

Komparasi Naïve Bayes, SVM, dan Decision Tree untuk Klasifikasi Komentar Provokatif pada Instagram Terkait Aksi Demonstrasi Agustus 2025

Muhamad Yusuf*, Erizal

¹ Fakultas Teknologi Industri dan Informatika, Program Studi Teknik Informatika, Universitas Muhammadiyah Prof. Dr. Hamka, Jakarta, Indonesia

² Fakultas Teknologi Industri dan Informatika, Program Studi Sistem dan Teknologi Informasi, Universitas Muhammadiyah Prof. Dr. Hamka, Jakarta, Indonesia

Email: ¹*my8486460@gmail.com, ²erizal@uhamka.ac.id

Email Penulis Korespondensi: my8486460@gmail.com

Submitted: 19/05/2026; Accepted: 22/05/2026; Published: 23/05/2026

Abstrak—Perkembangan media sosial menjadikan ruang digital sebagai tempat masyarakat menyampaikan pendapat terhadap isu sosial dan politik. Salah satu peristiwa yang banyak memunculkan komentar netizen adalah aksi demonstrasi 25–31 Agustus 2025. Banyak komentar yang mengandung ujaran kasar, hasutan, penghinaan, maupun ajakan konflik sehingga berpotensi memicu emosi negatif di media sosial. Penelitian ini dilakukan dengan tujuan mengklasifikasikan komentar provokatif dan non-provokatif pada Instagram menggunakan algoritma Naïve Bayes serta membandingkan performanya melalui pendekatan Support Vector Machine (SVM) dan Decision Tree. Data penelitian diperoleh dengan menerapkan teknik web scraping dan menghasilkan sebanyak 3396 komentar. Setelah dilakukan cleansing dan preprocessing, jumlah data menjadi 2490 komentar. Tahap preprocessing meliputi transform case, tokenizing, stopwords removal, filter tokens, dan stemming. Berikutnya dilakukan penghitungan bobot kata menggunakan metode TF-IDF serta implementasi algoritma menggunakan aplikasi RapidMiner dengan pembagian data 80:20. Berdasarkan hasil pelabelan manual diperoleh 1279 komentar provokatif dan 1211 komentar non-provokatif. Hasil evaluasi menunjukkan bahwa algoritma Naïve Bayes memperoleh accuracy sebesar 72,15%, SVM sebesar 69,44%, dan Decision Tree sebesar 72,91%. Penelitian ini berkontribusi dalam memberikan gambaran efektivitas algoritma Naïve Bayes, SVM, dan Decision Tree untuk klasifikasi komentar provokatif pada konteks demonstrasi politik di Instagram. Hasil penelitian menunjukkan bahwa Naïve Bayes memiliki performa yang lebih seimbang dalam mendeteksi komentar provokatif dan non-provokatif sehingga berpotensi mendukung proses moderasi konten di media sosial.

Kata Kunci: Komentar Provokatif; Naïve Bayes; SVM; Decision Tree; Text Mining

Abstract—The accelerating growth of social media has transformed digital platforms into spaces where people express opinions on social and political issues. One event that generated numerous public comments was the demonstration held on August 25–31, 2025. Many comments contained harsh language, provocation, insults, and calls for conflict that had the potential to trigger negative emotions in digital spaces. This study aims to classify provocative and non-provocative comments on Instagram using the Naïve Bayes algorithm and compare its performance with Support Vector Machine (SVM) and Decision Tree algorithms. The data were collected through a web scraping process, resulting in 3,396 comments. After the cleansing and preprocessing stages, the dataset was reduced to 2,490 comments. The preprocessing stages included transform case, tokenizing, stopwords removal, filter tokens, and stemming. Furthermore, word weighting was carried out using the TF-IDF method and implemented in RapidMiner with an 80:20 data split ratio. Based on manual labeling, 1,279 provocative comments and 1,211 non-provocative comments were obtained. The evaluation results showed that Naïve Bayes achieved an accuracy of 72.15%, SVM achieved 69.44%, and Decision Tree achieved 72.91%. Although Decision Tree produced a slightly higher accuracy, Naïve Bayes demonstrated a more balanced performance in detecting both comment classes, even though the accuracy value was still in the moderate category. The findings provide insights into the effectiveness of machine learning algorithms for identifying provocative comments and may support the development of automated content moderation on social media platforms.

