

HYBRID ENSEMBLE MODEL FOR FAKE NEWS DETECTION

M.S.H. BASRI¹, N.H. ABD RAHIM²

^{1,2} Faculty of Ocean Engineering Technology and Informatics, Universiti Malaysia Terengganu,

21030 Terengganu Malaysia

E-mail: ¹ p42724@umt.edu.my, ² noorhafhizah@umt.edu.my

ABSTRACT

The increasing number of fake news spread to give a negative impact to many parties including society. Hence, to counter the high volume of fake news dissemination problems, numerous machine learning techniques have been used for automated fake news detection. Thus, this study implements the machine learning algorithms by developing a hybrid ensemble model for the classification of fake news. The primary objective of this work is to increase the performance of the machine learning model by developing the hybrid ensemble model to assist in the automatic detection of fake news in classifying the labelled news. This paper focuses on implementing the model with unorganized data in textual forms of Malay language news. A few performance metrics such as recall, accuracy, precision, and f-score will be used to measure the performance of the proposed model. The model has successfully increased the accuracy of fake news detection with 75% score.

Keywords: *Automated Detection System, Fake News, Hybrid Ensemble Model, Machine Learning, Malay Language,*

1. INTRODUCTION

People access information from both online and offline sources, including newspapers, periodicals, television, and their online counterparts, which include social media and other news sites. With the introduction of social media, unfiltered information may be disseminated quickly. Various types of content can be shared with friends or followers. Hence, information could be transmitted via social media easily. The use of social media to consume news is a double-edged sword. On the one hand, social media's low cost, ease of use, and rapid transmission of information encourage consumers to seek and consume news through it. There are numerous social media applications commercially available, including Facebook, WhatsApp, and YouTube. According to the numbers obtained, about 2 million accounts were closed monthly to prevent the spread of fake news [1].

Based on research by [2], the concept of fake news is comprised of two components: authenticity and intent. Authenticity requires that fake news contain inaccurate material that could be validated as such, which excludes conspiracy theories, which are

notoriously tough to verify as factual or misleading in the majority of situations. The second component, intent, indicates that the incorrect material was created with the intention of deceiving the reader. Fake news is a notable issue since it has a significant impact on people and society. The widespread transmission of fake news has the possibility to provide a profoundly damaging effect on both individuals and communities. The rise of the fake news problem was believed from the US presidential election in 2016. The news about the election outperformed mainstream news outlets such as New York Times, Washington Post, Huffington Post, NBC News [3]. Since then, the term 'fake news' has passed on now [4].

In Malaysia, fake news became a significant part in Malaysia's political discourse in early 2017. One of the government initiatives to combat fake news is by constituting a party Anti-Fake News Act (AFNA) that defines fake news as news, information, data, and reports which are wholly or partly false. Fine and jail sentences will be included if offenders have been found guilty of creating, publishing, or disseminating any fake news or publishing anything that contains fake news. In recent years, another

government effort associated with combating fake news was through a web portal named *Sebenarnya.my*. The website *sebenarnya.my* has compiled a record of every fake news, rumours, and conspiracy theory in Malaysia [5]. The news collected on this website are from many online sources such as WhatsApp, Facebook, and Twitter. The source of this news brought up an idea to study the machine learning-based model to authenticate the news more importantly in Malay languages.

In this paper, we present a hybrid ensemble method for the classification model. The objective of this paper is to implement an advanced machine learning model consists of an ensemble model for hybridization of the model. The model is hoped to increase the performance of the classification of labelled news. We demonstrate this approach on textual features of news in the Malay language. The term "Textual Features" refers to information derived from the text corpus, such as the text body, the headline, and the news source's text message [6]. Moreover, the news source is obtained from *sebenarnya.my* website that collects news all over Malaysia. We extract 1000 news that is labelled as fake and others. Generally, a machine learning-based text classification system comprises preprocessing, feature extraction, feature selection, and classification. There are numerous machine learning techniques available for text classification. Although Natural Language Processing (NLP) resource for preprocessing our Malay datasets is poor, the Indonesian NLTK library can be implemented in preprocessing method as the language between Malay and Indonesia was alike. Contemporary Indonesian and Malay are descended from the identical language and hence exhibit striking similarities [7]. We try to hybridizing the different ensemble models within a single model. Finally, the single model with hybrid ensemble model.

