

Sentiment Analysis on Twitter Against IndiHome Providers Using Chi-Square and Ensemble Bagging Methods

Anisa Nur Aini¹, Jondri^{2*}, Widi Astuti³

School of Computing, Informatics Study Program, Telkom University, Bandung, Indonesia

Email: ¹anisanaini@student.telkomuniversity.ac.id, ^{2*}jondri@telkomuniversity.ac.id, ³widiwdu@telkomuniversity.ac.id

Email Author Correspondence: jondri@telkomuniversity.ac.id

Submitted: 25/07/2022; Accepted: 18/08/2022; Published: 30/09/2022

Abstract—During the Covid-19 pandemic, internet usage has increased rapidly. Now the internet is used as a means in the online teaching and learning process and work from home. One of the internet service providers is IndiHome. IndiHome is an internet service provider company that has a huge number of users. A large number of IndiHome users causes frequent problems, and this is one of the factors that IndiHome users provide various kinds of opinions or responses. Sentiment analysis is used to see the opinion or opinion given by someone on a particular object or problem. This study conducted a sentiment analysis using the Chi-square and the Ensemble Bagging method with three base classifier methods, namely K-Nearest Neighbor (K-NN), Support Vector Machine (SVM), and Naive Bayes (NB). Prediction results on labels obtained from each base classifier are combined using a hard majority vote. Tweet data collection was carried out in March 2022, and 6,962 tweets were collected. This study conducted two test scenarios. Scenario 1 is a scenario without oversampling with test results showing that Ensemble Bagging has the highest accuracy value of 83.32%, and in scenario 1 with hyperparameter tuning, Ensemble Bagging has the highest accuracy value of 83.93%. Scenario 2 is a scenario with oversampling, showing that Ensemble Bagging has the highest accuracy value of 84.51%, and scenario 2 with hyperparameter tuning also shows Ensemble Bagging has the highest accuracy value of 84.56%.

Keywords: Sentiment Analysis; IndiHome; Oversampling; Chi-Square; Ensemble Bagging

1. INTRODUCTION

During the Covid-19 pandemic, internet usage has increased rapidly. In 2021, 73.7% of the total population of Indonesia uses the internet [1], and 18.9% uses the internet for social media [2]. In Indonesia, one of the internet service providers is IndiHome. IndiHome is the most widely used internet service provider company by the people of Indonesia, where as many as 8.7% choose IndiHome as the internet service used [2]. With the large number of users who use IndiHome internet services, there are often obstacles such as unstable internet connections and the absence of an internet network, and this is one of the factors that people often give positive and negative opinions.

In 2019, according to data released by Twitter, Indonesia became the country with the largest growth in daily active Twitter users [3]. Twitter is one of the social media that allows users to share a message called a tweet [4]. Sentiment analysis is used to see the opinions or opinions given by someone on a particular object or problem. These opinions or opinions are usually channeled through customer care, one of which is through Twitter. Reviews or complaints expressed by users are subjective expressions and opinions that describe the user's feelings about the nature, event, or entity. Therefore, the opinions expressed by users on Twitter social media can be used as data to conduct sentiment analysis [5].

Previous research by Novan Dimas Pratama on sentiment analysis on consumer reviews used the Naive Bayes feature selection Chi-Square for recommendations for traditional food locations. This research aims to analyze traditional food consumers' opinion sentiments and provide location recommendations according to the required keywords. In this research, the method used is Naive Bayes and Chi-Square as a feature to assign a value to the features, which are then selected and sorted based on the percentage tested. The classification accuracy results with 25% feature selection is 81%, with 50% feature selection 80%, and with 77% feature selection 80%. From these results, it can be concluded that feature selection does not really affect the value of the accuracy results. It can be seen from the accuracy value between using and without using feature selection which is not too significant [6].