Keywords: Provocative Comments; Naïve Bayes; SVM; Decision Tree; Text Mining

1. PENDAHULUAN

Pesatnya perkembangan di bidang teknologi informasi dan komunikasi menyebabkan media sosial berkembang menjadi salah satu media unggulan masyarakat dalam mendapatkan serta menyebarkan informasi. Media sosial kini tidak sekadar digunakan sebagai media hiburan, melainkan juga sebagai ruang publik digital untuk menyalurkan opini mengenai berbagai permasalahan sosial dan politik [1]. Salah satu platform media sosial yang kerap diminati oleh kalangan masyarakat adalah Instagram. Instagram menyediakan fitur komentar yang memungkinkan pengguna memberikan tanggapan terhadap suatu unggahan dalam bentuk teks, emoji, maupun simbol tertentu [2].

Kemudahan dalam menyampaikan pendapat di media sosial menyebabkan munculnya berbagai jenis komentar, baik yang bersifat positif maupun negatif. Tidak sedikit komentar yang mengandung unsur provokatif, seperti ujaran kasar, hasutan, ajakan konflik, maupun pernyataan yang dapat memicu perpecahan di ruang digital. Komentar provokatif dapat menyebar dengan cepat dan berpotensi menimbulkan konflik sosial apabila tidak dikendalikan dengan baik. Menurut penelitian sebelumnya, ujaran provokatif di media sosial dapat mempercepat eskalasi konflik politik baik di ranah daring maupun luring [3].

Salah satu peristiwa yang banyak menimbulkan komentar netizen di media sosial adalah aksi demonstrasi yang berlangsung pada tanggal 25–31 Agustus 2025 di depan Gedung DPR RI. Demonstrasi tersebut menjadi perhatian publik karena memunculkan berbagai opini dan perdebatan di media sosial, khususnya Instagram. Banyak pengguna

Instagram memberikan komentar terhadap aksi demonstrasi tersebut, baik berupa dukungan, kritik, maupun komentar yang mengandung unsur provokatif. Komentar-komentar tersebut tidak hanya berasal dari masyarakat umum, tetapi juga dari berbagai kelompok dengan sudut pandang yang berbeda sehingga menimbulkan pro dan kontra di ruang digital.

Banyaknya tanggapan yang muncul di platform media sosial menyebabkan proses analisis secara manual menjadi kurang efektif. Oleh karena itu, diperlukan metode otomatis dalam rangka mengelompokkan komentar berdasarkan pengelompokan tertentu. Salah satu pendekatan yang bisa diterapkan adalah penambangan teks. Text mining merupakan proses pengolahan data teks untuk memperoleh informasi tertentu dari kumpulan data tidak terstruktur [4]. Dalam penerapannya, text mining sering digunakan pada analisis sentimen maupun klasifikasi teks di media sosial.

Dalam klasifikasi teks, salah satu model komputasi yang umum dipilih adalah Naïve Bayes. Pendekatan tersebut banyak digunakan karena memiliki proses perhitungan yang sederhana, mampu mengolah data teks dengan baik, serta memiliki performa yang cukup stabil dalam proses klasifikasi [5]. Di samping itu, algoritma Naïve Bayes juga dapat mengolah data berdimensi tinggi seperti teks serta memiliki proses komputasi yang lebih cepat dibandingkan metode klasifikasi lainnya [6].

Beberapa penelitian sebelumnya telah menerapkan algoritma Naïve Bayes pada proses klasifikasi komentar di media sosial. Penelitian yang dilakukan oleh Baehaqi dan Cahyono melalui penggunaan algoritma Naïve Bayes untuk menganalisis komentar cyberbullying pada Instagram dan memperoleh nilai accuracy sebesar 88% [7]. Penelitian lain yang dijalankan oleh Irsyad menggunakan pendekatan Multinomial Naïve Bayes untuk analisis sentimen komentar Instagram terkait cyberbullying dan memperoleh nilai accuracy sebesar 84% [8]. Karim juga menerapkan algoritma Naïve Bayes pada komentar Instagram layanan BPJS dan memperoleh nilai accuracy sebesar 86,67% [9]. Selanjutnya, penelitian yang dilakukan oleh Warraihan membandingkan metode Naïve Bayes dan KNN pada komentar pengguna layanan transportasi online Maxim di Instagram, di mana temuan pengujian membuktikan bahwa pendekatan Naïve Bayes memiliki tingkat performa yang jauh lebih stabil jika dibandingkan dengan KNN [10]. Kajian lain yang telah dilakukan oleh Mubarak et al. membahas ujaran kebencian pada kolom komentar Instagram dan menunjukkan bahwa bentuk komentar yang paling dominan adalah insult dan provokasi [11].

