

# Sentiment Analysis of Rohingya Refugees in Aceh using Support Vector Machine (SVM) and Multinomial Logistic Regression

<sup>1</sup>Gigih Army Buana Baliputra, <sup>2</sup>Slamet Kacung, <sup>3</sup>Budi Santoso  
<sup>1,2,3</sup>Informatics Department, Universitas Dr. Soetomo, Surabaya, Indonesia  
\*e-mail: [gigiharmy21@gmail.com](mailto:gigiharmy21@gmail.com)

(received: 14 March 2025, revised: 31 March 2025, accepted: 31 March 2025)

## Abstract

The rapid development of information technology affects the massive dissemination of information. Social media is one of them, and it contributes to communication and information technology. Information about Rohingya ethnic refugees in Aceh has spread widely on social media. This research aims to analyze public sentiment regarding ethnic Rohingya refugees in Aceh on X and YouTube, categorized into positive, neutral, and negative. This study aims to develop an application that uses the Support Vector Machine (SVM) and Multinomial Logistic Regression techniques to conduct sentiment analysis on public opinion with positive, neutral, and negative classifications regarding Rohingya refugees in Aceh. The 3683 comments collected through web crawling were categorized into positive, negative, and neutral sentiments. The analysis results show that 2112 data were classified as negative sentiments, 1400 as neutral sentiments, and 171 as positive sentiments. Based on the test results, the SVM and Multinomial Logistic Regression methods have similar accuracy of 83.18%. However the SVM method obtained 74.65% precision and 65.15% recall. Meanwhile, the Multinomial Logistic Regression method obtained 75.28% precision and 66.84% recall.

**Keywords :** rohingya, sentiment analysis, social media, support vector machine, multinomial logistic regression

## 1 Introduction

The rapid development of information technology is an important part of life in increasing the efficiency and flexibility of accessing the information needed [1]. The massive dissemination of information occurs on social media. People can easily obtain information from social media [2] especially the Rohingya ethnic issue in Aceh. The problem of Rohingya ethnic refugees originated from those who could not stand racist treatment by their government, the country of Myanmar [3]. The Ministry of Foreign Affairs of the Republic of Indonesia confirmed that the Aceh region and UNHCR (United Nations High Commissioner for Refugees) should be a temporary shelter for Rohingya refugees and other primary needs [3].

The influx of Rohingya ethnic refugees seeking asylum in Southeast Asia has sparked controversy among Indonesians, particularly the Acehese community [4]. Based on UNHCR data, Rohingya refugees are seeking safety not only in Indonesia but also in other countries, with more than 960,000 in Bangladesh, over 107,000 in Malaysia, and more than 22,000 in India. [4]. In December 2023, 1,445 Rohingya refugees landed on several beaches in Aceh. Most of the refugees consist of women and children more than 70%, as well as other vulnerable groups such as people with disabilities and the elderly [4]. Reports indicate tensions between Rohingya refugees and local Acehese stemming from cultural differences and different social expectations [5]. Therefore, sentiment analysis becomes essential to classify public responses objectively and understand the prevailing opinions regarding Rohingya refugees in Aceh.

Several studies have analyzed public sentiment on the issue of Rohingya refugees using sentiment analysis. Research [6] applied sentiment analysis on social media X, which utilized the Support Vector Machine (SVM) and Naive Bayes methods; the accuracy obtained by the SVM method was 76%, and the Naive Bayes method obtained 70%. The next research [7] also utilized the SVM method on 2000 YouTube social media comments, which resulted in an accuracy of 77%. Despite existing research, there remains a gap in analyzing public sentiment using multiple social

<http://sistemasi.ftik.unisi.ac.id>

media sources to enhance the reliability of sentiment classification. Furthermore, limited studies have compared different classification methods, particularly Support Vector Machine (SVM) and Multinomial Logistic Regression, in the context of sentiment analysis on Rohingya refugees in Aceh.

This research aims to create an application that can perform sentiment analysis on public opinion with positive, neutral, and negative classifications regarding Rohingya refugees in Aceh using the Support Vector Machine (SVM) and Multinomial Logistic Regression methods. This research utilizes the X Social Media dataset with the hashtag parameter #rohingyaaceh and the keyword "rohingya aceh". Moreover, the dataset was obtained from YouTube through the KOMPASTV channel. The novelty of this research is the utilization of two social media as a source of dataset and the use of a Multinomial Logistic Regression method as a comparison method. The results of this study can provide insights to UNHCR Indonesia regarding public sentiment in response to Rohingya refugees in Aceh so that it can be used as a basis for making decision. In addition, this study helps to understand the public response to Rohingya ethnic refugees in Aceh. This study plays a role in advancing sentiment analysis. and has a social impact on understanding public perception of humanitarian issues.

## **2 Literature Review**

Research using sentiment analysis has grown rapidly and is used in several fields. Previous research [6] utilized a dataset of Rohingya refugees in X by comparing the Support Vector Machine (SVM) and Naive Bayes for positive and negative sentiment classification. Both methods use SMOTE optimization, which results in an accuracy of 76% for SVM, while Naive Bayes is 70%. Researchers [6] said that the SVM model is better at predicting sentiment and has a low error rate than Naive Bayes model. Based on the information of researcher [6], the SVM model tends to be more accurate for negative classes than for positive classes.