The results show our accuracy prediction model improves the single classification model. We achieved 75% accuracy score of predicted label fake and others. On the other hand, as we implement Bag-of-Words and Term Frequency, Inverse Document Frequency, or TF-IDF to extract the important features, the result produced between two features shows TF-IDF yield a better result. We can conclude based on our study that the hybridization of the machine learning model can increase the performance of the classification model..

This paper contains five sections. This first section discusses briefly fake news in general. The

next section is related works. In this section, we review all works regarding fake news detection using machine learning algorithms. The following is the methodology. This section describes our proposed model starting from dataset collection and the model development. Next is the result and discussion section. In this section, we analyse the performance of the proposed model and factors that contribute to the result acquired. Last but not least is the conclusion section. For this section, we explained the overall flow of the proposed model and future works to improve the model and increase the performance of the result of fake news detection.

2. RELATED WORKS

Fake News detection studies mainly involve the source of information used as data, features extraction, and proposed method by the researchers. The result of the proposed approach would be the analysis which to find any improvement or drawback. The language of the future dataset would affect the selection of feature extraction and the proposed method in the research. columns.

As for [8], the objective of the research is to counter fake news on COVID-19 issues. They utilised 13 machine learning algorithms for text classification. The source of the collected dataset was a Twitter post related to COVID-19 issues in Arabic and English for the time being. They analyse the performance of different classification algorithms with various feature extraction techniques used in this research. columns.

Besides, [9] aimed to study the credibility of news disseminated in social for news languages other than English. The researcher proposed a language-independent approach for automatically distinguishing credible from fake news based on a rich feature set. The feature sets mentioned in this research were linguistic, credibility, and semantic. The dataset used in this study was collected from four online sources in the Bulgarian language. The researcher implements logistic regression with L-BFGS from [10] optimizer and elastic net regularization to evaluate these features' set performance.

Next, [11] studies spreading fake news and the deficiency of labelled datasets in the Portuguese language. Support Vector Machine (SVM) with Linear SCV in Scikit-Learn was implemented for the text classification technique. The researcher collected

the labelled dataset from 4 available news websites. All results of classification produced from various feature extraction were analysed. Then, [12] study fake news in Indonesia. The researcher proposed automatic hoax news detection. In order to achieve the objectives, this study implements machine learning classification such as Maximum Entropy, Naïve Bayes, and SVM for the dataset collected. The dataset was produced from a web crawler developed by the researcher.

Meanwhile, [13] also focus on the COVID-19 issue but in the Thai language. The researchers propose Thai COVID-19 fake news detection among word relations using transfer learning models. The author describes "transfer learning" as two crucial steps in word2vec deep neural training. By implementing the proposed model on the dataset that was consisted of Thai words/phrases gathered from social media, the result of the proposed techniques was considered and made a few more tuning.

Next, [14] propose a corpus of Spanish news language gathered from several websites. The corpus aimed to analyse the crucial aspect of fake news detection. The researchers also conduct experiments for automated fake news detection. The machine learning approach was selected to measure the credibility of the collected corpus. Support Vector Machine (SVM) with linear kernel, Logistic Regression (LR), Random Forest (RF), and boosting were used for the experiments.

As for [15], the researcher has worked with fake news detection in Portuguese. This research aimed to explore the hurdle in establishing automated fake news detection in Portuguese. The researchers utilise the text classification technique using machine learning algorithms. The dataset was gathered from nonpartisan and unbiased party news. The news in the dataset was labelled. Hence, the labels were verified by international organizations. Using several features and classification models such as Naïve Bayes and Support Vector Machine (SVM), the text classifications' performance using the novel dataset was analysed and recorded.

In addition, [16] introduces the German fake news dataset. This research's focal points would be the German language automated fake news detection. Techniques selected for fake news detection were Convolutional Neural Network (CNN) and Support Vector Machine. The origin of the dataset was

collected from a famous mainstream news sources and trustworthy parties. columns.

Next, [17] presents a novel corpus in the Urdu Language. This research also makes use of machine learning models for automated fake news detection. The researcher stated that they were focusing on style-based, which analyze writing style from three categories of fake news detection: knowledge-based context-based according to [18]. The proposed corpus was gathered from a variety of trustworthy sources between January 2018 to December 2018. Multiple n-grams feature extraction and various classifiers such as Naïve Bayes for text classification.

