

# Sentiment Analysis of TikTok Application on Twitter using The Naïve Bayes Classifier Algorithm

Denia Putri Fajrina, Syafriandi\*, Nonong Amalita, Admi Salma

Departemen Statistika, Universitas Negeri Padang, Padang, Indonesia

\*Corresponding author: [syafriandi\\_math@fmipa.unp.ac.id](mailto:syafriandi_math@fmipa.unp.ac.id)

Submitted : 19 September 2023

Revised : 25 Oktober 2023

Accepted : 27 Oktober 2023

## ABSTRACT

*TikTok is a popular social media platform that has gained a lot of attention lately. People of all ages are using this application to share short videos with their friends and followers. The content on TikTok is diverse and can be tailored to individual preferences, but there have been concerns about the presence of vulgar content that can be viewed by minors as there are no age restrictions. This has led to public scrutiny of the application on social media platforms like Twitter. To address this issue, sentiment analysis was conducted on reviews of the TikTok application to help users make informed decisions about its use. The aim of this analysis was to determine whether people's opinions about TikTok were positive or negative. To achieve this goal, researchers used the hashtag "TikTok Application". The results were classified into two categories positive and negative using the Naïve Bayes Classifier method. The analysis was carried out using 80% training data and 20% testing data, and the results showed an accuracy rate of 80.32%, with a recall value of 97% and a precision value of 78%. In general, positive feedback from Indonesians on the TikTok application is related to the invitation to download the TikTok application, while in negative feedback, information is obtained that the TikTok application is based on content that is inappropriate for TikTok users to download. This information can help users make informed decisions about using the TikTok application.*

**Keywords:** Naïve Bayes Classifier, Sentiment Analysis, TikTok, Twitter.



This is an open access article under the Creative Commons 4.0 Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited. ©2022 by author and Universitas Negeri Padang.

## I. PENDAHULUAN

Seiring dengan pesatnya perkembangan teknologi, masyarakat dapat menggunakan berbagai macam sarana untuk berkomunikasi, termasuk jejaring sosial. Jejaring sosial juga berperan sebagai platform yang dapat digunakan untuk berkomunikasi dan berbagi informasi dari berbagai kalangan masyarakat. Saat ini aplikasi yang sangat populer digunakan masyarakat untuk berkomunikasi adalah TikTok.

TikTok memberikan pengalaman yang cukup menyenangkan dari pada media sosial lainnya. Aplikasi ini berfokus pada hiburan dan kesenangan dengan konten yang di buat oleh pengguna TikTok. Menurut Zulqornain & Adikara (2021) penggunaan media sosial telah mengunduh aplikasi TikTok sebanyak 45,8 juta kali dan jumlah tersebut mengalahkan media sosial lainnya seperti WhatsApp, Youtube, Facebook, dan bahkan Instagram. DataIndonesia.id menyebutkan bahwa, Indonesia menduduki peringkat kedua dengan jumlah pengguna aktif sebanyak 109.9 juta pengguna di awal tahun 2023.

Beragam komentar dan opini yang diberikan oleh pengguna aplikasi TikTok baik komentar positif maupun komentar negatif. Komentar ini dapat dilihat pada media sosial Twitter. Twitter merupakan salah satu sumber data yang dapat digunakan untuk mengukur apakah aplikasi TikTok ini berdampak positif atau negatif dalam penggunaannya. Dukungan dan penolakan dalam suatu program dapat dihitung berdasarkan komentar dan opini publik di media sosial (Silitonga et al., 2023).

Analisis sentimen merupakan proses mengidentifikasi sentimen dalam teks dengan mengolah sebuah data teks yang tidak terstruktur menjadi sebuah teks yang terstruktur dan dapat memberikan informasi mengandung sentimen. Analisis sentimen dapat digunakan untuk menganalisis ulasan opini pada Twitter. Sebelum melakukan analisis sentimen, diperlukan *preprocessing* data untuk mengolah data teks agar siap untuk dianalisis, dengan tahapannya yaitu *cleaning, case folding, tokenizing, normalization, stopword removal, dan stemming* (Indriyani, et al. 2022).