Berdasarkan beberapa penelitian tersebut, dapat diketahui bahwa model klasifikasi Naïve Bayes cukup banyak digunakan dalam klasifikasi teks media sosial karena memiliki proses komputasi yang sederhana dan performa yang relatif stabil [5]. Namun, sebagian besar penelitian sebelumnya lebih banyak mengutamakan analisis sentimen yang bersifat umum, cyberbullying, maupun komentar layanan publik. Penelitian mengenai klasifikasi komentar provokatif pada konteks demonstrasi politik di media sosial Instagram masih sangat terbatas. Selain itu, komentar pada media sosial Instagram umumnya menggunakan bahasa informal, singkatan, slang, campuran bahasa, serta konteks kalimat yang ambigu sehingga proses klasifikasi teks menjadi lebih kompleks dibandingkan data teks formal [12]. Oleh karena itu, penelitian ini dilakukan untuk menganalisis klasifikasi komentar provokatif pada konteks demonstrasi politik di Instagram dengan karakteristik bahasa media sosial yang tidak terstruktur dan dinamis. Perbedaan nilai akurasi dibandingkan penelitian sebelumnya dapat disebabkan oleh perbedaan karakteristik data, kompleksitas konteks komentar demonstrasi, serta penggunaan dataset yang lebih beragam dan tidak terstruktur.

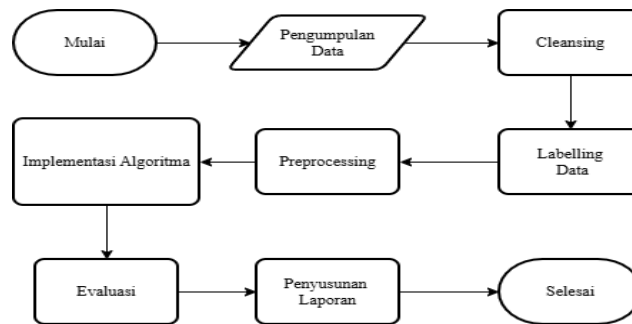
Pada penelitian ini dilakukan klasifikasi komentar provokatif dan non-provokatif pada media sosial Instagram terkait aksi demonstrasi 25–31 Agustus 2025 menerapkan metode Naïve Bayes. Selain hal tersebut, penelitian ini juga membandingkan performa algoritma Naïve Bayes dengan Support Vector Machine (SVM) dan Decision Tree guna mengetahui metode yang memiliki performa terbaik dalam rangka proses klasifikasi komentar. Data penelitian diperoleh melalui teknik web scraping pada komentar Instagram yang berkaitan dengan demonstrasi tersebut. Selanjutnya dilakukan tahap preprocessing yang terdiri dari transform case, tokenizing, stopwords removal, filter tokens, dan stemming. Setelah itu dilakukan pembobotan kata dengan menggunakan pendekatan TF-IDF sebelum proses klasifikasi dilakukan.

Penelitian ini bertujuan untuk mengetahui jumlah komentar provokatif dan non-provokatif serta mengukur serta membandingkan kemampuan algoritma Naïve Bayes, SVM, dan Decision Tree dalam melakukan klasifikasi komentar pada media sosial Instagram. Kontribusi penelitian ini terletak pada analisis klasifikasi komentar provokatif pada konteks demonstrasi politik di Instagram yang masih jarang dikaji dalam penelitian sebelumnya. Selain itu, penelitian ini menyajikan perbandingan performa Naïve Bayes, SVM, dan Decision Tree pada data komentar media sosial yang memiliki karakteristik bahasa informal, slang, singkatan, dan konteks ambigu sehingga dapat menjadi referensi dalam pengembangan sistem moderasi konten otomatis.

2. METODOLOGI PENELITIAN

2.1 Alur Penelitian

Data penelitian yang digunakan berasal dari komentar Instagram yang berkaitan dengan aksi demonstrasi 25–31 Agustus 2025 di Indonesia. Penelitian ini menerapkan pendekatan klasifikasi untuk mengidentifikasi komentar provokatif dan non-provokatif pada media sosial Instagram. Tahapan penelitian dilakukan secara sistematis mulai dari pengumpulan data, preprocessing, pembobotan kata menggunakan TF-IDF, implementasi algoritma, hingga proses evaluasi model. Tahapan atau alur penelitian yang digunakan oleh peneliti sebagaimana terlihat pada Gambar 1.



