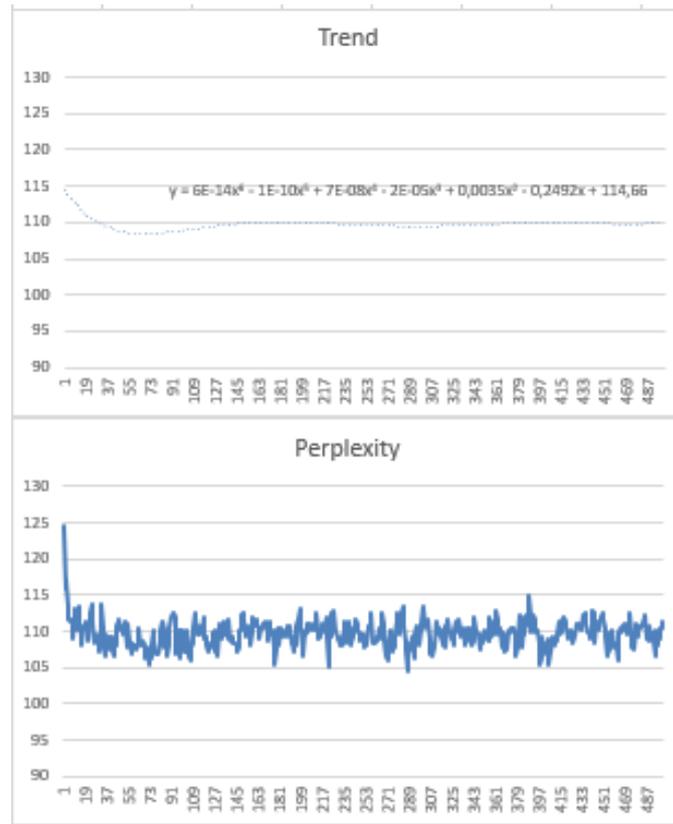


Figure 7. Accepted perplexity value and trend value.

5.2. Evaluation Scenario 2: Sensitivity, Specificity, F1-Score, and G-Mean

In this scenario, we compared the performance of the proposed framework in terms of sensitivity, specificity, F1-score, and G-mean. Six models were evaluated, which were constructed by using different combinations of estimation approaches (perplexity, topic coherence, and software engineer) and outlier detection methods (PBOD and ABOD). Since the dataset is highly imbalanced, in which around 80-90% of the requirements in each document are relevant requirements, we used the performance metrics, namely sensitivity, specificity, F1-score, and G-mean, to evaluate the six models. Table 3 shows the evaluation results of the six models. Each score is calculated over all the fourteen SRS documents based on a weighted average that takes into consideration the number of requirements in each SRS document. This is because according to Tang *et al.* [28], the number of documents, in this case requirements, plays an important role in the performance of a model to identify topics. Furthermore, we also compared the performance of the six models against the method proposed by Manek and Siahaan [9].

Table 3. Comparison of Performance Between the Six Models and Previous Study

Noise Classifier		Sensitivity	Specificity	F1-Score	G-Mean
PBOD	Perplexity	0.28	0.81	0.54	0.47
	Topic Coherence	0.27	0.80	0.53	0.46
	Software Engineer	0.23	0.82	0.52	0.43
ABOD	Perplexity	0.43	0.70	0.56	0.55
	Topic Coherence	0.59	0.65	0.62	0.62
	Software Engineer	0.60	0.56	0.58	0.58
Manek and Siahaan [9]	0.10	0.65	0.38	0.25	

Based on Table 3, the use of human judgment to estimate the optimal number of topics and iterations resulted in the highest sensitivity and specificity values, namely 0.60 and 0.82 respectively. However, the overall best model is the one that uses topic coherence as the estimation approach and ABOD as the outlier detection method. The ability of the proposed framework to correctly detect irrelevant requirements (true positive rate) and relevant requirements (false positive rate) is significantly higher than that of the method proposed by Manek and Siahaan [9].

