

Detection of Hoax Spread in The Whatsapp Group with Lexicon Based and Naive Bayes Classification



Jordan Andrean, Suharjito

Abstract: Spreading hoax through WhatsApp social media can lead to different beliefs and can cause disputes for those affected. This paper proposes a hybrid model for finding hoaxes in the WhatsApp group using a combination of knowledge-based and machine learning approaches. This Hybrid model combines two methods namely Lexicon based and Naive Bayes Classifier which will be applied to the WhatsApp monitoring application. This research focuses on two main aspects namely word weighting using the lexicon based method and data classification using the Naive Bayes Classifier and Decision tree-j48 methods. The dataset used is conversation data that is crossed from the WhatsApp group. Based on the experiments that have been carried out, it is obtained the results of classification using Naive Bayes classifier of 86.670% data conversation not indicated hoaxes and 13.330% indicated hoaxes. The average value of the percentage of truth obtained more than 75%. The average value of the classification performance evaluation results in a precision value of 0.771, a recall value of 0.754, an F-measure value of 0.773.

Keywords : WhatsApp monitoring, Hoax, Hybrid Approach, Lexicon based, Naive Bayes Classifier.

I. INTRODUCTION

Information diffusion about events that occur in the community quickly prominent in a variety of social media. These social media include Wikipedia, Facebook, Youtube, Twitter, Tumblr, BBM, WhatsApp, Instagram, and many more that can be used for social media [1]. This is an opportunity for parties who have a specific purpose to spread information about events that are not yet known to be true or hoaxes. Hoax distribution is often done through WhatsApp social media because it is considered easy for the public to deliver messages and is supported by group facilities that can accommodate more than 50 users. This causes the perpetrators of the spread of hoaxes to easily spread the hoax in the middle of the group. According to research from [2], Whatsapp application is chosen by many people (individuals, groups, organizations and even government) as a medium for delivering messages because it is considered more effective

and is a satisfaction when the information delivered is right on target. This research utilizes the Hybrid approach to the WhatsApp monitoring application to find the source of hoax in the WhatsApp group. The Hybrid approach used is a combination of a knowledge-based and machine learning approach. This hybrid model combines two methods namely lexicon based and Naive Bayes classifier. The lexicon based method is used in the word weighting process, which is a conversation dataset from preprocessing results compared to a dataset from the lexicon dictionary so that values are obtained for each word [3]. And then, the conversation dataset is classified using the naive bayes classifier method, which results in a percentage of conversation data that is indicated to be hoaxed or not. Data indicated hoax conversation if included in the negative sentiment category and not indicated hoax if included in the positive or neutral category. The negative sentiment category means that the conversation data contains hoax elements such as using emotional and provocative language. While the positive sentiment category means that the conversation data is not indicated as a hoax because it does not contain hoax elements.

II. RELATED WORK

This research was conducted by [3], proposes a lexicon-based approach to conducting entity level sentiment analysis on Twitter. Through the Chi-square test on the output, tweets containing opinions can be identified. A binary classifier is then trained to assign sentiment polarity to the tweet that has just been identified, the training data provided by the lexicon-based method. A study by [4], has proposed a Machine-Human (MH) model for detecting false news on social media. This model combines a human literacy news detection tool and a machine linguistic approach and a network-based approach. The model was stated to be able to improve the ability of humans to distinguish fake news with higher accuracy than when they did without using a model. The next interesting topic raised by [5], This paper proposes a hybrid approach that combines node embeddings and user-based features to enrich the detection of SOFNs on the Twitter social network.

This research shows knowledge extracted from social network graphs using node2vec is able to provide a general way to improve social networking embeddings and more helpful in detecting SOFNs.

Revised Manuscript Received on March, 28 2020.

* Correspondence Author

Jordan Andrean*, Magister Teknik Informatika, BINUS University, Jakarta, Indonesia

Suharjito, Magister Teknik Informatika, BINUS University, Jakarta, Indonesia.