Gambar 1. Alur Penelitian

Berikut merupakan penjelasan lebih rinci mengenai tahapan proses penelitian yang ditunjukkan pada Gambar 1, yang meliputi beberapa langkah sebagai berikut:

a. Pengumpulan Data (Web Scraping)

Merupakan proses pengambilan data komentar dari media sosial Instagram yang berkaitan dengan aksi demonstrasi 25–31 Agustus 2025. Proses ini dilakukan secara otomatis menggunakan teknologi web scraping dengan bantuan Node.js untuk memperoleh data berupa teks komentar, username, dan jumlah likes. Metode ini cukup populer digunakan dalam penelitian berbasis media sosial karena memiliki kemampuan mengumpulkan data dalam jumlah banyak secara efisien [13].

b. Cleansing

Cleansing merupakan proses pembersihan kumpulan data dari komponen yang tidak relevan seperti simbol, tanda-tanda baca, angka, emoji, dan juga mention. Tahapan ini bertujuan untuk menghilangkan noise dengan demikian data berubah menjadi lebih berkualitas dan telah siap untuk dimanfaatkan dalam proses analisis [14].

c. Labelling Data

Labelling data adalah proses pemberian label pada setiap komentar yang telah dikumpulkan. Pada penelitian ini, data digolongkan ke dalam dua kelas, yakni komentar provokatif dan non-provokatif. Pelabelan dilakukan secara manual oleh peneliti berdasarkan konteks kalimat dan indikator komentar provokatif, seperti ujaran kasar, hasutan, penghinaan, dan ajakan konflik, sehingga hasil klasifikasi dapat lebih memahami konteks data [15]. Untuk menjaga konsistensi pelabelan, peneliti menggunakan kriteria yang sama pada seluruh data komentar.

d. Preprocessing

Tahap preprocessing merupakan proses pengolahan data teks agar menjadi lebih terstruktur. Tahapan yang dilakukan meliputi transform case, tokenizing, stopwords removal, filter tokens, dan stemming. Proses tahapan ini bertujuan meningkatkan kualitas data sebelum proses klasifikasi [16].

e. Implementasi Algoritma

Proses ini dilaksanakan dengan menerapkan algoritma Naïve Bayes, Support Vector Machine (SVM), dan Decision Tree untuk mengklasifikasikan data komentar ke dalam kategori provokatif dan non-provokatif. Algoritma Naïve Bayes dikenal efektif dalam klasifikasi teks karena sederhana dan memiliki performa yang relatif stabil [5], sedangkan SVM dan Decision Tree banyak digunakan sebagai metode pembandingan dalam penelitian klasifikasi teks [17]. Implementasi algoritma SVM dilakukan menggunakan parameter bawaan RapidMiner dengan kernel type berupa dot, convergence epsilon sebesar 0.001, serta maximum iterations sebanyak 10000. Selain itu, fitur scaling diaktifkan untuk membantu proses normalisasi data. Implementasi algoritma Decision Tree menggunakan criterion gain ratio dengan maximal depth sebesar 10. Pada proses pembentukan pohon keputusan diterapkan pruning dan prepruning dengan confidence sebesar 0.1, minimal gain sebesar 0.01, minimal leaf size sebanyak 2, serta minimal size for split sebesar 4.

f. Evaluasi

Evaluasi dilakukan untuk mengetahui kinerja model klasifikasi yang telah dibangun. Metode evaluasi yang diterapkan berupa confusion matrix dengan parameter pengukuran meliputi accuracy, precision, recall, dan F1-score [18].

g. Penyusunan Laporan

Proses terakhir pada penelitian ini berupa penyusunan laporan penelitian yang berisi seluruh proses dimulai dari pengumpulan data hingga tahap evaluasi hasil, serta analisis terhadap hasil yang diperoleh.

2.2 Web Scraping

Pengumpulan data diperoleh dengan menerapkan teknik web scraping berbasis Node.js untuk memperoleh komentar Instagram yang berkaitan dengan aksi demonstrasi 25–31 Agustus 2025. Metode ini memungkinkan pengambilan data secara otomatis dan efisien dalam jumlah besar [13].

