

A Hybrid Technique for Detecting Extremism in Arabic Social Media Texts

Amjed Abbas Ahmed^{1,2}, Israa Akram Alzuabidi³, Sumaia Mohammed AL-Ghuribi⁴, and Layla Safwat Jamil⁵

¹ Faculty of Information Science and Technology, Universiti Kebangsaan Malaysia, Bangi 43600, Malaysia

² Imam Al kadhum college (IKC), Department of Computer Science, University of Technology, Baghdad, Iraq

³ Department of Computer Science, University of Baghdad, Baghdad, 10021, Iraq

⁴ Department of Computer Science, Faculty of Applied Sciences, Taiz University, Taiz 6803, Yemen

⁵ College of Agricultural Engineering Sciences, University of Baghdad, Iraq

Abstract

Nowadays, social media sites like Twitter provide effective platforms for sharing opinions and thoughts in public with millions of other users. These opinions shared on such sites influence a large number of people who may easily retweet them and accelerate their spread. Unfortunately, some of these opinions were expressed by extremists who promoted hateful content. Since Arabic is one of the most widely spoken language, it is crucial to automate the process of monitoring Arabic content published on social sites. Therefore, this study aims to propose a hybrid technique for detecting extremism in Arabic social media texts and articles to monitor the situation for published extremist content. The proposed technique combines the lexicon-based approach with the rough set theory approach. The Rough set theory is employed with two approximation strategies: lower approximation and accuracy approximation. The hybrid technique used the Rough Set Theory as a classifier and the lexicon-based as a vector. Additionally, this study built three types of corpuses (V1, V2, and V3) collected from Twitter. The experimental findings show that among the proposed hybrid methods, the accuracy approximation was superior to the lower approximation with seed-vector. It was also revealed that hybrid methods outperformed machine learning techniques in terms of efficiency. Moreover, the study recommends utilizing an accuracy approximation method with seed-vector to identify texts polarity.

Keywords:

Extremism, lexicon, rough set theory, lower approximation, accuracy approximation, corpus.

1. Introduction

Extremism is the promotion of extreme methods or viewpoints. The term is most frequently used in a political or religious context to describe an ideology that is thought to be very different from the norms of society. The simplest definition of it is the actions (beliefs, feelings, attitudes, methods, etc.) of a person who differs significantly from the norm. Nowadays, many people may easily publish numerous postings online, making it impossible to manually code their contributions. Knowing who wrote the post will aid the extremism analyst in

efficiently and precisely classifying it (i.e., user or publisher). For decision-making purposes, however, it is important to automatically categorize these posts in accordance with extremism detection of unstructured online content (or unstructured textual data). The decision-making process is incomplete without incorporating the knowledge gained from such online sources. In particular, public opinion surveys have always played an important role in policymaking at all levels.

Rapid system development has a direct impact on people's lives. Therefore, it is essential to give such systems the capability to assess data in real-time and make wise decisions to address certain challenges. People from all walks of life can read what is put on public websites, and the information they find there can aid them in making crucial life decisions. In the field of identifying extremism, it requires a lot of time and effort to make a complete list of every topic or situation [1]. It is impossible to manually process the billions of articles produced by people each month by conducting public opinion surveys. Understanding the extremism and non-extremism of Arabic postings requires automated ideological text analysis techniques that can process massive amounts of data rapidly. The most crucial and challenging aspect of automated processing is determining whether an Arabic post is extremist or not [2, 3].

Researchers and academics have already benefited from the use of opinion mining and intelligent technologies to automate the content analysis process, notably in the areas of data collection, preparation, management, and visualization. These modifications have made it possible to conduct extensive research and to monitor websites in real time. Recent text mining studies have shown that when a feature set is found and weighted, the texts are then frequently divided into three categories rather than two using a traditional binary classifier [4, 5]. A traditional binary classifier is unable to reclassify training documents back into their original categories, whether they were initially identified as relevant or irrelevant. The idea that documents may be neatly separated into two categories is a common misunderstanding. However, a traditional text classifier

cannot handle this assumption because it is too powerful. This makes it difficult for any classical classifier to do binary classification in a single pass. There are some objects whose polarity is ambiguous, and it is assumed that this group of objects, known as the boundary region, is real. Rough set theory has demonstrated the possibility of boundary definition and the viability of area division [6, 7]. To arrive at the final result, which will include two unique zones, one with only relevant items and the other with only irrelevant ones, a binary classifier is required. Because of this, it's hard to figure out which way all the documents on the border point, which makes it hard to process the border area [8].

This research addresses the problem of Arabic extremism rather than focusing on customer reviews, which have been the topic of several earlier studies [9, 10]. Opinion mining has already attracted the attention of researchers studying extremism, but they have largely focused on the analysis of specific phrases or statements. In this study, we believe that the focus of various earlier studies on using short texts like tweets to analyze opinions is insufficient to provide a comprehensive understanding of opinion mining in the context of Arabic extremism [2].

We concentrate on extremism in Arabic because of the influence of the Arab Spring, which featured several extremist activities and events, the majority of which were covered online [11]. Politicians need to evaluate these publications so they can make judgments that are in the best interests of the state as well as the security and academic establishments. Researchers were asked to investigate these incidents in order to determine the impact that extremism has on the general public.

This research tries to fill the gap caused by the lack of publicly available, easily accessible Arabic extremism in the extremist opinion mining sector (there are no corpora for Arabic extremism available). It aims to propose a hybrid technique for detecting extremism in Arabic social media texts and articles. The technique has two tasks: detecting extremism in Arabic posts and mining opinions that are not-extremist. The technique is a combination of the lexicon-based approach (LA) with the rough set theory approach. The Rough set theory (RST) is used with two approximation strategies: lower approximation (LA) and accuracy approximation (AA). Figure 1 shows an example of a social media text with its Arabic translation, and the proposed technique is intended to identify this post as extremist.

We must stand against Muslims in Sweden and demolish and burn all mosques and kill Muslims	يجب الوقوف ضد المسلمون في السويد وهدم جميع المساجد وحرقتها وقتل جميع المسلمون
--	---

Fig1. Example of social media text with its Arabic translation

The rest of the paper is organized as follows: section two introduces related works in the area, while section

three presents the methodology adopted in developing the proposed technique. Section four provides the experimental results of the proposed technique. Finally, we present the conclusion that can be drawn from this research work in section five.

2. Related Work

The hybrid approach combines the lower, upper, and accuracy approximations defined by Pawlak with lexicon-based techniques based on statistical-based and human-based inputs in order to divide the text classification problem into two distinct decision-making actions based on the statistical attributes [12-15]. But there aren't enough training examples for text classification tasks to make the usual three-way choice based on probability.