© The Authors. Published by Blue Eyes Intelligence Engineering and Sciences Publication (BEIESP). This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>)

The hybrid approach is also used for the message sentiment analysis on Twitter proposed by [6]. This approach combines two methods namely Lexicon based and Machine Learning based. Based on experiments that have been done, it is found that the use of the Lexicon based approach shows high precision but low memory causing performance problems. To improve performance, the two approaches are combined. Empirical evaluations are carried out with a variety of different training data sets to ensure that the proposed approach is very effective and better for Twitter sentiment analysis.

Table- I: Existing Approaches And Techniques

Authors	Approches and Techniques	Evaluation Result
[3]	Lexicon-based and Learning-based Methods	The ME method produces an accuracy of 0.756 with a value of Precision 0.170, Recall 0.202 and an F-score of 0.184. The FBS method produces an accuracy of 0.878 with a value of Precision 0.564, Recall 0.556 and an F-score of 0.560. The AFBS method produces an accuracy of 0.868 with a value of Precision 0.522, Recall 0.582 and an F-score of 0.569.
[4]	Combines the human literacy news detection tool and the machine linguistic and network-based approaches.	MH (Machine Human) = $\sum [A + B + C + D \dots + J] \leq 100$. If the MH results are ≤ 100 then the news is true and vice versa if $MH > 100$, then false news or hoaxes.
[5]	Twitter network analysis and machine learning.	Decission tree method produces 0.927 accuracy with a value of Precision 0.763, Recall 0.781 and F-score 0.928. The K-NN method produces an accuracy of 0.962 with a value of Precision 0.970, Recall 0.777 and F-score 0.956. The SVM method produces an accuracy of 0.980 with a value of Precision 0.976, Recall 0.765 and an F-score of 0.954.
[6]	Lexicon Based Approach and Machine Learning Approach.	Training data = 5000, testing data = 200 produces 60.53% unigram, 59.38% bigram and 57.13% trigram.

III. THE PROPOSED TECHNIQUE

Figure 1. Will show the workflow of the approach model that we propose.

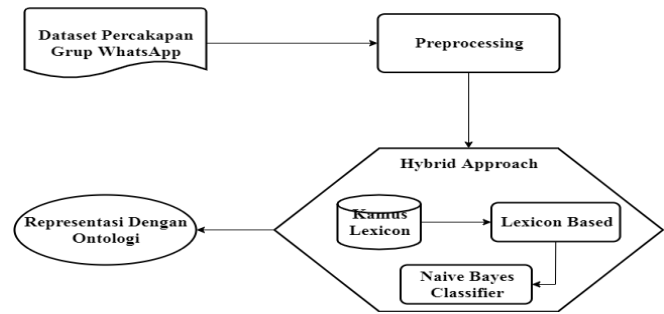


Fig. 1. Application workflow model

Fig 1. Shows the proposed application workflow model. The model includes the following steps.

STEP 1: DATASET

Conversation dataset is collected using crawling techniques through WhatsApp groups using Selenium and Redist. From the selenium tag class that initializes the group a dump process will be carried out to get the WhatsApp number, username and conversation information from that user. While Redist is used to process the session login from WhatsApp Web to generate WhatsApp Web Qrcode.

Table- II: Example of Conversation Dataset

No	1
User Account	+6283831286609
Conversation	Can be analyzed who did the offender who insulted our brother Papua, so that he could be tried quickly so that the problem can be resolved quickly just dismantle the account.
WhatsApp Groups Followed	<ul style="list-style-type: none"> Indonesia Bersatu WARTA TV POLRI NEWS JOKOWI SEKALI LAGI

Table- II shows the attributes of data that has been successfully crossed, namely the user account that contains the user's telephone number, user conversations and what groups the user has followed.

STEP 2: PREPROCESSING

The conversation data preprocessing stage consists of four stages:

1. Case Folding

The process for converting all uppercase contained in a conversation to lowercase.

