

# Journal of Dinda

Data Science, Information Technology, and Data Analytics

Vol. 4 No. 2 (2024) 123 - 134

E-ISSN: 2809-8064

## Post-Election Sentiment Analysis 2024 via Twitter (X) Using the Naïve Bayes Classifier Algorithm

Yessi Mayasari<sup>1\*</sup>, Yusuf Ramadhan Nasution<sup>2</sup>

<sup>1,2</sup>Computer Science, Science and Technology, Universitas Islam Negeri Sumatera Utara

<sup>1\*</sup>yesimayasari.2202@gmail.com, <sup>2</sup>ramadhannst@uinsu.ac.id

### Abstract

After the 2024 election, there have been several events that have attracted attention. One of them is the camp of Ganjar who had sued the election results to the Constitutional Court (MK). Likewise, the couple from the Anies Baswedan camp also did the same thing, namely suing the results of the presidential election. He assessed that there had been fraud in the results of the presidential election and proposed to hold a re-vote. However, the Constitutional Court rejected the lawsuit from the two camps. As an Indonesian society that has the right to have an opinion, of course, these things invite various opinions and comments. Platform X is the main place for this discussion where platform X is a means for the community to respond and provide opinions and comments. This research was conducted using the Naive Bayes Classifier approach to dig deeper into public attitudes on Twitter towards the 2024 presidential election. In this study, the data that will be analyzed is public opinion related to various post-election statements or comments, especially the 2024 presidential election. In this study, data analysis was carried out by collecting tweets using a twitter data crawling process with a tweet harvest library. From the research that has been carried out using the Naive Bayes algorithm using data of 1228 tweet data, then after the processing stage it became 1134 tweet data. With 907 data training and 227 data testing. In testing the model on the training data, the accuracy was 91%. Meanwhile, the model test on the test data obtained an accuracy of 95%. Based on the performance results produced, which is an accuracy value of 95%, this is of course a fairly accurate result

Keywords: *Sentiment Analysis, Naïve Bayes Classifier, Political Elections, Twitter Data Analysis*

© 2024 Journal of DINDA

### 1. Introduction

After the 2024 election, there will be many conflicts between public opinions on social media, one of which is on twitter. As for the post-election, especially the 2024 presidential election, which leaves various public opinions regarding the results of the presidential election. Where the final result of the presidential election invited various opinions and comments from the public, especially the supporters of the elected presidential pair and the supporters of the unelected presidential pair. After the 2024 election, there have been several events that have attracted attention. One of them is the camp of Ganjar who had sued the election results to the Constitutional Court (MK). Likewise, the couple from the Anies Baswedan camp also did the same thing, namely suing the results of the presidential election. He assessed that there had been fraud in the results of the presidential election and proposed to hold

a re-vote. However, the Constitutional Court rejected the lawsuit from the two camps.

Platform X is the main place for this discussion Where platform X is a means for the community to respond and provide opinions and comments. Twitter is a social media founded by Jack Dorsey in 2006. According to a twitter press release, more than 500 million tweets or tweets sent by users every day in 2019 have been used to post and various information about users, as well as content that can express feelings. Twitter is a website that can collect data on the opinions of people around the world[1]. Sentiment analysis, also known as opinion mining, is the process of extracting a person's opinion or opinion on a particular topic from a document.

Based on previous research by Hardi et al., it was concluded that the naïve bayes method is a fairly good method in classifying data mining or text mining. This is because the algorithm produces a fairly high accuracy

value, which is above 50%[3]. The same research was also conducted by Mahbubah et al. from the results of the study, it can be concluded that by using the naïve bayes method with 240 training data and 60 test data, a fairly high accuracy result of 73% [4]. Then the research using the naïve bayes method was also carried out by Kurniawan et al., from the results of the study also produced a very high accuracy with 100 training data and 150 test data obtained an accuracy of 93.35%.

In this study, the data that will be analyzed is public opinion related to various post-election statements or comments, especially the 2024 presidential election. The reason for choosing this topic is because the topic was hotly discussed by the people of Indonesia

The purpose of this research is to find out how to collect twitter data that is relevant to related topics, namely public opinion related to presidential candidates after the 2024 presidential debate. To find out the accuracy level of the Naïve Bayes Classifier method in the analysis of public sentiment. As well as how the Naive Bayes Classifier algorithm for sentiment analysis.

## 2. Research Methods

In this study, the system development method to be used, namely the KDD (Knowledge Discovery in Databases) method, is a general approach in data analysis that aims to extract useful knowledge from the available data.

