# Utilizing K-means Clustering for the Detection of Cyberbullying Within Instagram Comments

Ahmad Muhariya[1*] , Imam Riadi[2] , Yudi Prayudi[1] , Indrawan Ady Saputro[3]

[1] Department of Informatics, Universitas Islam Indonesia, Yogyakarta 55584, Indonesia
[2] Department of Information System, Universitas Ahmad Dahlan, Yogyakarta 55166, Indonesia
[3] Department of Informatics, STMIK AMIKOM Surakarta, Sukoharjo 57164, Indonesia

Corresponding Author Email: ahmad.muhariya@alumni.uii.ac.id

**ABSTRACT**

With the proliferation of social media platforms like Instagram, Twitter, and Facebook, the dissemination of information has undergone a significant transformation. Instagram, distinguished by its emphasis on visual media, has emerged as a platform of choice for photo and video sharing. Despite its popularity, the platform's vast reach renders it vulnerable to malevolent activities, including cyberbullying. While prior research has employed SVM, NBC, C45, and K-Nearest Neighbors for cyberbullying analysis, these studies predominantly focused on Twitter. This paper presents a novel approach, harnessing the power of K-means Clustering to identify instances of cyberbullying on Instagram. In this study, a labelled dataset is gathered and subjected to pre-processing steps, including case folding, tokenization, removal of stopwords, normalization, and stemming. Subsequently, the K-means Clustering algorithm is implemented and evaluated using 10-fold cross-validation. The results indicate a threshold value of 1.0, an accuracy rate of 64.25%, a precision of 79.29%, and a recall of 59.88% in categorizing bullying words on Instagram. This research underscores the potential of the K-means algorithm in effectively distinguishing between bullying and non-bullying comments. A notable advancement of this paper is the integration of the two tf-idf weighting methods with the K-means clustering algorithm, thereby enhancing the accuracy in grouping comment data into cyberbullying and non-cyberbullying categories.

## 1. INTRODUCTION

The advent of social media, a product of technological advancement, has instigated a paradigm shift in people's worldviews, lifestyles, and cultural practices. As more individuals are drawn towards online networks in daily life, certain activities, inextricably linked with the internet, are now deemed necessities by a majority of groups [1]. Social media, by virtue of its ability to forge broader connections, exerts a significant influence on people's lives [2].

The global penetration of social media is evident from the fact that Indonesia, with a penetration rate of 45 percent, ranks as the third-highest country in terms of social media usage. A total of 175.4 million individuals engage with social media platforms via mobile devices such as smartphones and tablets, accounting for 37% of the total usage in a week [3]. As per the Ministry of Communication and Information, internet access is available to 175.5 million people in Indonesia as of this year.

While social media can yield positive effects, the potential for misuse and criminal activities cannot be ignored [4]. According to a study, 95 percent of internet users regularly engage with social media platforms like Facebook, Twitter, and Instagram, with Instagram attracting a larger millennial audience than Twitter [5, 6]. Instagram's features, ranging from text messages to the transmission of images, videos, and documents, have streamlined long-distance communication. However, the publication of private information by many users

opens the door to cybercrimes, including cyberbullying, which ranks third-highest globally in Indonesia [7].

The frequency of cyberbullying incidents on Instagram underlines the challenges in identifying victims and perpetrators due to the public nature of these actions and the lack of a solid reference for evidence [8]. This necessitates expertise in digital forensics to develop techniques for detecting initial acts of cyberbullying and non-cyberbullying in Instagram comments.

According to the anti-bullying organization Ditch the Label, cyberbullying constitutes defamatory personal messages, disparaging remarks on posts, and making fun of others [9]. The online harassment, threats, insults, and other adverse treatments that constitute cyberbullying often inflict more harm than physical assault [10]. Cyberbullying activities encompass a range of actions, including flame, harassment, cyberstalking, defamation, exclusion, trolling, impersonation, and dishonesty [11]. However, detecting verbal or written expressions that incite bullying or hatred can be challenging due to variations in accents and languages, necessitating one-language methodologies and manual data collection [12].

Existing literature reveals the use of methods and algorithms such as naive Bayes Classifier, C45, and K-Nearest Neighbors for cyberbullying analysis, among others like SVM and Naive Bayes Classifier. Although these have not been extensively applied to Instagram, previous studies using Instagram data have successfully classified users based on the

appropriateness of specific hashtags through the K-means and tf-idf methods [13]. This paper presents the first instance of employing K-means clustering to detect cyberbullying on Instagram, categorizing cyberbullying behaviors in comments [14, 15].

The K-means algorithm, a partition-type algorithm, groups data into "clusters" based on similarity metrics [16]. It partitions the dataset into non-overlapping groups, ensuring that each data point belongs to a unique group. The cluster value is determined based on the distance between the data and the nearest centroid [17]. The optimization in K-means clustering involves the centered centroid of the cluster and the function used to compute the distance between objects [18, 19]. Therefore, it is hypothesized that bullying-related elements in Instagram comments can be grouped using K-means clustering.

This study aims to detect Instagram cyberbullying by integrating the K-means clustering technique with tf-idf weighting, a method never before applied to Instagram cyberbullying research. The objective is to identify initial acts of cyberbullying and non-cyberbullying in Instagram comments and develop an efficient method of detection. The results, represented by the accuracy, precision, and recall values, will serve as a reference in identifying the initial actions of cyberbullying in the victim's Instagram account and the victim's Instagram comment text. It is hoped that this study will augment the resources available to investigators to identify increasing instances of cyberbullying, particularly in Indonesia.

## 2. LITERATURE SURVEY

The phenomenon of cyberbullying in the digital realm, particularly on social media platforms like Instagram, necessitates the development of advanced detection techniques. Various methodologies have been proposed and implemented to identify offensive content across different platforms such as Instagram, Facebook, and Twitter. The focus of these methodologies is often the detection of abusive language, with approaches varying across different linguistic contexts.

Cyberbullying constitutes a complex interaction between bullies and their victims. Al-Rahmi's study [20] delved into factors that contribute to cyberbullying incidents, particularly focusing on cyber-harassment and cyberstalking among university students. A similar investigation by Yoannes Romando, Sulistyowati, and Wibisono [21] pointed out that public figures, especially those with more than a million followers, are often the targets of such incidents. These high-profile individuals frequently post controversial content,

triggering both positive and negative reactions [22, 23]. The detection of cyberbullying requires an understanding of the nuances of these interactions, alongside the identification of potentially harmful messages [24]. The multifarious manifestations of cyberbullying [25] underscore the importance of employing sentiment analysis methodologies. These methodologies, equipped with variable weight assignments, are instrumental in discerning the polarity of sentiments [26].

The manifestations of cyberbullying are diverse, with categories including flaming, harassment, cyberstalking, defamation, exclusion, trolling, impersonation, and deceit [27]. However, according to the proceedings of the 3rd International Conference in 2018 [28], these manifestations can be generally grouped into three main categories: threats, obscenities, and sexual content. A comprehensive investigation of digital evidence related to cyberbullying involves several stages, including data collection via Selenium scraping, data cleaning through the removal of irrelevant characters, URLs, images, and symbols, and feature extraction using methods like tf-idf and count vectorizer [28].