Table- III: Case Folding Stage

Conversation	Case Folding Results
Can be analyzed who did the offender who insulted our brother Papua, so that he could be tried quickly so that the problem can be resolved quickly just dismantle the account.	can be analyzed who did the offender who insulted our brother Papua, so that he could be tried quickly so that the problem can be resolved quickly just dismantle the account

Table- III shows the process of changing uppercase to lowercase letters in a user's conversation.

2. Tokenizing

The process of breaking a conversation into tokens uses delimiter spaces.

Table- IV: Tokenizing Stage

Unigram (one word)	can, analyzed, who, the culprit, who, maki2, brothers, us, papua, let, quickly, be tried, so, quickly, finished, problem, this, dismantle, alone, account, that
Bigram (two word)	can be analyzed, who is the culprit, who abuses us, our brothers, Papua so that we can be tried quickly, so quickly, problem is over, just dismantle, just account, analyzed who, the culprit, who insulted you, we Papua, so quickly, tried so quickly, this problem, just dismantle it, the account
Trigram (three word)	Can be analyzed who, the offender is insulting, our brothers and sisters in Papua, so that they can be tried quickly, so that it can be resolved quickly, this problem is dismantled, just that account

Table- IV shows the process of breaking words into three categories, namely Unigram, Bigram and Trigram. Solving words is done because a word can have different weights, for example one word entered in the unigram category and included in the bigram category will certainly have different weights.

3. Stopwordremoval (filtering)

The process of removing words that are considered not important, such as there are conjunctions, prepositions, pronouns, or words that have nothing to do with the sentiment of the analysis will be deleted.

Table- V: Stopwordremoval Stage

Conversation	Stopwordremoval Results
can, analyzed, who, the culprit, who, insulted, brothers, us, papua, let, quickly, be tried, so, quickly, finished, problem, this, dismantle, just, account, that.	can, analysis, who, the culprit, insulted, brother, us, papua, fast, fair, quick, finished, problem, dismantle, account.

Table- V shows the process of omitting words that are not important that can affect the results of sentiment analysis

4. Stemming

The process of converting infix or suffix-filled words into a basic word becomes more specific.

Table- VI: Stemming Stage

Conversation	Stemming Results
can, analysis, who, the culprit, insulted, brother, us, papua, fast, fair, quick, finished, problem, dismantle, account.	can, analysis, who, the offender, insulting, brothers, us, Papua, so, quickly, fair, fast, finished, problem, dismantle, account.

Table- VI shows the process of changing the words that affect the root words to make the process of weighting the words easier.

STEP 3: IMPLEMENTASI HYBRID APPROACH

The proposed hybrid approach utilizes a combination of the Lexicon-based approach and the machine learning-based approach. This hybrid model combines two methods namely lexicon based and Naive Bayes classifier. The lexicon based method is used in the word weighting process, which is after going through a preprocessing process, words from the stemming results are compared with the lexicon dictionary that has been made. Furthermore, the conversation dataset is classified using the naive bayes classifier method, which results in a percentage of conversation data that is indicated to be hoaxed or not. The classification process produces three categories of sentiments, namely positive, negative and neutral. Conversation data that falls into the category of positive and neutral sentiments is indicated that the conversation data does not contain hoax elements. While the conversation data that falls into the negative sentiment

category, it is indicated that the conversation data contains hoax elements.

A. Lexicon Dictionary

This is examples of some words in each lexicon dictionary used as a comparison in weighting words.

Table- VII: Emoticon Dictionary

Emoticon	Feeling	Quality
:) :-)	Happy	3
:(:-(Sad	-3
:D :-D	Very Happy	4
D: D=	Very Sad	-4
* * * * *	Interested	2
D:< D:	Afraid	-2
xD XD	Smile/Laughing	2

Table- VII is an example of some emoticons that already have values. Dictionary of emoticons can be enlarged in number. The more contents of the emoticon dictionary, the results of sentiment analysis will be better.