2.3 Data Preprocessing

Preprocessing merupakan tahapan pembuka dalam pengolahan data berbasis teks yang dilakukan guna menyaring serta mempersiapkan kumpulan data sebelum digunakan pada proses analisis. Tahap ini memiliki peran penting dalam

meningkatkan kualitas data dengan mengurangi noise dan menyederhanakan struktur teks [16]. Langkah-langkah preprocessing yang diterapkan pada penelitian ini terdiri dari:

- a. Transform Case (Case Folding) merupakan tahapan mentransformasi seluruh teks menjadi huruf kecil (lowercase) dengan tujuan menyeragamkan gaya penulisan serta menghindari adanya perbedaan makna yang disebabkan oleh penggunaan huruf kapital [17].
- b. Tokenizing adalah proses memecah teks menjadi satuan-satuan kata (token) sehingga mempermudah tahapan analisis lebih lanjut pada setiap kata [19].
- c. Stopwords Removal merupakan proses membuang kata-kata umum yang tidak memiliki nilai semantik yang signifikan, seperti kata hubung atau kata yang frekuensi kemunculannya tinggi dan tidak memberikan manfaat yang signifikan dalam proses analisis [20].
- d. Filter Tokens merupakan proses penyaringan kata berdasarkan jumlah karakter tertentu agar kata-kata yang kurang informatif dapat dihilangkan. Pada penelitian ini diterapkan batas minimal panjang kata sebanyak 4 karakter sehingga kata dengan jumlah karakter kurang dari 4 tidak digunakan dalam proses klasifikasi [12].
- e. Stemming merupakan proses mentransformasi kata berimbuhan menjadi bentuk dasar guna menyederhanakan variasi kata agar memiliki representasi yang seragam [21].

2.4 TF-IDF

Setelah proses pengolahan awal selesai dilaksanakan, tahap berikutnya adalah pembobotan kata melalui metode TF-IDF (Term Frequency–Inverse Document Frequency). Pendekatan ini diterapkan untuk menentukan tingkat signifikansi suatu kata dalam dokumen berdasarkan intensitas kemunculannya. TF-IDF membantu meningkatkan kualitas representasi fitur dengan mengalokasikan bobot yang lebih tinggi pada kata yang relevan dan lebih rendah pada kata yang umum [22].

2.5 Naïve Bayes

Naïve Bayes merupakan salah satu pendekatan pengklasifikasian berbasis peluang yang dibangun atas teorema Bayes dengan asumsi bahwa masing-masing fitur bersifat independen. Pendekatan ini kerap digunakan dalam text mining karena sederhana, efisien, serta mampu memproses data dengan dimensi tinggi seperti data teks [5]. Penelitian ini menerapkan metode Multinomial Naïve Bayes karena sangat sesuai diterapkan untuk pengolahan data teks berbasis frekuensi kata maupun bobot TF-IDF. Algoritma ini bekerja dengan mengukur peluang sebuah dokumen termasuk pada kelas tertentu berdasarkan frekuensi kemunculan kata-kata di dalamnya. Persamaan teorema Naïve Bayes yang digunakan adalah sebagai berikut:

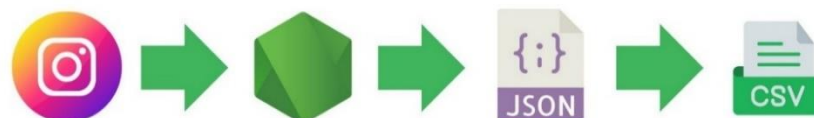
$$P(C|A) = \frac{P(A|C) \times P(C)}{P(A)} \quad (1)$$

Keterangan dari persamaan tersebut yaitu A merupakan data yang akan diklasifikasikan, C merupakan hipotesis atau kelas data, $P(C|A)$ merupakan probabilitas data A termasuk ke dalam kelas C (posterior probability), $P(C)$ merupakan peluang awal dari kelas C (prior probability), $P(A|C)$ merupakan probabilitas kemunculan data A pada kelas C (likelihood), sedangkan $P(A)$ merupakan probabilitas keseluruhan data.

3. HASIL DAN PEMBAHASAN

3.1 Pengambilan Data

Data penelitian diperoleh dari komentar pengguna Instagram yang berkaitan dengan aksi demonstrasi 25–31 Agustus 2025. Pengumpulan data dilakukan menggunakan teknik web scraping terhadap unggahan yang relevan dengan topik penelitian sehingga diperoleh sebanyak 3396 komentar. Komentar yang dikumpulkan mencerminkan berbagai respons masyarakat terhadap isu demonstrasi, baik berupa dukungan, kritik, maupun komentar yang mengandung unsur provokatif. Ilustrasi proses pengumpulan data komentar Instagram ditunjukkan pada Gambar 2.