Research [7] conducted sentiment analysis using the SVM method that utilizes YouTube social media to collect a dataset of 2000 comments. Split data is divided into three different training and testing data ratio scenarios: 70:30, 80:20, and 90:10. The highest accuracy result was obtained through the 90:10 scenario of 77%. However, the SVM model could not identify positive classes due to the imbalance of classes in the dataset. Research [4] significantly contributed to understanding public sentiment on the issue of Rohingya migration to Indonesia on social media using Decision Tree 5.0 and Naive Bayes methods. Research [4] utilized Rapidminer tools for sentiment class prediction and K-Fold Cross Validation evaluation where k is three to measure model performance three times; the Naive Bayes model has an accuracy value of 83% and an average error rate of 17%. Meanwhile, the Decision Tree 5.0 model obtained an accuracy value of 78% with an average error rate of 22%.

Based on these existing studies, SVM has been widely used and demonstrates strong performance in sentiment classification, particularly in identifying negative sentiments. Although SVM obtains high accuracy, it has limitations in recognizing positive classes due to data imbalance. Therefore, this study chose SVM for its ability to handle data with a clear hyperplane between sentiment classes [8]. Moreover, Multinomial Logistic Regression was selected to fill the gap in the literature for its simplicity and ability to provide clear interpretations of each feature's contribution to determining sentiment classes, making it a novel aspect of this study. In addition, this study utilizes two dataset sources from social media X and YouTube, in contrast to previous studies that only used one social media platform. Using two data sources, the results of this study can provide a more comprehensive analysis of public sentiment towards the Rohingya issue, and evaluate the performance of SVM and Multinomial Logistic Regression in handling more diverse datasets.

## **3 Research Method**

This study has several stages to analyze the sentiment towards Rohingya refugees in Aceh using SVM and Multinomial Logistic Regression methods. The sequence of stages in this research is illustrated in Figure (1).

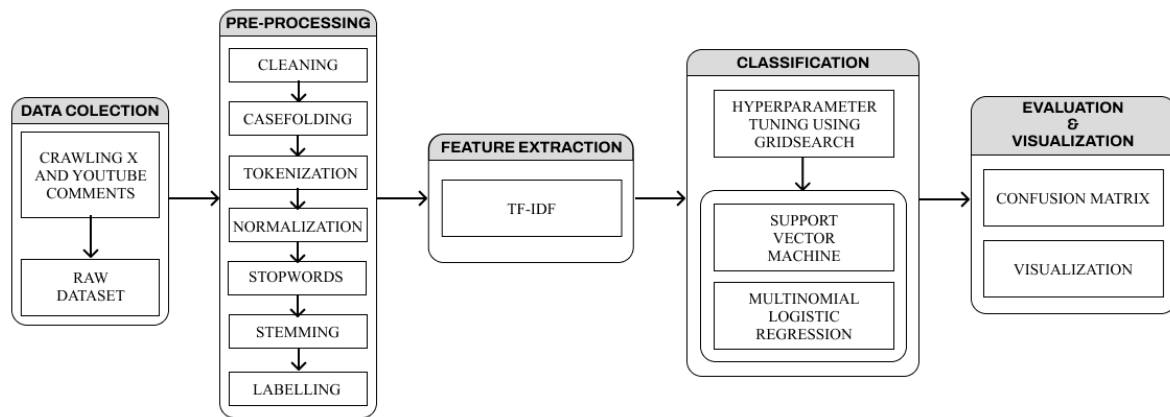


Figure 1. Research procedures

## Data Collection

The dataset was collected from X social media and YouTube, covering public opinion on Rohingya ethnic refugees in Aceh between November 2023 and May 2024. The X dataset collection stage utilized the keyword “rohingya aceh” and the hashtag #rohingyaaceh. The collection of YouTube datasets is taken from the KOMPASTV channel. Data retrieval in this system uses a crawling method that utilizes API tokens on social media X and YouTube. The successfully retrieved data is stored as a raw dataset in a CSV file format.

## Data Pre-processing

The raw dataset that has been stored needs to go through a data pre-processing stage because it contains a lot of noise and is still unstructured text [9]. This stage converts raw data into clean data by cleaning and homogenizing the data so that the data is ready for use [10]. The following are the stages of data pre-processing :

- Data Cleaning* : This process removes unnecessary elements, duplicate data and empty data [11].
- Case Folding* : Converts all letters in the text to lowercase to reduce data complexity [11].
- Tokenization* : Tokenization is the stage of breaking down text data into words or tokens [11].
- Normalization* : Normalization is the stage of converting words into formal forms based on dictionaries that have the same meaning [11].
- Stopwords* : The process of removing common words that provide meaningless information [11].
- Stemming* : Stemming is a pre-processing technique that reduces words to their basic form by cutting off affixes [11].
- Labeling* : Data labeling is part of the sentiment analysis process that labels text data based on the sentiment score calculation [12]. This study applies IndoBERT for automatic labeling. IndoBERT is a variant of the pre-trained BERT model specifically developed using the Indonesian language corpus [13].

## Feature Extration

The feature extraction stage utilizes one commonly used word weighting method, specifically the TF-IDF (Term Frequency-Inverse Document Frequency) method. TF-IDF is the process of assigning a value to each word in a document that aims to show the importance of a word in a particular context [14]. TF (Term Frequency) calculates how often a word appears in a document formulated in Equation (1). Equation (2) finds IDF (Inverse Document Frequency) by calculating the frequency of occurrence of terms in the corpus as a whole [14]. The TF-IDF formula can be seen in Equation (3) [15].

$$tf(t,d) = \frac{f(t,d)}{N(d)} \quad (1)$$

$$idf(t) = \log \frac{D}{D_t} \quad (2)$$

$$tfidf(t,d) = tf(t,d) * idf(t) \quad (3)$$

### Handling Imbalance Data using SMOTE

Class imbalance in a dataset occurs when the number of samples in a particular class exceeds other classes [16]. Imbalance data causes the machine learning models trained to focus on the majority class [16]. The Synthetic Minority Over-sampling Technique (SMOTE) method is suitable for overcoming class imbalance. SMOTE is intended to enhance the representation of the minority class by generating new synthetic samples [16]. Synthetic samples are generated by interpolating between minority data points [16].