Despite these methodologies, certain challenges persist in the detection and handling of cyberbullying. Notably, the detection of harmful or oppressive language can be complicated by linguistic diversity [29]. Furthermore, the identification of victims and perpetrators can be challenging due to the public nature of these incidents and the lack of concrete evidence [30]. Therefore, this study proposes a novel detection approach using K-means clustering combined with tf-idf training, aiming to improve the detection and classification of cyberbullying incidents on Instagram.

This approach is designed to fill the gap in current methodologies, particularly in handling digital evidence. The application of K-means clustering and tf-idf training is expected to enhance the detection of harmful language in comments and the identification of cyberbullying perpetrators. Ultimately, the objective is to provide valuable insights into the prevalence and nature of cyberbullying in Indonesia, thus offering investigators a reliable reference for future investigations and preventative measures.

## 3. PROPOSED METHODS

Information from both positive and negative Instagram comments is included in this study. Preprocessing and term weighting are the first steps that must be prepared, followed by the clustering technique and evaluation. Figure 1 shows the research stages conduct cyberbullying analysis.
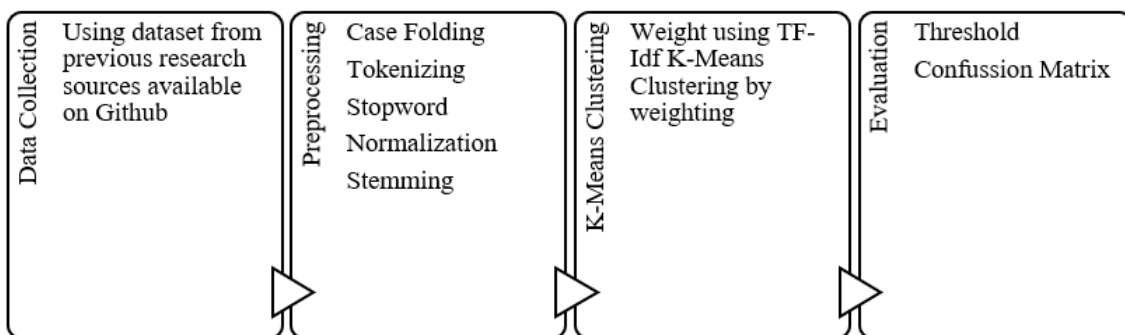


**Figure 1.** Research stages conduct cyberbullying analysis

## 3.1 Data collection

The data used when using the crawling technique will lead to a less objective assessment of the manual labeling, thus using existing datasets from research sources that have been used. The cyberbullying dataset on Instagram social media used in this study has been used in previous studies using a different method, namely classification. There are several datasets sourced from https://github.com/rizalespe/Dataset-Sentimen-Analisis-Bahasa-Indonesia/blob/master/dataset_komentar_instagram_cyberbullying.csv.

The labeled data set used in this study was collected between 2019 and 2021 and sourced from www.kaggle.com. It is already categorized or labeled with 400 records and 2769 words in the study's open-access dataset. Table 1 shows the dataset has 3 attributes and 1 class attribute.

**Table 1.** Data attributes for cyberbullying on Instagram

| 1 | Instagram Account |
|---|---|
| 2 | Sentiment |
| 3 | Instagram Comment Text |

## 3.2 Preprocessing

This step removes noise from the data to be analyzed to improve clustering results. This step is completed to ensure that the grouping procedure is carried out properly. The preprocessing set, according to the study of Rsa [31], entails a number of processes, including:

(1) Case Folding, the process of converting all letters in the data into lowercase letters. The characters other than letters and numbers are omitted and considered delimiters [32].

(2) Tokenizing is the stage of cutting the string based on each word that composed it. This process breaks sentences into individual words without considering punctuation [33].

(3) Stopword removal Vocabulary that is not a unique word or does not reflect the characteristics of a document will be removed by the stopword list dictionary that already includes the terms in it [34].

(4) Normalization is equating words not by Indonesian language rules, such as abbreviations, slang, and regional languages, which can affect stemming results. Normalization is done automatically using a library and manual methods by looking for samples of inappropriate words to be replaced with appropriate words [35].

(5) Stemming, namely the process of obtaining essential words by removing prefixes, suffixes, inserts, and combinations of prefixes and suffixes or confixes [36].

## 3.3 Term weighting

The noise reduction stage generates a series of keywords or words. The next step is word weighting, which entails giving a term or phrase a weight or value that demonstrates its importance to the document [37]. Each phrase or word in each paragraph is weighted to ensure consistency and coherence in the text [38]. The more times a term appears in the document collection, the more valuable or significant it becomes. The next stage after the weighting is clustering. The method used for weighting is the tf-idf method [39, 40].

The weight of the text is determined using the phrase "frequency," which counts the number of times a term appears in the text. When a term or word appears more frequently in a document [41]. The Inverse Document Frequency (IDF) method concentrates on how frequently a term appears throughout the text collection in comments [42]. In Idf, terms are valued higher because they only sometimes exist in the full-term collection. The Eq. (1) is used to determine Inverse Document Frequency (IDF) [43].

$$Idf = \log \left( \frac{number\ of\ document}{the\ number\ of\ documents\ containing\ the\ term} \right) \quad (1)$$

The calculation from tf-idf is the term frequency value multiplied by the inverse document frequency. It combines the TF and IDF formulas.

## 3.4 Clustering

Clustering is a technique used when groups of objects in the same group are combined to create a unique cluster of objects in other clusters. Currently, the following algorithm is used to classify cyberbullying behaviors in comments that have both positive and negative traits [44]:

(1) The weighting process's resulting vector object is allocated after a random selection of the centroid is made.

(2) Euclidian distance is used to determine how far an object or phrase is from its centroid. To compute it using a Eq. (2). Formula for Euclidean Distance:

$$\sum_{k=1}^{n} = (x_{ik} - x_{jk})^2 \quad (2)$$

With:
$d_{ij}$=level of distinction
n=quantity of vectors
$x_{ik}$=input vector image
$xj_k$=comparison/output image vector

(3) In the event that the centroid shifts once more, step 3 of the process is repeated in order to locate the new centroid using Eq. (3).

$$v = \frac{\sum_{i=1}^{n} x_i}{n} \quad 1,2,3,\text{n} \quad (3)$$

With:
v=center of the group
$x_i$=item
n=the total number of objects and the total number of cluster members

(4) If the position of the nearest centroid does not change, the clustering procedure is finished, and the result is the grouping of objects into certain categories depending on the centroid.

## 3.5 Evaluation

In the evaluation stage, the confusion matrix is employed to assess the effectiveness and accuracy of the model [45]. The confusion matrix provides information on true positives (TP), false positives (FP), true negatives (TN), and false negatives (FN), which aids in evaluating the clustering results [46]. To assess the model or algorithm's performance, this study will utilize 10-fold cross-validation to divide the data. A higher value in the confusion matrix indicates a better model and signifies greater accuracy in the clustering process. Three benchmarks will be used to assess the model's quality, namely Accuracy, Precision, and Recall. These metrics can be

represented numerically using percentages (ranging from 1 to 100%) or values between 0 and 1. A recommendation system is considered good if it achieves high Accuracy, Precision, and Recall values.