3. METHODOLOGY

Based on previous research, we learned more about machine learning algorithms for fake news detection. Thus, we developed the hybrid based on a machine learning algorithm. The process starts with dataset management, feature extraction, material method, and performance. The processes are shown in Figure 1.

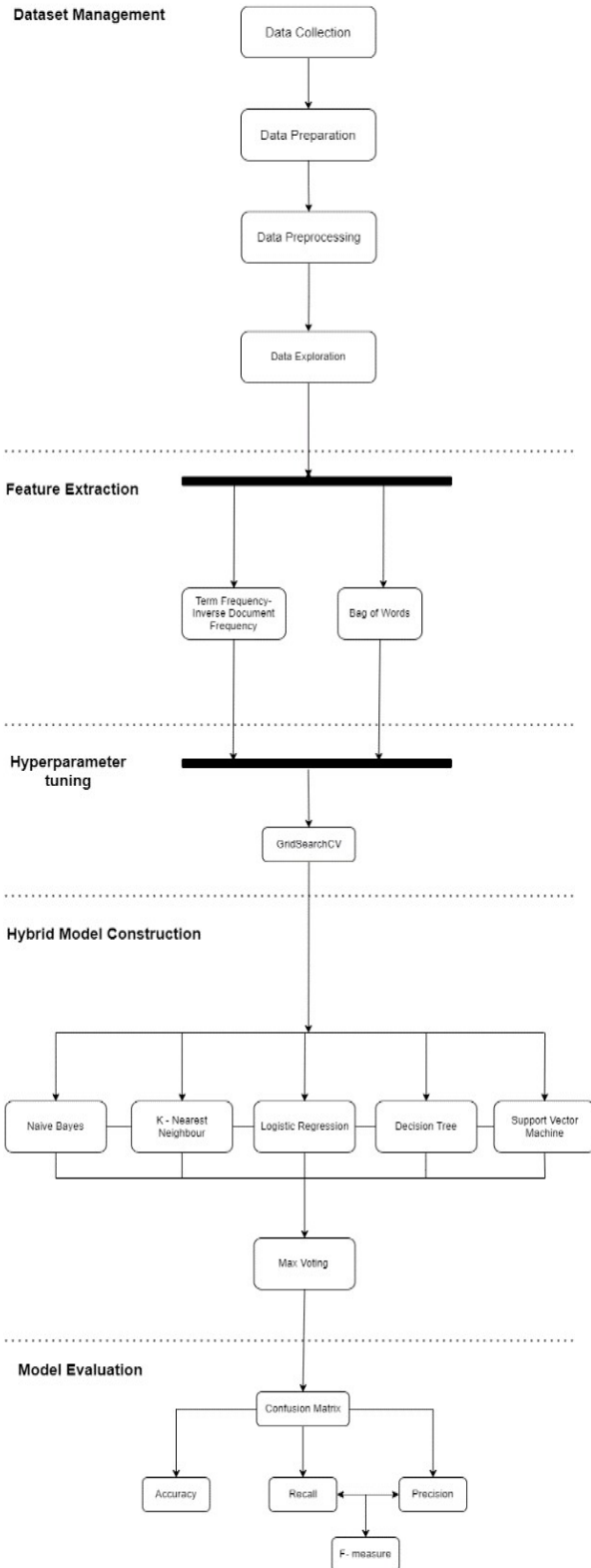


Figure 1: System Architecture of the Methodology

3.1 Dataset Preparation

Besides the hybrid model that we proposed, the novelty of this research could be found in our data collection. In America, numerous websites or blogs provide a bunch of news or information. All information is gathered from various trusted sources. However, some of the news probably comes from unreliable sources. In Malaysia, the same phenomenon happens. Thus, to ensure that society consumes only reliable information and filters out all the news transmitted, the local authorities launched a website called *Sebenarnya.my*. This website aims to clarify all types of news or buzz that go viral on social media and other online platforms. Most of the information was released online in an unstructured manner (for example, news, articles, videos, and audio) is rather challenging to discover and classify, as this requires solely human skill [19]. The statement is applicable on *Sebenarnya.my*. This site can still not classify news and rumours automatically via machine. Local authorities still rely on human fact-checks (e.g., local police authorities, constitutional law, legit annual reports, or verified press statements). Therefore, all the clarification and updates of the news are based on human verdicts.