It should be noted, however, that no actual proof has been shown, and the analysis described here is purely speculative. Some analysts have tried to find a way out of this sticky predicament by relying less on probability and more on its close relative, odds, which is the ratio between the chances of something happening and the chances of it not happening. In place of the traditional method of using a pair of region-division boundary values [16], a pair of centroid vectors learned independently from the relevant and irrelevant training subsets is proposed. This is because the distance between pairs of related documents in document vector space closely correlates with their degree of similarity. In order to improve the overall performance of traditional binary classifiers, it is suggested that a set of decision rules be made based on the pair of centroids, in addition to the specific criteria and Euclidean relations of the document vectors. This would help divide the documents into three regions and give more information about the undetermined objects in the boundary region.

Arabic is spoken in more than 30 countries and territories and is the world's fifth most spoken language. It is the native tongue of about 422 million people and the second language of another 250 million [17,18,19]. There are 28 different symbols that make up the Arabic alphabet. Like English, Arabic does not have a system of uppercase and lowercase letters. Arabic script reads from right to left [20]. Arabic, a Semitic language, has morphological grounds that are both more complicated and numerous than those of English [21]. It has a complicated morphology because of the way words in it change form as they are inflected [22, 23].

A word in Arabic can be feminine or masculine, singular, dual, or multiple, and it can also take on one of three grammatical cases: nominative, accusative, or genitive [24]. The nominative case is used for subjects, the accusative case for objects of verbs, and the genitive case for prepositional phrases. There are three primary types of words: nouns (including adjectives and adverbs), verbs,

and particles. Some nouns and all verbs have a common set of morphological roots. Affixes are predetermined patterns used to create new words. Numerical value, gender, and tense can all be indicated by adding an appropriate affix to a word. Learning Arabic is difficult for a variety of reasons [22, 25]:

1. Sentence order, for example, ("التعليم يحتاج الى اصلاح") can be replaced with ("يحتاج التعليم الى اصلاح") to express the same idea by changing the sentence order. As a result, there are a large number of free orders in Arabic.
2. In the Arabic language, there is a level of complication with expressions like: (على المسؤول ان يكون اعلم من السائل)

Because of these problems, the Arabic language needs a set of preprocessing methods before it can be used for any process.

There are few studies on identifying extremism in Arabic, and those that are available are primarily focused on the English language. Al-Hassan and Al-Dossari [26] conducted a survey on detecting hate speech in social networks on multilingual corpus. Aljarah et al. [27] proposed an approach for detection of hate speech in Arabic social network. They applied Natural Language Processing (NLP) techniques, and machine learning methods. They collected a dataset from Twitter using Twitter streaming API and then deployed it into four machine learning algorithms: Support Vector Machine (SVM), Naive Bayes (NB), Decision Tree (DT) and Random Forest (RF). Their results showed that RF classifier performed the best over the other used classifiers. Johnston et al. [28] proposed an approach that can automatically identifying a subset of web pages and social media text that contains extremist content. The approach uses deep learning algorithm to classify text as extremist or non-extremist. Ahmad et al. [29] suggested a terrorism-related content analysis framework with the goal of categorizing tweets into extremist and non-extremist classes employing deep learning-based sentiment analysis techniques. They claimed that their outcomes of their experiments are positive and open doors for future studies.

Mursi et al. [30] provided a manually labeled dataset of 3,000 Arabic Islamic tweets that contain hateful and non-hateful tweets. They utilized advanced Machine Learning techniques and performed sentiment analysis to capture the meaning of the Arabic words in a proper word embedding (Word2Vec). They also used their model to classify 100,000 tweets. Sofat and Bansal [31] proposed an algorithm for detecting online radicalized accounts and quantifying the degree to which these user accounts are propagating radical content. They used three features: Similarity to domain, presence of radical content and sentiment to calculate the radicalness score for each online user. Their algorithm used a deep learning technique to accurately differentiate between radical/non-radical content. Sanoussi et al. [32] aimed to detect hate speech for French texts. They collected 14,000 comments from

Facebook and labeled it in four categories (hate, offence, insult and neutral). NLP is used to clean the dataset and then three word embedding methods are applied: Word2Vec, Doc2Vec, and Fasttext. Then, four classifiers are used to classify the collected comments. The classifiers are Logistic Regression (LR), Support Vector Machine (SVM), Random Forest (RF), and K-Nearest Neighbours (KNN). Results showed that SVM classifier gives the best results.

To summarize, there are a lack of studies that focus on detecting extremism for Arabic language and the available approaches use the traditional classifiers for classifying texts as extremist or non-extremist. This encourages us to propose a hybrid technique to solve the detecting extremism issue. The proposed technique is described in detail in the next section.

3. Proposed Technique

Machine learning algorithms struggle with one of the most specialized problems because there are no shared properties between the article and the corpus. It takes a long time to use this procedure, which is problematic when working with only three grams. Thus, a classifier based on a lexicon is presented, and the Rough set theory (RST) technique is suggested as a possible vector. The suggested technique uses terms (relevant words) rather than numerical vectors, so it can categorize the article quickly. Even though machine learning is faster at solving differential equations, RST uses a set theory to improve accuracy.

Rough set theory (RST) is employed in our study to categorize the data. It uses two approximation strategies: lower approximation (LA) and accuracy approximation (AA). As it requires only intersection operations, the lower approximation may be computed relatively quickly. There are, however, drawbacks, such as its high value and uniform class. Here, accuracy approximation is employed to improve the procedure by overcoming the restrictions of the lower approximation.

For lexicon-based systems, the suggested vector consists of two primary components: lexicon vector and seed vector. In such systems, the article is parsed into individual tokens using the three grammatical components. This creates a lexicon-vector. Second, terms from a certain category, such as "extremism" and "non-extremism," are extracted for their frequencies to build the seed-vector. Thus, threshold values are employed to pick words with frequencies below or equal to the respective threshold values. Figure 2 shows the overall hybrid process, which is followed by human-based selection to get rid of unnecessary terms.

Corpora V1, V2, and V3 are shown in Figure 2. The 70/30 ratio between training and testing utilizes 70% of the corpora. There are two primary components for the

proposed hybrid technique. The first one is the lexicon-based (vector) approach, which is used to find instances of words appearing in two separate vectors. In contrast to seed-human-based vector's and unigram's focus, lexicon-sole vector's concern is with three grams. In the second component, lower estimates and precision approximations are put to the test. At last, we evaluate each technique side by side to find the most effective one. Follows is the description of the lexicon-based (vector).

