

Sentiment Analysis Of Public Opinion On Handling Stunting In Indonesia Using Random Forest

Ariska Fitriyana Ningrum⁽¹⁾, Ihsan Fathoni Amri⁽²⁾

^{1,2}Program Studi Sains Data, Universitas Muhammadiyah Semarang

Jln. Kedungmundu No 18, Kedungmundu, Kec. Tembalang, Kota Semarang

e-mail: ariskafitriyana@unimus.ac.id⁽¹⁾, ihsanfathoni@unimus.ac.id⁽²⁾

ABSTRAK

Masalah stunting penting untuk diselesaikan, karena berpotensi mengganggu potensi sumber daya manusia dan berkaitan dengan tingkat kesehatan, bahkan kematian anak. Pemerintah Indonesia menargetkan angka stunting turun menjadi 14 persen pada tahun 2024 melalui program percepatan penurunan stunting sebagai upaya meningkatkan status gizi masyarakat dan juga menurunkan prevalensi stunting atau balita pendek. Memahami sentimen publik terhadap inisiatif stunting sangat penting bagi para pembuat kebijakan dan pemangku kepentingan untuk merancang intervensi yang efektif dan mengalokasikan sumber daya secara efisien. Pada penelitian ini dilakukan klasifikasi pada sentiment positif dan negatif menggunakan algoritma random forest. Data yang digunakan adalah data komentar pada salah satu laman media sosial yaitu twitter mengenai sentiment masyarakat terhadap penanganan kasus stunting di Indonesia. Tahapan pertama pada penelitian ini setelah didapatkan sebuah data yaitu dilakukan preprocessing data. Tahapan preprocessing data dalam analisis sentimen berguna untuk membersihkan dan menormalkan teks, menghilangkan kata-kata tidak relevan, serta mempersiapkan data agar algoritma dapat menganalisis sentimen dengan lebih akurat dan efisien. Selanjutnya hasil data yang sudah di preprocessing diberikan label 0 untuk positif dan 1 untuk label negatif. Klasifikasi terhadap sentiment positif dan negatif ini dilakukan menggunakan random forest dan menghasilkan nilai akurasi sebesar 97,5%. Model ini sudah baik, namun kami menyarankan untuk mencoba algoritma lain dalam penelitian selanjutnya.

Kata kunci: Analisis Sentiment, Random Forest, Stunting

ABSTRACT

The problem of stunting is important to solve, as it has the potential to disrupt human resource potential and is linked to health outcomes and even child mortality. The Indonesian government targets the stunting rate to drop to 14 percent by 2024 through an accelerated stunting reduction program as an effort to improve the nutritional status of the community and also reduce the prevalence of stunting or short toddlers. Understanding public sentiment towards stunting initiatives is essential for policy makers and stakeholders to design effective interventions and allocate resources efficiently. In this research, classification of positive and negative sentiment is carried out using the random forest algorithm. The data used is comment data on one of the social media pages, namely Twitter, regarding public sentiment towards handling stunting cases in Indonesia. The first stage in this research after obtaining a data is data preprocessing. The data preprocessing stage in sentiment analysis is useful for cleaning and normalizing text, removing irrelevant words, and preparing data so that algorithms can analyze sentiment more accurately and efficiently. Furthermore, the results of the preprocessed data are labeled 0 for positive and 1 for negative labels. The classification of positive and negative sentiment was done using random forest and resulted in an accuracy value of 97.5%. This model is good, but we suggest trying other algorithms in future research.

Keywords: Sentiment Analyst, Random Forest, Stunting

INTRODUCTION

Stunting is a condition in children who experience growth disorders, so that the height and weight of children are not normal due to problems of nutritional deficiencies for a long time [1]. The problem of stunting in Indonesia is still quite large in the health sector today. According to the World Health Organization (WHO), as many as 22% or around 149.2 million children in the world under the age of five were recorded as stunted in 2020. Indonesia's position on the prevalence of stunting in the world is ranked 115 out of 151 countries. Meanwhile, in Southeast Asia, Indonesia is ranked second at 31.8% after Timor Leste at 48.8%. The third is Laos at 30.2%, the fourth is Cambodia at 29.9%, and the fifth is the Philippines at 29.9% [2] Stunting is caused by health problems, environmental factors and health services received by children. Genetic factors do not significantly affect stunting. Lack of nutrition in the fetus is the biggest cause of stunting in children. The first 1000 days of a child's life (1000 HPK) is the starting point for making important conclusions on long-term growth Thus, ineffective parenting and diet can increase the chance of stunting. Mental disorders and hypertension in mothers also affect the behavior and practices of nutrition in children. Limited access to health and sanitation services exacerbates the stunting conditions that occur in Indonesia such as lack of clean water, unclean latrines, and so on [3] Stunting in Indonesia is a deep-rooted problem. The problem of stunting is important to solve, because it has the potential to disrupt human resource potential and is related to health levels, even child mortality. In early 2021, the Indonesian government targeted the stunting rate to drop to 14 percent by 2024 through the accelerated stunting reduction program as an effort to improve the nutritional status of the community and also to reduce the prevalence of stunting or short toddlers [4].

