

Vol. 13 No. 3S1, pISSN: 2303-0577 eISSN: 2830-7062

http://dx.doi.org/10.23960/jitet.v13i3S1.8150

PERBANDINGAN ALGORITMA RANDOM FOREST DAN NAÏVE BAYES DALAM MENGANALISIS SENTIMEN ULASAN PADA PRODUK SKINCARE LOKAL DI MEDIA SOSIAL TIKTOK

Putri Widya Sari^{1*}, Firmansyah², Abdul Rahman Kadafi³

1,2,3 Universitas Bina Sarana Informatika, Teknik dan Informatika, Informatika

Keywords:

Sentiment Analysis, Random Forest, Naïve Bayes, Tiktok, Local *Skincare*, Web Scraping, Consumer Reviews, Classification

Corespondent Email: putriwidya730@gmail.com

Abstrak. Penelitian ini bertujuan untuk menganalisis sentimen ulasan konsumen terhadap produk skincare lokal yang dibagikan di media sosial, khususnya tiktok, dengan menggunakan algoritma random forest dan naïve bayes. Penelitian ini fokus pada pengklasifikasian sentimen konsumen menjadi kategori positif dan negatif untuk memberikan wawasan tentang preferensi konsumen. Data dikumpulkan melalui web scraping menggunakan skrip python dan diproses dengan teknik pre-processing standar seperti case folding,tokenisasi, penghapusan stopword, dan stemming. hasil analisis menunjukkan bahwa random forest mengungguli naïve bayes dalam hal akurasi, presisi, recall, dan fl-score. Temuan ini mengindikasikan bahwa random forest lebih efektif dalam menangani dataset yang kompleks dengan banyak fitur, sementara naïve bayes lebih cepat tetapi mungkin kesulitan dengan interaksi fitur yang lebih rumit. Selain itu, distribusi sentimen menunjukkan dominasi sentimen negatif yang sedikit lebih tinggi, menyoroti area yang perlu diperbaiki dalam produk skincare lokal. Penelitian ini memberikan wawasan yang berguna bagi konsumen dan produsen, membantu konsumen membuat keputusan pembelian yang lebih tepat dan membantu produsen dalam mengoptimalkan strategi pemasaran mereka.



Copyright © JITET (Jurnal Informatika dan Teknik Elektro Terapan). This article is an open access article distributed under terms and conditions of the Creative Commons Attribution (CC BY NC)

Abstract. This research aims to analyze the sentiment of consumer reviews on local skincare products shared on social media, particularly tiktok, using the random forest and naïve bayes algorithms. The study focuses on classifying consumer sentiments into positive and negative categories to provide insights into consumer preferences. The data was collected through web scraping using python scripts and processed with standard data cleaning and preprocessing techniques such as case folding, tokenization, stopword removal, and stemming. The results of the analysis showed that random forest outperformed naïve bayes in terms of accuracy, precision, recall, and flscore. The findings indicate that **random forest** is more effective in handling complex datasets with multiple features, while naïve bayes works faster but may struggle with more intricate feature interactions. Additionally, the sentiment distribution revealed a slight dominance of negative sentiments, highlighting areas for improvement in local skincare products. This study offers valuable insights for both consumers and producers, assisting consumers in making informed purchasing decisions and helping producers optimize their marketing strategies.

1. PENDAHULUAN

lokal Produk skincare semakin mendapat perhatian di pasar indonesia, seiring dengan meningkatnya kesadaran masyarakat pentingnya tentang merawat menggunakan bahan-bahan alami yang sesuai dengan kebutuhan kulit mereka. "Ketertarikan dan kesadaran masyarakat terhadap skincare ini menarik perhatian beberapa individual atau perusahaan lokal untuk membuat suatu brand yang menawarkan produk skincare yang mengandung banyak manfaat disertai dengan harga yang terjangkau"[1]. Meskipun produkproduk lokal ini sering kali lebih terjangkau, konsumen yang masih merasa banvak kebingungan dalam memilih produk yang tepat karena kurangnya informasi yang jelas mengenai kualitas dan efektivitas produk tersebut. Oleh karena itu, penelitian ini bertujuan untuk memberikan wawasan yang jelas bagi konsumen, agar mereka dapat membuat keputusan yang lebih tepat saat memilih produk skincare lokal berdasarkan analisis sentimen dari ulasan yang ada di media sosial.

Selain itu, peningkatan kesadaran tentang kualitas produk lokal menjadi hal yang penting. Banyak produk skincare lokal yang dianggap kurang berkualitas dibandingkan produk internasional, padahal banyak produk lokal yang memiliki bahan alami dan harga terjangkau. Dengan menggunakan perbandingan analisis sentimen berbasis random forest dan naïve bayes, konsumen dapat memperoleh gambaran yang lebih jelas tentang bagaimana produk lokal diterima di pasar, ulasan di media sosial melalui mengungkapkan sentimen positif, atau negatif, "data dari analisis ini dapat menjadi alat yang berguna sebagai pertimbangan bagi perusahaan dalam mengambil keputusan, serta memberikan wawasan kepada masyarakat yang berencana menggunakan produk tersebut"[2].

Saat ini, banyak perusahaan yang memanfaatkan platform media sosial untuk menjalankan strategi pemasaran mereka. "Media sosial berfungsi sebagai forum online yang menghubungkan pelanggan dan calon pelanggan, di mana mereka dapat berdiskusi dan memberikan komentar atau ulasan tentang suatu produk"[3]. Namun, seringkali ulasan yang beredar di media sosial tidak terstruktur dan bersifat subjektif, sehingga sulit untuk

mendapatkan gambaran yang jelas tentang kualitas produk. Oleh karena itu, analisis sentimen menjadi solusi yang tepat untuk mengklasifikasikan ulasan-ulasan tersebut ke dalam kategori positif atau negatif, serta memberikan pemahaman yang lebih jelas mengenai persepsi konsumen terhadap produk.