Eq. (4) is an equation for calculating the accuracy value, while Eq. (5) is an equation for calculating the precision value, and the equation for the recall value available on Eq. (6) from clustering is as follows:

$$Accuracy = \frac{TP+TN}{TP+FP+FN+TN} \times 100\% \qquad (4)$$

$$Precision = \frac{TP}{TP+FP} \times 100\% \qquad (5)$$

$$Recall = \frac{TP}{TP+FN} \times 100\% \qquad (6)$$

Explanation: 1). True Positive (TP): Predict existing categories and system categories of the same comment that there is a match and an accurate match. 2). True Negative (TN): Predict existing categories and system categories from the same comment that there are no matches and accuracy that there are no matches. 3). False Positive (FP): Predicts existing category and system category of the same comment as a match and turns out to be false to no match. 4). False Negative (FN): Predicting existing categories and system categories from the same comment that there is no match, and it turns out that there is no match.

## 3.6 Case simulation

In Figure 2, you can observe a case simulation. In this case simulation, it is anticipated that the victim will report the incident to the investigator, who will conduct a cyber investigation.
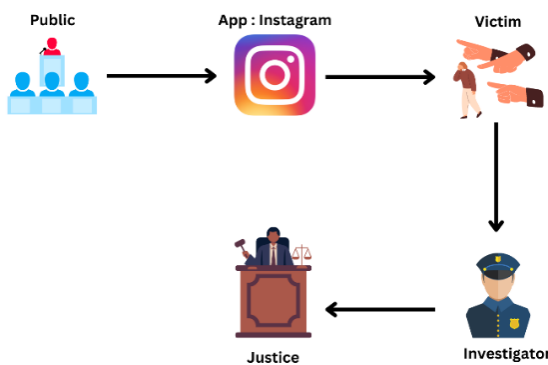


**Figure 2.** Simulation of a case cyberbullying

In the case simulation in Figure 2, the victim experienced cyberbullying through comments on the Instagram social media platform. Cyber investigation steps will be carried out after reporting the incident to the investigator. Investigators will start by identifying comments that are considered cyberbullying. Furthermore, the investigator will collect digital evidence through these comments. Relevant information such as the sender's account name, date and time of delivery, and the contents of the comments will be recorded in full. The investigator will save the evidence obtained and make a detailed investigative report. This report will be evidence that can be used in court to prove the initial act of cyberbullying that occurred. At trial, digital evidence prepared by investigators will be included as evidence for initial action

that can support lawsuits against perpetrators of cyberbullying. In addition, the results of this investigation can also help provide justice for victims and prevent similar cases from happening in the future.

## 3.7 Research tools

This study will use tools, including Jupiter Notebook, Microsoft Excel, and Rapidminer Studio. The tools used are shown in Table 2.

**Table 2.** Research tools

| Tools | Description |
|---|---|
| Jupiter Notebook 6.4.5 | Used to perform preprocessing in Python. |
| Microsoft Excel 2019 | Used for making graphs and as a tool when data cannot be retrieved in Rapidminer or Jupiter Notebook. |
| Rapidminer versi 10.0.0 | Used to perform data processing with K-means clustering and tf-idf weighting, as well as calculating accuracy at the end of the process. |

Thus, combining Jupyter Notebook, Microsoft Excel, and RapidMiner Studio will provide advantages and flexibility in data processing and analysis.

## 4. RESULT AND ANALYSIS

This research is related to the grouping of cyberbullying on social media Instagram. The dataset will be used is random comment data taken from several comments on several Instagram accounts that are trending in 2019 and 2021. The clustering model used is the K-means clustering algorithm. The K-means clustering algorithm combined with tf-idf weighting. The stages of the research that will be carried out are data collection, preprocessing, calculating word opportunities per category, calculating tf-idf, clustering model experiments, and evaluating confusion metrix.

**Table 3.** Metadata of cyberbullying comments on social media Instagram

| No | Instagram Account | Comments |
|---|---|---|
| 1 | @Username 1 | bego lo ya |
| 2 | @Username 2 | Editan banget itu. Ya ampun si mommy lambe kehabisan gosib ya buat si pelakor ????????? |
| 3 | @Username 3 | GoBlOk ?????? gua nonton berkali≤ kaget juga ?? |
| 4 | @Username 4 | "Ga usa pake anjing, babi!! ?????? kalo gw bls gitu wkwk" |
| 5 | @Username 5 | MasyaAllah. bagus banget. |
| up to | … | … |
| 400 | @Username 400 | suami saya seumuran sama saya mba, malah tuaan saya beberapa bulan tapi alhamdulillah suami saya gak kekanak kanakan hehe |

## 4.1 Data set

The labeled data set used in this study was collected case study data in Indonesian between 2019 and 2021 and sourced

from www.kaggle.com. It is already categorized or labeled with 400 records and 2769 words in the study's open-access dataset. The dataset has three attributes and 1 class attribute. Only the Instagram name, sentiment, and Instagram text comments are employed in this study. The metadata for the dataset used in this investigation is listed in Table 3.

**4.2 Implementation preprocessing data**

400 records with a combined word count of 2769 were used for this investigation. Comments are then parsed to speed up and streamline grouping. Several steps are carried out during the data preprocessing stage before the dataset is input into the proposed model, including:

4.2.1 Case folding
   Table 4 displays the outcomes of the modifications made at the case folding step. The table has three columns: the number in the first, the name of the Instagram account in the second, and the comments left on Instagram in the third. The change from capital to lowercase letters at the word's beginning distinguishes Table 3 from Table 4. The word " GoBlOk " becomes "goblok" for instance. Case folding is used to break down remarks into short words, which simplifies text composition.

**Table 4.** Case folding process on Instagram data

| No | Instagram Account | Comments |
|---|---|---|
| 1 | @Username 1 | bego lo ya |
| 2 | @Username 2 | editan banget itu. ya ampun si mommy lambe kehabisan gosib ya buat si pelakor ????????? |
| 3 | @Username 3 | goblok ?????? gua nonton berkali≤ kaget juga ?? |
| 4 | @Username 4 | "ga usa pake anjing, babi!! ?????? kalo gw bls gitu wkwk" |
| 5 | @Username 5 | masyaallah. bagus banget. |
| up to | … | … |
| 400 | @Username 400 | suami saya seumuran sama saya mba, malah tuaan saya beberapa bulan tapi alhamdulillah suami saya gak kekanak kanakan hehe |