Stunting in Indonesia is a deep-rooted problem. The problem of stunting is important to solve, because it has the potential to disrupt human resource potential and is related to health levels, even child mortality. In early 2021, the Indonesian government targeted the stunting rate to drop to 14 percent by 2024 through the accelerated stunting reduction program as an effort to improve the nutritional status of the community and also to reduce the prevalence of stunting or short toddlers[5] Understanding public sentiment towards stunting initiatives is crucial for policymakers and stakeholders to design effective interventions and allocate resources efficiently. Several previous studies have analyzed stunting predictions using the random forest algorithm which resulted in a classification accuracy value of 90.7%. [6]. Another study on social media analysis with the topic of stunting in Indonesia was conducted where the results showed that negative sentiment dominated by 60.6%, positive sentiment by 31.5%, and neutral by 7.9% [7]. In addition, this study shows that 'children', 'decline', 'numbers', 'prevention', and 'nutrition' are words that often appear in stunting [7]. Another study comparing SVM and random forest algorithms for the classification of stunting disease. the results show that the random forest algorithm provides higher accuracy of 88.2% compared to SVM of 65.6% [8].

The background of this study is based on the need to analyze public sentiment regarding the handling of stunting in Indonesia, which is a significant public health problem. The data used are the results of positive and negative reviews from the public on social media such as twitter regarding the handling of stunting cases in Indonesia. Some previous studies have also analyzed using comment data on twitter such as research on Sentiment Analysis of Twitter Netizens on the News of VAT on Basic Food and Education Services with Social Network Analysis and Naive

Bayes Classifier Approaches with data obtained as many as 4090 tweets [9]. While other studies have also conducted sentiment analysis of twitter users regarding online transportation service users [10].

The method used in this research is Random Forest. Random Forest was chosen because of its superior ability to handle complex and varied data, and provide accurate results in classification and prediction. This method is suitable for sentiment analysis because it is able to overcome overfitting, works well on large and irregular datasets, and provides a good interpretation of the features that affect the results. Random Forest is one of the state-of-the-art methods in machine learning that consists of a number of independently trained decision trees and the results are combined to produce more accurate and stable predictions. In the context of sentiment analysis, Random Forest can handle variations in language expression and identify relevant patterns from unstructured text data. In addition, Random Forest's ability to handle data with many features is very useful in sentiment analysis involving various emotional aspects and public opinions related to stunting treatment [11].

The state of the art of the Random Forest method shows that this technique has been successfully applied in various domains, including text and sentiment analysis, with satisfactory results. Previous studies have shown that Random Forest often outperforms other methods such as logistic regression and Support Vector Machines (SVM) in terms of accuracy and robustness to noise in the data. This makes Random Forest an appropriate choice for this research in an effort to understand and measure public sentiment towards stunting response efforts in Indonesia.

METHOD

In this study, a sentiment analysis about stunting in Indonesia was conducted. The remaining data will be analyzed using the random forest algorithm to check the accuracy of the sentiment results. The following are the stages of sentiment analysis research on stunting using the random forest algorithm

1. Data Collection

The data used in this study is the result of crawling Twitter data related to positive and negative responses to stunting conditions in Indonesia. The data collection process was carried out over a one-month period, starting from January 1, 2024 to January 31, 2024, to get a current picture of public sentiment on the issue of stunting. In the crawling process, certain filters and keywords related to infant stunting were used. The keywords used included "stunting", "child growth", and "child nutrition", as well as other keyword variations related to stunting and child health in Indonesia. In applying a series of filters, including language settings (Bahasa Indonesia) and geographical location (Indonesia) were used to ensure that the data captured was specifically related to responses to stunting conditions in Indonesia. This resulted in a total of 4601 comments showing both positive and negative opinions. Labeling the data into positive and negative categories was done through a manual process by the researcher, where each opinion was classified based on the sentiment expressed towards stunting conditions. This assessment is based on the context within each opinion, with the aim of gaining an accurate understanding of public sentiment.