"Perubahan tren memainkan peran penting dalam dinamika pasar perawatan kulit. telah Perkembangan waktu membawa pergeseran signifikan dalam preferensi dan perilaku konsumen terhadap produk perawatan kulit."[4]. Jika banyak ulasan positif, mereka bisa memaksimalkan kampanye pemasaran untuk menarik lebih banyak konsumen. sebaliknya, jika ada banyak ulasan negatif, mereka dapat segera mengidentifikasi masalah pada produk mereka dan melakukan perbaikan. Hal ini juga berkontribusi pada peningkatan kualitas dan variasi produk, karena produsen dapat menyesuaikan produk dengan kebutuhan pasar yang lebih tepat.

Penelitian ini diharapkan tidak hanya memberikan manfaat bagi konsumen yang ingin membuat keputusan yang lebih bijaksana dalam memilih produk skincare lokal, tetapi juga bagi para produsen dan penjual. Dengan mengelola media sosial secara efektif, hal ini bisa memberikan berbagai dampak mulai dari peningkatan perhatian hingga membangun kepercayaan publik terhadap suatu merek. "Berdasarkan deskripsi masalah di atas, diharapkan hasil penelitian ini meningkatkan wawasan dan pengetahuan praktis bagi para pelaku usaha yang sedang berusaha membangun merek melalui media sosial, khususnya pelaku usaha di bidang kosmetik dan Skincare"[5]. Sehingga mereka dapat memahami kekuatan dan kelemahan produk mereka serta merancang strategi pemasaran yang lebih efektif.

2. TINJAUAN PUSTAKA

2.1. Analasis Sentimen

Analisis sentimen adalah proses mengekstraksi opini, perasaan, dan emosi yang diekspresikan dalam teks untuk menentukan kecenderungan sentimen positif atau negatif [6]. Dalam konteks e-commerce dan media sosial, analisis sentimen membantu memahami persepsi konsumen terhadap produk dan layanan.

Menurut Liu, Analisis Sentimen adalah studi tentang menghitung opini, perasaan, dan emosi yang diekspresikan dalam teks. Menurut Hidayat, analisis sentimen diperlukan untuk mengklasifikasikan informasi tertentu dalam bentuk teks sebagai positif atau negatif. Analisis sentimen digunakan untuk menentukan kecenderungan suatu sentimen atau opini tertentu dan apakah itu merupakan opini positif atau negatif.[7]

2.2. Random Forest

Random Forest adalah algoritma ensemble learning yang membangun multiple decision trees dan menggabungkan hasilnya melalui voting untuk meningkatkan akurasi dan stabilitas prediksi [8]. Algoritma ini efektif menangani dataset kompleks dengan banyak fitur dan mengurangi risiko overfitting.

"Random Forest adalah Kumpulan beberapa pohon keputusan yang digunakan untuk membuat prediksi dengan membagi data menjadi beberapa kelas berdasarkan atribut tertentu dan membuat keputusan melalui perbandingan nilai-nilai spesifik." [9]

2.3. Naïve Bayes

Naïve Bayes adalah algoritma probabilistik berbasis Teorema Bayes dengan asumsi independensi antar fitur [10]. Meskipun asumsinya sederhana, algoritma ini efisien dan sering digunakan dalam klasifikasi teks karena kecepatan komputasinya.

2.4. Google Colab

Google Colab merupakan sebuah konsep pemrograman python dimana pemrosesan akan dilakukan oleh server google yang memiliki hardware berperforma tinggi pada sisi *software*, google colab telah menyediakan hampir semua *library* (pustaka) yang dibutuhkan [11]

2.5. Tiktok

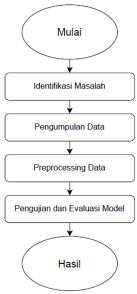
Tiktok adalah aplikasi yang menawarkan efek khusus yang unik dan menarik, mudah digunakan oleh siapa saja. Dengan Tiktok, pengguna bisa membuat video pendek yang luar biasa untuk dibagikan kepada teman-teman atau pengguna lainnya. Aplikasi video sosial ini juga dilengkapi dengan banyak pilihan musik, memungkinkan pengguna untuk menari, berimprovisasi, dan lebih banyak lagi. Ini semua mendorong kreativitas di antara penggunanya untuk menjadi pembuat konten yang hebat![12]

2.6. Penelitian Terkait

Beberapa penelitian terdahulu menunjukkan bahwa *Random Forest* umumnya menghasilkan akurasi lebih tinggi dibanding *Naïve Bayes* pada dataset kompleks [13], [14]. Namun, *Naïve Bayes* unggul dalam kecepatan pemrosesan pada dataset besar [15] dan contoh dalam penelitian Ryanizar [16], *Naïve Bayes* menunjukkan akurasi sebesar 76,00%.

3. METODE PENELITIAN

Metode penelitian ini mencakup proses dan tahapan-tahapan diambil untuk yang memastikan struktur penelitian terorganisir dengan baik dan keberhasilan untuk menyelesaikan tujuan penelitian ini dalam kinerja algoritma klasifikasi, yaitu random forest dan naïve bayes dalam menganalisis sentimen ulasan konsumen terhadap produk skincare lokal.



3.1.1. Identifikasi Masalah

Penelitian ini bertujuan untuk menganalisis sentimen dari ulasan konsumen tentang produk skincare lokal yang beredar di media sosial, terutama di TikTok. Penulis akan membandingkan efektivitas dua algoritma, yaitu random forest dan naïve bayes. Diharapkan, random forest dapat memberikan hasil yang lebih akurat dan stabil dengan menggabungkan informasi dari beberapa pohon keputusan. Sementara itu, naïve bayes akan diuji untuk melihat seberapa cepat dan efisien ia dapat mengklasifikasikan sentimen. Selain itu,

penelitian ini juga bertujuan untuk mengevaluasi perbandingan efektivitas kedua algoritma dalam mengklasifikasikan sentimen ulasan konsumen terhadap produk skincare lokal. Hasil analisis ini diharapkan dapat memberikan wawasan berharga bagi produsen untuk meningkatkan produk dan strategi pemasaran mereka.