Next, with the cooperation of the authorities, the responsible party who manages the website has exclusively provided all textual-format news from the year 2017-2020 with different types of domains. Our research focuses on textual format news in the Malay language without any other multimedia source includes. The dataset is readily compiled in the form of Excel files. The data consist of 1000 news with few features:

- Original website link source
- Date
- Textual information
- Label

The data is quite similar to the Politifact dataset used by [20] and the BuzzFeed dataset used by [21]. The news is broken down into four parts based on the label given: 'Penjelasan,' 'Waspada,' 'Palsu' and 'Penjelasan,' 'Sebenarnya'. We depend on ground-truth labelling because of the indefinite designate solution for detecting fake news, particularly in the Malay language. In addition, the local authorities such as the police department, responsible ministry, and government officials' statements are the ground truth source of the label of our dataset. The domain or topics discussed vary from each other. Based on Table 1, the distribution of label in the dataset is recorded.

Table 1: The Number of Occurrences in Dataset Label.

Label	Number of Occurrence
Palsu	585
Penjelasan	279
Waspada	133
Sebenarnya	3

Based on Table 1, there are only three (3) news recorded that is true which is labeled as 'Sebenarnya'. This is because the website mostly views the fake news that is disseminated through the internet. Most of the authorized news mostly appears in mainstream media such as television, newspaper, and radio.

3.2 Data Preprocessing

Data preprocessing is very crucial as it can remove any unnecessary information and terms for machine learning classification. The dataset's entire text was rendered to lower case letters, with punctuation removed. As data is taken from numerous sources, there is a tendency for emoticons and other punctuation to appear, which hinders our detection method. The dataset is searched for empty columns or rows. The date and website links are missing from some of the columns. The dataset was filtered to remove repeating letters, URL tags, dates, and stop words. Hashtags were also removed from the equation.

Following that, we go through Tokenization, Stemming, and Lemmatization. The tokenization would be useful for determining the frequency of words across the text. In preprocessing, next stemming is critical for locating the root of words. Word stems are stripped away from words. A word stem does not have to be the same root as a dictionary-based morphological root; it simply needs to be the same size as or smaller than the word. Finally, lemmatization is the process of searching the dictionary for the related word's form. From the preprocessing, we analyse the common words in our dataset after preprocessing step. The word frequency includes the words in both labels.

Table 2: Top 10 Words Frequency Appears In The Dataset.

words	Number of Occurrence
tular	741
sosial	522
dakwa	519
palsu	484
jelas	450
covid	316
kkm	312
rasmi	311
benar	311
maklum	297

3.3 Dataset Exploration

We explore the data to find out an in-depth understanding of the data. Based on (Sahoo et al., 2019), data should be assessed to obtain factual findings. By using exploratory data analysis results, decisive decisions can be made. We also followed the [22] method by utilising Python to explore the data. From our observation, we found that more than half of our data consists of fake news.

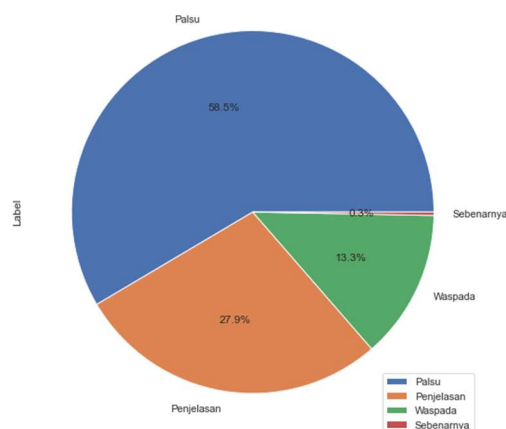


Figure 2: Dataset Label Distribution

Figure 2 shows the number of labels that appear according to the dataset's label. 'Penjelasan' is a piece of information that is true, but misunderstood and misrepresented by the public. 'Sebenarnya' label is the factual information that is once a rumour maybe because the information does not reach people. The above statistics show that the legitimate news reported on the website is very uncommon. Next, the 'waspada' label is the news that can be either true or false at the same time.