Pada penelitian ini dilakukan pengelompokan sentimen dengan menggunakan salah satu algoritma klasifikasi yaitu *Naïve Bayes Classifier* (NBC). Algoritma ini sangat baik dalam melakukan pengelompokkan pada data

teks(Herdhianto, 2020). Data teks yang digunakan berupa *tweet*. *Tweet* adalah komentar atau ulasan dari pengguna yang terdiri dari berbagai 280 karakter. Ada dua tahap dalam klasifikasi *tweet*, tahap pertama adalah training terhadap *tweet* yang telah diketahui kategorinya. Ada dua tahap dalam klasifikasi *tweet*, tahap pertama adalah training terhadap *tweet* yang telah diketahui kategorinya. Tahap kedua klasifikasi *tweet* yang belum diketahui kategorinya. Rahmadani, et al. (2022) mengusulkan untuk menerapkan algoritma NBC untuk menganalisis sentimen TikTok di media sosial menggunakan data Twitter. Studi tersebut menemukan bahwa algoritma NBC bekerja lebih baik dalam analisis sentimen aplikasi TikTok dengan data Twitter.

Tujuan dari penelitian ini adalah untuk mengetahui sentimen masyarakat terhadap aplikasi TikTok dan melihat akurasi algoritma NBC. Hasil yang diperoleh dapat bermanfaat bagi pengembang aplikasi TikTok dan dapat digunakan sebagai rujukan bagi masyarakat dan peneliti selanjutnya.

## II. METODE PENELITIAN

### A. Sumber Data dan Teknik Analisis Data

Periode pengumpulan data dari bulan Januari sampai bulan Juni 2023, data penelitian dikumpulkan melalui media sosial Twitter dengan Python. *Keyword* yang digunakan “Aplikasi TikTok” terkumpul sebanyak 1216 data.

Langkah- langkah yang dilakukan dalam analisis data yaitu sebagai berikut.

#### 1. Melakukan pelabelan data

Pada tahap ini dilakukan pelabelan data dengan menandai review aplikasi yang terkumpul berdasarkan opini positif dan negatif. Ketentuan rating atau rating positif adalah rating yang mendukung aplikasi dan penggunaannya, rating negatif adalah kritik atau keluhan terhadap penggunaan aplikasi TikTok. Pelabelan data dilakukan secara otomatis menggunakan *lexicon based*.

#### 2. Melakukan *preprocessing data*

Tahap *preprocessing data* merupakan tahap membersihkan data mentah menjadi data yang mudah dipahami. Tahap ini sangat perlu dilakukan dalam sebuah analisis sentimen. Ada beberapa tahapan yaitu (Irham & Wisesty, 2019)

:

- a. *Cleaning data* merupakan proses pembersihan teks dari karakter atau kata yang tidak perlu untuk mengurangi noise selama klasifikasi.
- b. *Case folding* merupakan peruses semua huruf dalam dokumen menjadi hurus kecil untuk menghindari perbedaan ejaan yang tidak perlu.
- c. *Tokenizing* merupakan proses pemotongan kalimat menjadi kata dengan menganalisis sekumpulan data dengan memisahkan kata dan menentukan struktur setiap kata.
- d. *Stopword* merupakan proses menghilangkan kata-kata yang tidak bermakna atau tidak relevan dalam analisis data, seperti kata penghubung.
- e. *Normalization* merupakan proses identifikasi penulisan kata berlebihan kemudian diganti dengan kata kamus KBBI.

#### 3. *Stemming* merupakan proses pencari kata dasar dari setiap kata hasil dari proses *normalization* Melakukan pembobotan data

Dalam tahap pembobotan data, digunakan metode *Term Frequency* dan *Invers Document Frequency* (TF-IDF). TF merupakan frekuensi kemunculan kata pada setiap dokumen. IDF yaitu frekuensi dokumen yang di invers. Dengan rumus TF-IDF dapat dilihat pada Persamaan 1 (Sammut, 2011).

$$W_{i,j} = t_{f_{i,j}} \times idf_i \quad (1)$$

$$idf_i = \log \left( \frac{N}{f_j} \right) \quad (2)$$

Keterangan :

- $W_{i,j}$  : Bobot dari kata ke-i pada dokumen ke-j  
 $N$  : Total seluruh dokumen  
 $t_{f_{i,j}}$  : Total kemunculan kata ke-i pada dokumen ke-j  
 $f_j$  : Total dokumen sampai ke-j  
 $idf_i$  : Total dokumen pada setiap kata ke-i