Table- VIII: Dictionary of Disclaimer

Word	Quality
Not yet	-2
Not	-3
Without	-3
No	-4
Abstinence	-4
Do not	-4
Never	-4
Arrogant	-3

Table- VIII is an example of a Lexicon dictionary containing several words of denial along with the value of each word.

Table- IX: Dictionary of Question Word

Word	Quality
Who	2
Where	3
When	3
Where	2
How	3
What	2
Why	3
Why	2

Table- IX is an example of a Lexicon dictionary which contains several question words along with the value of each word.

Table- X: Dictionary of Positive Word

Word	Quality
Good	4
Great	3
Clever	3
Fast	2
Honest	4
Can	4
True	4
Smart	4

Table- X is an example of a Lexicon dictionary that contains several positive words along with the value of each word.

Table- XI: Dictionary of Negative Word

Word	Quality
Lie	-4
Corruption	-3

Cruel	-4
Ugly	-4
Danger	-4
Take a part	-4
Disaster	-4
Clash	-4

Table- XI is an example of a lexicon dictionary containing several negative words along with the value of each word.

B. Word Qualiting Based on the Lexicon Dictionary

Qualiting the results of stemming is done by comparing words with the lexicon dictionary, each word will be matched with the lexicon dictionary to give value to each word. The results of word weighting can be seen in Table XII.

Table- XII: Word Qualiting Results

N	Word	Score
n1	Can	4
n2	Analysis	2
n3	Who	2
n4	Players	-3
n5	Insulted	-4
n6	Brothers	1
n7	Us	1
n8	Papua	4
n9	Fast	2
n10	Fair	4
n11	Fast	2
n12	Finish	2
n13	Problem	-2
n14	Take apart	-4
n15	Account	2

Table- XII shows the word weighting process, which is matching words from conversation data with words from the Lexicon dictionary to determine the value of each word in the conversation data.

C. Determination of Sentiment Value

Sentiment value search is performed on each word that has weight so that in one conversation will be known the total number of positive values (Spositive) and also negative values (Snegative) of each constituent word.

Look for total positive value :

$$\begin{aligned}
 &= \quad \quad \quad (1) \\
 &= + + + + + + + + + \\
 &= 4 + 2 + 2 + 1 + 1 + 4 + 2 + 4 + 2 + 2 + 2 \\
 &= 26
 \end{aligned}$$

Look for total negative value :

$$\begin{aligned}
 &= \quad \quad \quad (2) \\
 &= + + + + \\
 &= 4 + 2 + 2 + 1 + 1 + 4 + 2 + 4 + 2 + 2 + 2 \\
 &= 26
 \end{aligned}$$

After knowing the total positive and negative values, the next step is to determine the orientation of the sentiment by comparing the number of positive, negative and neutral values.

D. Sentiment Analysis Results

From the series of processes above it can be concluded that the conversation has positive sentiments.

Table- XIII: Conversational Sentiment Analysis Results

No	1
User Account	+6283831286609
Conversation	Can be analyzed who did the offender who insulted our brother Papua, so that he could be tried quickly so that the problem can be resolved quickly just

	dismantle the account
WhatsApp Groups Followed	<ul style="list-style-type: none"> Indonesia Bersatu WARTA TV POLRI NEWS JOKOWISEKALI LAGI
Sentiment	Positive

Table- XIII shows the results of a data sentiment analysis from a user account with the results of a positive sentiment analysis.

E. Conversation Data Classification

The classification process utilizes a machine learning based approach using the Naive Bayes classifier algorithm. The total amount of conversation data used for the classification process is 700 conversation data. The data is divided into two namely 500 conversation data as training data and 200 conversation data as testing data. Each training data has been labeled with three categories of sentiments, namely positive sentiment, negative sentiment and neutral sentiment.