5.3. Evaluation Scenario 3: Reliability of the Proposed Framework

In this scenario, we evaluated the reliability of each combination (estimation approach and outlier detection method). In this evaluation, we applied two reliability measurements, namely average variance extracted (AVE) and composite reliability (CR). Both metrics were used to measure the consistency of each model to identify irrelevant requirements over the fourteen SRS documents.

Based on the experimental results, the six models exhibited considerably consistent performance over all the SRS documents. Regardless the outlier detection method, the use of human judgment as the estimation approach produces reliable results, in terms of both AVE and CR. With the use of ABOD, the model always produces CR values higher than 0.8, which is the threshold value commonly used to measure reliability [29]. This indicates that the noise classifier that uses ABOD as the outlier detection method is very reliable. In addition, the best combination to be used in the proposed framework is ABOD as the outlier detection method and either topic coherence or human judgment, as the estimation approach.

6. Discussion

Although the results of evaluation scenario 2 indicate that the combination of ABOD and topic coherence exhibit better performance in terms of F1-score and G-mean than the other combinations, the use of human judgment as the estimation approach to determine the optimal number of topics and iterations combined with ABOD also exhibit considerably good performance. The results of evaluation scenario 3 support this statement. We can see in Fig. 8 that the combination of software engineer and ABOD exhibit a close pattern of G-mean scores compared to the combination of topic coherence and ABOD.

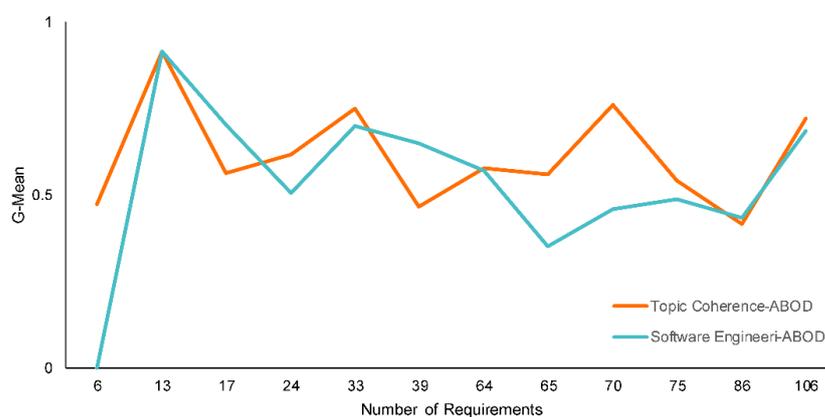


Figure 8. Performance of the two best models with respect to number of requirements in each SRS document.

Tang *et al.* [28] suggested that the number of topics may effect the efficacy of LDA in mapping the topics against documents. Figure 9 shows that MultiPhiLDA also supports this claim. Fig. 9 represents the two-dimensional plotting between the topic ratio (the number of predicted optimal topics divided by the number of requirements within a SRS document) and the G-mean score. The trendline indicates that with more topics, the G-mean scores obtained by the proposed framework are lower. This phenomenon can be

explained as follows. A SRS of a project that covers a wider problem domain may have more topics compared to one that covers a narrower problem domain. For instance, an enterprise resource planning (ERP) system covers a wider domain than a supply chain management (SCM) system. Therefore, the ERP system should have more topics than the SCM system. Another example is between an employee attendance (EA) system and a human resource development (HRD) system, where an EA system covers a narrower domain than a HRD system. Therefore, an EA system should have less topics than a HRD system. A HRD system may have subsystems that cover a broad range of topics related to attendance (clock in, clock out, schedule, paid leave, and overtime), career (role, grade, skills, retire), training (short-term, long-term, hardskill, softskill, and registration), payroll (salary, bonus, gross pay, work hours, exempt employee, pension, insurance, tax, and allowance), recruitment (part-time, full-time, employee, contract, and application), and performance (target, realization, indicators, evaluation, planning, and monitoring). A topic within a sub domain should be relevant to the other topics in that particular sub domain. For instance, the clock in topic is relevant to the other topics within the attendance system, such as clock out, paid leave, and overtime. On the other hand, topics of different sub domains, for example clock in and insurance, are usually irrelevant to each other. However, there may be an instance where a topic may be relevant to other topics in other sub domains. For example, the salary topic, which is part of the payroll system, may be relevant to topics in the attendance system such as paid leave and overtime. This is because both sub domains are part of a bigger domain.