### Support Vector Machine

Support Vector Machine (SVM) is a learning algorithm that uses a decision function to separate two classes and create a decision boundary (hyperplane) [17]. The decision function is chosen to maximize the distance to the training data on both sides of the surface. SVM represents the training data set in a multidimensional space, symbolising each data sample as a feature vector [18]. The separation of two classes with an optimal hyperplane margin is used to find the maximum point [17]. The data is mapped into a dimensional space, called the feature space, through a kernel transformation function. Figure (2) shows an illustration of the SVM hyperplane.

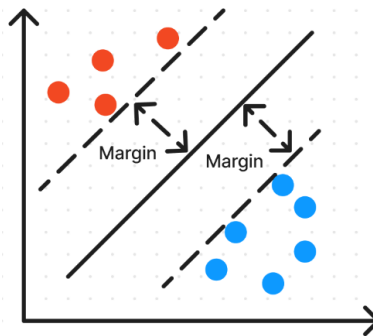


Figure 2. Hyperplane illustration

The calculation to find the hyperplane requires weights as a vector that determines the orientation of the hyperplane combined with the feature vector [18]. There is a bias in adjusting the decision boundary from the origin. The hyperplane equation in linear SVM can be seen in Equation (4) [18].

$$f(x) = wx + b \quad (4)$$

This research utilizes three kernels for sentiment classification, there are linear, polynomial, and Radial Basis Function (RBF) kernels. The following is a description of the three SVM kernels and their equations [19]:

- a) Linear Kernel : Linear measures the dot product of two vectors inserted into the original space without converting them to a higher feature space

$$K(x,y) = (x,y) \quad (5)$$

- b) Polynomial Kernel : This kernel measures the relationship between two input vectors in the original space.

$$K(x,y) = (x,y + c)^d \quad (6)$$

- c) Radial Basis Function (RBF) Kernel : RBF measures the similarity between two input vectors in feature space. There is a parameter  $\gamma$  that indicates the slope at which the output changes based on changes in the input.

$$K(x, x') = \exp(-\gamma(x, y)^2) \quad (7)$$

### Multinomial Logistic Regression

Multinomial Logistic Regression is a machine learning algorithm that can model dependent variables in more than two categories against independent variables [20]. Suppose the dependent variable  $Y$  has three categories coded as 0, 1, and 2. In the binary logistic regression model, the logit function comparing  $Y=1$  against  $Y=0$  acts as the dependent parameter. Meanwhile, in the multinomial logistic regression model, the logit function with  $Y=0$  as reference is used to compare  $Y=1$  and  $Y=2$ , so this model has two logit functions [21]. The formula of the multinomial logistic regression equation in the logit function of categories one and two with category zero can be seen in Equations (8) and (9) [21].

$$g_1(x) = \ln \left[ \frac{P(Y=1|x)}{P(Y=0|x)} \right] = \beta_{10} + \beta_{11}x_1 + \beta_{12}x_2 + \dots \beta_{1p}x_p \quad (8)$$

$$g_2(x) = \ln \left[ \frac{P(Y=2|x)}{P(Y=0|x)} \right] = \beta_{20} + \beta_{21}x_1 + \beta_{22}x_2 + \dots \beta_{2p}x_p \quad (9)$$

The form of the conditional probability equation for each category in Equations (10), (11), and (12)[21].

$$P(Y = 0|x) = \frac{1}{1 + e^{g_1(x)} + e^{g_2(x)}} \quad (10)$$

$$P(Y = 1|x) = \frac{e^{g_1(x)}}{1 + e^{g_1(x)} + e^{g_2(x)}} \quad (11)$$

$$P(Y = 2|x) = \frac{e^{g_2(x)}}{1 + e^{g_1(x)} + e^{g_2(x)}} \quad (12)$$

### Hyperparameter Tuning using Grid Search Cross Validation (GridSearchCV)

Machine learning algorithms have several parameter values that are possible to improve the accuracy of model performance. One method that can help improve model accuracy is the GridSearchCV method [22]. GridSearchCV is able to determine the best hyperparameters in machine learning models so as to produce the most optimal performance [22]. This study utilizes GridSearchCV with the help of pipeline library on SVM and Multinomial Logistic Regression classification models. The hyperparameters to be used from both classification models can be seen in Tables (1) and (2).

**Table 1. SVM hyperparameter tuning**

Kernel	Parameter	Value
Linear	C	0.01, 0.1, 1, 10, 100, 1000
Polynomial	C	0.01, 0.1, 1, 10, 100, 1000
	gamma degree	1, 0.1, 0.01, 0.001, 0.0001 1,2,3,4
RBF	C	0.01, 0.1, 1, 10, 100, 1000
	gamma	1, 0.1, 0.01, 0.001, 0.0001

**Table 2. Multinomial logistic regression hyperparameter tuning**

Parameter	Value
multi_class	'multinomial'
C	0, 0.001, 0.01, 0.1, 1, 10, 100, 1000
solver	'lbfgs', 'saga', 'newton-cg'
penalty	'l1', 'l2', 'elasticnet', 'none'
max_iter	100, 200, 500, 1000

## Confusion Matrix

The final step of this research evaluates the performance of the model using the Confusion Matrix method. This method has evaluation calculation metrics such as accuracy, precision, and recall. In Equation (13), the accuracy metric measures the proportion of correct predictions, both positive and negative, relative to the entire dataset [23]. The precision matrix determines the proportion of accurately predicted positive cases relative to the total predicted positives [23] in Equation (14). The recall matrix in Equation (15) compares the number of correctly predicted positive events to the total number of events that should be positive [23].

