

Indonesian Journal of Computer Science

ISSN 2549-7286 (*online*) Jln. Khatib Sulaiman Dalam No. 1, Padang, Indonesia Website: ijcs.stmikindonesia.ac.id | E-mail: ijcs@stmikindonesia.ac.id

Machine Learning on Opinion Mining of Netizen's Hate Speech

Mutiana Pratiwi^{1*}, Rima Liana Gema²

mutiana_pratiwi@upiyptk.ac.id, rimalianagema@upiyptk.ac.id ^{1,2}Universitas Putra Indonesia YPTK Padang

Article Information	Abstract
Submitted : 19 Dec 2023 Reviewed: 2 Jan 2024 Accepted : 17 Feb 2024	Netizen comments written in an online news portal through social media platforms, one of which is Instagram, can be used as material in the sentiment analysis process, which can be classified into positive, negative, or neutral sentiments. Sentiment analysis is part of the study of text mining, the
Keywords	science of discovering unknown knowledge by automatically extracting information from large volumes of unstructured text into useful information.
Machine learning Support Vactor Machine Opinion Mining	The resulting information is in the form of sentiment towards a topic, whether it tends to be positive, negative, or neutral. The classification method used in this research is Support Vector Machine (SVM) and TF-IDF
Hate Speech Netizen	data weighting to classify text. Stages to perform data analysis are pre- processing to clean data, word weighting, labeling data into positive, negative, or neutral classes, and classifying and visualizing data with graphs. Accuracy tests using 70:30 split data showed that the accuracy reached 98%. Tests with 80:20 and 90:10 split data also showed high accuracy of 98% and 99%.

A. Introduction

Based on the survey results of the Indonesian Internet Service Providers Association (APIII), internet users in Indonesia have reached 210 million people, with an internet penetration rate of 77.02%. This marks a significant increase compared to previous years, with a penetration rate of 64.80% in 2018 and 73.70% in 2019-2020. The rapid development of Internet technology has now given rise to many types of social media. Social media is a medium used by its users to exchange information and network via the internet. There are so many benefits of social media, some of the uses of social media are knowing sports information, business, tourism, and also for political affairs (Mahawardana, Imawati, and Dika 2022) (Iskandar Mulyana and Lutfianti 2023). This is also supported by developments in mobile phone technology so social media has become very popular because of its ease of communication (Putri Nur Lyrawati 2023)(Kowsari et al. 2019). Currently, the development of social media has entered the world of journalism where many online news portals have emerged. An online news portal is one of the press media that presents news, articles and features online (Anam et al. 2021). The development of online news portals in Indonesia is directly proportional to the development of audience development in the search for information sources (Arifin 2013).

Online news portals present news not only through websites and apps but also through social media platforms. These social media sites include Instagram, Facebook, Twitter, and YouTube. The current form of journalism also follows this trend to attract a wider audience. Digital platforms on online news portals enable audience involvement in the journalistic process. (Woro Harkandi Kencana, 2022). One of the social media applications that is widely used by the public is the Instagram application(Hashida, Tamura, and Sakai 2018). The Instagram application is often used as an efficient and effective means of campaigning, promoting, and socializing policies. In Instagram social media posts, there are many comments given by readers, so that they can generate public opinion in the form of both positive and negative comments (Fahlevvi 2022).

Based on the comments sent, a picture of public opinion can be seen, and then it will be extracted into structured information data to find out the polarity of an opinion and whether it is included in the positive, neutral, or negative classification (Wandani, Fauziah, and Andrianingsih 2021) (Darwis, Pratiwi, and Pasaribu 2020). The method that can be used for data to become structured information is the use of sentiment analysis, which is part of text mining (Alkaff et al. 2023). Text Mining is the process of discovering unknown knowledge through automatic extraction of information from large volumes of unstructured text (Sodik Pamungkas and Kharisudin 2021). Sentiment analysis also uses Natural Language Processing (NLP) to detect subjective information present in documents (Arsi and Waluyo 2021)(Idris, Mustofa, and Salihi 2023). The information used will be analyzed in several processes, starting with extracting data and then cleaning the data. The data-cleaning process is called pre-processing. After cleaning, the data is processed with the method that will be used. One of the methods used to classify sentiment analysis is Support Vector Machine (SVM) (Adhitya Karel Maulaya and Junadhi 2022) (Rahmaddeni et al. 2022).