Chi-Square is more often used for classification with quite a lot of documents compared to Information Gain. The Chi-Square can increase the F-Measure classification text [7]. In previous research, Armanda Eka Putra conducted research on predicting the labels obtained from each algorithm combined using vote type hard majority. From the results of this study, the authors get the accuracy results without using the Chi-Square, which is 73.33%, and the accuracy results with the Chi-Square, which is 93.33%. This proves that using the Chi-Square can affect the sentiment results obtained [8].

In 2019, Ainun Nisa, Eko Darwiyanto, and Ibnu Asror researched sentiment analysis using the Naive Bayes classifier with the Chi-Square selection feature for telecommunication service providers. The results showed that the accuracy value using Naive Bayes without the Chi-Square selection feature was 84.4%, and if using the Chi-Square selection feature, it was 85.5%. This study shows that the Chi-Square feature selection does not significantly differ from the Naive Bayes classification [9].

The social media used in this is Twitter. The data used is tweets addressed to customer care, namely @IndiHomeCare. Data retrieval is done with the help of the Twitter API. This study analyzes the sentiment of the IndiHome internet service provider through Twitter social media using the ensemble Bagging method and the Chi-Square selection feature. The base classifier used is K-Nearest Neighbor (K-NN), Support Vector Machine (SVM),

and Naive Bayes (NB). Adding the Chi-Square selection feature is expected to help improve system performance in conducting sentiment analysis. The ensemble method was chosen because it can predict better performance by combining several classification models to get the best classification model, rather than previous studies that only used one basic classification.

2. RESEARCH METHODOLOGY

2.1 Research Stages

At this stage, the author makes a research phase flow, which will be implemented in research documents. This study uses Chi-square and Ensemble Bagging methods with three base classifier methods K-Nearest Neighbor (K-NN), Support Vector Machine (SVM), and Naive Bayes (NB). The prediction results on the labels obtained from each base classifier are combined using the hard majority vote to increase the analysis results in determining user sentiment towards the IndiHome provider based on studies that have been done previously.

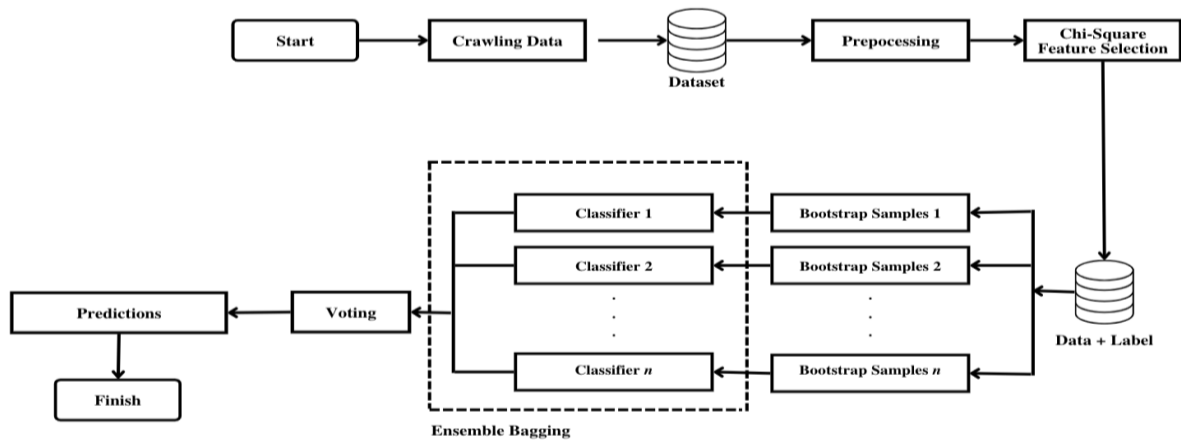


Figure 1. Research Phase Flow

Figure 1 is the flow of the research stage, the initial stage of this research is to collect data from Twitter. Tweet data obtained will be preprocessed so that the data is more structured and ready to be used. This study uses the Chi-Square selection feature and the Ensemble Bagging method with three base classifier methods, namely K-Nearest Neighbor (K-NN), Support Vector Machine (SVM), and Naive Bayes (NB). Prediction results on labels obtained from each base classifier are combined using the type of hard majority vote.