Next, the domain or issues of the news spread to the public were also analysed. This process was notable to understand the most controversial topic that could potentially be fake news. The graph below shows the domain extracted from the dataset and the ground-truth label provided by the responsible authorities

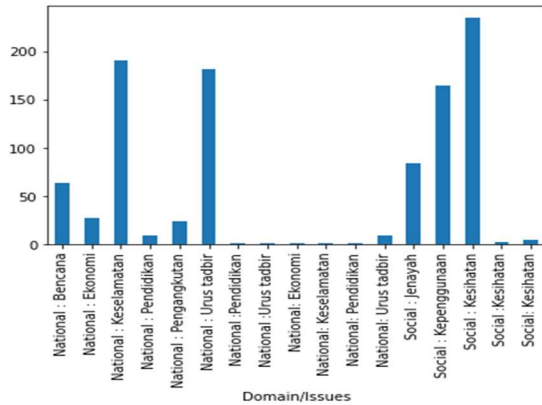


Figure 3: The Domain/ Issues In The Dataset

Figure 3 shows the issues or domains and its frequency that occur in the data. From our observation, Health issues are the most frequent topic that appears in our data. It may be because of the COVID-19 pandemic in recent years.

Next, we moved into the label of the dataset. Based on the collected dataset, we calculate and analyse the type of news based on the label. As the hybrid model focus on classification, the dataset is regrouped into fake and else. The 'fake' label comprises of half of all the collected datasets. The other half is regrouped into else

Label Distribution

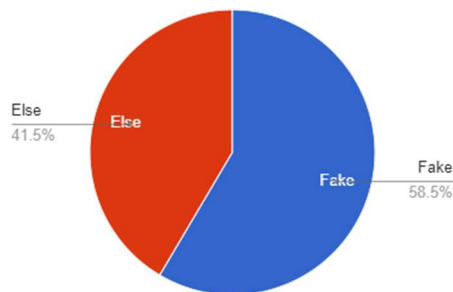


Figure 4: Number of Labels of Fake And Else After Regrouped

Based on Figure 4, from 1000 news obtained 585 news or 58.5% of them is fake and 415 or 41.5% of the dataset is else. In addition, the regrouping of else

was completed by merging the 'Waspada' and 'Penjelasan' and 'Sebenarnya' into a single class, and 'Palsu' into another class based on a study by [6]. Furthermore, the reason for uniting all three tags was to avoid imbalanced data. The imbalanced data occurs because of the representation of the dataset is not approximately even [23].

In addition, based on the obtained data, the majority of the label was 'palsu' or fake. Next, for the classification process, the value of fake could be assigned as 0 while everything else could be known as 1. Based on the dataset, we generate the word cloud to give a clear view on the frequency of words appearing in the dataset as can be shown in Figure 5.



Figure 5: Word Cloud for Fake News

Figure 5 shows the most frequent words that appear on fake labels.



Figure 6: Word Cloud for Else Label

Figure 6 shows the most frequent words that appear on else labels. There is a slight difference between the words that appear on both labels. The machine learning models could be used to classify the news.

3.4 Feature Extraction

The feature extraction procedure is a must to ensure the text in the dataset is parsed to evaluate words. We encrypt all words in integers values or floating-point before passing through into a machine learning algorithm for fake news identification. The "content" column has a considerable number of words for the vectorisation process. We extract textual features from the dataset by employing text representation models such as Count Vectorizer, Bag-

of-Words (BoW), and Term Frequency-Inverse Document Frequency (TF-IDF). TF-IDF and BoW text analysers are essential to make use of machine learning models. A vector enhances to get a better result in the classification process.

This technique measures the the training of the text document classification system

3.4.1 TF-IDF

This technique measures the relevancy of a word to a document in a collection of documents. The TF-IDF algorithm calculates values for each word in a document. The process begins with dividing the term's frequency in that document by the percentage of documents in which the word appears. This technique utilises multiplying two metrics: the number of word occurrences in a record and the term's inverse document frequency over a collection of documents. With the formula below, we can obtain TF and IDF values with:

$$1) TF_{ij} = \frac{n_{i,j}}{\sum_k n_{k,j}}$$

$$2) IDF = \log\left(\frac{N}{n}\right)$$

$$3) TF - IDF \text{ value of a terms} = (TF \cdot IDF)$$

3.4.2 Bag-Of-Words

Count vectoriser is the simplest way to represent text data numerically. This technique requires placing the words in data with a unique number. In addition, the count vectoriser ignores any specific information and focuses on the presence of terms in a document. The final look of this method is in the form of an encoded vector.