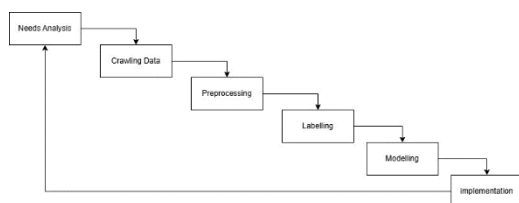


Figure 1. Research stages

In the research there are several stages, the first stage is the needs analysis where at this stage is the analysis of needs. Then in the second stage, namely the data *crawling* stage where this technique is used to collect tweet data on the X (Twitter) application. The third stage is data *preprocessing* where this stage is to prepare and clean raw data. The pre-processing used is case folding, emoji removal, cleaning, repetition character removal, word normalization, negation handling, stop words removal, *stemming*, and tokenization. Then the fourth stage is *labelling*. At this stage is data labelling, where *data labeling* in sentiment analysis refers to the process of labeling text data. These labels indicate the sentiment or opinion contained in the text, such as "positive," "negative," ". Furthermore, the fifth stage is *Modelling*. Where this process applies a sentiment analysis model using text data that has been preprocessed and labeled. Train the model to recognize sentiment by using *the naïve bayes classifier algorithm*. Then the sixth stage is

implementation, which is to implement a model that has been evaluated for use according to research needs.

### 2.1 Text Mining

The field of *text mining* is one of the branches of data *mining* that analyzes textual data. Text mining is a type of automated text analysis performed by a computer to extract high-quality information from unclear text types in documents [5]. Text Mining is widely used in knowledge-based organizations, as it is the process of examining large volumes of documents to find information and help answer specific research questions, identifying remaining facts, relationships, and statements [6]. To be more specific, text mining uses technologies including natural language processing, batch processing, information extraction, and data retrieval to help understand complex text analysis systems. Initially, text mining was used to provide intelligence to governments and security agencies to identify potentially dangerous activities and other security threats. To increase productivity, *text mining* uses text analysis and technology from external disciplines such as computer science, human resources, machine education, and statistics. Correction is essential for conducting practical research, and since then, related fields have used this technique extensively [7]. Text mining uses different methodologies for text processing, one of the most important being Natural Language Processing (NLP). Text classification is one of the most important tasks in natural language processing (NLP). In short, the statistics and summaries available are for the purpose of classifying broken texts [8]. Classification is a modeling process in which elements or data are classified into specific categories or labels based on existing features [9].

### 2.2 Sentiment Analysis

Sentiment analysis, also known as sentiment analysis, is typically used to express a positive or negative analysis of feelings. The method of sentiment analysis based on natural selection is based on measuring the sentiment of each word or phrase in a text using statistical analysis. The main strategy is based on the laws of nature, using a password defined by trial and error to determine its experiment on a specific number, which is an inaccurate method. This method is the most intuitive and is also similar to the process of analyzing human speech in text[10]. Sentiment analysis is the process of measuring the expression of emotions through language, which involves analyzing and interpreting emotional information, analyzing and interpreting data, and classifying the interpretation of written material. Emotion recognition, polarity detection and affective computing as the basis of sentiment analysis. Sentiment Analysis is a key technique in the field of natural language processing for sentiment analysis, which is

widely used in public opinion research, artificial intelligence, and business intelligence. There are three main types of text analysis techniques, namely machine-based text analysis, dictionary-based text analysis, and hybrid analysis [11].

### 2.3 X (Twitter)

In July 2023, Twitter changed to X. This applies not only to the name but also to the in-app features. Elon Musk, the new owner of Twitter, stated that this change is part of a larger Twitter revival. Musk's motive for buying Twitter is to turn it into a super app, similar to China's WeChat. Twitter began to undergo some changes, and eventually the name Twitter The project also involved rebranding by changing the logo old with an updated white X with a black background. The use of social media in life can also increase public participation in various government activities, where accountability is increasingly sensitive in the digital era and it proves that people are aware of negative issues in government policies. In addition, Twitter social media can also be a forum for exchanging information and news and forming opinions. To receive information about phenomena or problems that are being discussed by the public on Twitter[12].

### 2.4 Python

Python is a high-level programming language that supports object-oriented programming. Python has a difference from other programming languages, namely in syntax writing. In the python programming language, there are various libraries and frameworks used to perform data analysis[13]. Python is a high-level interpreted programming language, widely known for being easy to learn, but still able to harness the power of a system-level programming language when needed[14].

### 2.5 Preprocessing

The data preprocessing process is required after generating the target data to make it ready for use. In the next step, the ready-to-use data serves to analyze the data and produce some knowledge or results by applying several mining techniques. The initial processing of data includes five activities such as Data Cleaning, Data Optimization, Data Transformation, Data Integration and Data Conversion [15]. The pre-processing used is case folding, emoji removal, cleaning, repetition character removal, word normalization, negation handling, stopwords removal, stemming, and tokenization. [16]. Case folding is the process of converting data into the appropriate format. The purpose of folding this issue is to reduce the redundancy of information used in the classification process, reduce redundancy so that the calculation process is optimal.

Word normalization is the process of normalizing words, that is. correcting abbreviated words and converting them into words that have the same meaning based on KBBI. Stemming is a process carried out to find the root of a word. The procedure to be carried out is to remove all the corresponding suffixes, which consist of prefixes, suffixes, suffixes and combinations of prefixes, endings of derivative words. In the filtering process, it is the selection of important words and the elimination of unimportant words. The algorithm used in this process is a stop word. A stopword is a collection of words that are not related to the main topic (irrelevant), at this stage it is a stream of input obtained from a text file divided into sections. small portions. For example, breaking down sentences into words (tokens) [17].