2.2 Datasets

In this study, the data source used is a tweet from social media Twitter which is addressed to the @IndiHomeCare account. Data retrieval requires assistance from the Twitter API (Application Programming Interface). Tweet data collection was carried out in March 2022, with 6,962 tweets collected for the IndiHome dataset. After collecting data, the next step is to label each tweet with a sentiment label in a positive, neutral, or negative category. The final result of the sentiment label on each tweet is the joint decision of three people.

Table 1 explains examples of tweet data for each sentiment, namely positive, neutral and negative, that have been obtained.

Table 1. Example Dataset

Tweet	Sentiment
@IndiHomeCare definisi memudahkan, tengkyu min	Positive
@arisafitryy @IndiHomeCare Bermasalah Mulu, jangan pake IndiHome ajadeh bestie	Negative
@IndiHomeCare cek dm	Neutral

2.3 Pre-Processing

Pre-processing is used to process and clean up unstructured data. In the pre-processing process in this research, case folding, removing punctuation, tokenization, stopword, and stemming processes are carried out so that the data is more structured and clean.

Table 2 describes the number of tweets pre-processed in each scenario. Two scenarios are used: testing without oversampling and testing with oversampling.

Table 2. Number of tweets in each

Sentiment	Total Scenario 1 data without oversampling	Total Scenario 2 data with oversampling
Positive	579	2445

Negative	2155	2445
Neutral	2445	2445
Total	5179	7335

Table 3 is an example of the results of tweets that have gone through the pre-processing. The pre-processing stage using several stages, namely case folding to change capital letters to lowercase letters [10], remove punctuation to clean characters that cannot be read by the system [11], tokenization to separate each sentence into words [12], stopword to removing words that have no important meaning [13], and stemming to remove affixes that have unimportant meanings in each word [14].

Table 3. Example of pre-processing tweet

Tweet	Preprocessing Stage	Results
@IndiHomeCare Dear indihome. Ada gangguankah? Lampu di modem saya yg Pon merah? Tolong responnya'	Case folding	@indihomecare dear indihome.. ada gangguankah? lampu di modem saya yg pon merah? tolong responnya'
	Remove punctuation	indihomecare dear indihome ada gangguankah lampu di modem saya yg pon merah tolong responnya
	Tokenization	['indihomecare', 'dear', 'indihome', 'ada', 'gangguankah', 'lampu', 'di', 'modem', 'saya', 'yg', 'pon', 'merah', 'tolong', 'responnya']
	Stopword	['gangguankah', 'lampu', 'modem', 'pon', 'merah', 'responnya']
	Stemming	['gangguan', 'lampu', 'modem', 'pon', 'merah', 'respon']

2.4 Chi-Square Feature Selection

Chi Square is a feature selection method that looks at the dependence of terms with categories that have trial conditions that can be used as follows, namely, if the contingency table is in the form of 2x2, then there should be no cells with an expected frequency (Fn) of less than five or an expected count, there are no cells that have a reality frequency value (F0) which is 0 (zero) or an actual count, and if the table form is more than 2 x 2, then the number of cells with an expected frequency (Fh) which is less than five should not be more than 20% [15].

$$X^2 = \sum_{i=1}^n \frac{(O_i - E_i)^2}{E_i} \tag{1}$$

Equation (1) is the formula for Chi Square which shows O_i as the i -th observation value and E_i as the i -th expectation value.

2.5 Ensemble Bagging

Ensemble is a machine learning technique where several models will be trained to solve the same problem and then combined to get the best results. There are several ensemble learning methods, including the Ensemble Bagging method. Ensemble bagging is a method that can improve machine learning algorithms by combining several classification models [16]. Ensemble bagging uses a base classifier trained separately with a random sample from the training data by voting or averaging to have accurate and stable final results.