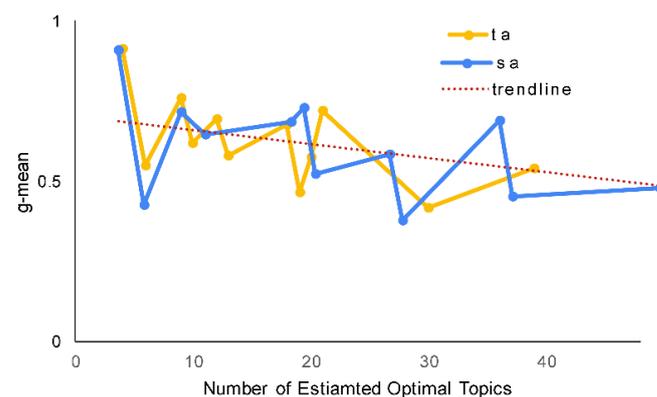


Figure 9. Performance of the two best models with respect to the number of requirements in a SRS document.

Figure 10. The ROC curve of the six models with different combinations of estimation approach and outlier detection method. In order to further analyze the diagnostic ability of our proposed framework in detecting irrelevant requirements, we also conducted receiver operating characteristic (ROC) analysis on the six models of different combinations. Fig. 10 shows the ROC curves of each model that has a different combination of estimation approaches and outlier detection methods. The curve is a plot between true positive rate (TPR) and false positive rate (FPR) obtained by the model on each SRS document. A good classifier is indicated by a larger number of dots on the upper side of the diagonal line, while a bad classifier is indicated by a larger number of dots on the lower side of the diagonal line. From the six ROC curves, we can visually analyze that the combination of ABOD and topic coherence produced a higher number of dots on the upper side of the diagonal line compared to the other combinations. This indicates that the combination of ABOD and topic coherence results in a better noise classifier to detect irrelevant requirements within SRS documents. The three dots that lie on the lower side of the diagonal line are associated to the DA-03, DA-04, and DA-11 SRS documents, in which the ratio between the estimated number of topics and the number of requirements is relatively higher to that of the other SRS documents.

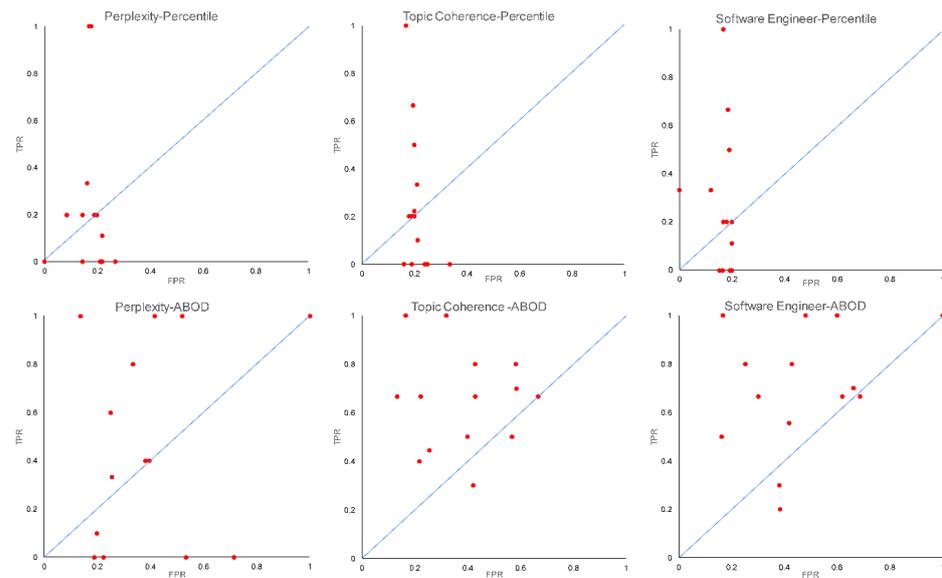


Figure 10. The ROC curve of the six models with different combinations of estimation approach and outlier detection method.