### 2.6 Algorithm Naive Bayes Classifier

Naive Bayes is the most popular classification method used with good accuracy. The Naive Bayes classification method is based on simple probability and is designed to be used assuming the independence of the explanatory variable. The purpose of this study is to analyze whether the emotional tone of the message is positive, negative[18]. Naive Bayes is a classification method that is often used in sentiment analysis because it is simple and easy to classify documents. Bayesian theory can generally be interpreted by the equation:

$$P(A|B) = \frac{P(B|A) \times P(A)}{P(B)} \quad (1)$$

Where:

A : Data hypothesis is a specific class.

B : Data with a class that is still unknown.

P(A|B) : Hypothesis probability based on conditions.

P(A) : Probabilitas hypothesis.

P(B|A) : Probability based on the conditions on the hypothesis.

P(B) : Probability B.

This algorithm emphasizes probability learning. The advantage of Naive Bayes' algorithm is that it has a lower error rate when used with large datasets. In addition, Naive Bayes' accuracy and speed are higher when applied to large datasets[19].

### 2.7 Feature Extraction

In this phase, words are extracted in vector form (as long as they retain the meaning of the words), because the computer can only recognize and process numbers, not words. *TF-IDF (Term Frequency - Inverse Document Frequency)* is the process of counting or extracting words into numbers in the form of vectors used to determine the weight of words in a document or corpus. This weight helps determine the importance of a word in the document. Basically, the calculation or formula of TF-IDF is divided into two parts: TF (Term Frequency) and IDF (Inverse Document Frequency), and there are different formulas and working methods, which are used at the end of the calculation between TF and IDF. Calculate as follows:

1. TF (Term Frequency) Counts the frequency of words in a document. On every document has different word frequencies then the TF value will be divided by the Document length:

$$TF = \frac{\text{Number of occurrences of words in the document}}{\text{Number of words on a document}} \quad (2)$$

Information:

TF = frequency of occurrence of words in a document

1. IDF (*Inverse Document Frequency*)

After completing the calculation of the TF value, the next step is to calculate the IDF value which is a value to measure how important a word is, by referring to the smaller the IDF value, the less important the word will be considered, and vice versa.

$$IDF = \log \frac{D + 1}{df + 1} + 1 \quad (3)$$

Information:

D : the frequency of words in D

df : many documents containing search glass

1. TF-IDF

After calculating the number of TF and IDF, the next thing is to calculate the value of TF-IDF by multiplying it:

$$TF - IDF = TF \times IDF \quad (4)$$

### 2.1.7 Confusion Matrix

In the field of data mining, a configuration matrix is used to calculate the accuracy of a method. Basically, a configuration matrix contains information that compares the classification results performed by the system with the classification results that should be. This confusion matrix performs calculations that produce 4 *outputs*, namely *recall*, *precision*, *accuracy*, and *error rate*.

$$\text{recall} = \frac{TP}{FN+TP} \times 100\% \quad (5)$$

$$\text{precision} = \frac{TP}{FP+TP} \times 100\% \quad (6)$$

The Operating Characteristic Receiver curve or ROC curve is often used as a visualization plot to measure the performance of binary classifiers. This curve is not a metric from the model, but rather a graphical representation of the True Positive Rate (TP) versus the False Positive (FP) at various classification thresholds from 0 to 1.

## 3. Results and Discussion

### 3.1 Data collection

In this study, the initial stage is to collect tweet data through a crawling process using the Tweet Harvest library in Python. This library is designed to make it easier to retrieve data from social media, including X. Before utilizing this library, Node.js installation is required to support its operations. The author uses Google Colab as a tool to execute code with the Python programming language. The data crawling process involves installing and importing the necessary data, which then allows for efficient retrieval of tweet data for further analysis.

```
# Import required Python package
!pip install pandas

# Install node.js (because tweet-harvest built using Node.js)
!sudo apt-get update
!sudo apt-get install -y ca-certificates curl gnupg
!curl -fsSL https://deb.nodesource.com/setup_20.x | sudo
!sudo -E bash -c 'curl -fsSL https://deb.nodesource.com/setup_20.x | sudo tee /etc/apt/sources.list.d/nodesource.list
!sudo apt-get update
!sudo apt-get install nodejs -y
node -v
```

Figure 2 Install Library

After successfully installing and importing the required library, what is needed is token authentication on account X. After having the `auth_token`, it is necessary to enter the desired data keywords and enter the desired amount of data limit as needed. After that, the data crawling process will run. Below is the source code for crawling data.

```
# Crawl Data
filename = 'pilpres2024.csv'
search_keyword = 'pilpres2024 lang:id'
limit = 1000

!topx -y tweet-harvest@0.1.0 -f "{filename}" -s "{search_keyword}" --tab "LATEST" -l {limit} --token {twitter_auth_token}
```