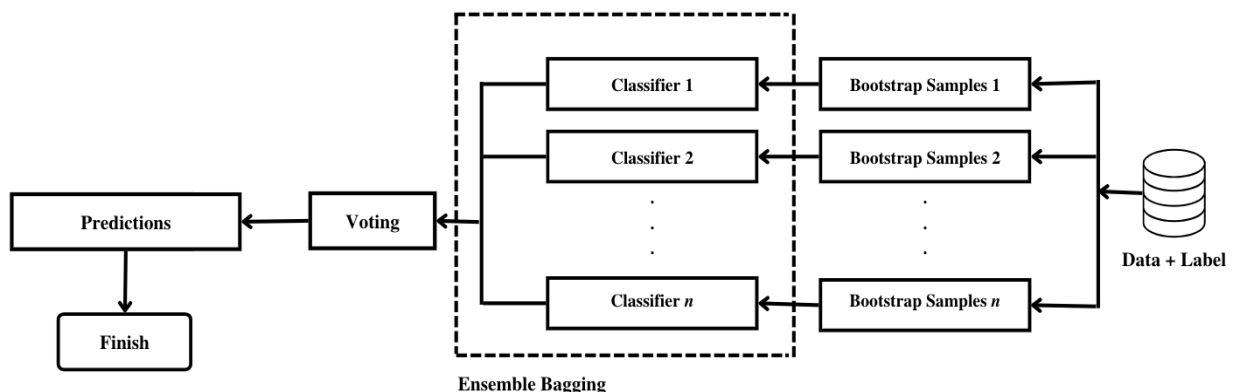


Figure 2. Ensemble Bagging Process Flow

Figure 2 describes the flow of the Ensemble Bagging, wherein in the first stage, the training data will be processed in bootstrap. Bootstrap is a data sampling technique to generate multiple random samples from training data. The classification process uses three classification algorithms, namely K-Nearest Neighbor (K-NN), Support

Vector Machine (SVM), and Naive Bayes (NB). Prediction results on labels obtained from each base classifier are combined using vote type hard majority. The way hard majority vote is by determining the most dominant prediction label among the three algorithms. This method was chosen because it is simple but gives maximum results. After obtaining the final label, the following process is to conduct testing using test data to get the best accuracy value.

3. RESULTS AND DISCUSSION

The results of the sentiment analysis of the IndiHome provider on Twitter social media using the Chi-Square classification method Ensemble Bagging are discussed in this section. The base classifier used is Support-Vector Machine, K-Nearest Neighbor, and Naive Bayes. In this study, two test scenarios were used: the first scenario of testing without oversampling and the second scenario with oversampling.

Table 4 is the result of 10-Fold Cross-Validations Scores of the three base classifier used: Support-Vector Machine, K-Nearest Neighbor, and Naive Bayes scenario oversampling.

Table 4. 10-Fold Cross-Validations Scores without oversampling

	Fold Score 1	Fold Score 2	Fold Score 3	Fold Score 4	Fold Score 5	Fold Score 6	Fold Score 7	Fold Score 8	Fold Score 9
SVM	0.804	0.768	0.732	0.719	0.827	0.791	0.798	0.829	0.791
KNN	0.676	0.632	0.568	0.588	0.636	0.610	0.610	0.626	0.641
Naive Bayes	0.809	0.817	0.755	0.820	0.837	0.809	0.817	0.837	0.842

Table 5 shows the results of the average prediction accuracy of the IndiHome dataset. The results obtained with the first scenario without oversampling, K-Nearest Neighbor has the lowest predictive accuracy value of 62.1%, and Naive Bayes has the highest predictive accuracy value of 81.7%.

Table 5. Cross Validation Accuracy without oversampling

Support-Vector Machine	K-Nearest Neighbor	Naive Bayes
0.785 +/- 0.034	0.621 +/- 0.028	0.817 +/- 0.023

Table 6 is the result of 10-Fold Cross-Validations Scores of the three base classifiers used, namely Support-Vector Machine, K-Nearest Neighbor and Naive Bayes scenario oversampling.