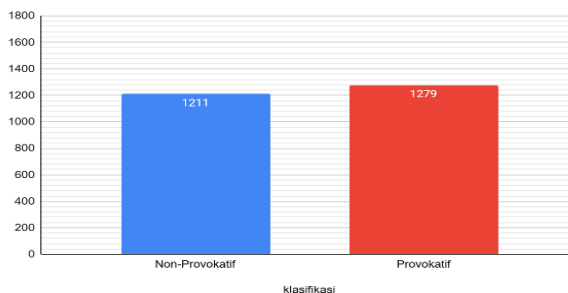
Gambar 2. Ilustrasi Pengumpulan Data

Hasil data komentar yang tersimpan dalam format JSON ditunjukkan pada Gambar 3.

```
{
  "username": "__sflxawr__",
  "text": "Bapak gua cuma lulusan sma, cuma pedagang,tapi kalo cuma ngitung nggak setolol dia bang 🤔",
  "likes": 0
},
```

Gambar 3. Data dalam file JSON

Berdasarkan hasil pelabelan diperoleh sebanyak 1279 komentar provokatif dan 1211 komentar non-provokatif. Visualisasi perbandingan jumlah data ditampilkan pada Gambar 6. Jumlah data pada kedua kategori terlihat relatif seimbang sehingga tidak terjadi ketimpangan kelas (class imbalance) yang signifikan. Kondisi ini dapat membantu meningkatkan performa model klasifikasi dan mengurangi bias dalam proses pembelajaran [18]. Visualisasi hasil pelabelan data komentar ditunjukkan pada Gambar 6.

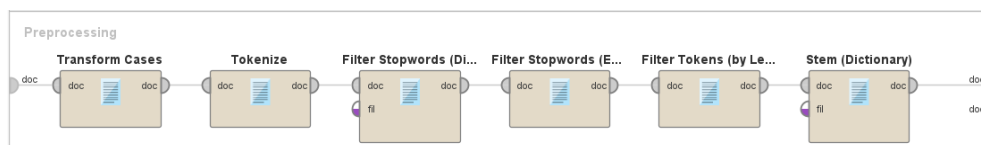


Gambar 6. Visualisasi Hasil Pelabelan Data secara Manual

Dominasi komentar provokatif dalam dataset menunjukkan bahwa isu demonstrasi cenderung memunculkan respons emosional dari pengguna platform media sosial. Hal tersebut menunjukkan bahwa media sosial telah berkembang menjadi ruang publik digital yang aktif digunakan masyarakat untuk menyampaikan opini, termasuk komentar yang bersifat provokatif.

3.4 Preprocessing

Tahap preprocessing dilakukan untuk mengolah data teks supaya siap digunakan dalam proses klasifikasi. Tahapan ini dilakukan menggunakan RapidMiner seperti terlihat pada Gambar 7.



Gambar 7. Proses Preprocessing

3.4.1 Transform Case

Transform case dilakukan dengan mentransformasi seluruh teks menjadi huruf kecil (lowercase) guna menyeragamkan format penulisan dan menghindari perbedaan makna akibat penggunaan huruf kapital [17]. Proses ini membantu meningkatkan konsistensi data sebelum memasuki tahap preprocessing selanjutnya.

3.4.2 Tokenize

Tokenizing merupakan tahap pemecahan teks menjadi unit-unit kata (token) untuk mempermudah proses analisis pada tingkat kata [20]. Tahap ini membantu sistem dalam mengenali serta mengolah setiap kata secara lebih terstruktur sebelum dilakukan proses klasifikasi.

3.4.3 Stopwords Removal

Stopwords removal diterapkan untuk membuang kata-kata umum yang tidak mengandung makna signifikan dalam proses analisis, seperti kata hubung dan kata yang sering muncul [20]. Pada penelitian ini, proses stopwords menggunakan daftar kata dalam bahasa Indonesia dan bahasa Inggris karena komentar Instagram yang digunakan mengandung campuran bahasa (code-mixing). Tahap ini bertujuan untuk meningkatkan kualitas data dengan menghilangkan kata-kata yang tidak memberikan pengaruh signifikan terhadap proses klasifikasi.

3.4.4 Filter Tokens

Filter tokens dilakukan untuk menyaring kata berdasarkan panjang karakter tertentu agar hanya kata-kata yang relevan digunakan dalam proses analisis [12]. Pada penelitian ini diterapkan batas minimal panjang kata sebanyak 4 karakter, sehingga kata dengan jumlah karakter kurang dari 4 akan dihapus. Tahap ini bertujuan untuk mengurangi kata-kata yang kurang informatif, seperti “yg”, “di”, dan “ke”, serta membantu meningkatkan kualitas fitur yang digunakan dalam proses klasifikasi.