3.5 Material Method

Due to the multifaceted structure of fake news, identifying it is not always straightforward. It is self-evident that a practical technique must include various perspectives to deal with the issue effectively. We proposed hybrid ensemble learning by training the machine learning algorithm for fake news classification. First and foremost, we try to implement ensemble learning because it is a technique that deals with classification and regression problems. This learning builds from the combination of the weak learner-type machine learning algorithm.

The term 'hybrid' is used in this model because we utilize the heterogenous-type algorithm for weak learners. In addition, the hybrid term is also valid because our model works by combining classification and classification models based on [24] explanation.

In this approach, we split our data into train and test. We choose the 80:20 ratio to train and test our data, respectively. First of all, we undergo classification techniques on single model classifiers. We pick the most popular homogeneous machine learning models for text classification. Next is to evaluate the performance of different machine learning models on the dataset. We selected the algorithm based on the learning technique consisting of Logistic Regression, k-Nearest Neighbors, Decision Tree., Support Vector Machine, and Naive Bayes.

Naïve Bayes, SVM, Decision Tree, and KNN are classifiers used in text mining issues. The classifiers are particularly popular, with Naive Bayes usually providing an accuracy baseline and SVMs outperforming in such tasks [25]. We observed several studies that utilised Logistic Regression for text mining cases. Hence, we try to train the model in our dataset. Next, for the second experiment of single model classifiers, we perform the default cross-validation with a value k is 10 and set the Seed value of 7 for the random state for all five models and observe the result.

Next, we undergo the second experiment by defining five machine learning models five times for each model. Then, set the estimator and implement Max Voting Classifiers to predict the final class of the ensemble. Finally. The voting machine learning model aims to improve model performance by combining the prediction from other models. Ideally, we want to achieve better performance than single model classification using ensemble. The same goes with the hybrid ensemble model, we set the same K cross-validation and put the same seed value in the second trial. columns.

After fitting all requirements to train all models, we analyse the best-performing models among the classifiers. Next, we consider the precision, accuracy, recall, F-Measure, and metrics for measuring the performance. Next, we analyse the classification of labelled data via a confusion matrix. For the hybrid ensemble model, we measure and compare the accuracy for all different settings we made.

4 RESULT AND DISCUSSION

We evaluate the performance of our models from the different settings. We found that the performance slightly increases after K-fold cross-validation. We assess the hybrid ensemble learning model only after

performing the cross-validation. We set 10 for k values to monitor the performance of the model. The K-Folds technique is widely used and intuitive; it typically results in a less unfair model than other methods. Since it assures that each direct observation in the training and test sets has a probability of occurring in the original dataset, this technique is one of the finest approaches if our input data is restricted.

Table 3: Performance Result Using TF-IDF features

Model	Precision	Recall	F-Score	Accuracy
SVM	0.72	0.88	0.79	0.73
K-Nearest Neighbour	0.69	0.85	0.76	0.71
Logistic Regression	0.71	0.87	0.78	0.71
Naïve Bayes	0.71	0.80	0.69	0.66
Decision Tree	0.62	0.70	0.66	0.59
Hybrid Ensemble	0.68	0.67	0.63	0.75

Table 4: Performance Result Using Bag-Of-Words features

Model	Precision	Recall	F-Score	Accuracy
SVM	0.6	0.52	0.59	0.71
K-Nearest Neighbour	0.55	0.48	0.51	0.64
Logistic Regression	0.73	0.54	0.62	0.70
Naïve Bayes	0.54	0.16	0.25	0.63
Decision Tree	0.59	0.49	0.53	0.63
Hybrid Ensemble	0.68	0.67	0.63	0.70

Based on Table 3 and 4, we can conclude the result of the TF-IDF features is much better compared to Bag-Of-Words. This shows that the hybrid model has improved the detection method to classify fake news based on the feature extraction selected. The special attribute of TF-IDF, which can differentiate the level of relevance of essential terms based on the document makes it better than Bag-Of-Words.

5 CONCLUSION

As a conclusion, we propose a hybrid ensemble model of machine learning classification for classifying labelled news. This model is an advanced method of utilising the standard machine learning algorithm from previous fake news detection. We use accuracy, precision, recall, and F1-score to compare

and measure the result of the single model and hybrid model. Based on the results, the hybrid model yields better results compared to the single classification model. The most possible reason for the result is that the hybrid model can achieve a better method with high flexibility. Moreover, this hybrid model which combines two or more models can have a high capability because it can have more units that use for prediction and optimization compared with single methods. Hybrid methods potentially can become a major preference for the classification of fake news as it has high potential and capability.