Table 6. 10-Fold Cross-Validations Scores with oversampling

	Fold Score 1	Fold Score 2	Fold Score 3	Fold Score 4	Fold Score 5	Fold Score 6	Fold Score 7	Fold Score 8	Fold Score 9
SVM	0.807	0.823	0.821	0.821	0.840	0.810	0.790	0.810	0.803
KNN	0.647	0.672	0.681	0.661	0.669	0.645	0.623	0.643	0.660
Naive Bayes	0.785	0.805	0.816	0.803	0.816	0.816	0.794	0.809	0.829

Table 7 shows the average prediction accuracy results from the IndiHome dataset. The results of the second scenario with oversampling are that K-Nearest Neighbor has the lowest predictive accuracy value of 65.8%, and Support-Vector Machine has the highest predictive accuracy value of 81.7%.

Table 7. Cross Validation Accuracy IndiHome with oversampling

Support-Vector Machine	K-Nearest Neighbor	Naive Bayes
0.817 +/- 0.015	0.658 +/- 0.017	+/- 0.012

3.1 Evaluation Model

At this stage, the author has conducted research using Chi-Square and Ensemble Bagging classification method. Three base classifiers Support-Vector Machine, K-Nearest Neighbor, and Naive Bayes. This study uses two test scenarios: the first scenario testing without oversampling and the second scenario testing with oversampling. The table below is the result of the performance of each test scenario on the IndiHome dataset.

3.1.1 First Scenario

This section discusses the results of the performance of the first scenario, namely a test scenario without oversampling. At this stage, testing is carried out without hyperparameter tuning and testing with hyperparameter tuning.

Table 8 shows the results of the classification performance of the first scenario, namely the test scenario without oversampling. The test results show that among the three base classifiers used, K-Nearest Neighbor with a value of k=3 has the lowest accuracy value of 66.79%, and Ensemble Bagging obtains the highest accuracy value of 83.32% because Ensemble Bagging uses several classification models. At once, combine them to get one classification model with the best results.

Table 8. Performance in the first scenario without oversampling

	SVM	K-Nearest Neighbor	Naive Bayes	Ensemble Bagging
Accuracy	0.828571	0.667954	0.768340	0.833205
Precision	0.830978	0.720117	0.812851	0.835533
Recall	0.828571	0.667954	0.768340	0.833205
F1 Score	0.827566	0.646472	0.771030	0.831846

Table 9 shows the results of the hyperparameter tuning in the first scenario with the best accuracy results. These results will be processed at the Ensemble Bagging stage.

Table 9. Hyperparameter tuning in the first scenario without oversampling

	Hyperparameter Tuning	Value
SVM	Kernel	Linear
	C	1
K-NN	n_neighbors	6
	leaf_size	1
Naive Bayes	Alpha	0.01

Table 10 shows the results of the classification performance with hyperparameter tuning in the first scenario. Hyperparameter tuning is used to give the optimal value to maximize the accuracy value. The test results show that among the three base classifiers used, Ensemble Bagging obtains the highest accuracy value of 83.93% because Ensemble Bagging uses several classification models at once and combines them to get one classification model with the best results.

Table 10. Performance hyperparameter tuning in the first scenario without oversampling

	SVM HT	K-NN HT	Naive Bayes HT	Ensemble Bagging HT
Accuracy	0.836293	0.715058	0.764479	0.839382
Precision	0.836089	0.722952	0.814428	0.841797
Recall	0.836293	0.715058	0.764479	0.839382
F1 Score	0.836081	0.713733	0.769377	0.839456

3.1.2 Second Scenario

This section discusses the results of the performance of the second scenario, namely the test scenario with oversampling. At this stage, testing without hyperparameter tuning and testing with hyperparameter tuning is also carried out.

Table 11 shows the results of the classification performance of the second scenario, namely the test scenario with oversampling. The test results show that among the three base classifiers used, K-Nearest Neighbor with a value of k=3 has the lowest accuracy value of 65.70%, and Ensemble Bagging obtains the highest accuracy value of 84.51% because Ensemble Bagging uses several classification models. At once and combine them to get one classification model with the best results.