Research on Sentiment Analysis on the use of the Shopee application using the Support Vector Machine (SVM) algorithm conducted (Dedi Darwis, Nery Siskawati, and Zaenal Abidin 2020). From the experimental results of his research, it produces quite good performance with an accuracy of 98% and an f1-score of 0.98 or 98%. Another study on the Application of the Support Vector Machine Method for Twitter User Sentiment Analysis was conducted in which the results of his research experiments resulted in SVM performance with an accuracy of 93% (Amin and Rainarli 2019). Likewise, in research on Sentiment Analysis of Online News Media Application Reviews on Google Play using the Support Vector Machines (SVM) and Naive Bayes Algorithm Methods and knowing the tendency of public opinion on Google Play about online news media applications where subjects were taken using scraping as many as 5615 opinions resulted in the conclusion that the accuracy of SVM (Ibrohim and Budi 2023).

Based on this, this research was conducted to determine the performance of the SVM method in classifying positive, neutral, and negative sentiments on infosumbar Instagram social media. By doing this research, it is hoped that it can provide an overview of how netizens react to each upload topic regarding positive, neutral, and negative comments on the account.

B. Research Method

In this section, the methods and techniques used for the classification of reviews on public opinion spaces and the steps taken during the experiments are classified. Figure 1 illustrates the stages of this research starting from the online review dataset until each review is classified into positive, negative, and neutral.



Fig.1. Research Framework

This research involves several stages to analyze public opinion reviews. The first stage is to conduct data collection to obtain a dataset of netizen comments related to city government policies. The dataset is then labeled with positive, negative, and neutral sentiments in the second stage. The labeled dataset is then processed by performing case-folding, stopword removal, tokenization, and stemming in the preprocessing stage. After preprocessing, TF-IDF analysis is performed to measure the importance of words in a document or text corpus. The next stage is text classification using the SVM method and model validation. a. Dataset

Indonesian dataset that will be used in this research. The amount of data obtained is a data sample from the scraping results of 1050 data and the amount is preprocessed to produce data selection. The following is the dataset taken using the scraping technique:

	Table 1. Dataset					
	No	Comment	Author			
0	1	Baguslah tidak bikin berantakan	van.mwessi			
1	2	Sebenarnya lebih baik direlokasi, karena	af_rina3			
		panta				
2	3	makanan yang layak dong yang dijual	dhiocwb			
		Taplau ini memang harus dirapikan, selain	pasar_seafood_online			
		tata				
		Mau tempat yang luas? Pindah ke by pass	verifirdaus			
1048	1049	Kebiasaan lama yang tidak bisa dirobah	ju_liar_di			
		Bisa				
1049	1050	Tetapkan menggunakan Pergub,, sehingga	dwi_ab_xv			
		ada dasa				

b. Labeling

This study first conducted an initial data selection to select data that fit the comment categories related to netizen comments on city government policies. After the selection process, comment data was obtained that could be classified as positive, negative, or neutral. Some examples of data selection results can be found in Table 2.

Table 2 Data Labeling

	Table 2. Data Labeling					
	No	Comment	Author	Sentiment		
0	1	Baguslah tidak bikin berantakan	van.mwessi	negative		
1	2	Sebenarnya lebih baik direlokasi,	af_rina3	possitive		
		karena panta				
2	3	makanan yang layak dong yang	dhiocwb	negative		
		dijual				
		Taplau ini memang harus	pasar_seafood_online	possitive		
		dirapikan, selain tata				
		Mau tempat yang luas? Pindah ke	verifirdaus	negative		
		by pass				
1048	1049	Kebiasaan lama yang tidak bisa	ju_liar_di	neutral		

		dirobah Bisa		
1049	1050	Tetapkan menggunakan Pergub,,	dwi_ab_xv	neutral
		sennigga ada dasa		

After the data labeling process, a comparison of the number of negative, positive, and neutral sentiments on the dataset can be seen in Figure 2.