Table- XIV: Training Data Classification Results

Percentage of correctness	Percentage of error
73,3333%	26,6667%

Based on Table- XIV, the percentage of correctness of the classification is 73.3333% and the percentage of errors is 26.6667%. The percentage of truth is the amount of labeling the right sentiment in the training data. Table-XV is a performance evaluation value from the initial classification results showing a precision value of 0.756, a recall value of 0.733 and an f-measure of 0.739.

Table- XV: Evaluation of training data classification

Precision	Recall	F-Measure
0.756	0.733	0.739

STEP 4: ONTOLOGY

Ontology is used to represent data that has been processed from Hybrid Approach to find out whether the perpetrators of the hoaxes are also included in the other Whtasapp groups.

IV. EXPERIMENTAL RESULT AND DISCUSSION

Based on the value of the percentage of errors that are still large and the value of performance evaluation that has not yet reached a value of 1, found factors that influence the results of the classification is there is a vocabulary in the conversation that is not in the Lexicon dictionary.

Based on these factors, the training data is improved by reviewing it by adding the word sentiment to the word dictionary. After the improvement process, the results of the sentiment analysis classification are tested again.

The results of the training data classification show the truth value of 98% and the percentage of misclassification of 2%. The value of the performance evaluation indicates a precision value of 0.987, a recall value of 0.973, an f-measure value of 0.975.

Table- XVI: Classification Results of Testing Data

Classification	Classification		
	Positive	Neutral	Negative
Hybrid Approach	34,330%	52,330%	13,330%
Naive Bayes Classifier	30,330%	44,330%	26,330%

Table- XVI shows the comparison of classification results between Hybrid Approach and Naive Bayes Classifier.



Table- XVII: Performance Evaluation of Testing Data Classification

Hybrid Approach	classification			
	Accuracy	Precision	Recall	F-measure
Hybrid Approach	75%	0,771	0,754	0,773
Naive Bayes Classifier	Accuracy	Precision	Recall	F-measure
Naive Bayes Classifier	68%	0,708	0,702	0,704

From the results of sentiment analysis of 200 test data, a confusion matrix table can be formed as in Table XVIII:

Table- XVIII: Matrix Confusion

Hybrid Approach	Actual	Prediction		
		Positive	Negative	Neutral
Hybrid Approach	Positive	58	0	10
Hybrid Approach	Negative	0	18	10
Hybrid Approach	Neutral	0	20	84
Naive Bayes Classifier	Actual	Prediction		
		Positive	Negative	Neutral
Naive Bayes Classifier	Positive	50	0	12
Naive Bayes Classifier	Negative	0	28	20
Naive Bayes Classifier	Neutral	0	25	75

Based on the results of the two experiments above, it can be concluded that the use of Hybrid Approach shows better results than those using only the Naive Bayes Classifier method. The application of Hybrid Approach has several advantages including having a higher level of accuracy and being able to do emotion sentiment.

This is the results of the WhatsApp monitoring system design that has been implemented:

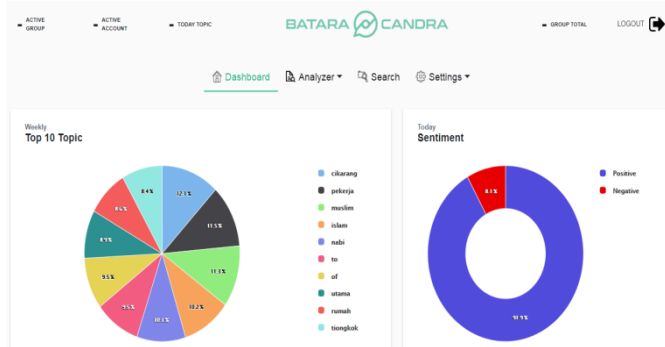


Fig. 2. The results of crawling WhatsApp group conversation topics.