Figure 3. Crawl Data On X

After the data crawl process is complete. So we need to save the crawl file using the source code as below.

```

import pandas as pd

# Specify the path to your CSV file
file_path = "tweets-data/(filename)"

# Read the CSV file into a pandas DataFrame
df = pd.read_csv(file_path, delimiter=",")

# Display the DataFrame
display(df)
    
```

Figure 4. Save Crawling Results

### 3.2 Preprocessing

The pre-processing process in this study is carried out systematically and sequentially. The steps taken include various data cleaning techniques, such as the removal of irrelevant characters, text normalization, as well as tokenization. Each step has an important role in preparing the data for further analysis, so that the data that has gone through this stage is ready to be used in effective and efficient classification.

The cleaning stage is a process used to clean text data by deleting irrelevant and inconsistent information. In this stage, characters that are considered unimportant, such as hashtags (#), numbers, usernames (@), URLs, punctuation, and emoticons, are omitted from the dataset. This data cleaning aims to improve the quality and consistency of data, so that the analysis carried out later can provide more accurate and reliable results

Table 1. Cleaning Data

Tweet	Cleaning
apa yang terjadi di MK Seperti Sebuah lagu Kau yang mulai kau yang mengakhiri #sengketapilpres2024 #ilc #mahkamahkonstitusi #pengamatpolitik #debatpublik <a href="https://t.co/F5IhWocLzK">https://t.co/F5IhWocLzK</a>	apa yang terjadi di MK Seperti Sebuah lagu Kau yang mulai kau yang mengakhiri
STOP Drama Omong kosong pada praktek Politik Sampah Dapur Ulang !! #ilc #pengamatpolitik #fyp #sengketapilpres2024 #faizalassegaf @tvonenews Lihat video <a href="https://t.co/bKvKvZ1H97!">https://t.co/bKvKvZ1H97!</a> #TikTok <a href="https://t.co/VZrcdyUJ8o">https://t.co/VZrcdyUJ8o</a>	STOP Drama Omong kosong pada praktek Politik Sampah Dapur Ulang

Pengadilan Tata Usaha Negara/PTUN : Tidak Akan Merubah Keputusan KPU RI @KPU\_ID & Keputusan MKRI @officialMKRI (Putusan Mahkamah Konstitusi bersifat Final & Mengikat) #Prediksi #Gugatan #PTUN #SengketaPilpres2024 #CakRasyid #SalamDamaiSejahtera <https://t.co/vwTzroJNDI>

The case folding stage is the process of equalizing the shape of letters in the data text by changing all uppercase letters to lowercase letters. This step is carried out on the tweet data so that there are no differences in the shape of the letters that can interfere with the sentiment analysis that will be carried out next. By implementing case folding, data consistency can be maintained, so that the analysis results will be more accurate and reliable

Table 2. Case Folding

Tweet	Case Folding
apa yang terjadi di MK Seperti Sebuah lagu Kau yang mulai kau yang mengakhiri	apa yang terjadi di mk seperti sebuah lagu kau yang dimulai kau yang mengakhiri
STOP Drama Omong kosong pada praktek Politik Sampah Dapur Ulang	stop drama bicara tidak berguna pada praktik politik tidak bernilai area kerja kembali
Pengadilan Tata Usaha Negara PTUN Tidak Akan Merubah Keputusan KPU RI	pengadilan tata usaha negara ptun tidak akan merubah keputusan kpu ri

Word normalization is a process that aims to simplify words in a text into a more basic and consistent form.

Table 3. Normalization

Tweet	Normalisasi kata
apa yang terjadi di mk seperti sebuah lagu kau yang dimulai kau yang mengakhiri	"apa yg terjadi di mahkamah konstitusi seperti satu lagu kamu yg mulai kamu yg mengakhiri"
stop drama omong kosong pada praktek politik sampah dapur ulang	"stop drama bicara tidak berguna pada praktik politik tidak bernilai daur kembali"
pengadilan tata usaha negara ptun tidak akan merubah keputusan kpu ri	"pengadilan tata usaha negara tidak akan mengubah keputusan komisi pemilihan umum republik indonesia"

Tokenization is a process that aims to break down a text document into smaller parts, or tokens, which are generally words. At this stage, each word is cut and lowercased, while other characters or punctuation marks are removed from the text.

Table 4. Tokenization

Tweet	Tokenization
apa yang terjadi di mk seperti sebuah lagu kau yang mengakhiri	['apa','yg','terjadi','di','mahkamah','konstitusi','seperti','satu','lagu','kamu','yg','mulai','kamu','yg','mengakhiri']
stop drama omong kosong pada praktek politik sampah dapur ulang	['stop','drama','bicara','tidak','berguna','pada','praktik','politik','tidak','bernilai','daur','kembali']
pengadilan tata usaha ptun tidak akan merubah keputusan kpu ri	['pengadilan','tata','usaha','negara','tidak','akan','mengubah','keputusan','komisi','pemilihan','umum','republik','indonesia']

Stopword removal is a stage in text processing that aims to remove words that are considered irrelevant or do not contribute significantly to the analysis of the topic of the document, such as connecting words and common words.