Selain itu, penelitian ini memiliki tujuan untuk menganalisis tingkat akurasi, presisi, recall, dan fl-score menggunakan algoritma tersebut dalam mengidentifikasi sentimen ulasan konsumen, serta untuk mengetahui apakah pandangan konsumen terhadap produk tersebut cenderung positif atau negatif. Data yang digunakan dalam penelitian ini merupakan data yang diperoleh melalui teknik web scraping pada ulasan konsumen di tiktok menggunakan skrip python yang dijalankan di Google Colaboratory. Data yang telah dikumpulkan akan disimpan dalam format csv untuk dianalisis lebih lanjut.

3.1.2. Pengumpulan Data

Data ulasan produk skincare dikumpulkan dari TikTok periode 4 Agustus 2024 hingga 24 Juni 2025 menggunakan web scraping dengan skrip Python di Google Colaboratory. TikTok dipilih karena tingginya volume interaksi dan ulasan konsumen terkait produk skincare lokal. Data yang dikumpulkan mencakup teks ulasan, waktu pengunggahan, dan informasi relevan lainnya, kemudian format CSV untuk disimpan dalam memudahkan analisis.

3.1.3. Preprocessing Data

- 1. Tahapan Preprocessing Meliputi:
 - a. Data Cleaning: Menghapus komponen seperti simbol, angka, tanda baca, url, dan karakter khusus yang tidak relevan.
 - b. Case Folding: Mengubah semua teks menjadi huruf kecil untuk menghindari perbedaan yang disebabkan oleh kapitalisasi.
 - c. Tokenization: Memecah teks menjadi kata-kata (token).
 - d. Stopwords Removal: Menghapus kata-kata umum seperti "dan", "yang", "adalah" yang tidak memiliki nilai informasi.
 - e. Stemming: Mengubah kata-kata menjadi bentuk dasar atau akar katanya.

3.1.4. Pengujian dan Evaluasi Model

Penguiian dan evaluasi ini dilakukan untuk memastikan bahwa model mampu memberikan hasil yang akurat dalam menentukan apakah sentimen dalam ulasan konsumen bersifat positif atau negatif. Tujuan dari pengujian dan adalah untuk memastikan evaluasi ini keakuratan model dalam menganalisis sentimen ulasan konsumen. Dalam konteks ini, penulis membandingkan algoritma random forest dengan naïve bayes, di mana kedua metode ini diuji untuk melihat mana yang lebih efektif dalam menganalisis sentimen dari ulasan.

1. Pengujian Model

Model yang telah dilatih menggunakan data uji untuk memprediksi sentimen ulasan konsumen, proses ini akan menghasilkan prediksi dalam bentuk kategori sentimen (positif dan negatif). random forest yang menggunakan beberapa pohon keputusan diharapkan memiliki kemampuan yang lebih baik dalam menangani data yang lebih kompleks, sementara *naïve bayes*, dengan asumsinya yang lebih sederhana tentang independensi antar fitur. mungkin menghasilkan kecepatan lebih tinggi tetapi dengan potensi akurasi yang lebih rendah pada data yang lebih rumit. Prediksi yang dihasilkan dari kedua algoritma akan dibandingkan dengan label yang sebenarnya dalam data uji untuk mengukur sejauh mana model mengklasifikasikan sentimen dengan benar.

2. Evaluasi Kinerja Model

Akurasi (Accuracy): Mengukur seberapa banyak prediksi yang benar dibandingkan dengan total data yang diuji. akurasi dihitung dengan rumus:

 $Akurasi = \frac{\textit{Jumlah prediksi benar}}{\textit{Jumlah total data}}$

Random Forest diharapkan menunjukkan akurasi yang lebih tinggi karena kemampuannya untuk menggabungkan hasil dari banyak pohon keputusan, sementara naïve bayes mungkin memiliki akurasi yang lebih rendah, terutama ketika ada ketergantungan

antar fitur yang tidak dapat ditangani dengan baik oleh model ini.

Untuk menilai kinerja model klasifikasi, peneliti akan menganalisis dan mengevaluasi efektivitas Model Klasifikasi Random Forest dan Naïve bayes . Pada tahap ini, digunakan confusion matrix, sebuah format matriks yang memvisualisasikan performa model dengan membandingkan prediksinya dengan data aktual. Penilaian kinerja model mencakup berbagai aspek, seperti akurasi (accuracy), daya ingat (recall), presisi (precision), dan f1-score. Selain itu, dilakukan juga analisis distribusi kelas sentimen untuk memahami proporsi setiap kategori sentimen dalam dataset. Berikut beberapa rumus perhitungan:

a. Rumus Confusion Matrix:

TN FP FN TP

Keterangan:

Confusion matrix adalah matriks yang menunjukkan jumlah prediksi yang benar dan salah untuk setiap kelas. Matriks ini memiliki 4 elemen dasar:

- 1) TP (True Positive): Jumlah data yang benar diprediksi positif.
- 2) TN (True Negative): Jumlah data yang benar diprediksi negatif.
- 3) FP (False Positive): Jumlah data yang salah diprediksi positif.
- 4) FN (False Negative): Jumlah data yang salah diprediksi negatif.
- b. Rumus Recall

$$Recall = \frac{TP}{TP + FN}$$

Keterangan:

Recall mengukur proporsi data positif yang berhasil diprediksi dengan benar, dihitung sebagai:

- 1) TP (True Positive) adalah jumlah prediksi positif yang benar.
- 2) FN (False Negative) adalah jumlah data positif yang salah diprediksi negatif.
- c. Rumus Precision

$$Precision = \frac{TP}{TP + FP}$$

Keterangan:

Precision mengukur proporsi prediksi positif yang benar, dihitung sebagai:

- 1) TP (True Positive) adalah jumlah prediksi positif yang benar.
- 2) FP (False Positive) adalah jumlah prediksi positif yang salah.
- d. Rumus F1-score.

$$F1 - Score = 2 \times \frac{Precision \times Recall}{Precision + Recall}$$

Keterangan:

F1-score adalah harmonic mean antara precision dan recall.