2. Preprocessing Data

The first stage of the system is preprocessing. This stage involves several processes including Case Folding, Tokenization, Normalization, and Stemming. Case Folding is a task of splitting review text into smaller units called tokens or terms [12]. For infant stunting cases, what is done before and after case folding is, for example, "Breastfeeding mothers must have good nutrition" becomes "breastfeeding mothers must have good nutrition". Next is Tokenizing, in this process the separation is carried out on each word that makes up a document. In general, each word is identified or separated from other words by space characters, so the tokenizing process relies on space characters in the document to perform word separation [13] In sentiment analysis of stunting cases, what is done is to present the number of tokens generated from a review or comment. For example, from the sentence "breastfeeding mothers should have good nutrition", the tokens generated are "mother", "breastfeeding", "should", "nutrition", "which", "good". Normalization (Stopword Removal) process Removes special characters, numbers, and stopwords (common words) from each token. In the case of sentiment analysis, it shows a list of stopwords used and examples of text before and after stopword removal. For example, from "mother", "breastfeeding", "should", "nutrition", "which", "good", after removing the stopwords "should", "which", then "mother", "breastfeeding", "nutrition", "good" remains. This research also uses Stemming techniques which aim to find the base word, by removing all affixes that are fused to the word.[14] In Indonesian, this usually involves the removal of prefixes, suffixes or infixes. As an example of words before and after the stemming process, for example, "menyusui" can be reduced to "susu".

3. Sentiment Analyst Using Random Forest

The last stage is sentiment classification. Each review will be classified into positive or negative category. In this study, we employ random forest for the classification task. Random forest algorithm is a supervised classification algorithm. It is an ensemble learning technique based on decision tree algorithm [15]. Random Forest Algorithm is the advancement of Classification and Regression Tree (CART) method with the implementation of bootstrap aggregating (bagging) and random feature selection. Procedure of random forest algorithm on the data of n observations and p predictor [16]

- a. Random samples of size n are drawn with the possibility of obtaining the same data (with replacement). This phase is called bootstrap.
- b. Using the bootstrap samples, the tree is grown until the maximum size is reached, which is done without pruning. At each node, the random feature selection is used to determine the split, which m number of variables randomly sampled as candidates at each split must be $m \ll p$, at which point, the best node will be chosen based on m number of variables available for splitting [17]
- c. Repeat stage 1 and 2 for k times to generate a forest that consists of k trees. Breiman and Cutler suggests to observe the error OOB when

$$m = \left(\frac{1}{2} \sqrt{p}, \sqrt{p}, 2\sqrt{p} \right) \quad (1)$$

where p is the total variable and the number of k is small, then m with the smallest error OOB will be chosen.

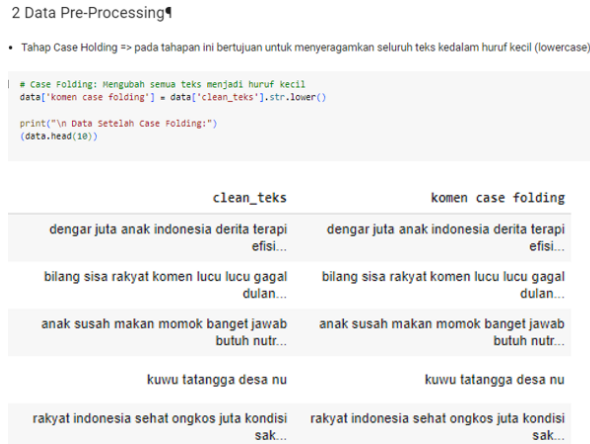


Figure 2. Case Holding Stages Result

Next, the tokenizing process is carried out. At this stage, the sentence in the comment will be broken down into words. The results of the tokenizing process are presented as in the following figure.

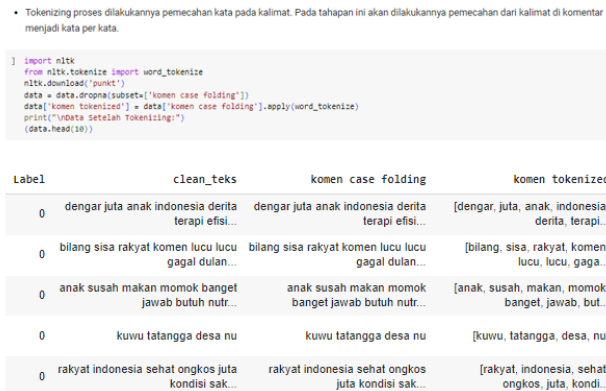


Figure 3. Tokenizing Stages Result

The next preprocessing is to perform normalization to change the values of a dataset so that they have a uniform scale. The main purpose is to ensure that variables with different value ranges have equal influence when used in the analysis. The results of data normalization are as follows.