Table 5. Stoword Removal

Tweet	Stopword
apa yang terjadi di mk seperti sebuah lagu kau yang mengakhiri	['terjadi','mahkamah','konstitusi','lagu','mulai','mengakhiri']
stop drama omong kosong pada praktek politik sampah dapur ulang	['stop','drama','bicara','berguna','praktik','politik','bernilai','daur','kembali']
pengadilan tata usaha ptun tidak akan merubah keputusan kpu ri	['pengadilan','tata','usaha','negara','mengubah','komisi','pemilihan','umum','republik','indonesia']

Stemming is a process in natural language processing that aims to transform words that have gone through the stage of stopword elimination to the basic form or root of the word by removing suffixes or affixes attached to the word.

Table 6. Stemming

Tweet	Stemming
apa yang terjadi di mk seperti sebuah lagu kau yang mengakhiri	['jadi','mahkamah','konstitusi','lagu','mulai','akhir']
stop drama omong kosong pada praktek politik sampah dapur ulang	['stop','drama','bicara','guna','praktik','politik','nilai','daur','kembali']
pengadilan tata usaha ptun tidak akan merubah keputusan kpu ri	['adil','tata','usaha','negara','ubah','komisi','pilih']

		Table 8. TF Weighting			
		Token	TF		
			D1	D2	D3
		jadi	1	0	0
		mahkamah	1	0	0
		konstitusi	1	0	0
		lagu	1	0	0
		mulai	1	0	0
		akhir	1	0	0
		stop	0	1	0
		drama	0	1	0
		bicara	0	1	0
		guna	0	1	0
		praktik	0	1	0
		politik	0	1	0
		nilai	0	1	0
		daur	0	1	0
		kembali	0	1	0
		adil	0	0	1
		tata	0	0	1
		usaha	0	0	1
		negara	0	0	1
		ubah	0	0	1
		komisi	0	0	1
		pilih	0	0	1
		umum	0	0	1
		republik	0	0	1
		indonesia	0	0	1

Table 7. Training Data			
	Sentiment Train	Sentiment Train	Class
D 1	apa yang terjadi di mk seperti sebuah lagu kau yang mengakhiri	'jadi','mahkamah','konstitusi', 'lagu', 'mulai', 'akhir']	Positiv e
D 2	stop drama omong kosong pada praktek politik sampah dapur ulang	['stop', 'drama', 'bicara', 'guna', 'praktik', 'politik', 'nilai', 'daur', 'kembali']	Negati ve
D 3	pengadilan tata usaha ptun tidak akan merubah keputusan kpu ri	['adil', 'tata', 'usaha', 'negara', 'ubah', 'komisi', 'pilih', 'umum', 'republik', 'indonesia']	Negati ve

D1 = Weight value of D1  
 D2 = Weight value of D2  
 D3 = Weight value of D3

After calculating the Term Frequency (TF) value, the next step is to calculate the Inverse Document Frequency (IDF), which requires determining the Document Frequency (DF) first the discussion is the basic explanation, relationship and generalization shown by the results.

Table 9. DF Weighting

Token	TF			DF
	D1	D2	D3	
jadi	1	0	0	1
mahkamah	1	0	0	1
konstitusi	1	0	0	1
lagu	1	0	0	1
mulai	1	0	0	1
akhir	1	0	0	1
stop	0	1	0	1
drama	0	1	0	1
bicara	0	1	0	1
guna	0	1	0	1
praktik	0	1	0	1
politik	0	1	0	1
nilai	0	1	0	1
daur	0	1	0	1
kembali	0	1	0	1
adil	0	0	1	1
tata	0	0	1	1
usaha	0	0	1	1
negara	0	0	1	1
ubah	0	0	1	1
komisi	0	0	1	1
pilih	0	0	1	1
umum	0	0	1	1
republik	0	0	1	1
indonesia	0	0	1	1

After calculating the Document Frequency (DF) value, the next step is to calculate the Inverse Document Frequency (IDF). In this example, we are using three documents or tweets, so D equals 3.

Table 10. IDF Weighting

Token	DF	IDF (Log(D/Df))
jadi	1	1.397
mahkamah	1	1.653
konstitusi	1	1.740
lagu	1	1.397
mulai	1	1.477
akhir	1	1.477
stop	1	1.397
drama	1	1.477
bicara	1	1.544
guna	1	1.397
praktik	1	1.602
politik	1	1.602
nilai	1	1.477
daur	1	1.397
kembali	1	1.602
adil	1	1.397
tata	1	1.397
usaha	1	1.477
negara	1	1.544
ubah	1	1.397
komisi	1	1.544
pilih	1	1.477
umum	1	1.397
republik	1	1.653
indonesia	1	1.698

After obtaining the Term Frequency (TF) and Inverse Document Frequency (IDF) values, the next step is to calculate the TF-IDF values used to give weight to each word in the document.