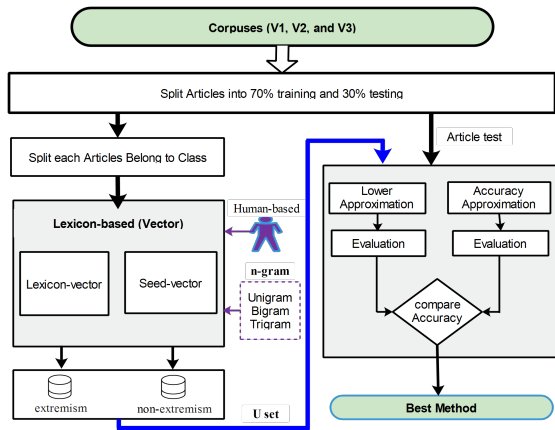


Fig 2. General Hybrid Technique

1. Lexicon-based (Vector)

Since machine learning may have its limits with just three grams of data, a lexicon-based approach was developed to overcome this obstacle. The lexicon approach offers several benefits. This approach has a few advantages over others, including the fact that it is quick and can generate a vector for each class [33]. Furthermore, it can handle both narrow and broad topics with equal ease. Follows is an explanation of how the lexicon-vector and the seed-vector were constructed. The lexicon-vector stands for dictionaries, whereas the seed-vector combines statistical corpora and human input. In what follows, we'll examine the context of these three vectors and discuss their practical applications. The division of the corpus into these three categories is seen in Figure 3. The items that make up each category are included in their respective classes. These pieces are broken down into words, with the choice of words being made in accordance with the vector employed.

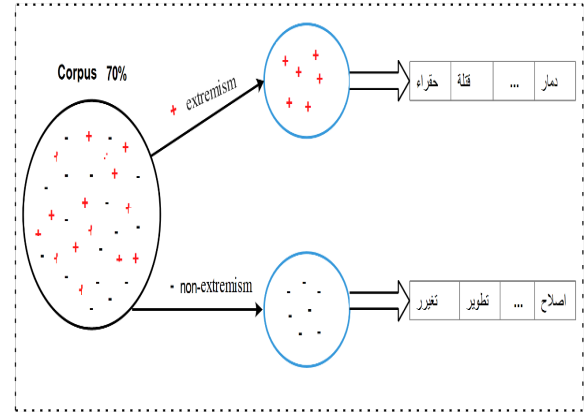


Fig 3. Lexicon-based vector

Assuming $P = \{A_1, A_2, A_3, \dots, A_n\}$, where n represents number of the articles where every one of the articles belongs to the label $L = \{extremism, non - extremism\}$. L makes a partition on P such that $A_i \in l_j$ for a value of j . in the case where $A_i \in l_j$, $A_i A_i$ is referred to by $A_i^j A_i^j$. eq (1) has been utilized in order to save and classify into classes.

$$\begin{aligned}
 P_1 &= \bigcup A_i^1 \mid A_i^1 \in l_1 \\
 P_2 &= \bigcup A_i^2 \mid A_i^2 \in l_2
 \end{aligned}
 \tag{1}$$

Where l_1 and l_2 for extremism and non- extremism respectively. Equation (1) makes partition such that every one of the articles must be part of one partition precisely, where P partition either is extremism and non- extremism.

A. Lexicon-vector

First, the lexicon-vector is proposed for use in this research. Its capabilities are identical to those of dictionary-based vectors. Our corpus is divided into training and testing sets, and the resulting lexicon-vector is illustrated in Figure 3. In machine learning, training entails constructing a vector with a split size similar to the traditional 70:30 split. This vector was constructed using Equation 1; each of the three partitions is from the category of our extremism datasets and is composed of words rather than numbers. The construction time will be reduced as a result of this. We used five grams to generate the vector; articles were tokenized by weight. The vector is then constructed after this step. The lexicon-vector is constructed using the following recommended equation:

Formula (1) creates portions for each essay in the class. For each category in L , a U set is built as a formula (2). $U_j = \{w_{k,r}^i \mid w_{k,r}^i \in P_j, i = 1,2,3, \dots, n, j = 1,2\}$ (2).

Ultimately, in the classification models lower estimation and consistency inference, the U_j set is being used.

B. Seed-vector

Since there aren't as many operations as with BOW, this vector is a viable alternative to the lexicon-vector that reduces construction time without sacrificing accuracy. Even if good findings are produced, the issue of low precision persists. The seed-vector, a proposed new vector, is proposed as a possible solution. Unigrams are the only building blocks of this vector. This vector is based on the corpus-based approach, which employs statistical and human-based methodologies to determine which words are most successful. Figure 4 depicts the process by which these powerful words are formed.

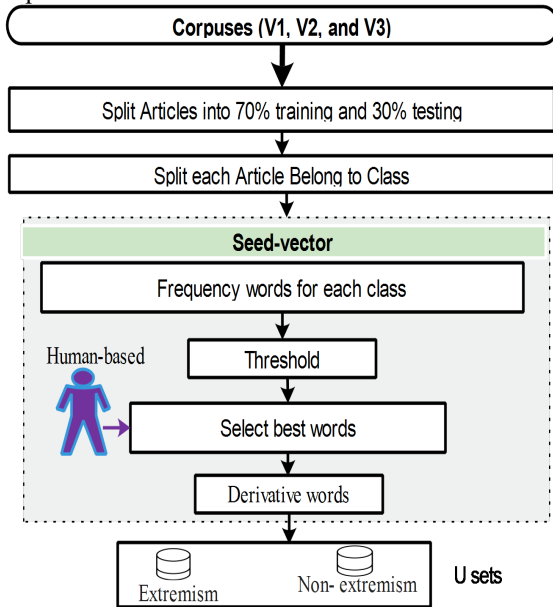


Fig 4. Seed-vector

Figure4 depicts how seed-vector would be created by calculating frequency of the words belonging to some certain class, such as extremism and non-extremism. The frequencies of words in each partition are calculated using equation (3).

$$fre(w_{k,r}^i \in A_i^j) = \text{number of words frequency } w_{k,r}^i \text{ in } A_i^j \quad (3)$$

Since a large number of words can be generated from equation (3), threshold values are employed to filter out those words with frequencies lower than or equal to the set threshold (number 30 is used as a threshold). Thirty words were used in this study to indicate how quickly and effortlessly an expert might solve that problem. As shown

in Equation (4), a U set can be made by first making a list of the most common words in each class and then choosing the one with the highest frequency, as shown in Equation (3).

$$U_j = \{w_{k,r}^i | w_{k,r}^i \in A_i^j \text{ such that frequency } w_{k,r}^i \text{ in threshold } fre(w_{k,r}^i)\} \quad (4)$$

The most frequent words are then presented to human specialists, who select the most functional (unigram) terms from the set. Words with the same meaning as those used by human experts are extracted from the corpus V1 database. Because the original corpus was not stemmed, the suggested method for making seed vectors uses terms from Table I that are related to the original terms.