$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN} \quad (13)$$

$$Precision = \frac{TP}{TP+FP} \quad (14)$$

$$Recall = \frac{TP}{TP+FN} \quad (15)$$

Based on the evaluation matrix equation above, the basic components in the confusion matrix are as follows [24]:

- TP (True Positive)* is the condition when the predicted data is positive which matches the actual data.
- FP (False Positive)* indicates that the predicted data is positive, but the actual data shows negative.
- TN (True Negative)* occurs when the predicted results with the actual data are exactly negative.
- FN (False Negative)* occurs when the predicted data is negative, but the actual data shows positive.

## 4 Results and Analysis

This section focuses on presenting the results of sentiment analysis on Rohingya refugees using SVM and Multinomial Logistic Regression methods. The results of this study were obtained through the stages of the research method described in the previous chapter. The evaluation results of the SVM and Multinomial Logistic Regression methods are visualized using heatmaps.

### Data Crawling

The raw dataset retrieved using the crawling technique contained as many as 3979 comments, 1157 of which were from the KOMPASTV YouTube channel, and 2822 tweets from social media X using the keyword “rohingya aceh” and hashtag #rohingyaaceh. The Indonesian-language raw datasets from the two social media were merged and given a dataset source category, which can be seen in Table (3).

**Table 3. Raw dataset**

No	Text	Dataset
1	<i>Perlu dijelaskan seberapa banyak dan sering pengungsi Rohingya mendarat di Aceh agar tak dibesar-besarkan demi alasan politik.</i>	X
2	<i>@sintingbuku Jujur kalo mereka ngusir dgn cara seperti itu malah takutnya jadi bumerang buat warga Aceh itu sendiri. Nanti malah jadi alasan buat sirewel dan kroco2nya buat ngeperangin rakyat Aceh dgn alasan tuh liat warga Aceh kejam sama Rohingya ayo balas balikkkk</i>	X
3	<i>Sedih tapi setuju. Di negara lain pun mereka jadi masalah.</i>	YouTube
4	<i>Kamu saja yg menampung !!!hampir semua negara yg seiman menolak !!! Kebanyakan ngga tahu diri ...&lt;br&gt;Di indonesia banyak yg lebih kasihan &lt;br&gt;Apa yg di Papua sana ngga kasihan karena ngga seiman??</i>	YouTube



## Data Pre-Processing

The collected raw datasets need to go through data pre-processing to remove irrelevant data rows and characters. The total data becomes 3683 comments after removing duplicate data. Furthermore, the raw data is converted to lowercase and then cleaned up unnecessary elements, tokenization, filtering, and stemming. The results of the pre-processing stage are in Table (4).

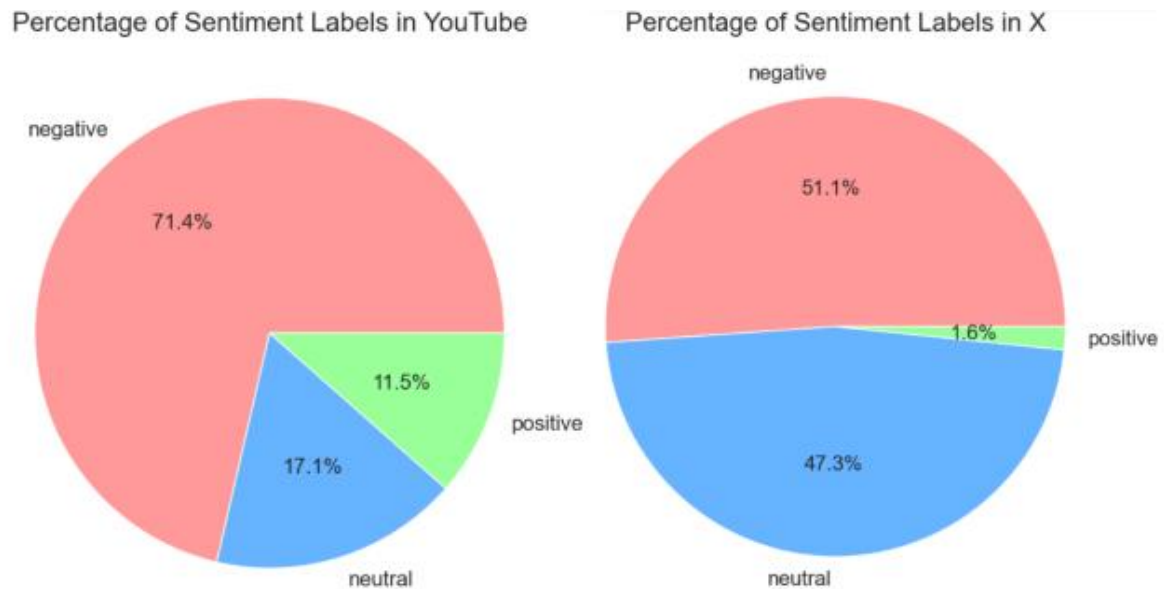
**Table 4. Data pre-processing results**