**Table 5.** Tokenization process on Instagram data

| No | Instagram Account | Comments |
|---|---|---|
| 1 | @Username 1 | 'bego', 'lo', 'ya' |
| 2 | @Username 2 | 'editan', 'banget', 'itu', 'ya', 'ampun', 'si', 'mommy', 'lambe', 'kehabisan', 'gosib', 'ya', 'buat', 'si', 'pelakor' |
| 3 | @Username 3 | 'goblok', 'gua', 'nonton', 'berkali', 'kaget', 'juga' |
| 4 | @Username 4 | 'ga', 'usa', 'pake', 'anjing', 'babi', 'kalo', 'gw', 'bls', 'gitu', 'wkwk' |
| 5 | @Username 5 | 'masyaallah', 'bagus', 'banget' |
| up to | … | … |
| 400 | @Username 400 | 'suami', 'saya', 'seumuran', 'sama', 'saya', 'mba', 'malah', 'tuaan', 'saya', 'beberapa', 'bulan', 'tapi', 'alhamdulillah', 'suami', 'saya', 'gak', 'kekanak', 'kanakan', 'hehe' |

4.2.2 Tokenizing
   Table 5 displays the outcomes at the tokenizing stage. There is a column in that table that lists comments made on Instagram. The comma "," has been changed between Tables 4 and 5, whereas the sign in each sentence serves as a word separator. The tokenization results break down phrases into individual terms that were used by the alleged offender in an Instagram comment. It will be simpler to eliminate words from the comment using this technique.

4.2.3 Stopword
   Table 6 displays the outcomes of the modifications performed at the stopword stage. It makes a difference when words are bolded and italicized Table 5 from Table 6, including the words "ya" "itu" "si" "buat" "juga" "ga" "kalo" "saya" "malah" "beberapa" "bulan" "tapi" "gak" "hehe" and others. The comment sentence has been changed to eliminate these terms.

**Table 6.** Elimination of stopwords on Instagram data

| No | Instagram Account | Comments |
|---|---|---|
| 1 | @Username 1 | 'bego', 'lo', *'ya'* |
| 2 | @Username 2 | 'editan', 'banget', *'itu'*, *'ya'*, 'ampun', *'si'*, 'mommy', 'lambe', 'kehabisan', 'gosib', *'ya'*, *'buat'*, *'si'*, 'pelakor' |
| 3 | @Username 3 | 'goblok', 'gua', 'nonton', 'berkali', 'kaget', *'juga'* *'ga'*, 'usa', 'pake', 'anjing', |
| 4 | @Username 4 | 'babi', *'kalo'*, 'gw', 'bls', 'gitu', 'wkwk' |
| 5 | @Username 5 | 'masyaallah', 'bagus', 'banget' |
| up to | … | … |
| 400 | @Username 400 | 'suami', *'saya'*, 'seumuran', *'sama'*, *'saya'*, 'mba', *'malah'*, 'tuaan', *'saya'*, *'beberapa'*, *'bulan'*, *'tapi'*, 'alhamdulillah', 'suami', *'saya'*, *'gak'*, 'kekanak', 'kanakan', *'hehe'* |

4.2.4 Normalization
   Table 7 displays the effects of the adjustments made during the normalization stage. The transition from non-standard to standard words, such as the word "yaa" becoming "iya," "lo" becoming "kamu", "mba" becoming "mbak", "gw" becoming "saya" and others, is what distinguishes Table 6 from Table 7. The comment's words will be modified.

**Table 7.** Instagram data normalization process

| No | Instagram Account | Comments |
|---|---|---|
| 1 | @Username 1 | 'bego', *'kamu'* |
| 2 | @Username 2 | 'editan', 'banget', 'ampun', 'mommy', 'lambe', 'kehabisan', 'gosib', 'pelakor' |
| 3 | @Username 3 | 'goblok', *'saya'*, 'nonton', 'berkali', 'kaget' |
| 4 | @Username 4 | 'usa', 'pakai', 'anjing', 'babi', *'saya'*, 'bls', 'gitu', 'wkwk' |
| 5 | @Username 5 | 'masyaallah', 'bagus', 'banget' |
| up to | … | … |
| 400 | @Username 400 | 'suami', 'seumuran', *'mbak'*, 'tuaan', 'alhamdulillah', 'suami', 'kekanak', 'kanakan' |

## 4.2.5 Stemming

As illustrated in Table 8, the final step is stemming, which makes use of the Sastrawi libraries. The italicized and bolded phrases "editan" "kehabisan" "berkali" "seumuran" "tuaan" "kenakan" and others are different in Tables 7 and 8. The comment's words will be modified.

**Table 8.** Stemming process on Instagram data

| No | Instagram Account | Comments |
|----|------------------|----------|
| 1 | @Username 1 | 'bego', 'kamu' |
| 2 | @Username 2 | ***'edit'***, 'banget', 'ampun', 'mommy', 'lambe', ***'habis'***, 'gosib', 'pelakor' |
| 3 | @Username 3 | 'goblok', 'saya', 'nonton', ***'kali'***, 'kaget' |
| 4 | @Username 4 | 'usa', 'pakai', 'anjing', 'babi', 'saya', 'bls', 'gitu', 'wkwk' |
| 5 | @Username 5 | 'masyaallah', 'bagus', 'banget' |
| up to | … | … |
| 400 | @Username 400 | 'suami', ***'umur'***, 'mbak', 'tua', 'alhamdulillah', 'suami', 'kekanak', ***'kana'*** |

The whole document that is produced after the preprocessing stage is 2109 words long.

### 4.3 Term weighting

After getting 400 data from preprocessing, the word is converted into vector data by multiplying tf*idf, resulting in 2109 syllables or words. As shown in Table 9, the following terms are obtained after the preprocessing results:

**Table 9.** The term "data" is the result of the preprocessing process

| No | Term |
|----|------|
| 1 | abang |
| 2 | acara |
| 3 | account |
| 4 | adat |
| 5 | adik |
| 6 | admin |
| 7 | adaptasi |
| 8 | Adek |
| 9 | Agama |
| 10 | Ajar |
| 11 | akhirat |
| 12 | akhlak |
| 13 | akun |
| 14 | Alam |
| … | … |
| 2109 | zina |

The technique of multiplying Term Frequency (Tf) by Inverse Document Frequency is the following step (Idf). Table 10 displays the Tf process results, and Table 11 displays the tf-idf results.

### 4.4 Clustering

After the preprocessing stage, which converts terms into vector data using Tf*Idf multiplication, K-means is used for clustering. Figure 3 and Figure 4 show the Output result from Cluster 0 and 1.

There are 363 data items in cluster 0, all of which are found in comments that do not indicate cyberbullying, as shown in Figure 3.