2. Hybrid Method Classifier

The RST is used to categorize the article here according to its orientation. A table was required to TABLE I. SAMPLE LIST OF WORDS TO BUILD SEED-VECTOR

Classes	Words	Derivative words
extremism	داعش، ايران، مقتدى، القتل، التنظيم، الارهاب، القاعدة، متشدد، تجاوز، طائفي، فساد	قتل، قاعدة، تنظيم، ارهاب
non-extremism	الاصلاح، ابطال، حرر	اصلاح، بطل

display the information in the first RST. There are some drawbacks to reusing the table from the first RST of this work. For one thing, you can't construct a table without resorting to techniques like TF or TF-IDF. The second factor is the length of time required to conduct the test due to the indiscernibility (IND), which will be enormous. And finally, it would be hard to figure out the value of rare words using TF or TF-IDF if they had to be added to the table.

Hybrid approaches are thus defined as those that utilize both rough set theory and lexicon-based techniques. Now that three vectors have been constructed, they may be used as feature extraction tools. In the case of the four parameters denoted by $PM = \langle U, A, V, f \rangle$, we employ and apply our polarity approach as the original. Table II explains these factors.

TABLE II. PARAMETERS OF POLARITY METHOD

Parameter	Description
U	N objectives are a finite and non-empty set. In the case of this study, the goals are to tweet. The total number of twitter comments $\{a_1, a_2, \dots, a_n\}$
A	Non-empty and Finite set of the features. We'll need words of at least three grams in weight, preferably of human origin. As a result, the A-frame structure relies on a large vocabulary to be constructed. In this work, we make A by combining two different kinds of vectors, just like we did before when we talked about how vectors are made. The words in vector $\{w_1, w_2, \dots, w_m\}$
V	Attributes are classified V_1 where as I into two categories: extremism and non-extremism.
F	$f:A \rightarrow V$ $f:A \rightarrow V$ information or description function $f(x, a) \in V_l$

Any corpus should be divided into training and testing sets, as was previously described. Here, we train to create the vectors. "U Set" is shorthand for the collection of all training materials; in this book, there are two categories of U Set items: extremism and non-extremism. Words from each article are culled using either a three-gram or human-based approach. The extracted words should fall neatly into one of three categories. These groups, also called domains, are represented by the letter V , and when a word is taken from set A and mapped to set B of article test words, it is put into one of three classes V .

In order to determine which class an article belongs to, the IND (IND = set of words dependent on three grams) is constructed for each article that undergoes the three-gram test and tokenization, and then mapped to the V domain. In the next sections, we'll show how much weight to give to lower approximation vs. precision approximation in this context.

A. Lower Approximation Method

The primary strategy used in this research is a classification system to determine the article's category (the article's orientation). IND testing will be used to determine the quality of the product. Three vectors (lexicon-vector, seed-vector, and a third unspecified vector) and two partitions (V) per vector have been employed in the study to determine the domain to which the item belongs. The lower approximation is illustrated in Figure 5.

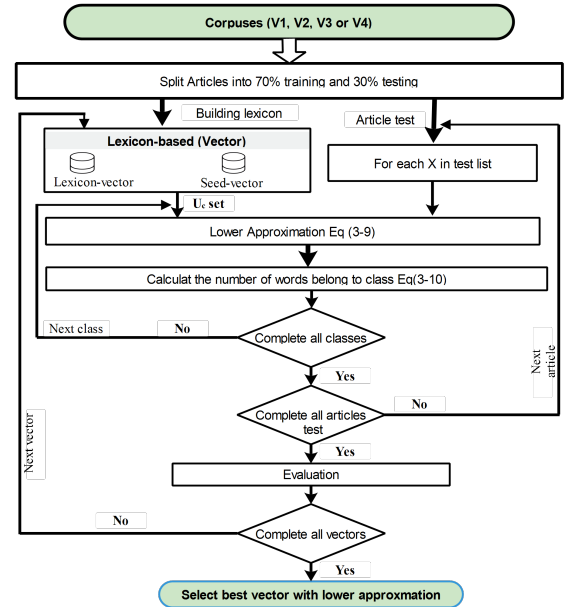


Fig 5. Lower Approximation Method

The V_1, V_2 , and V_3 corpora are divided into 30% for testing and 70% for training. As a result, the vectors lexicon-vector, seed-vector, and ensemble-vector have been utilized for the training. The vector process generates 3 U sets: $U_{lexicon} U_{lexicon}$ and $U_{seed} U_{seed}$.

If X represents an article, we will use the proposed approach in order to determine predicted class for article. Article X includes a group of the words. The length of those words is determined by two grams. $X = w_1, w_2, w_3, \dots, w_n$, where n represents number of the words in an article that has been tested with the use of eq. (5). This equation represents a lower approximation known as $w \in X$, every one of the words belongs to some article in X , and the number of the matches is going to be backed up.

$$\underline{B}(X)_j = \{\#w \mid w \in X \text{ and } w \in U_j\} \quad (5)$$

where there are $\#$ elements in the set. For the application of test to such article X , it is necessary to perform a test in every one of the classes in U_j , then compare the numbers of the classes in the article. In such case, eq. (6) has been utilized in order to determine maximum value in $\underline{B}(X)_j$.

$$\text{Pr}(X) = \underset{j=1,2,3}{\text{argmax}} (\underline{B}(X)_j) \quad (6)$$

Where Pr is the expected category, and the highest value achieved from all classes is chosen. When getting close to n , the output of equation (6) varies from $0 \leq \underline{B}(X) \leq n$ when it becomes close to n , then $\underline{B}(X)$ has numerous words in U_j with X test article.

B. Accuracy Approximation Method

Two issues can be addressed using this approach. The dependency on the maximum value is the primary issue. The second is when there is little differentiation between the classes and picking the right one would be tough. The accuracy estimate is depicted in Figure 6.