Fig. 2. Dataset distribution

C. Result and Discussion

This research is conducted in several stages, the first stage is data collection, the data collected is review data from social media and online news portals. The data used to test this research is Indonesian-language public opinion review data regarding city government policies. The dataset used in this research is 1050 reviews. In this research, we do pre-processing and then do TF IDF feature extraction. Then we classify the review text to find sentiment analysis with the ML model. Preprocessing the dataset needs to be done before the classification process, to eliminate words that are not suitable for research, equalize the shape of words or letters (lowercase), and reduce the number of words. By preprocessing the data, it can be ensured that the data used for training is better, cleaner, and more suitable for use in SVM analysis. After text preprocessing, it is necessary to weigh the words. Some examples of results from the text preprocessing process can be presented in Table 3.

Table 3. Data Pre-Proc	cessing (Cleaning)
------------------------	--------------------

No	Text	Result
1	Baguslah tidak bikin berantakan	Baguslah tidak bikin berantakan
2	Sebenarnya lebih baik direlokasi, karena	Sebenarnya lebih baik direlokasi karena
	pantai ini kan fasilitas umum siapa saja	pantai ini kan fasilitas umum siapa saja
	boleh menikmatisementara kalau ada	boleh menikmati sementara kalau ada
	pedangang dan kita tidak membeli	pedangang dan kita tidak membeli
	dagangannya kita nggak boleh duduk atau	dagangannya kita nggak boleh duduk
	parkir disana	atau parkir disana
3	makanan yang layak dong yang dijual	makanan yang layak dong yang dijual

	Table 4. Data Pre-Processing (Case Folding)				
No	Text	Result			
1	Baguslah tidak bikin berantakan	baguslah tidak bikin berantakan			
2	Sebenarnya lebih baik direlokasi, karena	sebenarnya lebih baik direlokasi karena			
	pantai ini kan fasilitas umum siapa saja	pantai ini kan fasilitas umum siapa saja			
	boleh menikmatisementara kalau ada	boleh menikmati sementara kalau ada			
	pedangang dan kita tidak membeli	pedangang dan kita tidak membeli			
	dagangannya kita nggak boleh duduk atau	dagangannya kita nggak boleh duduk			
	parkir disana	atau parkir disana			
3	makanan yang layak dong yang dijual	makanan yang layak dong yang dijual			

Table 4.	Data	Pre-Proc	cessing	(Case	Folding)
rabie n	Data	110 1100	- COULTE	Gabe	1 0101115

	No	Comment	Author	Sentiment
0	1	bagus bikin beranta	van.mwessi	0
1	2	relokasi pantai fasilitas menikmatisementara p	af_rina3	2
2	3	makan layak jual	dhiocwb	0
3	4	taplau rapi tata letak yg beranta sdm edukasi	pasar_seafood_online	2
4	5	luaspindah by pass	verifirdaus	0
1045	1046	amen preman sana tertib bikin unjung	mhdsyafiq08_	1
1046	1047	bikinkn menara pantau pakuntuk pantau unjung y	muhammadfikrii12	1
1047	1048	bahas muncul area dagang baru buat bijak siste	donny.aswin	1
1048	1049	biasa dirobah pantai padang sepi unjung ulah d	ju_liar_di	1
1049	1050	tetap pergubsehingga dasar kuat menindakbagi y	dwi_ab_xv	1

1050 rows × 4 columns

Fig. 3. Labeling Result

After going through the preprocessing process, a cleaned dataset is generated. To perform sentiment analysis, the cleaned data is converted into vectors by applying the TF-IDF method. Each model in the algorithm is evaluated in terms of accuracy. In the experiments conducted on the three Machine Learning algorithms using "TF-IDF" feature extraction, it was found that SVM provided the best training accuracy model with good classification. The following is a comparison of accuracy with three different ratios of training and testing data. The following is a comparison of the accuracy value of the SVM model

Table 5. Accuracy Split Data 70:30				
	precision	recall	f1-score	support
0	1.00	1.00	1.00	143
1	1.00	0.93	0.96	102
2	0.91	1.00	0.95	70
accuracy			0.98	315

macro avg	0.97	0.98	0.97	315
weighted avg	0.98	0.98	0.98	315

The results of this accuracy are presented in Figure 4 confusion matrix of SVM.