Table 11. Performance in second scenario with oversampling

	SVM	K-Nearest Neighbor	Naive Bayes	Ensemble Bagging
Accuracy	0.836423	0.657034	0.809706	0.845147
Precision	0.836130	0.717128	0.827020	0.845290
Recall	0.836423	0.657034	0.809706	0.845147
F1 Score	0.833736	0.621960	0.810993	0.842584

Table 12 shows the results of hyperparameter tuning in second scenario with oversampling of the best accuracy results. These results will be processed at the Ensemble Bagging stage.

Table 12. Hyperparameter tuning in the second scenario with oversampling

	Hyperparameter Tuning	Value
SVM	Kernel	Linear
	C	1
K-NN	n_neighbors	1
	Naive Bayes	Alpha

Table 13 shows the results of the classification performance with hyperparameter tuning in the second scenario, hyperparameter tuning is used to give an optimized value to maximize the accuracy value. The test results show that among the three base classifier used, Ensemble Bagging obtains the highest accuracy value of 84.56% because

Ensemble Bagging uses several classification models at once and combines them to get one classification model with the best results.

Table 13. Performance hyperparameter tuning in the second scenario with oversampling

	SVM HT	K-NN HT	Naive Bayes HT	Ensemble Bagging HT
Accuracy	0.834787	0.683751	0.822792	0.845692
Precision	0.835152	0.721818	0.834282	0.846907
Recall	0.834787	0.683751	0.822792	0.845692
F1 Score	0.831670	0.655680	0.823492	0.842209

4. CONCLUSION

In this study, sentiment analysis was conducted on the IndiHome internet service provider through Twitter social media using the Ensemble Bagging method and the Chi-Square feature selection with two scenarios. The base classifier used is K-Nearest Neighbor (K-NN), Support Vector Machine (SVM), and Naive Bayes (NB). The first scenario is a classification test without oversampling, and the second scenario is a classification test with oversampling. The results from the first scenario show that Ensemble Bagging has the highest accuracy value of 83.32%, and the first scenario with hyperparameter tuning shows Ensemble Bagging has the highest accuracy value of 83.93%. In the second scenario, Ensemble Bagging has the highest accuracy value of 84.51%, and the second scenario, with hyperparameter tuning, also shows Ensemble Bagging has the highest accuracy value of 84.56%. From the results of the two test scenarios obtained, it can be concluded that Ensemble Bagging gets the highest accuracy results compared to K-Nearest Neighbor (K-NN), Support Vector Machine (SVM), and Naive Bayes (NB). Ensemble Bagging gets the highest accuracy value because it uses several classification models and combines them to get one classification model with the best results.