Pre-Processing Stage	Result
Raw Dataset	@sintingbuku Jujur kalo mereka ngusir dgn cara seperti itu malah takutnya jadi bumerang buat warga Aceh itu sendiri. Nanti malah jadi alasan buat sirewel dan kroco2nya buat ngeperangin rakyat Aceh dgn alasan tuh liat warga Aceh kejam sama Rohingya ayo balas balikkkk
Case Folding	@sintingbuku jujur kalo mereka ngusir dgn cara seperti itu malah takutnya jadi bumerang buat warga aceh itu sendiri. nanti malah jadi alasan buat sirewel dan kroco2nya buat ngeperangin rakyat aceh dgn alasan tuh liat warga aceh kejam sama rohingya ayo balas balikkkk
Data Cleaning	jujur kalo mereka ngusir dgn cara seperti itu malah takutnya jadi bumerang buat warga aceh itu sendiri nanti malah jadi alasan buat sirewel dan kroconya buat ngeperangin rakyat aceh dgn alasan tuh liat warga aceh kejam sama rohingya ayo balas balikkkk
Tokenization	['jujur', 'kalo', 'mereka', 'ngusir', 'dgn', 'cara', 'seperti', 'itu', 'malah', 'takutnya', 'jadi', 'bumerang', 'buat', 'warga', 'aceh', 'itu', 'sendiri', 'nanti', 'malah', 'jadi', 'alasan', 'buat', 'sirewel', 'dan', 'kroconya', 'buat', 'ngeperangin', 'rakyat', 'aceh', 'dgn', 'alasan', 'tuh', 'liat', 'warga', 'aceh', 'kejam', 'sama', 'rohingya', 'ayo', 'balas', 'balikkkk']
Normalization	['jujur', 'kalau', 'mereka', 'mengusir', 'dengan', 'cara', 'seperti', 'itu', 'malah', 'takutnya', 'jadi', 'bumerang', 'buat', 'warga', 'aceh', 'itu', 'sendiri', 'nanti', 'malah', 'jadi', 'alasan', 'buat', 'sirewel', 'dan', 'kroconya', 'buat', 'ngeperangin', 'rakyat', 'aceh', 'dengan', 'alasan', 'itu', 'lihat', 'warga', 'aceh', 'kejam', 'sama', 'rohingya', 'ayo', 'balas', 'balikkkk']
Stopwords Removal	['jujur', 'mengusir', 'takutnya', 'bumerang', 'warga', 'aceh', 'alasan', 'sirewel', 'kroconya', 'ngeperangin', 'rakyat', 'aceh', 'alasan', 'lihat', 'warga', 'aceh', 'kejam', 'rohingya', 'ayo', 'balas', 'balikkkk']
Stemming	['jujur', 'usir', 'takut', 'bumerang', 'warga', 'aceh', 'alas', 'sirewel', 'kroco', 'ngeperangin', 'rakyat', 'aceh', 'alas', 'lihat', 'warga', 'aceh', 'kejam', 'rohingya', 'ayo', 'balas', 'balikkkk']

Raw datasets that have passed all the pre-processing stages above produce clean datasets. Furthermore, the clean dataset goes through a labeling process into three classes: positive, neutral, and negative. Table (5) shows the results of automatic labeling using the IndoBERT model, where the model has the confidence to determine how confident the model is with its predictions.

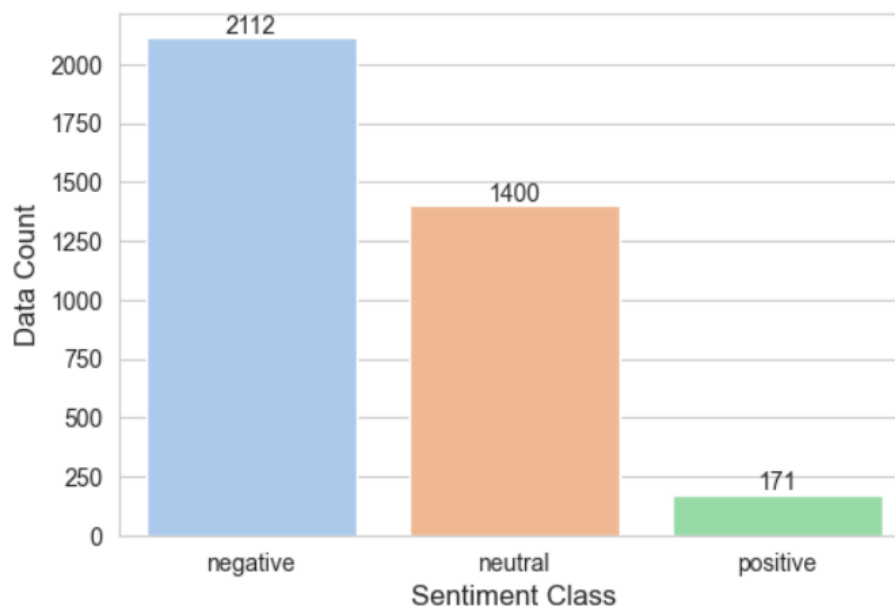
**Table 5. Labeling results using indobert**

No	Clean Text	Sentiment	Confidence
1	jujur usir takut bumerang warga aceh alas sirewel kroco ngeperangin rakyat aceh alas lihat warga aceh kejam rohingya ayo balas balikkkk	Negative	0.9981
2	ungsi rohingya darat aceh dibesarbesarkan alas politik	Neutral	0.8494
3	mending tamu bikin onar kriminal tangkap teman bantu nama	Negative	0.9941



**Figure 3. Visualization of sentiment proportion**

Based on the clean dataset that has been labeled with the model by IndoBERT, Figure (3) shows that the negative class is more than the other classes, where 51.1% on X and 71.4% on YouTube. This shows that the majority of people have negative opinions about Rohingya refugees in Aceh. In Figure (4), the total number of negative class data is 2112 out of 3683; there is an extreme data imbalance in the positive class, which only amounts to 171.



**Figure 4. Visualization of sentiment distribution**

### TF-IDF Weighting

After the sentiment class labeling is complete, the Term Frequency-Inverse Document Frequency (TF-IDF) calculation process is used on the clean dataset. TF-IDF method is used to calculate the weight value of certain words in the document. This method helps to improve accuracy in text analysis. Table (6) shows the average results of word weighting using the TF-IDF method.