**Table 10.** Terms for results frequency

| Doc | abang | acara | account | adat | adik | Up to | zina |
|-----|-------|-------|---------|------|------|-------|------|
| 1 | 0 | 0 | 0 | 0 | 0 | … | 0 |
| 2 | 0 | 0 | 0 | 0 | 0 | … | 0 |
| 3 | 0 | 0 | 0 | 0 | 0 | … | 0 |
| 4 | 0 | 0 | 0 | 0 | 0 | … | 0 |
| 5 | 0.316 | 0 | 0 | 0 | 0 | … | 0 |
| 6 | 0 | 0 | 0 | 0 | 0 | … | 0 |
| 7 | 0 | 0 | 0 | 0 | 0 | … | 0 |
| 8 | 0 | 0 | 0 | 0 | 0 | … | 0 |
| 9 | 0 | 0 | 0 | 0 | 0 | … | 0 |
| 10 | 0 | 0 | 0 | 0 | 0 | … | 0 |
| up to | … | … | … | … | … | … | 0 |
| 400 | 0 | 0 | 0 | 0 | 0 | … | 0 |

**Table 11.** Tf-idf results

| Doc | abang | acara | account | adat | adik | Up to | zina |
|-----|-------|-------|---------|------|------|-------|------|
| 1 | 0 | 0 | 0 | 0 | 0 | … | 0 |
| 2 | 0 | 0 | 0 | 0 | 0 | … | 0 |
| 3 | 0 | 0 | 0 | 0 | 0 | … | 0 |
| 4 | 0 | 0 | 0 | 0 | 0 | … | 0 |
| 5 | 0.349 | 0 | 0 | 0 | 0 | … | 0 |
| 6 | 0 | 0 | 0 | 0 | 0 | … | 0 |
| 7 | 0 | 0 | 0 | 0 | 0 | … | 0 |
| 8 | 0 | 0 | 0 | 0 | 0 | … | 0 |
| 9 | 0 | 0 | 0 | 0 | 0 | … | 0 |
| 10 | 0 | 0 | 0 | 0 | 0 | … | 0 |
| up to | … | … | … | … | … | … | 0 |
| 400 | 0 | 0 | 0 | 0 | 0 | … | 0 |

According to the data in Figure 4, cluster 1 comprises 37 records of data, all of which are in comments that contain aspects of cyberbullying.

The information on the clustering results obtained is shown in Figure 3 and Figure 4. There are three columns presented in the two figures. The first column displays the id, the second column displays clusters, and the third column displays Instagram text or comments. For example, in cluster 1, the data is in comments with cyberbullying elements. The results of the research that has been done are comments with id 16 and 17, which are indicated as cyberbullying comments. Can be seen in Figure 5.

The output generated in cluster 1 indicates that it is included in the bullying category. There needs to be a next stage where the results only display the ID. From the ID, it is then clarified with the existing dataset to find the username and comments of the perpetrator who carried out the cyberbullying action. These results can be used as initial digital evidence for trial purposes.

### 4.5 Evaluation

The accuracy of the clustering results generated by the system is evaluated by testing the confusion matrix model. Clustering with threshold values of 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, and 1.0 produced the data that will be analyzed. The information is utilized to contrast the clustering outcomes that the system found using various data sets. Data must have a label before being submitted into the dataset test. The following comparisons are made about how the system labeled and tested the cluster's 1.0 threshold value. The outcomes, when arranged by group, are shown in Table 12.

**Figure 3.** Output result from Cluster 0

| id | .. | .. | .. | .. | cluster ↑ | text |
|---|---|---|---|---|---|---|
| 2 | . | . | . | . | cluster_0 | geblek tatacowo banget bain balikanhadewwntar tinggal nyalahin cowopadahal kiten... |
| 3 | . | . | . | . | cluster_0 | kmrn mewek lengket duhhh labil banget mbak abege ajah kmrn cari sensasi markot... |
| 10 | . | . | . | . | cluster_0 | anyiennnnggg suara ancur banget merdu tukang goreng |
| 19 | . | . | . | . | cluster_0 | username ampun upil naruto cermin pecunpengalaman banget kasihan prihatin mat... |
| 39 | . | . | . | . | cluster_0 | jijik banget lihat tingkah upil mudah cepat tidur |
| 49 | . | . | . | . | cluster_0 | amanda muka tante umur muka boros banget dasar cabe jujur suka haltis |
| 51 | . | . | . | . | cluster_0 | yaelah laki laki jaman bejat ustad kayak kayak gila banget perempuan istri stju bang... |
| 55 | . | . | . | . | cluster_0 | username pikir manusia bodoh eehhh tanteibuomaempo otak picik banget manusi l... |
| 58 | . | . | . | . | cluster_0 | username pikir manusia bodoh eehhh tanteibuomaempo otak picik banget manusi l... |
| 60 | . | . | . | . | cluster_0 | username didik ortunya nagita bagus heran balajaer nagita rugi junjunganny benci b... |
| 68 | . | . | . | . | cluster_0 | munafik banget omong curhat sosmed kampung kitu keleusss wkkwkww mending ... |
| 82 | . | . | . | . | cluster_0 | najisapalagi philipine banget kalid banggajd usak moral bangganajis banget liatnya... |
| 97 | . | . | . | . | cluster_0 | dasar plagiat kreatipmemalukab banget |
| 113 | . | . | . | . | cluster_0 | edit banget ampun mommy lambe habis gosib pelakor |

**Figure 3.** Output result from Cluster 0

| id | .. | .. | .. | .. | cluster ↓ | text |
|---|---|---|---|---|---|---|
| 1 | . | . | . | . | cluster_1 | username tolol hubung gugur pakai hijab syar bayi panas dalem hubung woyyyy otak... |
| 4 | . | . | . | . | cluster_1 | inti jengkel gausah anak kasi kembang psikis anak orang tolol anak anak orang ben... |
| 5 | . | . | . | . | cluster_1 | hadewwwww permpuan lgsakit jiwaknp peran utama film hantu jeruk purutky khabis ... |
| 6 | . | . | . | . | cluster_1 | pantesan tinggalin laki laki mikir kali perempuan kayagni urus becus urus anak men... |
| 7 | . | . | . | . | cluster_1 | balajaer nyampah instagram artissuka instagram caption apakok balajaer heboh asi... |
| 8 | . | . | . | . | cluster_1 | rakyat indonesia bodoh beda buruk prilaku buruk sopan santun pilih panutan idola m... |
| 9 | . | . | . | . | cluster_1 | janda bego suami nikah |
| 11 | . | . | . | . | cluster_1 | syarat nikah agama islam saksi wali nikah kawinmahar perkara makeup manglingi ... |
| 12 | . | . | . | . | cluster_1 | cewek gatau ngerebut pacar orang abang ngedukung sifat jelek ckck |
| 13 | . | . | . | . | cluster_1 | username situ kids jaman micin ngina bego pelihara |
| 14 | . | . | . | . | cluster_1 | jijik liat suer minceeee |
| 15 | . | . | . | . | cluster_1 | tabu cium gitu eneg citra alim bijak citra alim suci serius |
| 16 | . | . | . | . | cluster_1 | cium orang bnyk malu kecuali mabok malu teman helooow teman libat cium bibir co... |
| 17 | . | . | . | . | cluster_1 | anjing ngamukhidup dmana musuhudh kliatan bukanorg musuh |

**Figure 4.** Output result from Cluster 1

| 16 | . | . | . | . | cluster_1 | cium orang bnyk malu kecuali mabok malu teman helooow teman libat cium bibir co... |
|---|---|---|---|---|---|---|
| 17 | . | . | . | . | cluster_1 | anjing ngamukhidup dmana musuhudh kliatan bukanorg musuh |

**Figure 5.** Example id indicated as a bullying comment that has been clarified with data

**Table 12.** Cluster results grouping (1.0 threshold)

| Cluster | Lots of Comments | Category |
|---|---|---|
| Cluster 0 | 363 Comments | Non-Cyberbullying |
| Cluster 1 | 37 Comments | Cyberbullying |

In this instance, the researcher must decide each cluster's category because the system cannot do so. The system produces only groups of each given comment. To make calculating the values of accuracy, precision, and recall easier, as shown in Table 13, relevant and irrelevant data are connected to the clustered data.