Fig. 4. Confusion Matrix Split Data 70:30

Table 6. Accuracy Split Data 80:20				
	precision	recall	f1-score	support
0	1.00	1.00	1.00	90
1	0.97	0.97	0.97	66
2	0.96	0.96	0.96	54
accuracy			0.98	210
macro avg	0.98	0.98	0.98	210
weighted avg	0.98	0.98	0.98	210

The results of this accuracy are presented in Figure 5 confusion matrix of SVM split data 80:20



Fig. 5. Co	nfusion	Matrix	Split	Data	80:20
------------	---------	--------	-------	------	-------

	Table. 7. Accuracy Score Split Data 90:10				
	precision	recall	f1-score	support	
0	1.00	1.00	1.00	46	
1	0.97	1.00	0.99	36	
2	1.00	0.96	0.98	23	
accuracy			0.99	105	
macro avg	0.99	0.99	0.99	105	
weighted avg	0.99	0.99	0.99	105	



Fig. 6. Confusion Matrix Split Data 90:10

D. Conclusion

This research uses the SVM method and TF-IDF data weighting to classify text. Through accuracy testing using 70:30 split data, it was found that the accuracy reached 98%. Tests with 80:20 and 90:10 split data also resulted in high accuracy of 98% and 99% respectively. In this research, it is important to define technical

abbreviations to ensure clear understanding. The text uses objective, neutral language, and is free from embellishment. A formal structure was followed and argumentation was balanced, avoiding bias. Appropriate vocabulary selection was made, and grammar and formatting were consistent. However, there were differences between testing directly on the model and testing using the comparison chart. This offers opportunities for future research using other Machine Learning algorithms and methods, as well as expanding the number of datasets. In carrying out intent sentiment analysis, the number and variety of datasets can be increased.

E. Acknowledgment

Thanks are due to the Chairperson of the Padang Computer College Foundation Universitas Putra Indonesia YPTK Padang who has provided full support in this research.

F. References

- [1] Adhitya Karel Maulaya, and Junadhi. 2022. "Analisis Sentimen Menggunakan Support Vector Machine Masyarakat Indonesia Di Twitter Terkait Bjorka." Jurnal CoSciTech (Computer Science and Information Technology) 3(3): 495–500.
- [2] Alkaff, Muhammad et al. 2023. "Hate Speech Detection for Banjarese Languages on Instagram Using Machine Learning Methods." *MATRIK : Jurnal Manajemen, Teknik Informatika dan Rekayasa Komputer* 22(3): 495–504.
- [3] Amin, Nur Daniar, and Ednawati Rainarli. 2019. "Klasifikasi Konten Instagram Berdsarkan Komentar Menggunakan Support Vector Machine." : 23–120. https://elibrary.unikom.ac.id/id/eprint/1112/.
- [4] Anam, M. Khairul et al. 2021. "Penerapan Na[¬]ive Bayes Classifier, K-Nearest Neighbor (KNN) Dan.Pdf." *Jurnal Manajemen, Teknik Informatika, dan Rekayasa Komputer* 21(1): 139~150.
- [5] Arifin, Pupung. 2013. "Persaingan Tujuh Portal Berita Online Indonesia Berdasarkan Analisis Uses and Gratifications." Jurnal ILMU KOMUNIKASI 10(2): 195–211.
- [6] Arsi, Primandani, and Retno Waluyo. 2021. "Analisis Sentimen Wacana Pemindahan Ibu Kota Indonesia Menggunakan Algoritma Support Vector Machine (SVM)." Jurnal Teknologi Informasi dan Ilmu Komputer 8(1): 147.
- [7] Darwis, Dedi, Eka Shintya Pratiwi, and A Ferico Octaviansyah Pasaribu.
 2020. "Penerapan Algoritma Svm Untuk Analisis Sentimen Pada Data Twitter Komisi Pemberantasan Korupsi Republik Indonesia." *Edutic -Scientific Journal of Informatics Education* 7(1): 1–11.
- [8] Dedi Darwis, Nery Siskawati, and Zaenal Abidin. 2020. "Penerapan Algoritma Naive Bayes Untuk Analisis Sentimen Review Data Twitter BMKG Nasional." Jurnal TEKNO KOMPAK 15(1): 131–45.
- [9] Fahlevvi, Mohammad Rezza. 2022. "Analisis Sentimen Terhadap Ulasan Aplikasi Pejabat Pengelola Informasi Dan Dokumentasi Kementerian Dalam Negeri Republik Indonesia Di Google Playstore Menggunakan Metode