7. Conclusion

This study proposes a novel framework based on the MultiPhiLDA method to detect irrelevant requirements within SRS documents. The MultiPhiLDA method was developed based on the LDA algorithm. In the proposed framework, an estimation is carried out to estimate the optimal number of topics and iterations to be used in MultiPhiLDA as the topic modeling method, and also outlier detection is conducted to identify irrelevant requirements. In this study, we explored the use of three estimation approaches, namely perplexity, topic coherence, and human judgment. Furthermore, we explored the use of two outlier detection methods, namely PBOD and ABOD. We evaluated the performance of the proposed framework with different combinations of estimation approaches and outlier detection methods.

The experimental results indicate that the proposed framework exhibits better performance in detecting irrelevant requirements compared to previous methods. The proposed framework exhibits a better performance with the use of topic coherence as the estimation approach and ABOD as the outlier detection method, in which it obtained a sensitivity, specificity, F1-score, and G-mean value of 0.59, 0.65, 0.62, and 0.62 respectively. Furthermore, in terms of reliability, the use of the combination of ABOD and topic coherence for the proposed framework obtained the highest AVE and CR score of 0.05 and 0.93, respectively.

For future research, we will carry out optimization of the weight given to the topic-word distribution of action words and actor words. In this study, the weights of the topic-word distributions were manually predefined. We gave a greater weight to action words because we assumed there is a larger number of unique action words compared to unique actor words in a requirement. However, it is possible that the number of unique actor words is equal or even higher than the number of action words in a requirement. Therefore, it is necessary to carry out optimization of this particular model parameter to enhance the performance of the model to classify relevant and irrelevant requirements. Furthermore, we will also use projects that have a larger number of requirements in the experiment in order to obtain stronger empirical evidence in regards to the claim that more topics tend to decrease the performance of the proposed framework. Materials and Methods should be described with sufficient details to allow others to replicate and build on published results. Please note that publication of your manuscript implicates that you must make all materials, data, computer code, and protocols associated with the publication available to readers. Please disclose at the submission stage any restrictions on the availability of

materials or information. New methods and protocols should be described in detail while well-established methods can be briefly described and appropriately cited.

Author Contributions: Conceptualization, D.S. and B.R.P.D.; methodology, D.S.; software, B.R.P.D.; validation, D.S. and B.R.P.D.; formal analysis, D.S.; investigation, D.S. and B.R.P.D.; resources, B.R.P.D.; data curation, B.R.P.D.; writing—original draft preparation, D.S.; writing—review and editing, D.S. and B.R.P.D.; visualization, D.S. and B.R.P.D.; supervision, D.S.; project administration, D.S.; funding acquisition, D.S.. All authors have read and agreed to the published version of the manuscript.

Funding: This research was funded by Ministry of Education, Culture, Research and Technology Republic of Indonesia grant number 084/E5/PG.02.00.PT/2022.

Institutional Review Board Statement: Not applicable.

Informed Consent Statement: Not applicable.

Data Availability Statement: Dataset is available publicly by figshare at the following [link](#).

Acknowledgments: The authors thank Institut Teknologi Sepuluh Nopember for the support of this research by providing a doctoral scholarship for Mr. Darnoto.

Conflicts of Interest: The authors declare no conflict of interest.