Fig. 2. shows conversational topics that have been successfully crossed from all conversations in the WhatsApp group based on conversations that often appear and are used.

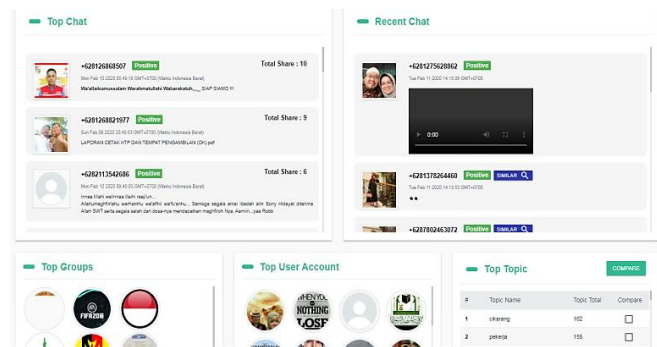


Fig. 3. The results of crawling conversations on WhatsApp groups.

Fig. 3. Shows the results of crawling conversation data consisting of Top Chat, Recent Chat, Top Groups, Top User Accounts and Top Topic.

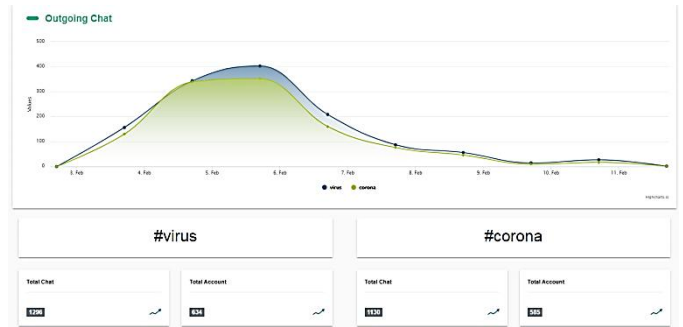


Fig. 4. Graph of conversation by topic.

Fig. 4. shows a graph of conversations by user accounts based on topic compare consisting of total chat and total account.

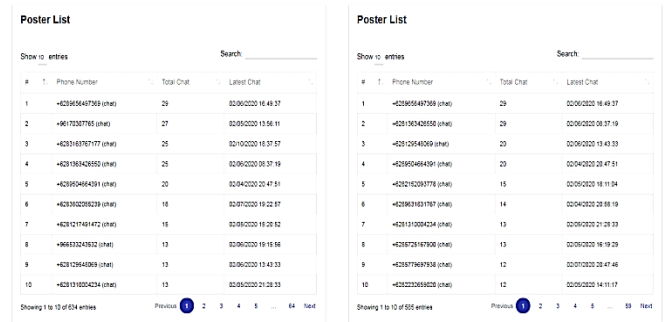


Fig. 5. User account data based on topic.

Fig. 5. Shows user account data that addresses a particular topic based on topic compare.

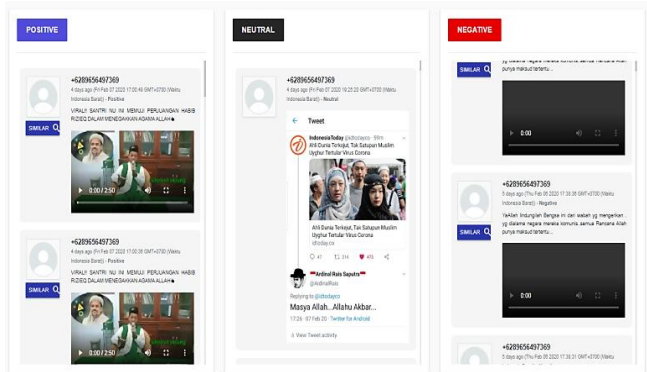


Fig. 6. The results of sentiment analysis on user accounts based on topic.

Fig. 6. shows the results of sentiment analysis on user accounts that discuss certain topics based on topic compare. The results of this analysis will be used to determine the existence of hoax indications in every conversation.