3. Interpretasi Hasil:

Hasil evaluasi menunjukkan bahwa random forest diharapkan memberikan akurasi yang lebih tinggi dibandingkan naïve bayes, karena model ini menggunakan banyak pohon keputusan yang membuatnya lebih stabil dan andal dalam mengklasifikasikan data. Naïve Bayes, meskipun lebih cepat dan sederhana, dapat lebih rentan terhadap ketergantungan antar fitur yang tidak dapat ditangani dengan kedua baik. Hasil dari model akan menunjukkan bagaimana sentimen positif dan negatif dapat diidentifikasi dengan akurasi yang berbeda oleh masing-masing algoritma.

3.1 Metode Pengolahan Dan Analisis Data

Proses berikut digunakan untuk mengolah dan menganalisis data:

1. Pengolahan Data:

Data ulasan yang diperoleh melalui web scraping diproses terlebih dahulu untuk memastikan kualitas dan relevansinya dengan tujuan analisis. Dalam proses ini, digunakan metode *tf-idf* (*term frequency-inverse document* frequency) untuk memberikan bobot pada Teknik setiap kata dalam ulasan. ini menghitung pentingnya sebuah kata berdasarkan seberapa sering kata tersebut muncul dalam ulasan tertentu dan seberapa jarang kata itu muncul di seluruh kumpulan data, sehingga kata-kata yang lebih signifikan akan mendapatkan nilai yang lebih tinggi.

2. Penerapan Algoritma Machine Learning

Random Forest sendiri merupukan algoritma yang meningkatkan stabilitas dan akurasi prediksi dengan menggabungkan hasil

klasifikasi dari beberapa pohon keputusan yang berbeda, sehingga menghasilkan model yang lebih kuat dan andal.

Di sisi lain, *naïve bayes* adalah algoritma probabilistik yang bekerja dengan menghitung probabilitas kata-kata dalam data dan mengasumsikan independensi antar fitur. Meskipun *naïve bayes* lebih sederhana dan cepat dalam proses perhitungan, terutama pada dataset besar, model ini bisa kurang efektif dalam menangani interaksi kompleks antar fitur seperti yang bisa dilakukan oleh *random forest*.

3. Evaluasi Model

Untuk menilai seberapa baik kinerja model, kita biasanya menggunakan beberapa metrik evaluasi, seperti akurasi, presisi, recall, dan fl-score. Akurasi sendiri mengukur seberapa banyak prediksi yang benar dibandingkan dengan total prediksi yang dibuat oleh model. Dengan kata lain, metrik ini memberikan gambaran umum tentang seberapa akurat model dalam memprediksi sentimen secara keseluruhan.

Selain itu, presisi mengukur seberapa akurat prediksi positif yang dibuat oleh model, yaitu sejauh mana prediksi positif benar-benar sesuai dengan kenyataan. Sedangkan recall berfokus pada seberapa efektif model dalam mengidentifikasi semua data positif yang relevan dari keseluruhan data yang tersedia, mengindikasikan kemampuan model untuk menangkap semua kasus positif. Terakhir, f1score merupakan ukuran yang menggabungkan presisi dan recall, memberikan gambaran yang lebih menyeluruh tentang kinerja model secara keseluruhan. terutama ketika ada ketidakseimbangan antara presisi dan recall.

4. Visualisasi Data

Untuk mempermudah pemahaman hasil analisis, hasil-hasil yang diperoleh akan divisualisasikan dalam bentuk grafik atau diagram, seperti grafik yang menggambarkan distribusi sentimen ulasan, yaitu positif dan negatif, grafik yang menunjukkan kinerja algoritma berdasarkan metrik evaluasi yang telah diterapkan.

4. HASIL DAN PEMBAHASAN

4.1 Hasil Penelitian

Pengumpulan data dalam penelitian ini dilakukan dengan mengumpulkan ulasan dari pengguna tiktok mengenai produk *skincare* lokal. data yang digunakan mencakup ulasan yang dipublikasikan oleh pengguna dalam periode waktu 04 agustus 2024 - 24 juni 2025 dan menghasilkan 1464 ulasan.

4.1.1 Data Set

Dataset ini diambil dari platform ulasan tiktok, berdasarkan atribut seperti username pengguna, ulasan, dan juga tanggal. Scraping Data dilakukan menggunakan tools web scraping yang mendukung pengambilan data dari website berbasis visual. jumlah data dalam dataset ini adalah 1464 ulasan.

Ulasan	T T	
	Username	Waktu
		2/3/2025, 2:24:10
7 8 88	alakkk	PM
labore ga si? punya paragon itu		2/3/2025, 10:22:03
gong banget	whatuneedd	PM
		2/3/2025, 5:16:27
	nbiii21	PM
aubreeeee, beneran gue suka bgt		2/3/2025, 6:08:19
	lwirahmaaah	PM
TERATUUUUUU, plisss		
sunscreen nya bagus bangettt,		
mana ada sensasi dinginnya lagii		2/4/2025, 10:13:43
luvvv???????	callmemputtt	AM
		2/3/2025, 4:40:35
	niken_larashati	PM
tulus skin paling the best buat		
kulit menenangkan dan		
melembabkan sih, hanya hargaya		2/12/2025,
	ovelycreamy_21	11:09:12 AM
The Aubree, Avoskin, Bhumi, fss,		
dll. Banyak sebenernya, mana		2/5/2025, 10:58:49
produknya gong semua lagi y	akalipusing	AM
jarang ada yg bahas skinouru,		
padahal cleanser sm moistnya		2/8/2025, 1:50:20
cocok polll p	antiesushi	PM
		2/3/2025, 6:44:17
	whiphiphuraa	PM
aku jg baru kenal sama brand		
lokal chelvy, asli sih ini keknya		
emang ga overclaim bagus bgt		2/3/2025, 3:07:29
	ee.woona0	PM
VIVAAAAAAAAAA PLISS		3/5/2025, 1:36:02
	ecaccttaa	PM
setujuu, forebie bagus bagus bgt		2/3/2025, 2:46:53
	lexnaw	PM
npure kak ga kaleng kaleng bagus		
bgt cocok buat kulitku yg oily		4/27/2025, 2:30:36
	inaaaaash	PM
Hanasui lokal ga si? ,,, bagus bngt		
sih vit c nya ?? iseng2 bngt beli	1 1 2	4/24/2025, 5:47:43
,, eh dipakai 3hr ada perubahan p	palepaleeeee_2	AM
	,	4/24/2025,
	ouelcare	10:55:57 PM
eiem beauty serum niacinamide		4/8/2025, 1:22:03
	acun.fitrinj	PM
kakkk coba review make up		0/10/2005 0 50 5:
remover mireya itu bagus tp blm		2/10/2025, 2:59:51
banyak yg tauuu k	spgirl	PM

Dataset ini berisi informasi tentang beberapa data mentah dari ulasan pengguna produk *skincare* lokal yang di perlukan sebelum kita mendapatkan hasil perbandingan.