Figure 4. Normalization Stages Result

The last step in data preprocessing is stemming. This aims to find the base word, by removing all affixes that are fused to the word. The results of stemming performed for sentiment analysis are as follows.

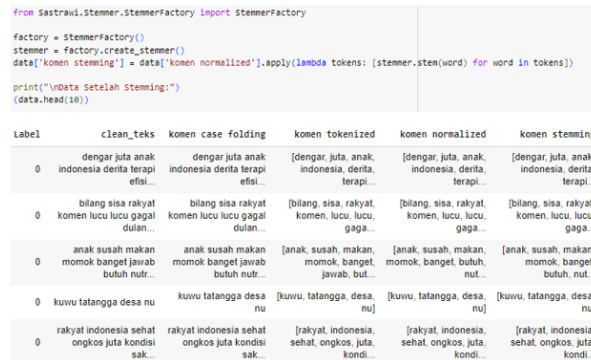


Figure 5 Stemming Performed for sentiment analysis

Furthermore, the results of data preprocessing that have been carried out can be seen visually regarding positive and negative opinions. Visualization aims to display the words that appear most or most often in a sentiment. Wordcloud this time describes each sentiment, the more often a word is used when giving a review, the larger the size of the word displayed on the wordcloud visualization. The following figure shows the visualization results for positive and negative sentiments



Figure 6. Visualization Positive Sentiment

Based on the figure above, it can be seen that in the positive sentiment there are several words that stand out such as, "balanced nutrition," "helps reduce", "Welfare Growth" and several other words which indicate that the public's response to handling stunting in Indonesia has helped reduce stunting rates, provide balanced nutrition to children and can foster community welfare.

The performance generated by the random forest algorithm provides considerable accuracy, which is 97.50%, indicating that this model can classify data has a very good indication, and produces precision on Label 0 (negative comments) of 97% and recall of 100%, the results obtained are very high, and F1 score of 99%, indicating a high balance of precision and Recall. Meanwhile for precision on Label 1 (Positive comments) of 100%, and recall of 18% and the result for f1-score is 30%.

CONCLUSION

This study shows that public sentiment towards the handling of stunting cases in Indonesia can be divided into positive and negative based on the analysis of 4601 comments from Twitter social media. The results show that positive responses include the view that the handling of stunting has succeeded in reducing stunting rates, providing balanced nutrition to children, and potentially improving the general welfare of society. On the other hand, negative responses include dissatisfaction with the effectiveness of stunting handling, the existence of unresolved social inequalities, and the lack of effort from the government in handling stunting cases in the community.

This research continued with the classification of comment data based on sentiment using TF-IDF feature extraction. This method is important because it converts text into a numerical vector representation, where the TF-IDF weight of each word gives an idea of the importance of the word in determining positive or negative sentiment. Through this classification, it is possible to identify and categorize the sentiments present in the text data, enabling a deeper understanding of the public's views on stunting in Indonesia.

We performed sentiment analysis using random forest algorithm and achieved about 97.5% accuracy. I would recommend you to try using some other machine learning algorithms such as LSTM or KNN and see if you can get better results.

ACKNOWLEDGMENTS

This work was supported by Data Science Program, Faculty of Science and Agricultural Teknologi, University of Muhammadiyah Semarang.

REFERENCE

- [1] H. Rahman, M. Rahmah, and Nur Saribulan, "UPAYA PENANGANAN STUNTING DI INDONESIA," *Jurnal Ilmu Pemerintahan Suara Khatulistiwa (JIPSK)*, vol. VIII, no. 01, pp. 44–59, 2023.
- [2] W. H. S. W. H. Organization, "Monitoring Health for the SDGs, World Health Organization," 2022.
- [3] A. Muhaimin *et al.*, "Social Media Analysis and Topic Modeling: Case Study of Stunting in Indonesia Analisis Sosial Media dan Pemodelan Topik: Kasus Studi Stunting di Indonesia," *Jurnal Informatika dan Teknologi Informasi*, vol. 20, no. 3, pp. 406–415, 2023, doi: 10.31515/telematika.v20i3.10797.
- [4] R. Rizqi Robbi Arisandi, B. Warsito, and A. Rachman Hakim, "Aplikasi Naive Bayes Classifier (NBC) Pada Klasifikasi Status Gizi Balita Stunting Dengan Pengujian K-Fold Cross Validation," *Jurnal Gaussian*, vol. 11, no. 1, pp. 130–139, 2022, [Online]. Available: <https://ejournal3.undip.ac.id/index.php/gaussian/>