REFERENCES

- [1] DOI 10.47065/bits.v4i2.1901[1] “Peluang penetrasi Internet dan Tantangan Regulasi Daerah : Asosiasi Penyelenggara Jasa Internet Indonesia,” 2021.
- [2] “Penetrasi & Profil Perilaku Pengguna Internet Indonesia: Asosiasi Penyelenggara Jasa Internet Indonesia,” 2018.
- [3] Bill Clinton, “Pengguna Aktif Harian Twitter Indonesia Diklaim Terbanyak,” Kompas.com, Aug. 2019. <https://tekno.kompas.com/read/2019/10/30/16062477/pengguna-aktif-harian-twitter-indonesia-diklaim-terbanyak> (accessed Jul. 22, 2022).
- [4] C. S. Rao, G. S. Prasad, and V. V. Rao, “Prediction and Analysis of Sentiments on Twitter Data using Hybrid Naive Bayes Approach,” *International Journal of Innovative Technology and Exploring Engineering (IJITEE)*, vol. 8, no. 8, pp. 2278–3075, 2019.
- [5] S. Farah Aliyah, H. Yasin, B. Warsito, T. Widiaris, D. Statistika, and F. Sains dan Matematika, “Analisis Sentimen PT TIKI Jalur Nugraha Ekakurir (PT TIKI JNE) Pada Media Sosial Twitter Menggunakan Model Feed Forward Neural Network,” *Statistika*, vol. 8, no. 2, pp. 103–113, 2020.
- [6] N. Dimas Pratama, Y. A. Sari, and P. P. Adikara, “Analisis Sentimen Pada ReviewKonsumen Menggunakan Metode Naive Bayes Dengan Seleksi Fitur Chi Square Untuk Rekomendasi Lokasi Makanan Tradisional,” *Jurnal Pengembangan Teknologi Informasi dan Ilmu Komputer*, vol. 2, no. 9, pp. 2982–2988, 2018, [Online]. Available: <http://j-ptiik.ub.ac.id>
- [7] D. B. Satmoko, P. Sukarno, and E. M. Jadied, “Peningkatan Akurasi Pendeteksian Serangan DDoS Menggunakan Multiclassifier Ensemble Learning dan Chi-Square,” *e-Proceeding of Engineering*, vol. 5, no. 3, pp. 7877–7985, 2018.
- [8] J. Ling, I. Putu, E. N. Kencana, and T. B. Oka, “Analisis Sentimen Menggunakan Metode Naive Bayes Classifier dengan Seleksi Fitur Chi Square,” *E-Jurnal Matematika*, vol. 3, no. 3, pp. 92–99, 2014.
- [9] A. Nisa, E. Darwiyanto, and I. Asror, “Analisis Sentimen Menggunakan Naive Bayes Classifier dengan Chi-Square Feature Selection Terhadap Penyedia Layanan Telekomunikasi,” *e-Proceeding of Engineering*, vol. 6, pp. 8650–8658, 2019.
- [10] A. Susanto, M. Atho’il Maula, I. Utomo, W. Mulyono, and K. Sarker, “Sentiment Analysis on Indonesia Twitter Data Using Naive Bayes and K-Means Method,” *Journal of Applied Intelligent System*, vol. 6, no. 1, pp. 40–45, 2021, [Online]. Available: <http://kateglo.com>.
- [11] B. Pahwa, S. Taruna, and N. Kasliwal, “Sentiment Analysis- Strategy for Text Pre-Processing,” *International Journal of Computer Applications*, vol. 180, no. 34, pp. 15–18, Apr. 2018, doi: 10.5120/ijca2018916865.
- [12] C. Toraman, E. H. Yilmaz, F. Şahinuç, and O. Özcelik, “Impact of Tokenization on Language Models: An Analysis for Turkish,” Apr. 2022, [Online]. Available: <http://arxiv.org/abs/2204.08832>
- [13] A. W. Pradana and M. Hayaty, “The Effect of Stemming and Removal of Stopwords on the Accuracy of Sentiment Analysis on Indonesian-language Texts,” *Kinetik: Game Technology, Information System, Computer Network, Computing, Electronics, and Control*, vol. 4, pp. 375–380, Oct. 2019, doi: 10.22219/kinetik.v4i4.912.
- [14] S. Sharma and M. Bansal, “Stemming and lemmatization of tweets for sentiment analysis using R,” *International Journal of Recent Technology and Engineering*, vol. 8, no. 2, pp. 2038–2040, Jul. 2019, doi: 10.35940/ijrte.B2157.078219.
- [15] I. C. Negara and A. Prabowo, “Penggunaan Uji Chi-Square untuk Mengetahui Pengaruh Tingkat Pendidikan dan Umur terhadap Pengetahuan Penasun Mengenai HIV-AIDS di Provinsi DKI Jakarta,” *FMIPA Unsoed Purwokerto*, pp. 1–8, 2018.
- [16] Zhi-Hua Zhou, “Ensemble Methods Foundations and Algorithms,” in *Chapman & Hall/CRC Machine Learning & Pattern Recognition Series*, 2012, pp. 47–66. Accessed: May 28, 2022. [Online]. Available: <https://tjzhifei.github.io/links/EMFA.pdf>