4.1.2 Hasil Pengujian

Hasil pengujian digunakan untuk mengevaluasi kemampuan algoritma dari random forest dan naïve bayes dalam mengklasifikasikan sentimen ulasan sekaligus membandingkan kinerja dari masing-masing algoritma. Pengujian Ini menggunakan dataset yang telah melalui proses preprocessing.

1. Pembersihan Data

Data yang tidak lengkap atau tidak relevan dihapus untuk memastikan kualitas data yang digunakan. Gambar tersebut menampilkan tabel yang berisi tahapan-tahapan dalam preprocessing data teks. Preprocess Text berikut merupakan preprocess text yang akan meliputi penggunaan case folding, cleaning, tokenization, stopwords removal dan stemming.

Stenned_Ulasan	Tokenized_Ulasan_Nithout_Stopwords	Tokenized_Ulasan	Cleaned_Ulasan	Processed_Illasan	liaktu	Usernane	Vlasan
teratu file aubree gongg sitti	Beats, the autree googs shift]	Jeratu the autree gongo, shift]	leratu the aubree gongg shint	ferati, the authree gongg sitht	23/2025; 2/24/10 PM	BBWK_	teratu, the autree gongg sithh
latore ga si paragon gong banget	(atore ga si paragon goog bangel)	(abore ga si, punya, paragon, itu, gong, ba .	labore ga si punya paragon itu gong banget	labore ga si? punya paragon itu gong bangat	23/2025; 10/22/03 PM	_whatuneedd	labore ga si? punya paragon itu gong bangal
wartah crystal secret best byt	[wardah, crystal, secret, best, byf]	jwardah, crystal, secret, best bgfj	wardah orjstal secret best byt	wardah crystal secret best bgt 77	2/3/2025; 5:16:27 PM	mbii21	wardah crystal secret best byt 77
aubressee beneran gue suka bgi serum centella	Jaubreeeee, beneran, gue, suka, byd senum, ce	jaubreesse, beneran, gue, suka, byt, serum, ce.	aubreeese beneran gue suka bgi serum centella	aubreeece, beneran gue suka byl serum centella	23/2025, 6:08:19 PM	dwishneech	aubreeese, beneran gue suka bgl serum cerdella
teratuuuuuu pisss sunscreen nya bagus banget.	Berahuuuuu, pisse, sunsceen, ma, bagus,	(feratuuuuu, pisss, sunscreen, nya, bagus,	teratuuuuuu pisss sunschen nya bagus tanget	teratuuuuuu, pisss sunscreen nya bagus banga	242025 _, 10:13:43:AM	_calmemputt	TERATUUUUUU pisss sunscreen nja bagus bange

Preprocessing data dilakukan melalui beberapa tahapan sistematis menggunakan library NLTK dan Sastrawi. Pertama, data CSV dimuat menggunakan pd.read_csv() dan diubah menjadi DataFrame untuk memudahkan manipulasi data.

Tahapan preprocessing meliputi:

- a. Case folding menggunakan str.lower() untuk mengubah semua teks menjadi huruf kecil agar kata seperti "Apel" dan "apel" diperlakukan sama.
- b. Data cleaning menggunakan regex untuk menghapus karakter tidak relevan seperti angka, tanda baca, dan spasi ekstra.

- c. Tokenization dengan word_tokenize() dari NLTK untuk memecah teks menjadi kata-kata terpisah.
- d. Stopwords removal menggunakan daftar stopwords NLTK untuk menghapus kata umum seperti "dan" dan "atau" yang tidak memberikan nilai informasi.
- e. Stemming menggunakan Sastrawi Stemmer untuk mengubah kata ke bentuk dasar, seperti "berlari" menjadi "lari".

Hasil preprocessing ditampilkan menggunakan df.head() untuk verifikasi, menghasilkan data bersih yang siap untuk analisis sentimen lebih lanjut.

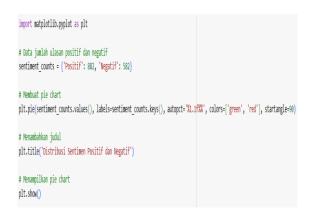
2. Analisis Sentimen

	Vlasan	Usernane	Waktu	Processed_Ulasan	Cleaned_Ulasan	Sentiment	Sentiment_Category
0	teratu, the aubree gongg sithih	lalaldk	2/3/2025, 2:24:10 PM	teratu, the aubree gongg sithh	teratu the aubree gongg sihih	1	Positi
1	labore ga s?? punya paragon itu gong bangel	_whatureedd	2/3/2025, 10:22:03 PM	labore ga si? punya paragon ilu gong bangel	labore ga si punya paragon itu gong bangel	1	Positi
2	wardah crystal secret best bgt 97	mbii21	2/3/2025, 5:16:27 PM	wardah crystal secret best bgt ??	wardah crystal secret best bgt	1	Positi
3	aubreeeee, beneran gue suka byt serum centella	dwrahmaaah	2/3/2025, 6:08:19 PM	aubreeeee, beneran gue suka bgt serum centella	aubreeeee beneran gue suka bgt serum centella 	1	Positi
4	TERATUUUUUU, pisss sunscreen nya bagus bange	_calmemputt	242025, 10:13:43 AM	teratuuuuuu, plisss sunscreen nya bagus bange	teratuuuuuu plisss sunscreen nya bagus banget	1	Positi
-	-	-					1
63	menurulku brand lokal harganya tu mahal isinya	bingungdeh_	1/10/2025, 12:16:05 AM	menuruku trand lokal harganya tu mahal isinya	menurutku brand lokal harganya tu mahal isinya	0	Negati
64	kebanyakan skincare lokal tuh kurang ini marke	nyonyahpuput25	1/12/2025, 7:08:07 PM	kebanyakan skincare lokal tuh kurang ini marke	kebanyakan skincare lokal tuh kurang ini marke	0	Negati
65	aku make dan suka banget moisturizer tamanu ru	arayrayrayraysu	1/13/2025, 7:35:44 AM	aku make dan suka bangel moisturizer lamanu ru	aku make dan suka bangel moisturzer tamanu ru	1	Positi
66	ak dulu salah1 distributor delmoraskin priya su_	iceleechytea	1/8/2025, 11:06:10 PM	ak dulu salah1 distributor delmoraskin pnya su	ak dulu salah distributor delmoraskin pnya sua	1	Positi
67	dokkkk//////naku secinta itu sama tamanu t	farrahfafa_	1/10/2025, 10:00:34 PM	dokkkki/7777777/naku secinta itu sama tamanu t	dokkkik aku secirta itu sama tamaru tapi semen	1	Positi