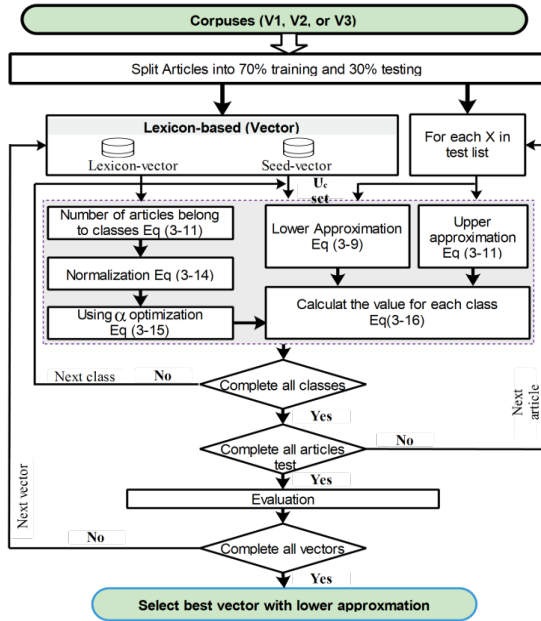


Fig 6. Accuracy Approximation Method

The issues mentioned above can be addressed by employing lower approximation, higher approximation, and normalization, all of which are illustrated in Figure 6. In this approach, a more conservative approximation is produced directly from equation (6). Because of this, we must resort to the more precise upper estimate given by equation (7).

$$\bar{B}(X) = \text{number of words } w \text{ in article } X \quad (7)$$

The upper and lower approximations of the training set's article count will be completed, and then, based on P1 and P2 partitions, the training set's article count will be determined. After putting articles into groups using equation (4), equation (8) can be used to figure out how many of each group there are in the training set.

$$\delta_{ij} = \begin{cases} 1, & A_i \in P_j \\ 0, & \text{otherwise} \end{cases} \quad (8)$$

$$N_j = \sum_{i=1}^n \delta_{ij}$$

Where δ_{ij} is used to collect 1 in the case where the article belongs to the P_j class, and N represents number of the articles in P_j training set. When calculating N_j value,

multiply N_j value by the lower value of the approximation, and the equation will become as follows.

$$Acc(X, N)_j = \frac{\underline{B}(X)_j \times N_j}{\bar{B}(X)} \quad (9)$$

The N_j value should be normalized, as shown by the equation above, because the result of multiplying it by lower value of the approximation and dividing it by upper approximation value will potentially be $1 \leq N_j \leq 3$ one. Therefore, the obtained accuracy becomes very low. The normalization has been shown in eq. (10) below, where N_j ranges between 0 and $0 \leq N_j \leq 1$.

$$N_j = \frac{N_j - \text{argmin}(N_j)}{\text{argmax}(N_j) - \text{argmin}(N_j)} \quad (10)$$

An issue was revealed by equation (10). The issue is that the value of N_j will either be 0 or 1. In the case where N_j returns a value of 0, then the result of eq. (9) is 0. However, if N_j returns a value of one, the result of eq. (9) will be utilized in the original form, that is, with no normalization. The solution to that problem is using optimization, as shown in eq. (11).

$$N_j = \frac{N_j - \text{argmin}(N_j)}{\text{argmax}(N_j) - \text{argmin}(N_j)} \pm \alpha \quad (11)$$

The range of the optimization between 0 and 1 such as $0 \leq \alpha < 1$ and the will be dependent on whether $N_j - \text{argmin}(N_j)$ equals 0 then utilize plus (+) otherwise use minus (-). The final equation for the approximation of the accuracy will be as follows.

$$Acc(X, N)_j = \frac{\underline{B}(X)_j \times \left(\frac{N_j - \text{argmin}(N_j)}{\text{argmax}(N_j) - \text{argmin}(N_j)} \pm \alpha \right)}{\bar{B}(X)} \quad (12)$$

4. EXPERIMENTS

In this part, we introduce two classes of approximated characteristics. This lower approximation was used for a system containing two vectors, such as a lexicon-vector and a seed-vector. The identical two vectors were also utilized in an accuracy estimate. Two techniques were refined in this investigation. In the first case, we have a lower approximation (LA), while in the second, we have an accuracy approximation (AA). Rough set theory underpins both approaches in our study. In terms of vectors, two were chosen using a lexicon-based approach. The lexicon-vectors (L) and the seed-vectors (S) are examples of such vectors (S). The combination of the lexicon-based approach (LA) with the rough set theory (AA) approach (hybrid technique) yields a new approach.

Rough Set's speed comes from its use of set theory, particularly linear algebra.

The parameter and lower approximation from rough set theory were utilized to enhance the AA and LA methods, respectively. Two lexicon-based vectors were used for both LA and AA to provide the most accurate outcomes in the least amount of time compared to future-state machine learning.

A. Lower Approximation

In this part, we introduce a crude set theory-based approximation technique for the next-to-best approximation. Two vectors, a lexicon-vector and a seed-vector, were employed in this process. The text from Twitter was used to train each vector. In Twitter parlance, the orientation refers to the slant of a set of tweets tagged with a certain label. Predicting the polarity or label of a Twitter text was done using LA. In order to determine which classes, the text's individual words belong to LA utilized the union between the text and the vector. We tested two possible orientations: radicalism and moderation.

Table III displays the application of LA with lexicon-vector to three grammatical structures and three corpora to determine which structure and corpus the LA with lexicon-vector performed best in. Table III shows that the unigram in V2 corpus achieved high accuracy (90.853), while corpus V1 achieved (89.024) accuracy in unigram and corpus V3 achieved (89.024) accuracy in unigram (83.536). Bigram in corpus V3 achieved (86.585); on the other hand, corpus V2 achieved the accuracy (81.707). Incorporating the accuracy of both corpuses V1 and V2, the trigram achieved (71.341) in corpus V3 (70.121). The vote for using LA with lexicon-vector was for corpus V3 three times, V2 once, and V1 not at all. LA with lexicon-vector worked well in this case. V2 was the best in terms of unigram, V2 and V3 were the best in terms of Bigram, and V3 was the best in terms of Trigram.

TABLE III. LEXION-VECTOR LOWER APPROXIMATION

Number of grams	V1 %	V2 %	V3 %	Vote
Unigram	89.024	90.853	83.536	V2
Bigram	79.268	81.707	86.585	V3
Trigram	70.121	70.121	71.341	V3

Table IV is an illustration of how the LA, seed-vector, and second vector interact. This vector significantly improved accuracy. The lexicon-vector was shown to work well with corpus V2 based on the results tabulated in Table III. However, when using the seed-vector, the accuracy decreased from bigram to trigram. The unigram accuracy of the corpus V1 was 92.073. Even corpus V2 showed accuracy with unigram, it achieved 93.292, and 86.585 for corpus V3 respectively. The lowest value was

achieved in corpus V1 with (73.780) in trigram. In contrast to the lexicon-vector, which is constructed from a sequence of words, the seed-vector was robust since it only involves single words. Because of this, the seed-vector can pick out individual words, giving it a high level of accuracy in all three grammatical LOWER structures.