Table 11. TF-IDF Weighting

Token	TFIDF		
	D1	D2	D3
jadi	1.397	0	0
mahkamah	1.653	0	0
konstitusi	1.740	0	0



lagu	1.397	0	0	<b>Data Latih</b>	stop drama omong kosong pada praktek politik sampah dapur ulang	['stop', 'drama', 'bicara', 'guna', 'praktik', 'politik', 'nilai', 'daur', 'kembali']	<b>Negative</b>
mulai	1.477	0	0				
akhir	1.477	0	0				
stop	0	1.397	0				
drama	0	1.477	0				
bicara	0	1.544	0		pengadilan tata usaha ptun tidak akan merubah keputusan kpu ri	['adil', 'tata', 'usaha', 'negara', 'ubah', 'komisi', 'pilih', 'umum', 'republik', 'indonesia']	<b>Negative</b>
guna	0	1.397	0				
praktik	0	1.602	0				
politik	0	1.602	0	<b>Data Uji</b>	sdh terbaca gak bakal menang. dah males politik indonesia"	["baca", "menang", "males", "politik", "indonesia"]	?
nilai	0	1.477	0				
daur	0	1.397	0				
kembali	0	1.602	0				
adil	0	0	1.397				
tata	0	0	1.397				
usaha	0	0	1.477				
negara	0	0	1.544				
ubah	0	0	1.397				
komisi	0	0	1.544				
pilih	0	0	1.477				
umum	0	0	1.397				
republik	0	0	1.653				
indonesia	0	0	1.698				

At the class classification stage, it begins with calculating the prior probability

$$P(\text{Sentiment Class}) = \frac{\text{Number of Class } X}{\text{Total sentiment}} \quad (7)$$

Using the equation above, the probability of each class on sentiment is

- a.  $P(\text{Positive} | \text{Sentiment}) = \frac{1}{3} = 0,3$
- b.  $P(\text{Negative} | \text{Sentiment}) = \frac{2}{3} = 0,6$

Calculation of conditional probability value

$$P(A|B) = \frac{P(B|A) \times P(A)}{P(B)} \quad (8)$$

With the equation of the formula above, the probability of the term in each class on sentiment can be calculated.

- a. Positive
  - $P(\text{baca} | \text{Positive}) = (0+1)/6+25 = 0,032$
  - $P(\text{menang} | \text{Positive}) = (0+1)/6+25 = 0,032$
  - $P(\text{males} | \text{Positive}) = (0+1)/6+25 = 0,032$
  - $P(\text{politik} | \text{Positive}) = (0+1)/6+25 = 0,032$
  - $P(\text{indonesia} | \text{Positive}) = (0+1)/6+25 = 0,032$
- b. Negative
  - $P(\text{baca} | \text{Negative}) = (0+1)/19+25 = 0,022$
  - $P(\text{menang} | \text{Negative}) = (0+1)/19+25 = 0,022$
  - $P(\text{males} | \text{Negative}) = (0+1)/19+25 = 0,022$
  - $P(\text{politik} | \text{Negative}) = (1+1)/19+25 = 0,045$
  - $P(\text{indonesia} | \text{Negative}) = (1+1)/19+25 = 0,045$

Positive Class

### 3.4 Naive Bayes Classifier

Datasets that have gone through preprocessing and feature extraction will then be followed by a learning process using the naive bayes classification method. After being given training data as % to the system, the system can then conduct tests using test data with the aim of testing the accuracy of a system in classifying data. In the classification process, it is necessary to prepare a library that will be used, namely the sklearn or scikit learn library, just like the feature extraction process in using the library. As for the library sklearn.

Table 12. Training Data & Test Data

apa yang terjadi di mk seperti sebuah lagu kau yang mengakhiri	['jadi', 'mahkamah', 'konstitusi', 'lagu', 'mulai', 'akhir']	<b>Positif</b>
--	--	----------------

$P(\text{positive}) = 0,3 \times 0,032 \times 0,032 \times 0,032 \times 0,032 = 1,006$

$P(\text{Negative}) = 0,6 \times 0,022 \times 0,022 \times 0,022 \times 0,045 \times 0,045 = 2,845$

From the calculation above, the results were obtained 2,845 for the negative class and 1,006 for the positive class. Because the probability value in the negative class is greater than the probability value of the positive class, the test data on tweets "baca", "menang", "males", "politik", "Indonesia" are included in the negative class

Model evaluation needs to be done to determine the performance of the model. After sentiment testing using the naïve Bayes algorithm, sentiment classification results in the form of sentiment labels will be obtained. The resulting classification labels will be compared with the actual labels so that the accuracy, precision, recall, and f1-score values of the model used on the dataset will be known.