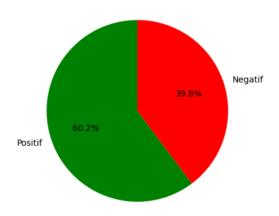
Klasifikasi sentimen dilakukan menggunakan fungsi sentiment label yang mengidentifikasi sentimen berdasarkan kata kunci. Daftar positive keywords mencakup kata seperti "bagus", "cocok", "senang", dan negative keywords sedangkan mencakup "jelek", "breakout", "bermasalah", dan "tidak puas". Fungsi ini diterapkan pada kolom 'Cleaned Ulasan' menggunakan apply() untuk memberikan label sentimen pada setiap ulasan. Jika teks mengandung kata positif, diberi label positif; jika mengandung kata negatif, diberi label negatif; dan jika tidak mengandung keduanya, label default adalah negatif. Label numerik kemudian dikonversi menjadi kategori menggunakan

sentiment_to_category dan disimpan dalam kolom 'Sentiment Category'.

Hasil analisis menunjukkan 882 ulasan positif (60,2%) dan 582 ulasan negatif (39,8%), mengindikasikan distribusi data tidak seimbang (imbalanced data) dengan kelas positif sebagai mayoritas. Kondisi ini dapat memengaruhi performa model karena cenderung lebih baik mengenali kelas mayoritas. Statistik distribusi sentimen divisualisasikan dalam diagram batang dan pie chart untuk memperjelas proporsi kedua kategori.

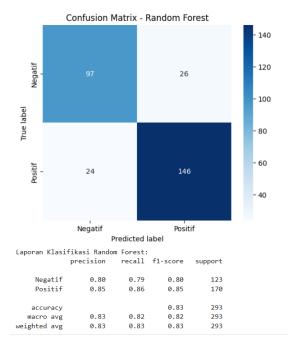


Distribusi Sentimen Positif dan Negatif



Distribusi sentimen divisualisasikan menggunakan diagram pie dengan fungsi plt.pie(). Data disimpan dalam dictionary sentiment_counts yang berisi 882 ulasan positif dan 582 ulasan negatif. Diagram menggunakan parameter autopct='%1.1f%%' untuk menampilkan persentase dengan satu desimal, warna hijau untuk sentimen positif dan merah untuk negatif, serta startangle=90 untuk rotasi estetik. Visualisasi ini memperjelas proporsi distribusi sentimen dalam dataset secara visual.

3. Hasil Klasifikasi

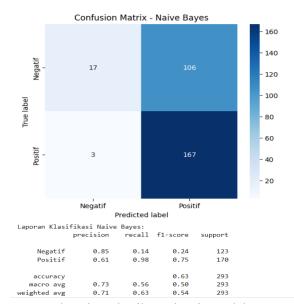


Berdasarkan hasil evaluasi model *Random* Forest menggunakan Confusion Matrix dan Classification Report, diperoleh performa yang baik dalam klasifikasi sentimen ulasan. Confusion Matrix menunjukkan 97 True Negative, 146 True Positive, 24 False Negative, dan 26 False Positive.

Classification Report memberikan rincian metrik sebagai berikut:

- Sentimen Negatif: *precision* 0.80, *recall* 0.79, dan *f1-score* 0.80.
- Sentimen Positif: *precision* 0.85, *recall* 0.86, dan *f1-score* 0.85.

Secara keseluruhan, model mencapai akurasi sebesar 0.83, dengan nilai *macro average* dan *weighted average* untuk *fl-score* juga sebesar 0.83. Hasil ini menunjukkan bahwa model memiliki kemampuan yang solid dalam mengklasifikasikan sentimen, dengan performa yang sedikit lebih unggul pada kategori sentimen positif.



Berdasarkan hasil evaluasi model *Naive Bayes* menggunakan *Confusion Matrix* dan *Laporan Klasifikasi*, ditemukan performa yang kurang seimbang dalam klasifikasi sentimen ulasan. *Confusion Matrix* menunjukkan 17 *True Negative*, 167 *True Positive*, 3 *False Negative*, dan 106 *False Positive*.

Laporan Klasifikasi memberikan rincian metrik sebagai berikut:

- Sentimen Negatif: *precision* 0.85, *recall* 0.14, dan *f1-score* 0.24.
- Sentimen Positif: *precision* 0.61, *recall* 0.98, dan *fl-score* 0.75.

Secara keseluruhan, model ini mencapai akurasi sebesar 0.63, dengan nilai *macro average f1-score* sebesar 0.50. Hasil ini mengindikasikan bahwa model *Naive Bayes* memiliki performa yang tidak seimbang. Model ini menunjukkan kemampuan yang sangat tinggi dalam mengenali sentimen positif (*recall* 0.98), namun sangat lemah dalam mengidentifikasi sentimen negatif (*recall* 0.14), yang menandakan adanya bias prediksi yang signifikan terhadap kelas positif.