**Table 6. TF-IDF weighting result**

Word	TF-IDF (average)
<i>ungsi</i>	0.02498
<i>selundup</i>	0.00237
<i>rohingya</i>	0.03756
<i>warga</i>	0.01703
<i>aceh</i>	0.03908
<i>merdeka</i>	0.00113
<i>negara</i>	0.00994
<i>tolak</i>	0.01453

### SVM and Multinomial Logistic Regression Modeling

Before modeling, the data is divided into two parts: 80% training data and 20% testing data. The SMOTE technique is used in the pipeline to improve model performance on unbalanced data. To find the best parameters, SVM and Multinomial Logistic Regression models were evaluated using Grid Search with cv=5 (cross-validation). Table (7) shows the hyperparameter tuning results of both models.

**Table 7. Model hyperparameter tuning results**

Classification Method	Parameter	Best Parameter Value	Best Score
SVM	Linear	C	10
	Polynomial	C	0.01
		gamma	0.01
		degree	1
	RBF	C	100
Multinomial Logistic Regression		gamma	0.01
		max_iter	500
		C	10
		penalty	L2
		solver	lbfgs

The parameter tuning results in Table (7) show that in the SVM model, the RBF kernel achieved a score of 80.41%, superior to the linear kernel of 80.14% and the polynomial kernel of 68.06%. Meanwhile, the Multinomial Logistic Regression model produced a score of 80.65%. Based on these results, the Multinomial Logistic Regression model is slightly superior to the SVM kernel RBF model.

### Evaluation

The final stage of this study involves evaluating the performance of the Support Vector Machine (SVM) and Multinomial Logistic Regression models. This evaluation allows a detailed assessment of the classification capabilities of the models through metrics such as accuracy, precision, and recall. The evaluation results of the SVM and Multinomial Logistic Regression methods can be seen in Tables (8) and (9).

**Table 8. SVM evaluation result**

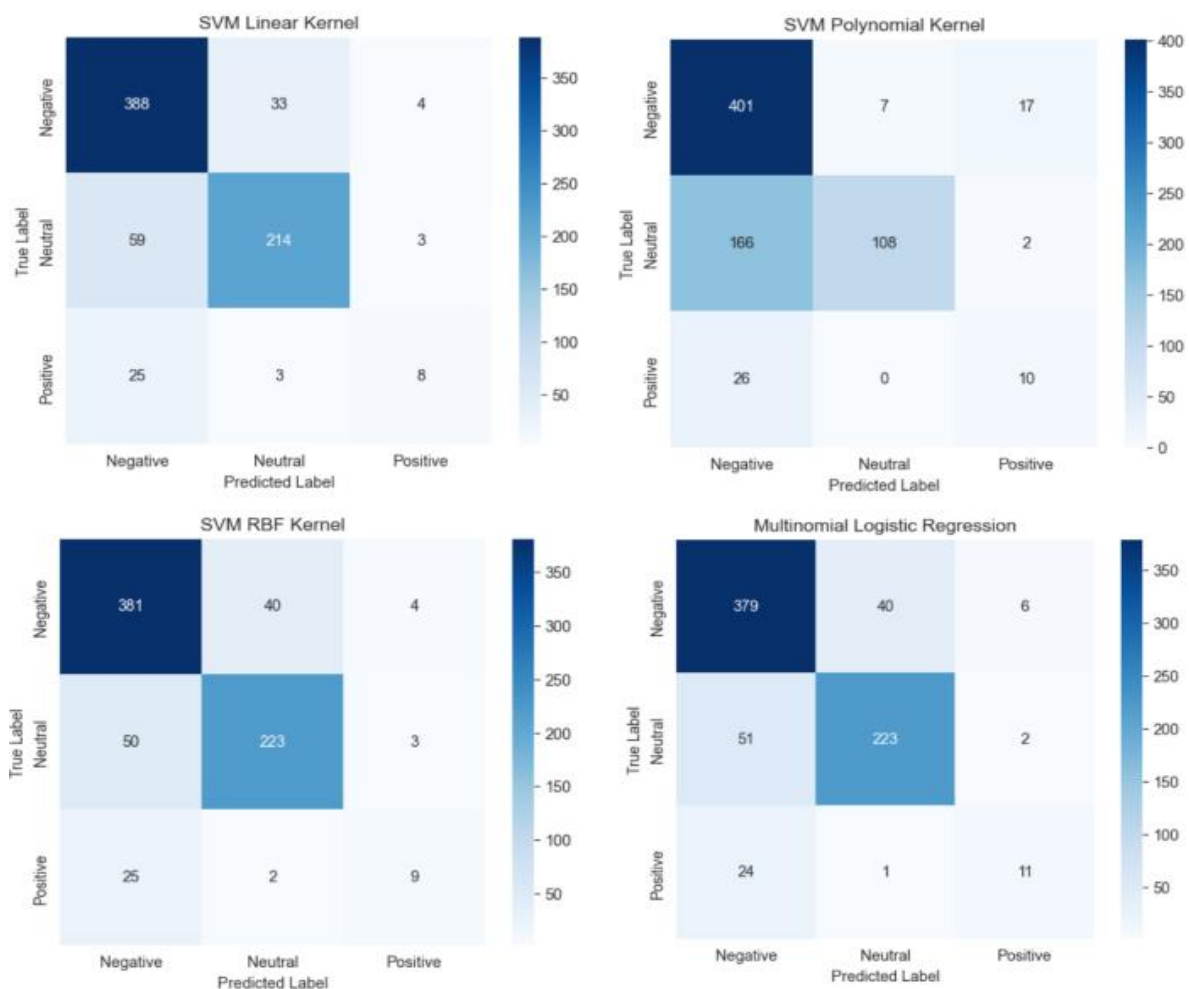
Matrix	Kernel		
	Linear	Polynomial	RBF
Accuracy	82.77%	70.42%	83.18%