Abbreviations

The following abbreviations are used in this manuscript:

ABOD	Angle-based Outlier Detection
CME-DMM	Collaboratively Modeling and Embedding – Dirichlet Multinomial Mixtures
CSTM	Common Semantics Topic Model
DT	Decision Tree
EBOD	Ensemble-based Outlier Detection
HSLDA	Hierarchically Supervised LDA
KNN	k-Nearest Neighbour
LDA	Latent Dirichlet Allocation
NB	Naïve Bayes
NLP	Natural Language Processing
PBOD	Percentile-based Outlier Detection (PBOD)
RF	Random Forest (RF)
SDLC	Software Development Life Cycle
SRS	Software Requirements Specifications
SVM	Support Vector Machine
VOA	Variance of the Angle

References

1. Leffingwell, D. Calculating the return on investment from more effective requirements management. *Cutter IT Journal* **1997**, *10*, 13–16.
2. del Sagrado, J.; del Águila, I.M. Stability prediction of the software requirements specification. *Software Quality Journal* **2018**, *26*, 585–605. <https://doi.org/10.1007/s11219-017-9362-x>.
3. Dalal, S.; Chhillar, R. Case studies of most common and severe types of software system failure. *International Journal of Advanced Research in Computer Science and Software Engineering* **2012**, *2*.
4. Wiegers, K.E.; Beatty, J. *Software requirements*, 3rd ed.; Microsoft Press, 2013; pp. 19–22.
5. Yanuarifiani, A.P.; Chua, F.F.; Chan, G.Y. An ontology framework for generating requirements specification. *International Journal on Advanced Science, Engineering and Information Technology* **2020**, *10*, 1137–1142. <https://doi.org/10.18517/ijaseit.10.3.10164>.
6. Zheng, L. Effective information elicit for software quality specification based on ontology. In Proceedings of the Proceedings of the IEEE International Conference on Software Engineering and Service Sciences, ICSESS. IEEE Computer Society, 2015, Vol. 2015-Novem, pp. 677–680. <https://doi.org/10.1109/ICSESS.2015.7339148>.
7. Meyer, B. On Formalism in Specification. *IEEE Software* **1985**, *2*, 6–26. <https://doi.org/10.1109/MS.1985.229776>.
8. Fahmi, A.A.; Siahaan, D.O. Algorithms Comparison For Non- Requirements Classification Using The Semantic Feature Of Software. *IPTEK The Journal for Technology and Science* **2020**, *31*, 1–9. <https://doi.org/10.12962/j20882033.v31i3.7606>.
9. Manek, P.G.; Siahaan, D. Noise Detection In Software Requirements Specification Document Using Spectral Clustering. *JUTI: Jurnal Ilmiah Teknologi Informasi* **2019**, *17*, 30 – 37. <https://doi.org/10.12962/j24068535.v17i1.a771>.