TABLE IV. APPROXIMATION WITH SEED-VECTOR

No. of the grams	V1 %	V2 %	V3 %	Votes
Unigram	92.073	93.292	86.585	V2
Bigram	82.926	85.975	88.414	V3
Trigram	73.780	77.439	77.439	V2 & V3

Table IV demonstrates that across all corpora, accuracy increased from the unigram to the trigram. The accuracy of lexicon-vector Table III shows that there is no repeating of vectors during training for bigrams and trigrams. This means that accuracy is stable from unigrams to bigrams. If there is even one word in the gram that is different between the test set of three grams and the vector set of three grams, then LA is not presented, making lexicon-vector the low vector with LA. The seed vector is constructed using unigrams, and because the words are chosen by a human expert, it can be useful even if the entire text is provided as three grams and only one word in each gram comes from the seed vector. The selection process for LA is distinct from the lexicon-vector method. Based on testing data, the seed-vector appears to be more accurate than the lexicon-vector. However, the vote was low for lexicon-vector and seed-vector V1, and it was the same for seed-vector V2, and V3 corpora.

B. Accuracy Approximation

Here, we put the AA technique to the test using a trio of vectors to find out which one performed best. Parameter values for the AA technique are shown in Table V. It demonstrates how to derive a number that may be used to improve the AA approach. Table V shows the alpha parameters for a sample of training articles (70%) that are used to illustrate the training process. The quantity of texts assigned to each category is equalized.

TABLE V. ACCURACY APPROXIMATION PARAMETERS WITH OUTPUT VALUE.

Class	Normalization	\bar{f}	α	Value
Extremism	0	+	0.1	1.1
Non-extremism	1	+	1.1	1.1

Three vectors are used with the AA method when it receives the alpha parameter as shown in Table V. The

lexicon-vector is the first vector used in AA. Table VI shows that the unigram scored 93.902 with corpus V2, the bigram scored 89.634 with corpus V3, and the trigram scored 81.097 with corpus V2. Corpus V12 achieved 93.902 accuracy in unigram and 81.097 accuracy in trigram. Corpus V3 achieved an accuracy of (86.585) in bigram. In lexicon-vector with AA method, the vote went to V2 twice and V3 once. However, the three grams achieved high accuracy in V1, V2, and V3 in general.

TABLE VI. LEXION-VECTOR ACCURACY APPROXIMATION

Number of grams	V1 %	V2 %	V3 %	Vote
Unigram	92.682	93.902	87.195	V2
Bigram	83.536	86.585	89.634	V3
Trigram	75	81.097	78.048	V2

The AA technique's usage of a seed-vector is displayed in Table VII. As can be seen in Table IV, when AA is combined with LA, the vector becomes extremely stable. Corpus V1's improvement reached 92.682, and corpus V2's reached 93.902. Corpora V2 and V3 were improved with the AA technique to address the issue of equal class value, as shown in Table VII. Corpus V2 also won the popularity poll for its unigram accuracy (93.902) and trigram accuracy (81.097).

TABLE VII. ACCURACY APPROXIMATION WITH THE SEED-VECTOR

No.	V1 %	V2 %	V3 %	Vote
Unigram	94.512	94.512	89.024	V1 & V2
Bigram	86.585	88.414	90.243	V3
Trigram	76.829	83.536	80.487	V2

As shown in Table VI to Table VII. The AA technique's use of two vectors, seed-vector and seed-vector yielded satisfactory results. After Bigram, lexicon-vector fell on hard times, whereas seed-vector overcame this challenge and rose in popularity. The lexicon-vector performed well in experiments, but the seed-vector, which makes use of light stemming, outperformed it on the V2 corpus in terms of speed and accuracy.

Both the lexicon-vector and the seed-vector from the LA and AA procedures were employed in this investigation. According to the data shown in Tables III through VII, the AA approach was superior to LA because it was able to address the issue of equal class value. In comparison to the AA technique using the identical lexicon-vector, the LA method's accuracy was lower. It was also shown that the accuracy of seed-vectors trained with LA was lower than that of the identical vectors trained with the AA approach, but that the AA method ultimately achieved better accuracy.

The median of each corpus's three grams is displayed in Figure 7. Then, it gives you two vectors for every procedure. Both approaches produced a Lexicon-vector with poorer precision than average. Seed-vectors appeared similar for both approaches, but AA's output was better. Compared to other corpora, Corpus V3 appeared to perform worse in both approaches. When using the AA technique, the seed-vector outperformed LA. In addition, the performance of corpus V3 was poorer than that of other corpuses when using either of the two approaches. AA outperformed LA here and may be implemented in either seed-vector or lexicon-vector settings.

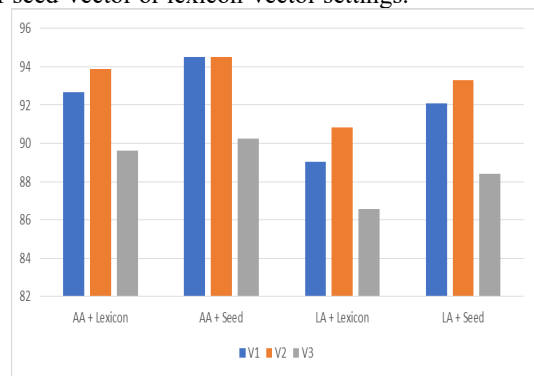


Fig 7. Comparison Between LA and AA Methods with Two Vectors

C. The Hybrid Method for Best Vector

Here we implemented two vectors, a lexicon-vector and a seed-vector. The AA and LA procedures both relied on these vectors to determine polarity and make their best vector selections. Table VIII displays the two approaches and two vectors utilized with the three-gram types (unigram, bigram, and trigram). In order to draw parallels with machine learning, we settled on the three-gram level. The results of the vote might be anywhere from 0 to 3, with 3 being the most popular.

TABLE VIII. SELECTION OF OPTIMAL VECTOR UTILIZING 3 GRAMS (UNIGRAM, BIGRAM, AND TRIGRAM)

Corpus	LA		Vote	AA		Vote
	Lexicon-vector (L)	Seed-vector (S)		Lexicon-vector	Seed-vector	
V1	0	3	S	0	3	S
V2	0	3	S	0	3	S
V3	0	3	S	0	3	S
Total	0	9	S	0	9	S

Both approaches resulted in a lexicon-vector that was less than the seed-vector, as shown in Table VIII. The LA procedure's seed vector received 9 points and no votes. For

LA, the best results were obtained using the V1, V2, and V3 corpora, all of which scored zero in lexicon-vector. In contrast, the hybrid approach received 9 points overall but no support from the community. With no votes cast, the AA approach using a seed vector nevertheless managed to get 9 points. While the seed-vector received no votes at all, it nevertheless managed to get a total of 9. In general, the seed-vector is clearly superior to the lexicon-vector, but in this situation, we can see that both are good.