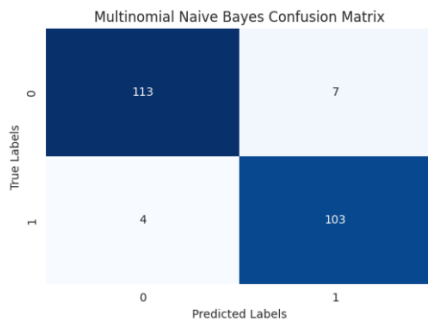


Figure 5. Confusion Matrix on Test Data

Then look for accuracy, precision, recall, and f1-score values in python language programming. So that the results of the code are obtained from the performance results of the naïve bayes method from each class through precision, recall, and f1-score values. The result of this value is in the form of decimal numbers with a range of values ranging from 0 to 1. Where it is said that the higher the number obtained, the better the results obtained.

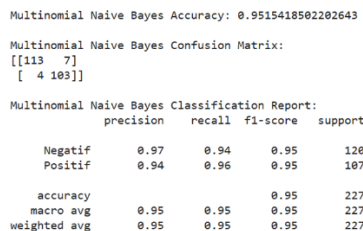


Figure 6. Accuracy Results on Test Data

$$Accuracy = \frac{113 + 103}{113 + 7 + 4 + 103} \times 100\% = 95\%$$

From the results of the calculation above, it can be seen that the number of test data is 227 data and an accuracy value of 95% is produced

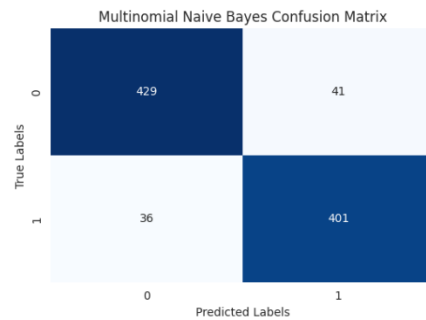


Figure 7. Confusion Matrix on Training Data

Meanwhile, in the training data, the model was tested and the results were obtained as below

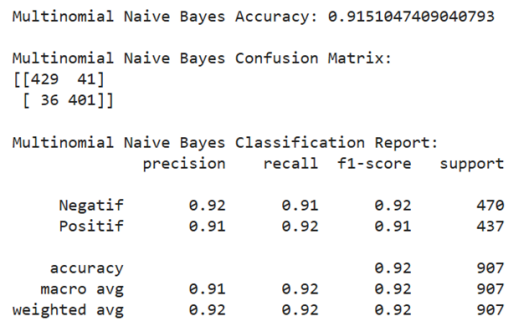


Figure 8. Accuracy Results on Training Data

$$Accuracy = \frac{429 + 41}{429 + 41 + 36 + 401} \times 100\% = 91\%$$

From the results of the calculation above, it can be seen that the number of test data is 907 data and an accuracy value of 91% is produced.



Figure 9. Wordcloud Sentiment

From wordcloud figure 9. The dominating words are the words that most often appear in tweets made by users on social media X. The dominating words are tweets related to Prabowo, Ganjar, Anies, the Constitutional Court, etc. This means that Prabowo gets a lot of negative sentiment from users on platform X.

#### 4. Conclusion

In this study, it can be concluded that public sentiment related to the post-election, especially the 2024 presidential election, is dominated by negative sentiment. However, not a few people also gave positive sentiments related to this. From the research that has been carried out, using data as many as 1228 tweet data, then it became 1134 tweet data, after the pre-processing stage was carried out. With 907 data training and 227 data testing. In testing the model on the training data, the accuracy was 91%. Meanwhile, the model test on the test data obtained an accuracy of 95%. Based on the performance results produced, which is an accuracy value of 95%, this is of course a fairly accurate result.