On the other hand, there are a few limitations of this research that can be addressed. Firstly, we focus only on a machine learning algorithm. Next, the number of datasets in this research may be a little low because the dataset was collected only in the Malay language. We are unable to use this type of feature because the Malay language does not have the standard pre-trained corpus, especially in 2021, where COVID-19 pandemics affect the final result of the model.

In the near future, there are many improvements and a diversity of research methods to produce a high-quality result in detecting fake news. Moreover, hybrid models based on deep learning can improve the overall performance of the classification model. One of the main ideas for future improvement in our work is to utilize deep learning algorithms for hybridizing the classification model. Previous research already proved the correlation between deep learning. Deep learning algorithms such as Recurrent Neural Network, Artificial Neural networks can be used to process high amount of data for classification tasks. Furthermore, few improvements can be made in extracting important features in our future works. Numerous technique such as word embedding with fastText or spacy can be used. Numerous technique such as word embedding with fastText or spacy can be used.

ACKNOWLEDGEMENT

This research was supported by the fundamental research grant scheme for research acculturation of early career researchers (FRGS-RACER) with the reference code of RACER/1/2019/ICT02/UMT//2 and vote number 59547 under the Malaysia Ministry of Higher Education (RACER 2019-1).

REFERENCES:

- [1] J. Vasandani, "I Built a Fake News Detector Using Natural Language Processing and Classification Models," *Towards Data Science*, 2019. <https://towardsdatascience.com/i-built-a-fake-news-detector-using-natural-language-processing-and-classification-models-da180338860e> (accessed Sep. 21, 2021).
- [2] S. Lorent and A. Itoo, "Fake News Detection Using Machine Learning," *Master Thesis*, p. 91, 2019.
- [3] S. Craig, "This Analysis Shows How Viral Fake Election News Stories Outperformed Real News On Facebook," *BuzzFeed News*, Sep. 16, 2016. <https://www.buzzfeednews.com/article/craigsilverman/viral-fake-election-news-outperformed-real-news-on-facebook> (accessed Apr. 07, 2021).
- [4] H. Allcott and M. Gentzkow, "Social media and fake news in the 2016 election," *J. Econ. Perspect.*, vol. 31, no. 2, pp. 211–236, May 2017, doi: 10.1257/jep.31.2.211.
- [5] H. Wahed, "Misinformation and Disinformation During Covid-19: the Effects and the Relevant Laws in Malaysia," *Int. J. Law, Gov. Commun.*, vol. 5, no. 21, pp. 202–209, 2020, doi: 10.35631/ijlgc.5210015.
- [6] J. C. S. Reis, A. Correia, F. Murai, A. Veloso, F. Benevenuto, and E. Cambria, "Supervised Learning for Fake News Detection," *IEEE Intell. Syst.*, vol. 34, no. 2, pp. 76–81, 2019, doi: 10.1109/MIS.2019.2899143.
- [7] N. Lin, S. Fu, S. Jiang, G. Zhu, and Y. Hou, "Exploring Lexical Differences Between Indonesian and Malay," *Proc. 2018 Int. Conf. Asian Lang. Process. IALP 2018*, pp. 178–183, 2019, doi: 10.1109/IALP.2018.8629131.
- [8] M. K. Elhadad, K. F. Li, and F. Gebali, *COVID-19-FAKES: A Twitter (Arabic/English) dataset for detecting misleading information on COVID-19*, vol. 1263 AISC. Springer International Publishing, 2021.
- [9] M. Hardalov, I. Koychev, and P. Nakov, "In Search of Credible News," in *Artificial Intelligence: Methodology, Systems, and Applications*, 2016, pp. 172–180.
- [10] D. C. Liu and J. Nocedal, "On the limited memory BFGS method for large scale optimization," *Math. Program.*, vol. 45, no. 1, pp. 503–528, Aug. 1989, doi: 10.1007/BF01589116.
- [11] R. A. Monteiro, R. L. S. Santos, T. A. S. Pardo, T. A. de Almeida, E. E. S. Ruiz, and O. A. Vale, "Contributions to the Study of Fake News in Portuguese: New Corpus and Automatic Detection Results," *Lect. Notes Comput. Sci. (including Subser. Lect. Notes Artif. Intell. Lect. Notes Bioinformatics)*, vol. 11122 LNAI, pp. 324–334, 2018, doi: 10.1007/978-3-319-99722-3_33.
- [12] I. Y. Risca Pratiwi, A. A. Rosa, and F. Rahutomo, "Study of Hoax News Detection using Naive Bayes Classifiers," *2017 11th Int. Conf. Inf. Commun. Technol. Syst.*, pp. 73–78, 2017, doi: 10.1109/ICTS.2017.8265649.
- [13] P. Mookdarsanit and L. Mookdarsanit, "The COVID-19 fake news detection in Thai social texts," *Bull. Electr. Eng. Informatics*, vol. 10, no. 2, pp. 988–998, 2021, doi: 10.11591/eei.v10i2.2745.
- [14] J. P. Posadas-Durán *et al.*, "Detection of fake news in a new corpus for the Spanish language," *J. Intell. Fuzzy Syst.*, vol. 36, no. 5, pp. 4868–4876, May 2019, doi: 10.3233/JIFS-179034.
- [15] C. Ruschel, M. Dias, and P. Alegre, "Towards fake news detection in Portuguese: New dataset and a claim-based approach for automated detection," 2019.
- [16] I. Vogel and P. Jiang, "Fake News Detection with the New German Dataset 'GermanFakeNC,'" in *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 2019, vol. 11799 LNCS, pp. 288–295, doi: 10.1007/978-3-030-30760-8_25.
- [17] M. Amjad, G. Sidorov, A. Zhila, H. Gómez-Adorno, I. Voronkov, and A. Gelbukh, "'Bend the truth': Benchmark dataset for fake news detection in Urdu language and its evaluation," *J. Intell. Fuzzy Syst.*, vol. 39, no. 2, pp. 2457–2469, 2020, doi: 10.3233/JIFS-179905.
- [18] M. Potthast, J. Kiesel, K. Reinartz, J. Bevendorff, and B. Stein, "A stylometric inquiry into hyperpartisan and fake news," in *ACL 2018 - 56th Annual Meeting of the Association for Computational Linguistics, Proceedings of the Conference (Long Papers)*, Feb. 2018, vol. 1, pp. 231–240, doi: 10.18653/v1/p18-1022.
- [19] I. Ahmad, M. Yousaf, S. Yousaf, M. Ovais Ahmad, and M. O. Ahmad, "Fake News Detection Using Machine Learning Ensemble Methods," *Complexity*, vol. 2020, 2020, doi: 10.1155/2020/8885861.