Fig. 7. Ontology representation on user accounts based on topic.

Fig. 7. Shows the ontology representation on the user account by topic, within the ontology there is a node that denotes the WhatsApp group and the WhatsApp account node that addresses the topic.

V. CONCLUSION AND FUTURE WORK

In this paper we propose the WhatsApp monitoring application to look for sources of hoaxes on the WhatsApp group. The experimental results show that the proposed hybrid model for detecting hoaxes is superior to hoax detection which uses only one method, the Naive Bayes Classifier.

This is also evidenced by the results of comparisons between two methods which show that Hybrid Approach has better results in terms of accuracy as well as others. The use of Hybrid Approach also has advantages such as being able to process emotion sentiment which will facilitate the detection of hoaxes from conversation data in the WhatsApp group. In further developments, the use of the Hybrid approach can be applied to other social media platforms to detect hoaxes.

REFERENCES

1. A. S. Cahyono, "Pengaruh Media Sosial Terhadap Perubahan Sosial Masyarakat di Indonesia," pp. 140–157, 2016.
2. Trisnani, "Pemanfaatan WhatsApp Sebagai Media Komunikasi Dan Kepuasan Dalam Penyampaian Pesan Dikalangan Tokoh Masyarakat," vol. 6, 2017.
3. L. Zhang, R. Ghosh, M. Dekhil, M. Hsu, and B. Liu, "Combining lexicon-based and learning-based methods for twitter sentiment analysis," HP Laboratories Technical Report, no. 89, 2011.
4. E. M. Okoro, B. A. Abara, A. O. Umagba, A. A. Ajonye, Z. S. Isa, and N. A. I. Nstitute, "A Hybrid Approach to Fake News Detection on Social Media," vol. 37, no. 2, pp. 454–462, 2018.
5. T. Hamdi, H. Slimi, I. Bounhas, and Y. Slimani, "Hybrid Approach Based on Graph Embedding and Users Features to Detect Source of Fake News in Twitter Social Network Using Machine Learning."
6. V. Nandi and S. Agrawal, "Sentiment Analysis using Hybrid Approach," Int. Res. J. Eng. Technol., pp. 1621–1627, 2016.
7. B. Kaur and N. Kumari, "A Hybrid Approach to Sentiment Analysis of Technical Article Reviews," Int. J. Educ. Manag. Eng., vol. 6, no. 6, pp. 1–11, 2016.
8. A. Altaher, "Hybrid approach for sentiment analysis of Arabic tweets based on deep learning model and features weighting," Int. J. Adv. Appl. Sci., vol. 4, no. 8, pp. 43–49, 2017.
9. J. Z. Pan, S. Pavlova, C. Li, N. Li, Y. Li, and J. Liu, "Content based fake news detection using knowledge graphs," Lect. Notes Comput. Sci. (including Subser. Lect. Notes Artif. Intell. Lect. Notes Bioinformatics), vol. 11136 LNCS, pp. 669–683, 2018.
10. Z. Zhou, H. Guan, M. Bhat, and J. Moorthy Hsu, "Fake News Detection via NLP is Vulnerable to Adversarial Attacks," 2015.
11. T. Granskogen and J. A. Gulla, "Automatic Detection of Fake News in Social Media using Contextual Information," 2018.
12. P. Kussa Laksana Utama, "Identifikasi Hoax pada Media Sosial dengan Pendekatan Machine Learning," Widya Duta J. Ilm. Ilmu Agama dan Ilmu Sos. Budaya, vol. 13, no. 1, pp. 69–76, 2018.
13. C. Boididou, S. Papadopoulos, M. Zampoglou, L. Apostolidis, O. Papadopoulou, and Y. Kompatsiaris, "Detection and visualization of misleading content on Twitter," Int. J. Multimed. Inf. Retr., vol. 7, no. 1, pp. 71–86, 2018.
14. M. Cs and A. Justitia, "Sistem Deteksi Hoax Dengan Menggunakan Algoritma Naive Bayes," pp. 94–95, 2018.
15. Sani M Isa, "Sentiment Analysis Approaches and Methods," BINUS, 2017. [Online]. Available: <https://mti.binus.ac.id/2017/10/04/1900/>. [Accessed: 29-Jul-2019].
16. S. Lorent and A. Itoo, Fake News Detection Using Machine Learning. 2019.
17. J. C. S. Reis, A. Correia, F. Murai, A. Veloso, and F. Benevenuto, "Explainable machine learning for fake news detection," WebSci 2019 - Proc. 11th ACM Conf. Web Sci., pp. 17–26, 2019.
18. A. Driif, Z. F. Hamida, and S. Giordano, "Fake News Detection Method Based on Text-Features," no. c, pp. 26–31, 2019.
19. P. Kaur, R. S. Boparai, and D. Singh, "A Review on Detecting Fake News through Text Classification," pp. 393–406.
20. K. Shu, A. Sliva, S. Wang, J. Tang, and H. Liu, "Fake News Detection on Social Media," ACM SIGKDD Explor. Newsl., vol. 19, no. 1, pp. 22–36, 2017.
21. R. J. Poovaraghan, M. V. K. Priya, P. V. S. S. Vamsi, M. Mewara, and S. Loganathan, "Fake News Accuracy using Naive Bayes Classifier," Int. J. Recent Technol. Eng., vol. 8, no. 1C2, pp. 962–964, 2019.