Support Vector Machine." Jurnal Teknologi dan Komunikasi Pemerintahan 4(1): 1–13.

- [10] Hashida, Shuichi, Keiichi Tamura, and Tatsuhiro Sakai. 2018. "Classifying Sightseeing Tweets Using Convolutional Neural Networks with Multi-Channel Distributed Representation." *Proceedings - 2018 IEEE International Conference on Systems, Man, and Cybernetics, SMC 2018* (February): 178–83.
- [11] Ibrohim, Muhammad Okky, and Indra Budi. 2023. "Hate Speech and Abusive Language Detection in Indonesian Social Media: Progress and Challenges." *Heliyon* 9(8): e18647. https://doi.org/10.1016/j.heliyon.2023.e18647.
- [12] Idris, Irma Surya Kumala, Yasin Aril Mustofa, and Irvan Abraham Salihi. 2023. "Analisis Sentimen Terhadap Penggunaan Aplikasi Shopee Mengunakan Algoritma Support Vector Machine (SVM)." *Jambura Journal of Electrical and Electronics Engineering* 5(1): 32–35.
- [13] Iskandar Mulyana, Dadang, and Nesti Lutfianti. 2023. "Penerapan Sentimen Analisis Dengan Algoritma SVM Dalam Tanggapan Netizen Terhadap Berita Resesi 2023." *Sisfotenika* 13(1): 53–64.
- [14] Kowsari, Kamran et al. 2019. "Text Classification Algorithms: A Survey." *Information (Switzerland)* 10(4): 1–68.
- [15] Mahawardana, Putu Pasek Okta, Ida Ayu Putu Febri Imawati, and I Wayan Dika. 2022. "Analisis Sentimen Berdasarkan Opini Dari Media Sosial Twitter Terhadap 'Figure Pemimpin' Menggunakan Python." Jurnal Manajemen dan Teknologi Informasi 12(2): 50–56. https://ojs.mahadewa.ac.id/index.php/jmti/article/view/2111.
- [16] Putri Nur Lyrawati, Dayang. 2023. "Hate Speech Detection on Twitter Approaching The Indonesian Election Using Machine Learning." The Journal on Machine Learning and Computational Intelligence (JMLCI) (2808-974X): 26–31.
- [17] Rahmaddeni, Rahmaddeni et al. 2022. "Comparison of Support Vector Machine and XGBSVM in Analyzing Public Opinion on Covid-19 Vaccination." *ILKOM Jurnal Ilmiah* 14(1): 32–38.
- [18] Sodik Pamungkas, Fajar, and Iqbal Kharisudin. 2021. "Analisis Sentimen Dengan SVM, NAIVE BAYES Dan KNN Untuk Studi Tanggapan Masyarakat Indonesia Terhadap Pandemi Covid-19 Pada Media Sosial Twitter." *Prisma* 4: 628–34.
- [19] Wandani, Aprilia, Fauziah, and Andrianingsih. 2021. "Sentimen Analisis Pengguna Twitter Pada Event Flash Sale Menggunakan Algoritma K-NN, Random Forest, Dan Naive Bayes." *Jurnal Sains Komputer & Informatika (J-SAKTI* 5(2): 651–65.