**Table 13.** Results of a cluster

| | | Actual (Real) | | Total Prediction |
|---|---|---|---|---|
| | | C0 | C1 | |
| Prediction | C0 | 221 | 142 | 363 |
| (system) | C1 | 1 | 36 | 37 |
| **Actual Total** | | **200** | **200** | **N=400** |

The accuracy numbers obtained are as follows based on Table 13 that has been completed: accuracy $= \frac{257}{400} \times 100\% = 64.25\%$. Meanwhile, to calculate the value of multi-class precision clustering is as follows:

$$P\,(Non - Cyberbullying) = \frac{221}{221 + 142}x100\% = 60.88\%$$
$$P\,(Cyberbullying) = \frac{36}{36 + 1}x100\% = 97.29\%$$
$$All\ precision = \frac{0.6088 + 0.9729}{2}x100\% = 79.29\%$$

Meanwhile, to calculate the recall value of multi-class clustering is as follows:

$$P\,(Non - Cyberbullying) = \frac{221}{221 + 1}x100\% = 99.54\%$$
$$P\,(Cyberbullying) = \frac{36}{36 + 142}x100\% = 20.22\%$$
$$All\ recall = \frac{0.6088 + 0.9729}{2}x100\% = 59.88\%$$

Table 14 displays the precision and recall values for each cluster.

**Table 14.** Results for precision and recall

| Cluster | Category | Precision | Recall |
|---------|----------|-----------|--------|
| Cluster 0 | Non-Cyberbullying | 60.88% | 99.54% |
| Cluster 1 | Cyberbullying | 97.29% | 20.22% |
| | **Average** | **79.29%** | **59.88%** |

It is clear from Table 14 above that cluster 1 has a good degree of precision. This indicates that more accurate data were obtained within a single group than inaccurate data. Comparatively speaking, cluster 0 is less precise than cluster 1. Table 15 compares the accuracy, precision, and recall values with thresholds of 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0, 8, 0.9, and 1.0.

**Table 15.** Values for recall, precision, and accuracy are compared

| No | Threshold | Accuracy | Precision | Recall |
|----|-----------|----------|-----------|--------|
| 1 | 0.1 | 50.00% | 50.00% | 50.00% |
| 2 | 0.2 | 52.50% | 53.91% | 52.50% |
| 3 | 0.3 | 54.17% | 58.57% | 54.17% |
| 4 | 0.4 | 55.62% | 59.93% | 55.63% |
| 5 | 0.5 | 61.50% | 65.55% | 61.50% |
| 6 | 0.6 | 55.00% | 62.94% | 55.00% |
| 7 | 0.7 | 53.57% | 57.62% | 53.57% |
| 8 | 0.8 | 62.50% | 71.06% | 62.51% |
| 9 | 0.9 | 58.33% | 69.42% | 58.34% |
| 10 | 1.0 | 64.25% | 79.29% | 59.88% |

Table 14 and Figure 6 show that the level of assessment using the confusion matrix with a threshold of 1.0 is obtained, with accuracy values of 64.25 percent and precision values of 79.29 percent, respectively. From the other data tests, the recall is 59.88 percent.

The application of the K-means Algorithm can work properly using as many as 400 data. One of the suggested algorithms and a popular one for analyzing Instagram social media is K-means. The number of clusters that have been identified is used in this study to calculate K = 2. Rapidminer software is used to process the data. The study findings reveal

that a threshold value of 1.0 results in a cluster of 0, which has 363 records. cluster 1 that was created, meanwhile, has 37 records in it. Another information discovered is the deal of accuracy, average precision, and recall. With accuracy values of 64.25 percent and precision values of 79.29 percent, respectively. From the other data tests, the recall is 59.88 percent. These findings came from testing several pieces of information. In earlier studies [47], the results of testing for cyberbullying on Instagram comments using the Support Vector Machine Classification Method with 400 data records generated the best accuracy rate of 90%, a precision of 94.44%, and a recall of 85%. Therefore, research using the K-means clustering approach yields results with a lesser level of accuracy than research using the Support Vector Machine method in the past. Because there are fewer combinations of data obtained, fewer datasets are used, and more different data features are employed, the accuracy of the K-means method is lower than that of the SVM (Support Vector Machine) algorithm.
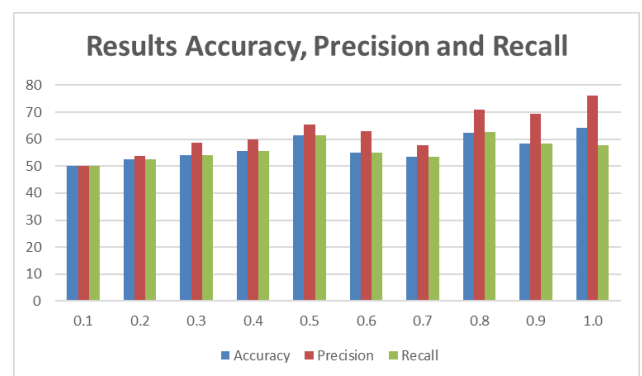


**Figure 6.** Results from the confusion matrix

This research faces several challenges and limitations. One of them is the limited amount of data used in this study, which is only 400 data. This limitation can affect the representativeness and accuracy of the analysis results. In addition, this research also needs help in interpreting Indonesian language texts. The lack of access to the complete Indonesian language corpus and the use of regional languages and slang in Instagram comments hinder accurate processing and understanding of the text. As a result, the results of data mining analysis have yet to reach the expected level.

To address this challenge, future research could expand the amount of data used, increase the variety of data characteristics, and assemble a richer corpus of Indonesian. In addition, developing more sophisticated natural language processing techniques to understand regional languages and slang can improve the accuracy of the analysis. A combination approach between the K-means clustering method and other methods, such as SVM, can also be explored to improve the quality of the analysis results. Cross-disciplinary collaboration and more advanced technology will be the key to overcoming this challenge and producing a more accurate and comprehensive social media analysis in the Indonesian context.

## 5. CONCLUSION

The findings and recommendations of this study indicate that in the many data groups analyzed, using a threshold value of 1.0, the highest levels of accuracy, precision, and recall

were obtained compared to other data groups, especially with an accuracy of 64.25%, a precision of 79.29% and a recall of 59.88 % However, it should be noted that the test results in previous studies using the Support Vector Machine Classification Method produced the best accuracy rate of 90%, precision of 94.44%, and recall of 85%. Thus, this study using the K-means clustering method has a lower accuracy level than previous studies using the Support Vector Machine method. The results of a cyberbullying comment detection simulation that generates user IDs so that it is possible to retrieve usernames and comments from perpetrators can serve as initial digital evidence that is relevant for trial purposes.