#### References

- [1] A. P. Nardilasari, A. L. Hananto, S. S. Hilabi, T. Tukino, dan B. Priyatna, "Analisis Sentimen Calon Presiden 2024 Menggunakan Algoritma SVM Pada Media Sosial Twitter," *JOINTECS (Journal Inf. Technol. Comput. Sci.*, vol. 8, no. 1, hal. 11, 2023.
- [2] I. Kurniawan dan A. Susanto, "Implementasi Metode K-Means dan Naïve Bayes Classifier untuk Analisis Sentimen Pemilihan Presiden (Pilpres) 2019," *Eksplora Inform.*, vol. 9, no. 1, hal. 1–10, 2019.
- [3] N. Hardi, Y. Alkahfi, P. Handayani, W. Gata, dan M. R. Firdaus, "Analisis Sentimen Physical Distancing pada Twitter Menggunakan Text Mining dengan Algoritma Naive Bayes Classifier," *Sistemasi*, vol. 10, no. 1, hal. 131, 2021.
- [4] L. D. Mahbubah dan E. Zuliarso, "Analisa Sentimen Twitter Pada Pilpres 2019 Menggunakan Algoritma Naive Bayes," *Sintak*, hal. 194–195, 2019, [Daring]. Tersedia pada: <https://www.unisbank.ac.id/ojs/index.php/sintak/article/view/7585>
- [5] A. Nugraha, Y. H. Chrisnanto, dan R. Yuniarti, "Prediksi Sentimen Pada Sosial Media Twitter Mengenai Produk Smartphone Menggunakan Algoritma K-NN Classification," *Sensasi*, hal. 251–258, 2019.
- [6] L. Saidaliyeva, J. Uzokova, dan J. Juzjasarova, "European Journal of Interdisciplinary Research and Development Volume-12 Website : [www.ejird.journalspark.org](http://www.ejird.journalspark.org) ISSN ( E ): 2720-5746 European Journal of Interdisciplinary Research and Development Volume-12
- [7] P. Wang *et al.*, "Classification of Proactive Personality: Text Mining Based on Weibo Text and Short-Answer Questions Text," *IEEE Access*, vol. 8, hal. 97370–97382, 2020.
- [8] I. Athiyyah Rahma dan L. Hulliyyatus Suadaa, "Penerapan Text Augmentation Untuk Mengatasi Data Yang Tidak Seimbang Pada Klasifikasi Teks Berbahasa Indonesia Studi Kasus: Deteksi Judul Clickbait Dan Komentar Hate Speech Pada Berita Online," *J. Teknol. Inf. dan Ilmu Komput.*, vol. 10, no. 6, hal. 1329–1340, 2023.
- [9] S. Fauziah, D. D. Saputra, R. L. Pratiwi, dan M. R. Kusumayudha, "Komparasi Metode Feature Selection Text Mining Pada Permasalahan Klasifikasi Keluhan Pelanggan Industri Telekomunikasi Menggunakan Smote Dan Naïve Bayes," *IJIS - Indones. J. Inf. Syst.*, vol. 8, no. 2, hal. 174, 2023.
- [10] H. Liu, X. Chen, dan X. Liu, "A Study of the Application of Weight Distributing Method Combining Sentiment Dictionary and TF-IDF for Text Sentiment Analysis," *IEEE Access*, vol. 10, hal. 32280–32289, 2022.
- [11] Z. Li, R. Li, dan G. Jin, "Sentiment analysis of danmaku videos based on naïve bayes and sentiment dictionary," *IEEE Access*, vol. 8, hal. 75073–75084, 2020.
- [12] M. E. Atmojo dan V. P. Pratiwi, *Media Sosial Twitter sebagai Platform Media Informasi Digital dalam Penerapan New Normal*, vol. 1, no. November. 2021.
- [13] D. Duei Putri, G. F. Nama, dan W. E. Sulistiono, "Analisis Sentimen Kinerja Dewan Perwakilan Rakyat (DPR) Pada Twitter Menggunakan Metode Naive Bayes Classifier," *J. Inform. dan Tek. Elektro Terap.*, vol. 10, no. 1, hal. 34–40, 2022.
- [14] S. Raschka, J. Patterson, dan C. Nolet, "Machine learning in python: Main developments and technology trends in data science, machine learning, and artificial intelligence," *Inf.*, vol. 11, no. 4, 2020.
- [15] A. P. Joshi dan B. V. Patel, "Data Preprocessing: The Techniques for Preparing Clean and Quality
- Website : [www.ejird.journalspark.org](http://www.ejird.journalspark.org) ISSN ( E ): 2," no. 2013, hal. 74–78, 2023.

- Data for Data Analytics Process,” *Orient. J. Comput. Sci. Technol.*, vol. 13, no. 0203, hal. 78–81, 2021. [18]
- [16] N. V. Pusean, N. Charibaldi, dan B. Santosa, “Comparison of Scenario Pre-processing Performance on Support Vector Machine and Naïve Bayes Algorithms for Sentiment Analysis,” *Inf. J. Ilm. Bid. Teknol. Inf. dan Komun.*, vol. 8, no. 1, hal. 57–63, 2023.
- [17] H. Jurnal, P. W. Setyaningsih, A. Witanti, dan I. Susilawati, “Jurnal Teknik Informatika Data Preperation Untuk Automatic Summarization Video To Text,” vol. 10, no. 2, 2022.
- [19] N. R. Siahaan, R. Y. Tiffany, and S. R. E. Sinaga, “Analisis Sentimen Ulasan Aplikasi Media Sosial Whatsapp Menggunakan Metode Naive Bayes Classifier,” *J. Ilm. ...*, no. 02, hal. 343–354, 2023, [Daring]. Tersedia pada: <https://ejournal.pppmitpa.or.id/index.php/betrik/article/view/104%0Ahttps://ejournal.pppmitpa.or.id/index.php/betrik/article/download/104/76>
- L. A. Waskito, K. M. Lhaksmana, dan D. T. Murdiansyah, “Analisis Sentimen Terhadap Pemilihan Presiden Indonesia 2019 Pada Media Sosial Twitter Menggunakan Metode Naïve Bayes,” *eProceedings Eng.*, vol. 6, no. 2, hal. 9753–9765, 2019.