- [20] K. Shu, D. Mahudeswaran, S. Wang, D. Lee, and H. Liu, "FakeNewsNet: A Data Repository with News Content, Social Context, and Spatiotemporal Information for Studying Fake News on Social Media," *Big Data*, vol. 8, no. 3, pp. 171–188, 2020, doi: 10.1089/big.2020.0062.
- [21] C. Buntain and J. Golbeck, "Automatically Identifying Fake News in Popular Twitter Threads," *Proc. - 2nd IEEE Int. Conf. Smart Cloud, SmartCloud 2017*, pp. 208–215, 2017, doi: 10.1109/SmartCloud.2017.40.
- [22] K. Sahoo, A. K. Samal, J. Pramanik, and S. K. Pani, "Exploratory data analysis using python," *Int. J. Innov. Technol. Explor. Eng.*, vol. 8, no. 12, pp. 4727–4735, 2019, doi: 10.35940/ijitee.L3591.1081219.
- [23] N. V. Chawla, "Data Mining for Imbalanced Datasets: An Overview," in *Data Mining and Knowledge Discovery Handbook*, O. Maimon and L. Rokach, Eds. Boston, MA: Springer US, 2005, pp. 853–867.
- [24] C. F. Tsai and M. L. Chen, "Credit rating by hybrid machine learning techniques," *Appl. Soft Comput. J.*, vol. 10, no. 2, pp. 374–380, Mar. 2010, doi: 10.1016/j.asoc.2009.08.003.
- [25] G. Gravanis, A. Vakali, K. Diamantaras, and P. Karadais, "Behind the cues: A benchmarking study for fake news detection," *Expert Syst. Appl.*, vol. 128, pp. 201–213, 2019, doi: 10.1016/j.eswa.2019.03.036.