Recommendations for further research are to expand the amount of data used, increase the variety of data characteristics, and collect a richer corpus of Indonesian. In addition, the development of more sophisticated natural language processing techniques to understand regional languages and slang also needs attention to increase the accuracy of the analysis. A combination approach between the K-means clustering method and other methods, such as SVM, can also be explored to improve the quality of the analysis results. Cross-disciplinary collaboration and more advanced technology will be vital in overcoming this challenge and producing a more accurate and comprehensive social media analysis in the Indonesian context.

## REFERENCES

[1] Wijayanto, A., Riadi, I., Prayudi, Y., Sudinugraha, T. (2022). Network forensics against address resolution protocol spoofing attacks using trigger, acquire, analysis, report, action method. Register: Jurnal Ilmiah Teknologi Sistem Informasi, 8(2): 156-169. http://doi.org/10.26594/register.v8i2.2953

[2] Riadi, I., Siregar, N.H. (2022). Mobile forensic analysis of signal messenger application on android using Digital Forensic Research Workshop (DFRWS) Framework. Ingénierie des Systèmes d'Information, 27(6): 903-913. https://doi.org/ https://doi.org/10.18280/isi.270606

[3] Ministry of Communications and Information Technolog. (2021). Networking is increasing, indonesia needs to increase cultural values on the internet. https://aptika.kominfo.go.id/2021/09/Warganet-Meningkat-Indonesia-Perlu-Tingkatkan-Nilai-Budaya-Di-Internet/.

[4] Noh, C.H.C., Ibrahim, M.Y. (2014). Kajian penerokaan buli siber dalam kalangan pelajar UMT. Procedia-Social and Behavioral Sciences, 134: 323-329. https://doi.org/10.1016/j.sbspro.2014.04.255

[5] Herman, Riadi, I., Rafiq, I.A. (2022). Forensic mobile analysis on social media using national institute standard of technology method. International Journal of Safety & Security Engineering, 12(6): 707-713. https://doi.org/10.18280/ijsse.120606

[6] Riadi, I., Herman, Rafiq, I.A. (2022). Mobile forensic investigation of fake news cases on instagram applications with digital forensics research workshop framework. International Journal of Artificial Intelligence Research, 6(2). https://doi.org/10.29099/ijair.v6i2.311

[7] Riadi, I., Sunardi, Widiandana, P. (2022). Cyberbullying detection on instant messaging services using rocchio and digital forensics research workshop framework.

Journal of Engineering Science and Technology, 17(2): 1408-1421.

[8] Ardi, Z., Putri, S.A. (2020). The analysis of the social media impact on the millennial generation behavior and social interactions. Southeast Asian Journal of Technology and Science, 1(2): 70-77. https://doi.org/10.29210/81065100

[9] Chan, T.K.H., Cheung, C.M.K., Lee, Z.W.Y. (2021). Cyberbullying on social networking sites: A literature review and future research directions. Information Management, 58(2): 103411. https://doi.org/10.1016/j.im.2020.103411

[10] Fazry, L., Apsari, N.C. (2021). The influence of social media on cyberbullying behavior among teenager. Jurnal Penelitian dan Pengabdian Kepada Masyarakat, 2(1): 28-36. https://doi.org/10.24198/jppm.v2i1.33435

[11] Pillai, T.M.R., Vasudevan, H. (2020). Workplace bullying and management of mistreated behaviour: A case study in the banking sector. IIUM Journal of Case Studies in Management, 11(2): 15-22.

[12] Rahman, A., Zaman, N., Asyhari, A. T., Sadat, S.M.N., Pillai, P., Abdullah, R. (2021). Ad hoc networks spy-bot: Machine learning-enabled post filtering for social network-integrated industrial internet of things. Ad Hoc Networks, 121: 102588. https://doi.org/10.1016/j.adhoc.2021.102588

[13] Habibi, M., Cahyo, P.W. (2019). Clustering user characteristics based on the influence of hashtags on the instagram platform. IJCCS (Indonesian Journal of Computing and Cybernetics Systems), 13(4): 399. https://doi.org/10.22146/ijccs.50574

[14] Vankayalapati, R., Ghutugade, K.B., Vannapuram, R., Prasanna, B.P.S. (2021). K-means algorithm for clustering of learners performance levels using machine learning techniques. Revue d'Intelligence Artificielle, 35(1): 99-104. https://doi.org/10.18280/ria.350112

[15] Dewi, I.C., Gautama, B.Y., Mertasana, P.A. (2017). Analysis of clustering for grouping of productive industry by k-medoid method. International Journal of Engineering and Emerging Technology, 2(1): 26. https://doi.org/10.24843/ijeet.2017.v02.i01.p06

[16] Jarrah, A., Amri, S. (2022). Optimized fpga-based implementation of brain tumor detection by combining K-means and grey wolf optimization algorithms. Traitement du Signal, 39(6): 1879. https://doi.org/10.18280/ts.390601

[17] Khan, K.B., Khaliq, A.A., Shahid, M., Khan, S. (2016). An efficient technique for retinal vessel segmentation and denoising using modified ISODATA and CLAHE. IIUM Engineering Journal, 17(2): 31-46. https://doi.org/10.31436/iiumej.v17i2.611

[18] Daoudi, S., Zouaoui, C.M.A., El-Mezouar, M.C., Taleb, N. (2021). Parallelization of the K-means++ clustering algorithm. Ingénierie des Systèmes d'Information, 26(1): 59-66. https://doi.org/10.18280/isi.260106

[19] Gustientiedina, G., Adiya, M.H., Desnelita, Y. (2019). Application of the K-means algorithm for clustering drug data. Jurnal Nasional Teknologi Dan Sistem Informasi, 5(1): 17-24. https://doi.org/10.25077/teknosi.v5i1.2019.17-24

[20] Al-Rahmi, W.M., Yahaya, N., Alamri, M.M., Aljarboa, N.A., Bin Kamin, Y., Moafa, F.A. (2019). A model of factors affecting cyber bullying behaviors among university students. IEEE Access, 7: 2978-2985.

https://doi.org/10.1109/access.2018.2881292

[21] Romando, Y., Sulistyowati, R., Wibisono, I.S. (2019). Identification of negative comments in Indonesian on Instagram using K-means. Multimatrix, II(1): 6-8.