Matrix	Kernel		
	Linear	Polynomial	RBF
Precision	73.71%	65.34%	74.65%
Recall	63.68%	53.75%	65.15%

**Table 9. Multinomial logistic regression evaluation result**

Matrix	Score
Accuracy	83.18%
Precision	75.28%
Recall	66.84%

In Table (8), the SVM kernel RBF model test results have the best performance, with the highest accuracy of 83.18%. This indicates that the model has the ability to identify most samples better than the linear and polynomial kernels. In Table (9), the Multinomial Logistic Regression model test results have the same accuracy as the SVM kernel RBF model in Table (8). The precision and recall values of the Multinomial Logistic Regression model are slightly superior to the SVM kernel RBF model. The confusion matrix visualization can be used to measure the effectiveness of the model in predicting the testing data. Multinomial Logistic Regression is slightly better at predicting the positive class which can be seen in Figure (5).



**Figure 5. Confusion matrix visualization**

## 5 Conclusion

Based on the results of research that has been conducted, 3683 datasets of public sentiment on the issue of Rohingya refugees in Aceh on social media X and YouTube, It is discovered that the majority of people have a poor opinion of Rohingya refugees in Aceh, classified by 2112 negative comments. In contrast, 1400 neutral comments and 171 positive comments were found. The same accuracy is obtained based on the test results using SVM and Multinomial Logistic Regression methods with the help of SMOTE and GridSearchCV techniques. The SVM kernel RBF method and Multinomial Logistic Regression obtained 83.18% accuracy. Although both methods have the same accuracy, Multinomial Logistic Regression has a precision value of 75.28% and recall of 66.84%, slightly superior to SVM kernel RBF, which has a precision value of 74.65% and recall of 65.15%. This study has a drawback where the model has difficulty classifying positive classes due to extreme data imbalance despite applying the SMOTE oversampling technique.

## Reference

- [1] O. Manullang, C. Prianto, and N. E. Harani, "Analisis Sentimen untuk memprediksi Hasil Calon Pemilu Presiden menggunakan Lexicon based dan Random Forest," *J. Ilm. ...*, no. 54, pp. 1–11, 2023, [Online]. Available: <https://forum.upbatam.ac.id/index.php/jif/article/view/7987%0Ahttps://forum.upbatam.ac.id/index.php/jif/article/download/7987/3319>
- [2] A. N. Ramadhan, E. Utami, and A. D. Hartanto, "Analisis Sentimen Opini Publik menggunakan Metode BiLSTM pada Media Sosial Twitter," *Semiot. (Seminar Nas. Teknol. Inf. dan Mat.*, vol. 2, no. 1, pp. 1–12, 2023.
- [3] Y. Azhari and Wilopo, "Pencegahan Potensi Konflik antara Pengungsi Rohingya dan Masyarakat Lokal Indonesia," *JPM (Jurnal Pengabd. Mandiri*, vol. 1, no. 3, pp. 475–488, 2022.
- [4] U. Kurniasih and A. T. Suseno, "Analisis Sentimen Masyarakat terhadap Isu Migrasi Rohingya ke Indonesia," *J. Teknol. Dan Sist. Inf. Bisnis*, vol. 7, no. 1, pp. 199–207, 2025, doi: <https://doi.org/10.47233/jteksis.v7i1.1815>.
- [5] A. R. Usman, A. Sulaiman, M. Muslim, and T. Zulyadi, "Conflict and Cultural Adaptation of the Aceh Rohingya Refugees in Media Opinion," *Profetik J. Komun.*, vol. 16, no. 1, pp. 107–122, 2023, doi: 10.14421/pjk.v16i1.2491.
- [6] D. Ananda and R. R. Suryono, "Analisis Sentimen Publik terhadap Pengungsi Rohingya di Indonesia dengan Metode Support Vector Machine dan Naïve Bayes," *J. Media Inform. Budidarma*, vol. 8, no. 2, pp. 748–757, 2024, doi: 10.30865/mib.v8i2.7517.
- [7] H. Hidayat, F. Santoso, and L. F. Lidimillah, "Analisis Sentimen Pengguna YouTube tentang Rohingya menggunakan Algoritma SVM (Support Vector Machine)," *G-Tech J. Teknol. Terap.*, vol. 8, no. 3, pp. 1729–1738, 2024, doi: 10.33379/gtech.v8i3.4497.
- [8] A. Novantika and Sugiman, "Analisis Sentimen Ulasan Pengguna Aplikasi Video Conference Google Meet menggunakan Metode SVM dan Logistic Regression," *Prism. Pros. Semin. Nas. Mat.*, vol. 5, pp. 808–813, 2022, [Online]. Available: <https://journal.unnes.ac.id/sju/index.php/prisma/>
- [9] N. K. Putri, A. V. Vitianingsih, S. Kacung, A. L. Maukar, and V. Yasin, "Sentiment Analysis of Brand Ambassador Influence on Product Buyer Interest using KNN and SVM," *Indones. J. Artif. Intell. Data Min.*, vol. 7, no. 2, pp. 327–336, 2024, doi: 10.24014/ijaidm.v7i2.29469.
- [10] T. S. Sabrila, Y. Azhar, and C. S. K. Aditya, "Analisis Sentimen Tweet tentang UU Cipta Kerja menggunakan Algoritma SVM berbasis PSO," *JISKA (Jurnal Inform. Sunan Kalijaga)*, vol. 7, no. 1, pp. 10–19, 2022, doi: 10.14421/jiska.2022.7.1.10-19.
- [11] K. H. Yuniur, A. V. Vitianingsih, S. Kacung, A. Lidya Maukar, and A. Dwi Arumsari, "Sentiment Analysis of Cyberbullying Detection on Social Networks using the Sentistrength Method," *Sistemasi*, vol. 13, no. 4, pp. 1587–1596, 2024, doi: 10.32520/stmsi.v13i4.4226.
- [12] R. I. Putra Selian, A. V. Vitianingsih, S. Kacung, A. Lidya Maukar, and J. Febrian Rusdi, "Sentiment Analysis of Public Responses on Social Media to Satire Joke using Naive Bayes and KNN," *Sinkron*, vol. 8, no. 3, pp. 1443–1451, 2024, doi: 10.33395/sinkron.v8i3.13721.
- [13] P. Sayarizki, Hasmawati, and H. Nurrahmi, "Implementation of IndoBERT for Sentiment