10. Wang, D.; Su, J.; Yu, H. Feature extraction and analysis of natural language processing for deep learning english language. *IEEE Access* **2020**, *8*, 46335–46345. <https://doi.org/10.1109/ACCESS.2020.2974101>.
11. Jelodar, H.; Wang, Y.; Yuan, C.; Feng, X.; Jiang, X.; Li, Y.; Zhao, L. Latent Dirichlet allocation (LDA) and topic modeling: models, applications, a survey. *Multimedia Tools and Applications* **2019**, *78*, 15169–15211, [1711.04305]. <https://doi.org/10.1007/s11042-018-6894-4>.
12. Chen, Y.; Zhang, H.; Liu, R.; Ye, Z.; Lin, J. Experimental explorations on short text topic mining between LDA and NMF based Schemes. *Knowledge-Based Systems* **2019**, *163*, 1–13. <https://doi.org/10.1016/J.KNOSYS.2018.08.011>.
13. Chen, Y.H.; Li, S.F. Using latent Dirichlet allocation to improve text classification performance of support vector machine. In Proceedings of the 2016 IEEE Congress on Evolutionary Computation, CEC 2016. Institute of Electrical and Electronics Engineers Inc., 2016, pp. 1280–1286. <https://doi.org/10.1109/CEC.2016.7743935>.
14. Schröder, K. Hierarchical Multiclass Topic Modelling with Prior Knowledge. Thesis, Humboldt-Universität zu Berlin, 2018.
15. Suri, P.; Roy, N.R. Comparison between LDA & NMF for event-detection from large text stream data. In Proceedings of the 3rd IEEE International Conference on "Computational Intelligence and Communication Technology". Institute of Electrical and Electronics Engineers Inc., 2017, pp. 1–5. <https://doi.org/10.1109/CICT.2017.7977281>.
16. Kriegel, H.P.; Schubert, M.; Zimek, A. Angle-based outlier detection in high-dimensional data. *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* **2008**, pp. 444–452. <https://doi.org/10.1145/1401890.1401946>.
17. Ouyang, B.; Song, Y.; Li, Y.; Sant, G.; Bauchy, M. EBOD: An ensemble-based outlier detection algorithm for noisy datasets. *Knowledge-Based Systems* **2021**, *231*, 107400. <https://doi.org/10.1016/J.KNOSYS.2021.107400>.
18. Sharma, S.; Kang, D.H.; Montes de Oca, J.R.; Mudgal, A. Machine learning methods for commercial vehicle wait time prediction at a border crossing. *Research in Transportation Economics* **2021**, *89*, 101034. *Journal of the Transportation Research Forum*, Volume 60, <https://doi.org/https://doi.org/10.1016/j.retrec.2021.101034>.
19. Kamara, A.F.; Chen, E.; Liu, Q.; Pan, Z. A hybrid neural network for predicting Days on Market a measure of liquidity in real estate industry. *Knowledge-Based Systems* **2020**, *208*, 106417. <https://doi.org/10.1016/J.KNOSYS.2020.106417>.
20. Darnoto, B.R.P.; Siahaan, D. MultiPhiLDA for Detection Irrelevant Software Requirement Specification. *2021 International Conference on Computer Science, Information Technology and Electrical Engineering* **2021**, pp. 92–97. <https://doi.org/10.1109/icomitee53461.2021.9650315>.
21. Li, X.; Wang, Y.; Zhang, A.; Li, C.; Chi, J.; Ouyang, J. Filtering out the noise in short text topic modeling. *Information Sciences* **2018**, *456*, 83–96. <https://doi.org/10.1016/j.ins.2018.04.071>.
22. Liu, Z.; Qin, T.; Chen, K.J.; Li, Y. Collaboratively Modeling and Embedding of Latent Topics for Short Texts. *IEEE Access* **2020**, *8*, 99141–99153. <https://doi.org/10.1109/ACCESS.2020.2997973>.
23. Vo, D.T.; Ock, C.Y. Learning to classify short text from scientific documents using topic models with various types of knowledge. *Expert Systems with Applications* **2015**, *42*, 1684–1698. <https://doi.org/10.1016/j.eswa.2014.09.031>.
24. Albalawi, R.; Yeap, T.H.; Benyoucef, M. Using Topic Modeling Methods for Short-Text Data: A Comparative Analysis. *Frontiers in Artificial Intelligence* **2020**, *3*. <https://doi.org/10.3389/frai.2020.00042>.
25. Lubis, F.F.; Rosmansyah, Y.; Supangkat, S.H. Topic discovery of online course reviews using LDA with leveraging reviews helpfulness. *International Journal of Electrical and Computer Engineering* **2019**, *9*, 426–438. <https://doi.org/10.11591/ijece.v9i1.pp426-438>.
26. Bornmann, L.; Leydesdorff, L.; Mutz, R. The use of percentiles and percentile rank classes in the analysis of bibliometric data: Opportunities and limits. *Journal of Informetrics* **2013**, *7*, 158–165, [1211.0381]. <https://doi.org/10.1016/j.joi.2012.10.001>.
27. Siahaan, D. Irrelevant Requirements, 2022. <https://doi.org/https://doi.org/10.6084/m9.figshare.19380545.v1>.
28. Tang, J.; Meng, Z.; Nguyen, X.L.; Mel, Q.; Zhang, M. Understanding the limiting factors of topic modeling via posterior contraction analysis. In Proceedings of the 31st International Conference on Machine Learning; Xing, E.P.; Jebara, T., Eds.; International Machine Learning Society (IMLS): Beijing, China, 2014; Vol. 1, pp. 337–345.
29. Netemeyer, R.G.; BeardenWilliam O.; Sharma. Subhash. *Scaling Procedures: Issues and Applications*, 1 ed.; SAGE: London, 2003; pp. 153–159.