3.4.5 Stemming

Proses stemming dilaksanakan untuk menyederhanakan kata berimbuhan menjadi bentuk dasar sehingga keragaman kata memiliki representasi yang lebih seragam [21]. Tahap ini membantu meningkatkan konsistensi data serta kualitas fitur yang digunakan dalam proses klasifikasi. Contoh hasil stemming ditunjukkan pada Tabel 3.

accuracy: 72.15%

	true Provokatif	true Non-Provokatif	class precision
pred. Provokatif	714	245	74.45%
pred. Non-Provokatif	310	724	70.02%
class recall	69.73%	74.72%	

Gambar 9. Confusion Matrix

Evaluasi model dilakukan menggunakan confusion matrix untuk mengetahui performa algoritma dalam mengklasifikasikan komentar provokatif dan non-provokatif. Berdasarkan hasil pengujian, algoritma Naïve Bayes memperoleh accuracy sebesar 72,15%, precision kelas provokatif sebesar 74,45%, recall sebesar 69,73%, precision kelas non-provokatif sebesar 70,02%, serta recall sebesar 74,72%. Hasil tersebut menunjukkan bahwa algoritma Naïve Bayes mampu melakukan klasifikasi komentar dengan performa yang cukup baik dan relatif seimbang pada kedua kelas. Nilai precision yang cukup tinggi pada kelas provokatif menunjukkan bahwa model mampu mengidentifikasi komentar provokatif dengan tingkat ketepatan yang baik. Sementara itu, nilai recall menunjukkan bahwa sebagian besar komentar provokatif berhasil dikenali oleh model meskipun masih terdapat beberapa kesalahan klasifikasi. Kesalahan klasifikasi umumnya dipengaruhi oleh karakteristik komentar media sosial yang menggunakan bahasa informal, singkatan, kata ambigu, serta konteks kalimat yang sulit dipahami oleh model. Beberapa komentar netral terkadang mengandung kata kasar namun tidak bermaksud provokatif, sehingga dapat menyebabkan model salah dalam menentukan kategori komentar. Kondisi tersebut menunjukkan bahwa klasifikasi teks pada media sosial memiliki tingkat kompleksitas yang cukup tinggi karena dipengaruhi oleh konteks bahasa dan makna kalimat. Secara keseluruhan, hasil penelitian menunjukkan bahwa algoritma Naïve Bayes cukup efektif digunakan dalam klasifikasi komentar provokatif pada media sosial Instagram, terutama karena mampu menghasilkan performa yang relatif stabil pada kedua kelas data.

3.8 Komparasi Algoritma

3.8.1 SVM

Hasil evaluasi algoritma SVM menggunakan confusion matrix ditunjukkan pada Gambar 10.

accuracy: 69.44%

	true Provokatif	true Non-Provokatif	class precision
pred. Provokatif	699	284	71.11%
pred. Non-Provokatif	325	685	67.82%
class recall	68.26%	70.69%	

Gambar 10. Hasil SVM

Berdasarkan hasil pengujian, algoritma SVM memperoleh accuracy sebesar 69,44%, precision kelas provokatif sebesar 71,11%, recall sebesar 68,26%, precision kelas non-provokatif sebesar 67,82%, dan recall sebesar 70,69%. Nilai tersebut menunjukkan bahwa performa SVM masih berada di bawah Naïve Bayes dalam proses klasifikasi komentar provokatif dan non-provokatif. Performa yang lebih rendah ini diduga dipengaruhi oleh karakteristik komentar Instagram yang sangat beragam, tidak terstruktur, serta banyak mengandung bahasa informal, singkatan, slang, campuran bahasa, dan konteks yang ambigu. Kondisi tersebut menyebabkan proses pemisahan kelas pada algoritma SVM menjadi kurang optimal. Meskipun SVM secara teoritis cukup baik dalam menangani data berdimensi tinggi seperti teks, algoritma ini cenderung sensitif terhadap distribusi fitur dan pemisahan kelas yang tidak konsisten. Selain itu, penggunaan parameter default pada RapidMiner juga diduga memengaruhi performa model karena belum dilakukan optimasi hyperparameter secara khusus. Meskipun demikian, algoritma SVM masih mampu menghasilkan performa klasifikasi yang cukup baik dalam mengenali pola data teks, namun pada penelitian ini performanya belum mampu melampaui algoritma Naïve Bayes.