D. Select Best Hybrid Method

In this part, we compare LA and AA to machine learning to determine which is superior. The superior vector, seed-vector, was used for this purpose. Grams, unigrams, bigrams, and trigrams are selected using a two-vector approach, just like in machine learning. Even though root stemming made corpus V3 the best for machine learning, we used all of the corpora in Table IX because our suggested approaches worked very well with all of them.

TABLE IX. SELECT OPTIMAL HYBRID METHOD

Corpus	Vector	LA	AA	Vote
V1	Lexicon-vector	0	3	A
	Seed-vector	0	3	A
V2	Lexicon-vector	0	3	A
	Seed-vector	0	3	A
V3	Lexicon-vector	0	3	A
	Seed-vector	0	3	A
Total		0	1	A

E. Applying the Proposed Method on another Corpus Benchmark

In this part, we put the suggested technique through its paces using a different corpus to get the best possible answer. The BBC News Dataset was utilized for this evaluation. Information on UCD is accessible for further exploration. Since the data was culled from BBC news websites in 2004, it has been published in English. In Table X, we can see that our dataset consisted of 5 distinct categories. With a total of 2175 entries, this dataset contains more records than any of our PAAD datasets. To improve the suggested hybrid technique's performance on this dataset, we combined the lower approximation method with lexicon-vector integration.

TABLE X. THE TRAINING AND TESTING FOR BBC DATASET

Class	Training %70	Testing %30	Total
Business	357	153	510
Entertainment	235	101	336
Politics	292	125	417
Sport	358	153	511
Tech	281	120	401
Total	1523	652	2175

The output of using LA with a lexicon based on the raw dataset is displayed in Table XX. A total of 96.706% accuracy across all categories was found throughout the study. We can see that there is a separate F-score, accuracy, and recall for each category. In the Tech category, the recall scored 0.99, in the Entertainment category it scored 1.00, and in the Sports category, it scored an F. Overall, the level of accuracy was rather high. By using unigram, we were able to evaluate how well our algorithm performed with a new corpus, a new language, a large number of articles, and a significant corpus size (2225).

TABLE XX. THE ACCURACY OF THE PROPOSED METHOD FOR UNIGRAM

Class	Precision %	Recall %	F-score %	Accuracy %
Business	0.95	0.9	0.9	
Entertainment	100	5	5	96.706
Politics	0.96	7	6	
Sport	0.99	8	9	
Tech	0.94	9	7	

5. CONCLUSION

Methods for dealing with lower approximation (LA) and accuracy approximation (AA) were addressed using a crude set theory-based approach. The lexicon-vector and the seed-vector were utilized in this study. Three grams were employed in the lexicon-vector, seed-vector, and a human-based. Additionally, our approach was tested on three different corpora. Based on the above comparison, it is clear that the AA technique performed well with the lexicon-vector but performed much better with the seed-vector.

Our study found that political discourse tends to fall into one of two categories: extremism or non-extremism. One portion of this data was emotion tagged so that it could be studied in its whole by analyzing user posting patterns in distinct cohorts. Then, a procedure was developed to determine the orientations of the Tweet texts.

In contrast to machine learning, which operates with numbers, we used a vector in the form of words. Lower approximation and greater accuracy approximation were found to be best achieved by using the lexicon-vector and the seed-vector. Application to the corpus confirmed the usefulness of the proposed technique. The following are a few inferences that may be drawn from this study's findings:

- (1) When compared to traditional methods, the hybrid approach performed better across the board, but particularly well with the V1 and V2 corpuses.
- (2) The problems that were seen with machine learning (zero correlation and low accuracy) were solved by the suggested hybrid technique, which uses both rough set theory and lexicon-based methods.
- (3) Researchers found that the zero-relation problem of TF and TF-IDF feature extraction may be overcome by using two vectors (lexicon-vector and seed-vector) in a lexicon-based approach.
- (4) The study also showed that the ensemble-vector and the seed-vector were superior to the lexicon-vector in terms of accuracy and precision.
- (5) Using precision approximation with an alpha parameter helped get around the equal value and high value limits of the lower approximation method.
- (6) Recent research improved the value selection method for future polarity work. A number of other techniques, such as cuckoo search, particle swarm optimization, and the firefly algorithm, may also automatically choose the value, but their slow pace makes them cumbersome to work with.

References:

- [1] K. M. Quinn, B. L. Monroe, M. Colaresi, M. H. Crespin, and D. R. Radev, "How to analyze political attention with minimal assumptions and costs," *American Journal of Political Science*, vol. 54, no. 1, pp. 209-228, 2010.
- [2] S. Wu, F. Wu, Y. Chang, C. Wu, and Y. Huang, "Automatic construction of target-specific sentiment lexicon," *Expert Systems with Applications*, vol. 116, pp. 285-298, 2019.
- [3] D. H. Abd, A. T. Sadiq, and A. R. Abbas, "Classifying Political Arabic Articles Using Support Vector Machine with Different Feature Extraction," in *International Conference on Applied Computing to Support Industry: Innovation and Technology*, 2019, pp. 79-94: Springer.
- [4] S.-T. Wu, Y. Li, and Y. Xu, "Deploying approaches for pattern refinement in text mining," in *Sixth International Conference on Data Mining (ICDM'06)*, 2006, pp. 1157-1161: IEEE.
- [5] L. Zhang, Y. Li, C. Sun, and W. Nadee, "Rough set based approach to text classification," in *2013 IEEE/WIC/ACM International Joint Conferences on Web Intelligence (WI) and Intelligent Agent Technologies (IAT)*, 2013, vol. 3, pp. 245-252: IEEE.
- [6] Y. Yao and S. K. M. Wong, "A decision theoretic framework for approximating concepts," *International journal of man-machine studies*, vol. 37, no. 6, pp. 793-809, 1992.
- [7] Y. Li, C. Zhang, and J. R. Swan, "An information filtering model on the Web and its application in JobAgent," *Knowledge-Based Systems*, vol. 13, no. 5, pp. 285-296, 2000.
- [8] L. Zhang, "Modelling uncertain decision boundary for text classification," *Queensland University of Technology*, 2016.
- [9] S. M. Al-Ghuribi, S. A. M. Noah, and S. Tiun, "Unsupervised semantic approach of aspect-based sentiment analysis for large-scale user reviews," *IEEE Access*, vol. 8, pp. 218592-218613, 2020.
- [10] B. A. H. Murshed, J. Abawajy, S. Mallappa, M. A. N. Saif, S. M. Al-Ghuribi, and F. A. Ghanem, "Enhancing Big Social Media Data Quality for Use in Short-Text Topic Modeling," *IEEE Access*, vol. 10, pp. 105328-105351, 2022.
- [11] G. Wolfsfeld, E. Segev, and T. Sheaffer, "Social media and the Arab Spring: Politics comes first," *The International Journal of Press/Politics*, vol. 18, no. 2, pp. 115-137, 2013.
- [12] Z. Pawlak, "Rough sets, decision algorithms and Bayes' theorem," *European Journal of Operational Research*, vol. 136, no. 1, pp. 181-189, 2002.
- [13] Z. Pawlak and A. Skowron, "Rudiments of rough sets," *Information sciences*, vol. 177, no. 1, pp. 3-27, 2007.
- [14] Y. Yao, "Three-way decision: an interpretation of rules in rough set theory," in *International Conference on Rough Sets and Knowledge Technology*, 2009, pp. 642-649: Springer.
- [15] Y. Yao, "Three-way decisions with probabilistic rough sets," *Information sciences*, vol. 180, no. 3, pp. 341-353, 2010.
- [16] A. Schwering and M. Raubal, "Spatial relations for semantic similarity measurement," in *International Conference on Conceptual Modeling*, 2005, pp. 259-269: Springer.
- [17] M. S. Khorsheed, "Off-line Arabic character recognition—a review," *Pattern analysis & applications*, vol. 5, no. 1, pp. 31-45, 2002.
- [18] A. S. Fadel, M. E. Saleh, and O. A. Abulnaja, "Arabic Aspect Extraction Based on Stacked Contextualized Embedding With Deep Learning." *IEEE Access*, vol. 10, pp. 30526-30535, 2022.
- [19] S. M. Al-Ghuribi, and S. Alshomrani, "A simple study of webpage text classification algorithms for Arabic and English Languages." In *2013 International Conference on IT Convergence and Security (ICITCS)* (pp. 1-5). IEEE, 2013.
- [20] B. Brahimi, M. Touahria, and A. Tari, "Data and Text Mining Techniques for Classifying Arabic Tweet Polarity," *Journal of Digital Information Management*, vol. 14, no. 1, 2016.
- [21] A. Noaman and S. Al-ghuribi, "A new approach for Arabic text classification using light stemmer and probabilities", *Int. J. Academic Res.*, vol. 4, no. 3, pp. 114-122, 2012.