#### 4. Melakukan klasifikasi menggunakan algoritma *Naive Bayes Classifier*

Pada tahap klasifikasi digunakan algoritma NBC. Algoritma NBC merupakan algoritma dengan menggunakan nilai probabilitas dalam melakukan klasifikasi data uji pada kategori yang paling tepat. Metode NBC salah satu metode

sederhana tetapi memiliki nilai akurasi yang tinggi. Metode NBC sangat baik dalam melakukan klasifikasi pada data teks berupa *tweet* (Herdhianto, 2020). Pada saat klasifikasi algoritma akan mencari probabilitas tertinggi dari semua kategori/kelas dokumen yang diujikan ( $Y_{map}$ ). Nilai  $P(Y_i)$  dan  $P(X_i|Y_i)$  dihitung pada saat *training*, didapat dengan rumus sebagai berikut (Herdhianto, 2020).

$$P(Y_i) = \frac{|doc\ i|}{|training|} \tag{3}$$

$$P(X_i|Y_i) = \frac{n_i + 1}{|n + kosakata|} \tag{4}$$

Adapun persamaan  $Y_{map}$  adalah sebagai berikut (Herdhianto, 2020).

$$Y_{map} = \underset{Y_i = Y}{\arg \max} P(Y_i) \prod_i P(X_i|Y_i) \tag{5}$$

Keterangan:

- $Y_{map}$  : Semua kategori/kelas untuk data *test*
- $Y_i$  : Kategori/kelas dokumen yaitu positif dan negatif
- $P(X_i|Y_i)$  : Probabilitas  $X_i$  pada kategori
- $P(Y_i)$  : Probabilitas dari  $Y_i$
- $|doc\ i|$  : Total dokumen pada kelas ke-i dalam *training*
- $|training|$  : Total dokumen pada data *training*
- $n_i$  : Total kemunculan kata pada data ke-i
- $n$  : Banyaknya seluruh kata dalam dokumen pada kelas  $Y_i$
- $|kosakata|$  : Banyaknya kata dalam *training*

5. Melakukan evaluasi kinerja model

Tahap selanjutnya yaitu mengevaluasi kinerja model menggunakan *confusion matrix*. Menurut Flach (2007) pada pengukuran kinerja menggunakan *confusion matrix*, terdapat 4 istilah sebagai representasi hasil proses klasifikasi. Pada *confusion matrix* menghasilkan nilai yaitu *True Positive* (TP), *True Negative* (TN), *False Positive* (FP), dan *False Negative* (FN). Nilai-nilai yang diperoleh dapat disajikan seperti pada Tabel 1.

**Tabel 1.** *Confusion Matrix*

	Nilai Prediksi		
	Kelas	Positif	Negatif
Nilai Aktual	Positif	TP	FN
	Negatif	FP	TN

$$Akurasi = \frac{TP + TN}{TP + TN + FP + FN} \times 100\% \tag{6}$$

$$Presisi = \frac{TP}{TP + FP} \times 100\% \tag{7}$$

$$Sensitivitas = \frac{TP}{TP + FN} \times 100\% \tag{8}$$

Agar sistem dapat memberikan nilai informasi dan tingkat keakuratan kepada peneliti, maka terbagi tiga informasi yaitu akurasi, presisi dan sensitivitas. Akurasi didefinisikan sebagai tingkat kesesuaian antara nilai yang diproyeksikan dengan nilai sebenarnya dengan membandingkan data yang telah diklasifikasikan dengan benar ke seluruh data set. Sedangkan presisi adalah jumlah data teks berkategori positif yang diklasifikasikan dengan benar dibagi dengan seluruh data yang digolongkan positif. Sensitivitas adalah jumlah persentase data kategori positif yang dapat diklasifikasikan oleh sistem secara akurat untuk menentukan seberapa suksesnya pemulihan informasi.

6. Melakukan visualisasi kata dengan *word cloud*

Pada tahap visualisasi dilakukan dengan *word cloud*. *Word cloud* merupakan salah satu grafik sederhana yang mampu menunjukkan hubungan antara frekuensi kata dengan melihat ukuran kata yang lebih sering digunakan dengan

cepat (Yanuarti, 2021). Kata-kata yang ditampilkan dengan ukuran besar dan lebih mencolok lebih mudah menarik perhatian pembaca dan mempercepat dalam mengambil kesimpulan.