4. Akurasi Model

```
import matplotlib.pyplot as plt

# Data Akurasi Model
models = ['Naive Bayes', 'Random Forest']
accuracies = [0.6280, 0.6284] # Akurasi dari Naive Bayes dan Random Forest

# Membuat Bar Chart
plt.figure(figsize=(8, 5))
plt.bar(models, accuracies, color=['skyblue', 'pink'])

# Menambahkan label dan judul
plt.title('Perbandingan Akurasi Model Naive Bayes dan Random Forest')
plt.xlabel('Nkouel')
plt.ylabel('Nkouesi')

# Menambahkan nilai di atas bar
for i, acc in enumerate(accuracies):
    plt.text(i, acc + 0.005, f'(acc:.4f)', ha='center', va='bottom')

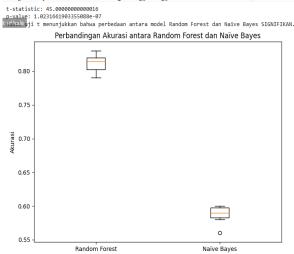
# Menampilkan grafik
plt.tight_layout()
plt.show()
```



Untuk memvisualisasikan perbandingan kinerja, dibuat sebuah diagram batang yang menampilkan akurasi dari model *Naive Bayes* dan *Random Forest*. Diagram ini menggunakan nilai akurasi yang telah dihitung sebelumnya, yaitu 0.6280 untuk *Naive Bayes* dan 0.8294 untuk *Random Forest*.

Setiap batang pada diagram diberi anotasi nilai akurasinya untuk memberikan informasi yang presisi. Grafik ini juga dilengkapi dengan judul dan label sumbu yang jelas untuk mempermudah interpretasi.

Secara keseluruhan, visualisasi ini secara efektif mengilustrasikan bahwa model *Random Forest* memiliki performa akurasi yang jauh lebih unggul dibandingkan dengan model *Naive Bayes* pada dataset yang digunakan.



Sebuah analisis perbandingan dilakukan untuk mengevaluasi kinerja model *Random Forest* dan *Naive Bayes* pada dataset yang sama, dengan fokus pada metrik akurasi. Hasil observasi menunjukkan bahwa model *Random Forest* secara konsisten mencapai akurasi yang lebih tinggi, dengan rentang nilai antara 0.79 hingga 0.82, sementara model *Naive Bayes* mencatatkan akurasi pada rentang 0.58 hingga 0.62.

Untuk memvalidasi signifikansi perbedaan kinerja ini secara statistik, dilakukan uji t-berpasangan (*paired t-test*). Uji ini bertujuan untuk menentukan apakah selisih akurasi antara kedua model tersebut signifikan dengan membandingkan p-value yang dihasilkan dengan tingkat signifikansi ($\alpha = 0.05$).

Sebagai pelengkap, perbandingan kinerja divisualisasikan menggunakan box plot. Visualisasi ini memberikan gambaran yang jelas mengenai sebaran, median, dan variasi data akurasi dari masing-masing model, sehingga mempermudah interpretasi perbandingan performa keduanya secara visual.

4.2 Pembahasan

Penelitian ini membandingkan kinerja model *Random Forest* dan *Naïve Bayes* untuk tugas klasifikasi sentimen pada ulasan produk *skincare* lokal. Hasil perbandingan menunjukkan bahwa model Random Forest secara signifikan lebih unggul daripada *Naïve Bayes*.

Model *Random Forest* mencapai akurasi sebesar 82%, sementara *Naïve Bayes* hanya mencapai 62%. Keunggulan ini disebabkan oleh kemampuan *Random Forest* dalam menangani dataset yang kompleks dan interaksi antar fitur, berbeda dengan *Naïve Bayes* yang mengasumsikan independensi fitur.

Analisis lebih lanjut menggunakan confusion matrix mengonfirmasi superioritas Random Forest, yang mencatatkan jumlah True Positives dan True Negatives lebih tinggi. Selain itu, metrik precision, recall, dan F1-score untuk model Random Forest juga secara konsisten lebih tinggi, yang mengindikasikan bahwa model ini lebih akurat dan efisien dalam meminimalkan kesalahan klasifikasi untuk kedua kategori sentimen.

5. KESIMPULAN

Berdasarkan hasil analisis dan pembahasan yang telah dilakukan dalam penelitian ini mengenai perbandingan algoritma *random forest* dan *naïve bayes* dalam analisis sentimen ulasan konsumen produk *skincare* lokal di media sosial tiktok, dapat ditarik beberapa kesimpulan sebagai berikut:

- Kinerja Algoritma: Algoritma random forest menunjukkan performa yang lebih dibandingkan naïve bayes dalam melakukan klasifikasi sentimen ulasan konsumen. Random Forest mencapai akurasi keseluruhan sebesar 82%, dengan precision, recall, dan f1-score yang konsisten di angka 81%-85% untuk sentimen positif maupun negatif. Hal ini mengindikasikan kemampuan random forest dalam menangani kompleksitas data dan interaksi antar fitur dengan lebih baik. Sementara itu, naïve bayes mencapai akurasi 62%, dengan kinerja yang lebih bervariasi antar kelas sentimen.
 - 2. Efektivitas *Random Forest*: terlihat dari keunggulannya dalam menggabungkan banyak pohon keputusan. Ini mengurangi risiko overfitting dan sekaligus meningkatkan stabilitas serta akurasi prediksi. Jadi, metode ini jadi pilihan yang lebih kuat untuk analisis sentimen pada dataset ulasan konsumen yang bervariasi.
 - 3. Distribusi Sentimen Konsumen: Dari total 1464 ulasan, sentimen negatif mencapai (39.8%)dibandingkan sentimen positif yang lebih unggu mencapai (60,2%). Distribusi ini memberikan gambaran umum bahwa meskipun produk skincare lokal mendapatkan ulasan positif, masih ada bagi produsen untuk ruang meningkatkan kepuasan konsumen dan mengatasi isu-isu yang menyebabkan sentimen negatif.
 - 4. Wawasan Preferensi Konsumen:
 Analisis Tf-Idf terhadap kata-kata
 penting dalam ulasan memberikan
 wawasan berharga mengenai preferensi
 dan kekhawatiran konsumen. kata-kata
 seperti "bagus" dan "banget" sering
 mengindikasikan aspek positif produk,
 sementara kata "ga" (tidak) dan
 "produk" sering muncul dalam konteks

- evaluasi, baik positif maupun negatif, menyoroti aspek-aspek yang menjadi perhatian konsumen.
- 5. Manfaat Penelitian: Hasil penelitian ini dapat dimanfaatkan oleh konsumen panduan dalam memilih sebagai produk skincare lokal yang sesuai berdasarkan sentimen mayoritas. Bagi produsen dan penjual, analisis sentimen menjadi alat penting untuk memahami persepsi pasar, mengidentifikasi kekuatan dan kelemahan produk, serta merumuskan strategi pemasaran dan pengembangan produk yang lebih tepat sasaran di platform media sosial seperti tiktok.