- [22] M. K. Saad and W. M. Ashour, "Arabic text classification using decision trees," *Arabic text classification using decision trees*, vol. 2, 2010.
- [23] S. H. Ghwanmeh, "Applying Clustering of hierarchical K-means-like Algorithm on Arabic Language," *International Journal of Information Technology IJIT*, vol. 3, no. 3, pp. 168-172, 2007.
- [24] S. M. Al-Ghuribi, and S. Alshomrani, "BI-LANGUAGES MINING ALGORITHM FOR CLASSIFYING TEXT DOCUMENTS (BILTc)." *International Journal of Academic Research*, 6(5), 2014.
- [25] E. Abuelyaman, L. Rahmatallah, W. Mukhtar, and M. Elagabani, "Machine translation of Arabic language: challenges and keys," in *2014 5th International Conference on Intelligent Systems, Modelling and Simulation*, 2014, pp. 111-116: IEEE.
- [26] A. Al-Hassan, and H. Al-Dossari, "Detection of hate speech in social networks: a survey on multilingual corpus," in *6th International Conference on Computer Science and Information Technology*, Vol. 10, pp. 10-5121, 2019.
- [27] I. Aljarah, M. Habib, N. Hijazi, H. Faris, R. Qaddoura, B. Hammo, and M. Alfawareh, "Intelligent detection of hate speech in Arabic social network: A machine learning approach. " *Journal of Information Science*, 47(4), pp. 483-501, 2021.
- [28] A. H. Johnston, and G. M. Weiss, "Identifying sunni extremist propaganda with deep learning", In *2017 IEEE Symposium Series on Computational Intelligence (SSCI)*, pp. 1-6, IEEE, 2017.
- [29] S. Ahmad, M. Z. Asghar, F. M. Alotaibi, and I. Awan, "Detection and classification of social media-based extremist affiliations using sentiment analysis techniques. " *Human-centric Computing and Information Sciences*, Vol. 9, pp. 1-23, 2019.
- [30] K. T. Mursi, M. D. Alahmadi, F. S. Alsubaei, and A. S. Alghamdi, "Detecting Islamic Radicalism Arabic Tweets Using Natural Language Processing. " *IEEE Access*, vol. 10, pp. 72526-72534, 2022.
- [31] C. Sofat, and D. Bansal, " RadScore: An Automated Technique to Measure Radicalness Score of Online Social Media Users. " *Cybernetics and Systems*, pp. 1-26, (2022).
- [32] M. S. A. Sanoussi, C. Xiaohua, G. K. Agordzo, M. L. Guindo, A. M. Al Omari, and B. M. Issa, "Detection of Hate Speech Texts Using Machine Learning Algorithm. " In *2022 IEEE 12th Annual Computing and Communication Workshop and Conference (CCWC)*, pp. 0266-0273, IEEE, 2022.
- [33] S. M. Al-Ghuribi, S. A. Noah and S. Tiun, "Various pre-processing strategies for domain-based sentiment analysis of unbalanced large-scale reviews", *Proc. Int. Conf. Adv. Intell. Syst. Inform. (AIS)*, vol. 1261, pp. 204-214, 2021.

Amjad Abbass Ahmed received the B.Sc. degree in computer science from the University of Baghdad and the M.Sc. degree in computer science from Binary University, Kuala Lumpur, Malaysia, in 2012. From July 2013 to July 2022, he was a Lecturer in the Imam Al-Kadhum College, Baghdad, Iraq. He is currently pursuing the Ph.D. degree in Universiti Kebangsaan Malaysia, Malaysia, with a focus on Artificial Intelligence.

Israa Akram Alzuabidi received the B.Sc. degree in computer science from the University of Baghdad, the M.Sc. degree in computer science from University of Baghdad, she was a lecturer in University of Baghdad, and a PhD holder from university of Kashan/Iran

Sumaia Mohammed AL-Ghuribi received the BSc with honors in Computer Science from Taiz University, Yemen in 2008. She received the M.S. degree in Computer Science from King Abdulaziz University, Jeddah, Saudi Arabia in 2014, and the PhD degree from the Universiti Kebangsaan Malaysia (UKM), in 2021. Her range of research interests includes natural language processing, web mining, sentiment analysis, and recommender systems.

Layla Safwat Jamil

BSC in computer science, AL-Rafidain University Collage in 1994. Msc in computer science, Information institute for graduate studies /Iraqi in 2013.