### III. HASIL DAN PEMBAHASAN

#### A. Pelabelan Data

Proses pelabelan data dilakukan secara otomatis, dengan cara mengekstrak kalimat *tweet*. Metode yang digunakan untuk mengekstrak yaitu lexicon based. Berikut ini hasil proses pelabelan data *tweet*. Dapat ditunjukkan pada Tabel 2.

**Tabel 2.** Pelabelan Data

Kelas Sentimen	Jumlah
Positif	768
Negatif	449

Berdasarkan Tabel 2, data yang diperoleh dalam hasil pelabelan data untuk dilakukan *preprocessing* dan sudah dikelompokkan menjadi 2 kategori yaitu positif dan negatif. Total data yang diperoleh sebesar 1217 data, dimana data positif sebesar 768 data dan data negatif sebesar 449 data.

#### B. Preprocessing Data

Tahap *preprocessing* merupakan tahapan setelah tersedianya satu atau lebih dataset hasil dari tahapan pengumpulan data yang dilakukan untuk analisis menggunakan algoritma NBC. Tahapan pada *preprocessing* disajikan pada Tabel 3.

**Tabel 3.** Preprocessing Data pada Analisis sentimen

Nama Proses	Output Kalimat
Data Set	<a href="https://t.co/9MNUOw7L9t">@kenjthana @tanyakanrl</a> Tapi banyak yg make. Mau dihina segimanapun, ga usah denial, tiktok adalah aplikasi paling sering digunakan untuk sekarang.
Cleaning Data	Tapi banyak yg make Mau dihina segimanapun ga usah denial tiktok adalah aplikasi paling sering digunakan untuk sekarang
Case Folding	tapi banyak yg make mau dihina segimanapun ga usah denial tiktok aplikasi paling sering digunakan untuk sekarang
Tokenization	['make', 'dihina', 'segimanapun', 'ga', 'denial', 'tiktok', 'aplikasi']
Stopword	make dihina segimanapun ga denial tiktok aplikasi
Normalization	pakai dihina segimanapun tidak denial tiktok aplikasi
Stemming	pakai: pakai dihina: hina segimanapun : segimanapun tidak : tidak denial:denial tiktok : tiktok aplikasi :aplikasi

Berdasarkan Tabel 3, diperoleh hasil dari *crawling* data dari Twitter. Data yang dikumpulkan berupa data *tweet* yang tidak terstruktur yang perlu dilakukan *preprocessing* data. Langkah pertama yaitu *tweet* dibersihkan dari simbol, angka, url, dan tag yang tidak memberikan makna disebut proses *cleaning* data. Langkah kedua dengan proses menyetarakan bentuk huruf seperti “Tapi” menjadi “tapi”, “Mau” menjadi “mau” disebut proses *case folding*. Setelah bentuk huruf sama tidak ada huruf kapital lagi, langkah ketiga yaitu membuat kalimat terpisah oleh spasi menjadi token disebut dengan proses *tokenization*. Kalimat yang sudah berbentuk token dilakukan penghapusan pada kata yang tidak bermakna proses ini disebut *stopword*. Kemudian agar kata-kata yang diperoleh menjadi kata baku menurut KBBI dilakukan *normalization* seperti “make” menjadi “pakai”, “ga” menjadi “tidak”. Setelah proses normalisasi kata langkah terakhir yaitu proses *stemming*, *stemming* ini akan mengubah kata berimbuhan menjadi kata dasar seperti kata “dihina” menjadi “hina”.

**C. Pembobotan**

Pada tahap pembobotan dilakukan menggunakan TF-IDF. TF-IDF bertujuan pembobotan pada kata dan mengatasi masalah dalam mengklasifikasikan data ke dalam sentimen, diperoleh hasil seperti Tabel 4.

**Tabel 4.** Bobot Kata menggunakan TF-IDF

Kata No.Dok	abal	abang	abis	aborsi	...	video	videographer
1.	0	0	0	0	...	0.257827	0
2.	0	0	0	0	...	0	0
...	...	...	...	...	...	...	...