UCAPAN TERIMA KASIH

Penulis mengucapkan terima kasih yang tulus kepada semua pihak yang telah memberikan bantuan dan dukungan dalam penyelesaian jurnal penelitian ini.

DAFTAR PUSTAKA

- [1] K. Sriwenda Putri, R. Setiawan, and A. Pambudi, "Analisis Sentimen Terhadap Brand Skincare Lokal Menggunakan Naïve Bayes Classifier," 2023.
- [2] H. Harnelia, "Analisis Sentimen Review Skincare Skintific Dengan Algoritma Support Vector Machine (Svm)," *Jurnal Informatika dan Teknik Elektro Terapan*, vol. 12, no. 2, Apr. 2024, doi: 10.23960/jitet.v12i2.4095.
- [3] D. Prihartini and R. Damastuti, "Pengaruh e-WOM terhadap Minat Beli Skincare Lokal pada Followers Twitter @ohmybeautybank," 2022.
- [4] N. Nawiyah, R. C. Kaemong, M. A. Ilham, and F. Muhammad, "Penyebab Pengaruhnya Pertumbuhan Pasar Indonesia Terhadap Produk Skin Care Lokal Pada Tahun 2022," *Armada: Jurnal Penelitian Multidisiplin*, vol. 1, no. 12, pp. 1390–1396, Dec. 2023, doi: 10.55681/armada.v1i12.1060.
- [5] R. Damastuti, "Membedah Feeds Instagram Produk Skincare Lokal (Analisis Isi Kuantitatif Akun Instagram Avoskin) Discovering Local Skincare Product Instagram Feeds (Quantitative Content Analysis Instagram Account Avoskin),"

 Des, vol. 5, no. 2, pp. 189–199, 2021.
- [6] M. Z. Rahman, Y. A. Sari, and N. Yudistira, "Analisis Sentimen Tweet COVID-19

- menggunakan Word Embedding dan Metode Long Short-Term Memory (LSTM)," 2021. [Online]. Available: http://j-ptiik.ub.ac.id
- [7] M. N. Muttaqin and I. Kharisudin, "Analisis Sentimen Pada Ulasan Aplikasi Gojek Menggunakan Metode Support Vector Machine dan K Nearest Neighbor," *UNNES Journal of Mathematics*, vol. 10, no. 2, pp. 22–27, 2021, [Online]. Available: http://journal.unnes.ac.id/sju/index.php/ujm
- [8] G. E. Pratiwi and A. Nugroho, "Implementasi Metode Random Forest Untuk Klasifikasi Penjualan Produk Sabun Paling Laris," *Jurnal Teknik Informasi dan Komputer (Tekinkom)*, vol. 7, no. 2, p. 541, Dec. 2024, doi: 10.37600/tekinkom.v7i2.1610.
- [9] Intan Permata and Esther Sorta Mauli Nababan, "Application Of Game Theory In Determining Optimum Marketing Strategy In Marketplace," *Jurnal Riset Rumpun Matematika Dan Ilmu Pengetahuan Alam*, vol. 2, no. 2, pp. 65–71, Jul. 2023, doi: 10.55606/jurrimipa.v2i2.1336.
- [10] A. Nugroho and Y. Religia, "Analisis Optimasi Algoritma Klasifikasi Naive Bayes menggunakan Genetic Algorithm dan Bagging," *Jurnal RESTI (Rekayasa Sistem dan Teknologi Informasi)*, vol. 5, no. 3, pp. 504–510, Jun. 2021, doi: 10.29207/resti.v5i3.3067.
- [11] R. Gelar Guntara, "Pemanfaatan Google Colab Untuk Aplikasi Pendeteksian Masker Wajah Menggunakan Algoritma Deep Learning YOLOv7," *Jurnal Teknologi Dan Sistem Informasi Bisnis*, vol. 5, no. 1, pp. 55–60, Feb. 2023, doi: 10.47233/jteksis.v5i1.750.
- [12] N. Ilahin, "Pengaruh Pengunaan Media Sosial Tik-Tok Terhadap Karakter Siswa Kelas V Madrasah Ibtidaiyah", doi: 10.37850/ibtida.
- [13] S. Ubaidillah Royan, N. Suarna, I. Ali, and D. Solihudin, "ANALISIS Sentimen Ulasan Produk Skincare Di Shopee Untuk Meningkatkan Kualitas Produk Menggunakan Metode Support Vector Machine," *Jurnal informasi dan Komputer*, vol. 13, no. 1, p. 2025.
- [14] H. Harnelia, "Analisis Sentimen Review Skincare Skintific Dengan Algoritma Support Vector Machine (Svm)," *Jurnal Informatika dan Teknik Elektro Terapan*, vol. 12, no. 2, Apr. 2024, doi: 10.23960/jitet.v12i2.4095.
- [15] A. F. Setyaningsih, D. Septiyani, and S. R. Widiasari, "Implementasi Algoritma Naïve

- Bayes untuk Analisis Sentimen Masyarakat pada Twitter mengenai Kepopuleran Produk Skincare di Indonesia," *Jurnal Teknologi Informatika dan Komputer*, vol. 9, no. 1, pp. 224–235, May 2023, doi: 10.37012/jtik.v9i1.1409.
- [16] P. Mixue, D. Metode, N. Bayes, and C. Dan, "Skripsi Analisis Sentimen Pada Media Sosial Twitter Terhadap."