[22] Andriansyah, M., Akbar, A., Ahwan, A., Gilani, N.A., Nugraha, A.R., Sari, R.N., Senjaya, R. (2017). Cyberbullying comment classification on Indonesian selebgram using support vector machine method. In 2017 Second International Conference on Informatics and Computing (ICIC), Jayapura, Indonesia, pp. 1-5. https://doi.org/10.1109/IAC.2017.8280617

[23] Nurrahmi, H., Nurjanah, D. (2018). Indonesian twitter cyberbullying detection using text classification and user credibility. In 2018 International Conference on Information and Communications Technology (ICOIACT), Yogyakarta, Indonesia, pp. 543-548. https://doi.org/10.1109/ICOIACT.2018.8350758

[24] Zhao, R., Mao, K. (2017). Cyberbullying detection based on semantic-enhanced marginalized denoising auto-encoder. IEEE Transactions on Affective Computing, 8(3): 328-339. https://doi.org/10.1109/TAFFC.2016.2531682

[25] Imam Riadi, P.W., Sunardi. (2021). Investigation of cyberbullying on WhatsApp using digital forensics. Rekayasa Sistem dan Teknologi Informasi, 1(10): 730-735. https://doi.org/10.29207/resti.v4i4.2161

[26] Riadi, I., Sunardi, S., Widiandana, P. (2020). Mobile forensics for cyberbullying detection using term frequency-inverse document frequency (tf-idf). Jurnal Ilmiah Teknik Elektro Komputer dan Informatika, 5(2): 68. http://dx.doi.org/10.26555/jiteki.v5i2.14510

[27] Hang, O.C., Dahlan, H.M. (2019). Cyberbullying lexicon for social media. 2019 6th International Conference on Research and Innovation in Information Systems (ICRIIS), Johor Bahru, Malaysia, pp. 1-6. https://doi.org/10.1109/ICRIIS48246.2019.9073679

[28] Gorro, K.D., Sabellano, M.J.G., Gorro, K., Maderazo, C., Capao, K. (2018). Classification of cyberbullying in facebook using selenium and SVM. In 2018 3rd International Conference on Computer and Communication Systems (ICCCS), Nagoya, Japan, pp. 183-186. https://doi.org/10.1109/CCOMS.2018.8463326

[29] Amali, H.I., Jayalal, S. (2020). Classification of cyberbullying Sinhala language comments on social media. In 2020 Moratuwa Engineering Research Conference (MERCon), Moratuwa, Sri Lanka, pp. 266-271. https://doi.org/10.1109/MERCon50084.2020.9185209

[30] Tapia, F., Aguinaga, C. (2018). Detección de patrones de comportamiento a través de redes sociales como twitter, utilizando técnicas de minería de datos como método para detectar el acoso cibernético. 2018 7th International Conference on Software Process Improvement (CIMPS), pp. 111-118. https://doi.org/10.1109/CIMPS.2018.8625625

[31] Rsa. (2016). Current State of Cybercrime. Available: Https://www.Rsa.com/Content/Dam/Rsa/Pdf/2016/05/2016-current-state-of-cybercrime.pdf.

[32] Christopher, H.S., Manning, D., Raghavan, P. (2009)., Introduction to modern information retrieval. Cambridge University Press Cambridge, England.

[33] Ruthven, I., Lalmas, M. (2003). A survey on the use of relevance feedback for information access systems. Knowledge Engineering Review, 18(2): 95-145. https://doi.org/10.1017/S0269888903000638.

[34] Sahu, S.K., Sarangi, S., Jena, S.K. (2014). A detail analysis on intrusion detection datasets. In 2014 IEEE International Advance Computing Conference (IACC), Gurgaon, India, pp. 1348-1353. https://doi.org/10.1109/Iadcc.2014.6779523

[35] Riyadduloh, R., Romadhony, A. (2021). Normalization of Indonesian text based on a slang dictionary: A case study of gadget product tweets on Twitter. eProceedings of Engineering, 8(4): 4216-4228. Available: https://openlibrarypublications.telkomuniversity.ac.id/index.php/engineering/article/view/15246/14969.

[36] Selberg, E.W. (1997). Information retrieval advances using relevance feedback. UW Dept CSE Gen. Exam. Retrieved from http://www.cs.rpi.edu/~chapaa/userskill/paper/generals.pdf.

[37] Suwija Putra, I.M., Adiwinata, Y., Singgih Putri, D.P., Sutramiani, N.P. (2021). Extractive text summarization of student essay assignment using sentence weight features and fuzzy C-Means. International Journal of Artificial Intelligence Research, 5(1): 13-24. https://doi.org/10.29099/ijair.v5i1.187

[38] Yugianus, P., Dachlan, H.S., Hasanah, R.N. (2013). Development of a library catalog search system using the Rocchio relevance feedback method. Jurnal EECCIS (Electrics, Electronics, Communications, Controls, Informatics, Systems), 7(1): 47-52. https://doi.org/10.21776/jeeccis.v7i1.201

[39] Jumeilah, F.S. (2017). Application of Support Vector Machine (SVM) for research categorization. Jurnal RESTI (Rekayasa Sistem Dan Teknologi Informasi), 1(1): 19-25. https://doi.org/10.29207/resti.v1i1.11

[40] Ridho Lubis, A., Nasution, M.K.M., Salim Sitompul, O., Muisa Zamzami, E. (2021). The effect of the tf-idf algorithm in times series in forecasting word on social media. Indonesian Journal of Electrical Engineering and Computer Science, 22(2): 976-984. https://doi.org/10.11591/ijeecs.v22.i2.pp976-984

[41] Shehzad, F., Rehman, A., Javed, K., Alnowibet, K.A., Babri, H.A., Rauf, H.T. (2022). Binned term count: An alternative to term frequency for text categorization. Mathematics, 10(21): 4124. https://doi.org/10.3390/math10214124

[42] Muhariya, A., Riadi, I., Prayudi, Y. (2022). Cyberbullying analysis on Instagram using K-means clustering. JUITA: Jurnal Informatika, 10(2): 261-271. https://doi.org/10.30595/juita.v10i2.14490

[43] Widyasanti, N.K., Putra, I.D., Rusjayanthi, N.D. (2018). Word weight feature selection using the TFIDF method for Indonesian summaries. Jurnal Ilmiah Merpati (Menara Penelitian Akademika Teknologi Informasi), 6(2): 119. https://doi.org/10.24843/jim.2018.v06.i02.p06

[44] Malhotra, A., Jindal, R. (2022). Deep learning techniques for suicide and depression detection from online social media: A scoping review. Applied Soft Computing, 109713. https://doi.org/10.1016/j.asoc.2022.109713

[45] Wirasto, A., Nisa, K. (2021). Comparison of machine learning algorithms for classification of drug groups. Jurnal Ilmiah Sistem Informasi dan Teknik Informatika, 11(2): 196-207.

https://dx.doi.org/10.30700/jst.v11i2.1134

[46] Dinata, R.K., Safwandi, S., Hasdyna, N., Azizah, N. (2020). A K-means clustering analysis on motorcycle data. INFORMAL: Informatics Journal, 5(1): 10-17. https://doi.org/https://doi.org/10.19184/isj.v5i1.17071

[47] Luqyana, W.A., Cholissodin, I., Perdana, R.S. (2018). Sentiment analysis of cyberbullying in Instagram comments using the Support Vector Machine classification method. Jurnal Pengembangan Teknologi Informasi dan Ilmu Komputer, 2(11): 4704-4713.