Berdasarkan Tabel 4 diperoleh hasil bahwa pada kata “abal”, “abang”, “abis”, “aborsi” diperoleh 0, yang artinya tidak ada kata “abal”, “abang”, “abis”, “aborsi” berada pada dokumen 1 dan tidak memiliki nilai. Sementara pada kata “video” diperoleh hasil 0.257827 yang artinya kata “video” ada pada dokumen 1 dan memiliki nilai 0.257827 pada dokumen 1.

**D. Klasifikasi menggunakan Naïve Bayes Classifier**

Pada tahapan pengujian data uji digunakan algoritma NBC pada proses pengujian algoritma akan menghasilkan dua klasifikasi yaitu positif dan negatif. Hasil pengujian menggunakan algoritma NBC dapat Tabel 5.

**Tabel 5.** Probabilitas Klasifikasi pada Aplikasi TikTok

No Dok	Ulasan	Negatif	Positif	Sentiment Label
1	pengaruh aplikasi tiktok terhadap gaya hidup generasi milenial tapi balik lagi ke kamu ya kalo bahannya or kamu cocoknya opsi pertama ga ada masalah juga topiknya sama bagus	0.245891	0.754109	Positif
2	Ngga punya aplikasi lain selain Twitter sama M-banking, bosan banget. Kapan ujiannya selesai. Mau download TikTok sama Instagram lagi :(((	0.932961	0.067039	Negatif

Berdasarkan hasil klasifikasi pada Tabel 5 diperoleh bahwa pada dokumen 1 didapatkan nilai probabilitas negatif sebesar 0.245891 dan probabilitas positif sebesar 0.754109, Hasil yang diperoleh nilai probabilitas positif lebih besar dari probabilitas negatif maka dinyatakan sentimen positif. Hal ini terlihat bahwa pada dokumen 1 memiliki kecenderungan pada ulasan positif dan seterusnya pada dokumen berikutnya. Sementara ada dokumen 2 didapatkan nilai probabilitas negatif lebih besar dari pada nilai probabilitas positif, maka dokumen ke dua berlabel negatif.

**E. Confusion Matrix**

Tahapan ini dilakukan evaluasi model, penganalisaan serta menemukan tingkat akurasi dari sistem yang dikembangkan. Metode yang digunakan yaitu algoritma NBC untuk melakukan klasifikasi untuk data uji. Dalam pengujian akurasi data latih dan data uji dibagi menjadi 80% dan 20%, hasil yang diperoleh seperti pada Tabel 6.

**Tabel 6.** Hasil Evaluasi Model

Algoritma	Akurasi	Presisi	Sensitivitas
NBC	80.32%	78%	97%

Berdasarkan hasil pada Tabel 6, proses klasifikasi menggunakan algoritma NBC menghasilkan nilai akurasi sebesar 80.32%, nilai presisi sebesar 78% dan nilai sensitivitas 97%. Hasil dari *confusion matrix* ini memiliki ketepatan klasifikasi pada data sentimen Aplikasi TikTok secara keseluruhan baik.

#### F. Word Cloud

Visualisasi mengenai ulasan Aplikasi TikTok di Twitter disajikan dalam bentuk *word cloud*. Word Cloud berfungsi untuk memberikan gambaran kata-kata yang sering muncul dimana besar atau kecilnya suatu kata bergantung pada banyaknya kemunculan kata tersebut dalam keseluruhan data. Didapatkan hasil seperti Gambar 1.



Gambar 1. (a) Word Cloud tanggapan positif (b) Word Cloud tanggapan negatif

Berdasarkan pada Gambar 1 (a) adalah *word cloud* dari keseluruhan tanggapan positif masyarakat Indonesia mengenai Aplikasi TikTok yang diperoleh dari bulan Januari sampai Juni 2023 melalui media sosial Twitter, sedangkan (b) adalah *word cloud* dari tanggapan negatif masyarakat Indonesia mengenai Aplikasi TikTok yang diperoleh dari bulan Januari sampai Juni 2023 melalui media sosial Twitter. Berdasarkan *word cloud* tersebut dapat diketahui bahwa perbincangan positif masyarakat Indonesia mengenai aplikasi TikTok banyak dikaitkan dengan kata “tiktok”, “aplikasi”, “tidak”, “mengunduh”, “menonton”. Hal ini menunjukkan adanya ajakan untuk mengunduh aplikasi TikTok dan menonton video di Aplikasi TikTok. Sedangkan pada tanggapan negatif perbincangan sering dikaitkan dengan “tiktok”, “aplikasi”, “banget”, “kalau”, “tidak”, “sudah”, “saya”. Hal itu menunjukkan bahwa adanya Aplikasi TikTok dengan konten tidak sesuai bagi pengguna yang sudah mengunduh Aplikasi TikTok.