- [5] U. R. Gurning, S. F. Octavia, D. R. Andriyani, N. Nurainun, and I. Permana, "Prediksi Risiko Stunting pada Keluarga Menggunakan Naïve Bayes Classifier dan Chi-Square," *MALCOM: Indonesian Journal of Machine Learning and Computer Science*, vol. 4, no. 1, pp. 172–180, Jan. 2024, doi: 10.57152/malcom.v4i1.1074.
- [6] A. A. Reza1 and M. S. Rohman2, "JITE (Journal of Informatics and Telecommunication Engineering) Prediction Stunting Analysis Using Random Forest Algorithm and Random Search Optimization," *JITE*, vol. 7, no. 2, 2024, doi: 10.31289/jite.v7i2.10628.
- [7] A. Muhaimin *et al.*, "Social Media Analysis and Topic Modeling: Case Study of Stunting in Indonesia Analisis Sosial Media dan Pemodelan Topik: Kasus Studi Stunting di Indonesia," *Jurnal Informatika dan Teknologi Informasi*, vol. 20, no. 3, pp. 406–415, 2023, doi: 10.31515/telematika.v20i3.10797.
- [8] M. Banurea, D. Betaria Hutagaol, and O. Sihombing, "KLASIFIKASI PENYAKIT STUNTING DENGAN MENGGUNAKAN ALGORITMA SUPPORT VECTOR MACHINE DAN RANDOM FOREST," *Jurnal TEKINKOM*, vol. 6, no. 2, 2023, doi: 10.37600/tekinkom.v6i2.927.
- [9] J. A. Nursiyono and C. Chotimah, "Analisis Sentimen Netizen Twitter terhadap Pemberitaan PPN Sembako dan Jasa Pendidikan dengan Pendekatan Social Network Analysis dan Naive Bayes Classifier," *Jurnal Ilmiah Teori dan Aplikasi Statistika*, vol. 14, pp. 52–58, 2021.
- [10] "Penerapan Text Mining pada Analisis Sentimen Pengguna Twitter Layanan Transportasi Online Menggunakan Metode Density Based Spatial Clustering of Applications With Noise (DBSCAN) dan K-Means," *Jurnal Ilmiah Teori dan Aplikasi Statistika*, vol. 15, pp. 184–194, 2022.
- [11] Stephenie, B. Warsito, and A. Prahutama, "Sentiment Analysis on Tokopedia Product Online Reviews Using Random Forest Method," in *E3S Web of Conferences*, EDP Sciences, Nov. 2020. doi: 10.1051/e3sconf/202020216006.
- [12] M. Ali Fauzi, A. Z. Arifin, S. C. Gosaria, and I. S. Prabowo, "Indonesian news classification using naïve bayes and two-phase feature selection model," *Indonesian Journal of Electrical Engineering and Computer Science*, vol. 8, no. 3, pp. 610–615, Dec. 2017, doi: 10.11591/ijeecs.v8.i3.pp610-615.
- [13] Bahrawi, "Sentiment Analysis Using Random Forest Algorithm Online Social Media Based," Makassar, Dec. 2019.
- [14] G. Khanvilkar and D. Vora, "Sentiment analysis for product recommendation using random forest," *International Journal of Engineering and Technology(UAE)*, vol. 7, no. 3, pp. 87–89, 2018, doi: 10.14419/ijet.v7i3.3.14492.
- [15] L. Breiman, "Random Forests," 2001.
- [16] Leo Breiman and A. Cutler, "Manual on Setting Up, Using, and Understanding Random Forest 4.0," 2003.
- [17] M. A. Fauzi, "Random forest approach fo sentiment analysis in Indonesian language," *Indonesian Journal of Electrical Engineering and Computer Science*, vol. 12, no. 1, pp. 46–50, Oct. 2018, doi: 10.11591/ijeecs.v12.i1.pp46-50.
- [18] I. Afdhal *et al.*, "Penerapan Algoritma Random Forest Untuk Analisis Sentimen Komentar Di YouTube Tentang Islamofobia," *Jurnal Nasional Komputasi dan Teknologi Informasi*, vol. 5, no. 1, 2022.