22. F. A. Ozbay and B. Alatas, "A Novel Approach for Detection of Fake News on Social Media Using Metaheuristic Optimization Algorithms," Elektron. IR ELEKTROTEHNIKA, vol. 25, no. 4, pp. 62–67, 2019.
23. H. Ahmed, I. Traore, and S. Saad, "Detection of Online Fake News Using N-Gram Analysis and Machine Learning Techniques," Springer Int. Publ., no. LNCS 10618, pp. 127–138, 2017.
24. O. Lloyd and C. Nilsson, "How to Build a Web Scraper for Social Media." 2019.
25. F. Sommar and M. Wielondek, "Combining Lexicon- and Learning-based Approaches for Improved Performance and Convenience in Sentiment Classification." 2015.
26. M. Taboada, J. Brooke, M. Tofiloski, K. Voll, and M. Stede, "Lexicon-Based Methods for Sentiment Analysis," Comput. Linguist., vol. 37, no. 2, pp. 267–307, 2011.
27. S. Baccianella, A. Esuli, and F. Sebastiani, "SENTIWORDNET 3.0 : An Enhanced Lexical Resource for Sentiment Analysis and Opinion Mining," vol. 0, pp. 2200–2204, 2008.
28. Z. Zhou, X. Zhang, and M. Sanderson, "Sentiment Analysis on Twitter through Topic-Based Lexicon Expansion," pp. 98–99, 2014.
29. O. Kolchyna, P. C. Treleaven, and T. Aste, "Twitter Sentiment Analysis: Lexicon Method, Machine Learning Method and Their Combination." .

AUTHORS PROFILE



Jordan Andrian is an informatics activist and is developing a startup in the field of big data processing and is pursuing an information technology education program at the S2 level at BINUS University.



Suharjo is the Head of Computer Science Department in Binus Online Learning Program of Bina Nusantara University. He received under graduated degree in mathematics from The Faculty of Mathematics and Natural Science in Gadjah Mada University, Yogyakarta, Indonesia in 1994. He received master degree in information technology engineering from Sepuluh November Institute of Technology, Surabaya, Indonesia in 2000. He received the PhD degree in system engineering from the Bogor Agricultural University (IPB), Bogor, Indonesia in 2011. His research interests are intelligent system, Fuzzy system, image processing and software engineering