#### IV. KESIMPULAN

Berdasarkan analisis yang telah dilakukan menggunakan dataset ulasan dari twitter dan algoritma yang diusulkan, maka dapat disimpulkan berdasarkan 1217 *tweet*, arah pandangan (sentimen) masyarakat Indonesia terhadap aplikasi TikTok dengan pengujian menggunakan algoritma *Naïve Bayes Classifier* ada di angka 80,32% dan mendapatkan hasil *recall* 97%, dan *presisi* 78%. Secara umum *feedback* positif dan negatif masyarakat Indonesia terhadap Aplikasi TikTok berkaitan dengan ajakan mendownload Aplikasi TikTok, sedangkan pada *feedback* negatif informasi yang menarik yaitu Aplikasi TikTok dengan konten tidak sesuai dengan pengguna yang mengunduh Aplikasi TikTok. Dengan proses yang tepat sangat penting untuk mencapai hasil optimal untuk proses selanjutnya, aplikasi ini akan menjadi lebih baik jika data latih memiliki klasifikasi sentimen positif dan negatif dengan jumlah yang sama.

#### DAFTAR PUSTAKA

- Flach, P. (2007). Performance Evaluation in Machine Learning: The Good, The Bad, The Ugly and The Way Forward. *The Alan Turing Institute*.
- Herdhianto, A. (2020). *Sentiment Analysis Menggunakan Naïve Bayes Classifier (NBC) Pada Tweet Tentang Zakat* [Universitas Islam Negeri Syarif Hidayatullah]. <http://repository.uinjkt.ac.id/dspace/handle/123456789/53661>
- Indriyani, E. R., Paradise, P., & Wibowo, M. (2022). Perbandingan Metode Naïve Bayes dan Support Vector Machine Untuk Analisis Sentimen Terhadap Vaksin Astrazeneca di Twitter. *Jurnal Media Informatika Budidarma*, 6(3), 1545. <https://doi.org/10.30865/mib.v6i3.4220>
- Rahmadani, P. S., Tampubolon, F. C., Jannah, A. N., Hutabarat, N. L. H., & Simarmata, A. M. (2022). Tiktok Social Media Sentiment Analysis Using the Nave Bayes Classifier Algorithm. *Sinkron*, 7(3), 995–999. <https://doi.org/10.33395/sinkron.v7i3.11579>
- Sammut (Ed.). (2011). *TF-IDF*. In: Sammut, C., Webb, G.I. (eds) *Encyclopedia of Machine Learning*. Springer, Boston, MA. [https://doi.org/https://doi.org/10.1007/978-0-387-30164-8\\_832](https://doi.org/https://doi.org/10.1007/978-0-387-30164-8_832)

- Silitonga, P. D. P., Hasibuan, M., Situmorang, Z., & Purba, D. (2023). Comparison of Tiktok User Sentiment Analysis Accuracy with Naïve Bayes and Support Vector Machine. *International Journal of Advanced Trends in Computer Science and Engineering*, 12(1), 11–15. <https://doi.org/10.30534/ijatcse/2023/031212023>
- Yanuarti, R. (2021). Analisis Media Sosial Twitter Terhadap Topik Vaksinasi Covid-19. *JUSTINDO (Jurnal Sistem Dan Teknologi Informasi Indonesia)*, 6(2), 121–130. <https://doi.org/10.32528/justindo.v6i2.5503>
- Zulqornain, J. A., & Adikara, P. P. (2021). Analisis Sentimen Tanggapan Masyarakat Aplikasi Tiktok Menggunakan Metode Naïve Bayes dan Categorical Propotional Difference ( CPD ). 5(7), 2886–2890.