<http://sistemasi.ftik.unisi.ac.id>

- Analysis of Indonesian Presidential Candidates,” *Indones. J. Comput.*, vol. 9, no. 2, pp. 61–72, 2024, doi: 10.34818/indojc.2024.9.2.934.
- [14] Y. Nurdiansyah, F. Rahman, P. Pandunata, and A. Infantono, “Analisis Sentimen Opini Publik terhadap Undang-Undang Cipta Kerja pada Twitter menggunakan Metode Naive Bayes Classifier,” *Pros. Semin. Nas. Sains Teknol. dan Inov. Indones.*, vol. 3, no. November, pp. 201–212, 2021, doi: 10.54706/senastindo.v3.2021.158.
- [15] F. Rizal, A. Wijaya, and F. Hasyim, “Analisis Sentimen Masyarakat Indonesia terhadap Aplikasi TikTok menggunakan Algoritma Logistic Regression,” *AKIRATECH J. Comput. Electr. Eng.*, vol. 1, no. 2, pp. 57–65, 2024, [Online]. Available: <https://journal.ajbnews.com/index.php/akiratech>
- [16] M. P. Pulungan, A. Purnomo, and A. Kurniasih, “Penerapan SMOTE untuk mengatasi Imbalance Class dalam Klasifikasi Kepribadian MBTI menggunakan Naive Bayes Classifier,” *J. Teknol. Inf. dan Ilmu Komput.*, vol. 11, no. 5, p. 1033=1042, 2024, doi: 10.25126/jtiik.1077989.
- [17] N. S. Ramadan and D. Darwis, “Perbandingan Metode Naïve Bayes dan SVM untuk Sentimen Analisis Masyarakat terhadap Serangan Ransomware pada Data KIP-K,” *J. Sist. Inf. dan Inform.*, vol. 8, no. 1, pp. 12–23, 2025, doi: <https://doi.org/10.47080/simika.v8i1.3621>.
- [18] R. Fajriah and D. Kurniawan, “Optimalisasi Model Klasifikasi Naive Bayes dan Support Vector Machine dengan Fast Text dan Chi Square,” *Fakt. Exacta*, vol. 17, no. 4, pp. 334–345, 2025, doi: 10.30998/faktorexacta.v17i4.24751.
- [19] S. Rabbani, D. Safitri, N. Rahmadhani, A. A. F. Sani, and M. K. Anam, “Perbandingan Evaluasi Kernel SVM untuk Klasifikasi Sentimen dalam Analisis Kenaikan Harga BBM,” *MALCOM Indones. J. Mach. Learn. Comput. Sci.*, vol. 3, no. 2, pp. 153–160, 2023, doi: 10.57152/malcom.v3i2.897.
- [20] B. L. Fauzan, T. Agustin, and A. M. H. Mahmudah, “Prediksi Klasifikasi Kecelakaan Lalu Lintas di Kota Surakarta dengan menggunakan Metode Regresi Logistik Multinomial,” *Sustain. Civ. Build. Manag. Eng. J.*, vol. 1, no. 4, pp. 1–9, 2024, doi: 10.47134/scbmej.v1i4.3159.
- [21] S. S. Dewi, R. Resmawan, and L. O. Nashar, “Analisis Regresi Logistik Multinomial dengan Metode Bayes untuk Identifikasi Faktor-Faktor terjadinya Diabetes Melitus,” *J. Math. Theory Appl.*, vol. 5, no. 2, pp. 51–60, 2023, doi: 10.31605/jomta.v5i2.2520.
- [22] Z. M. E. Darmawan and A. F. Dianta, “Implementasi Optimasi Hyperparameter GridSearchCV pada Sistem Prediksi Serangan Jantung menggunakan SVM Implementation of GridSearchCV Hyperparameter Optimization in Heart Attack Prediction System using SVM,” *Unipdu*, vol. 13, no. 1, pp. 8–15, 2023, doi: <https://doi.org/10.26594/teknologi.v13i1.3098>.
- [23] E. Damayanti, A. V. Vitianingsih, S. Kacung, D. Cahyono, and A. L. Maukar, “Sentiment Analysis of Alfagift Application User Reviews using Long Short-Term Memory (LSTM) and Support Vector Machine (SVM) Methods,” *Decod. J. Pendidik. Teknol. Inf.*, vol. 4, no. 2, pp. 509–521, 2024, doi: 10.51454/decode.v4i2.478.
- [24] B. Ayuningrum, H. H. M. Haqq, S. M. P. Lestari, and M. Al Haris, “Analisis Sentimen Survei Regsosek pada Twitter menggunakan Algoritma K-Nearest Neighbor (K-NN),” *ESTIMASI J. Stat. Its Appl.*, vol. 4, no. 2, pp. 196–207, 2023, doi: <https://doi.org/10.20956/ejsa.v4i2.25522>.