3.8.2 Decision Tree

Hasil evaluasi algoritma Decision Tree menggunakan confusion matrix ditunjukkan pada Gambar 11.

accuracy: 72.91%

	true Provokatif	true Non-Provokatif	class precision
pred. Provokatif	504	20	96.18%
pred. Non-Provokatif	520	949	64.60%
class recall	49.22%	97.94%	

Gambar 11. Hasil Decision Tree

Berdasarkan hasil evaluasi, algoritma Decision Tree memperoleh accuracy sebesar 72,91%, precision kelas provokatif sebesar 96,18%, recall sebesar 49,22%, precision kelas non-provokatif sebesar 64,60%, dan recall sebesar 97,94%. Nilai accuracy tersebut merupakan yang tertinggi dibandingkan algoritma lainnya. Namun demikian, performa model belum dapat dikatakan optimal karena nilai recall pada kelas provokatif masih tergolong rendah. Nilai recall sebesar 49,22% menunjukkan bahwa masih banyak komentar provokatif yang gagal dikenali oleh model. Kondisi tersebut mengindikasikan bahwa Decision Tree cenderung lebih dominan dalam mengklasifikasikan komentar ke dalam kategori non-provokatif. Hal ini dapat terjadi karena karakteristik data komentar media sosial memiliki variasi kata dan konteks yang cukup kompleks sehingga proses pembentukan pohon keputusan menjadi kurang stabil. Meskipun Decision Tree memperoleh accuracy tertinggi, algoritma Naïve Bayes menunjukkan performa yang lebih seimbang antara precision dan recall pada kedua kelas. Dengan demikian, Naïve Bayes dinilai lebih efektif dalam mendeteksi komentar provokatif dan non-provokatif secara konsisten dibandingkan algoritma pembanding lainnya. Rendahnya nilai recall pada kelas provokatif menunjukkan bahwa masih banyak komentar provokatif yang diklasifikasikan sebagai non-provokatif (false negative). Kondisi ini menjadi kelemahan penting karena tujuan utama penelitian adalah mendeteksi komentar provokatif pada media sosial. Tingginya nilai precision menunjukkan bahwa Decision Tree sangat selektif dalam memberikan label provokatif, namun selektivitas tersebut menyebabkan model gagal mengenali sebagian besar komentar provokatif yang memiliki konteks ambigu atau pola bahasa yang bervariasi. Untuk mengurangi false negative, penelitian selanjutnya dapat menerapkan optimasi parameter, penambahan jumlah dataset, maupun penggunaan metode ensemble agar proses klasifikasi menjadi lebih stabil. Perbandingan hasil evaluasi algoritma Naïve Bayes, SVM, dan Decision Tree ditunjukkan pada Tabel 4.

Tabel 4. Perbandingan Hasil Evaluasi Algoritma

Algoritma	Accuracy	Precision Provokatif	Recall Provokatif	Precision Non-Provokatif	Recall Non-Provokatif	F1-Score
Naive Bayes	72,15%	74,45%	69,73%	70,02%	74,72%	71,00%
SVM	69,44%	71,11%	68,26%	67,82%	70,69%	69,66%
Decision Tree	72,91%	96,18%	49,22%	64,60%	97,94%	65,11%

3.9 Pembahasan

Berdasarkan hasil penelitian, algoritma Naïve Bayes memperoleh accuracy sebesar 72,15%, SVM sebesar 69,44%, dan Decision Tree sebesar 72,91%. Meskipun Decision Tree menghasilkan accuracy tertinggi, algoritma Naïve Bayes menunjukkan performa yang lebih seimbang berdasarkan nilai precision, recall, dan F1-score pada kedua kelas komentar. Hasil penelitian ini sejalan dengan penelitian Baehaqi dan Cahyono [7] yang menunjukkan bahwa algoritma Naïve Bayes mampu menghasilkan performa yang baik dalam klasifikasi komentar cyberbullying pada Instagram dengan accuracy sebesar 88%. Temuan serupa juga dilaporkan oleh Stephanie dan Irsyad [8] yang memperoleh accuracy sebesar 84% menggunakan Multinomial Naïve Bayes pada klasifikasi komentar cyberbullying. Selain itu, penelitian Warraihan et al. [10] menunjukkan bahwa Naïve Bayes memiliki performa yang lebih stabil dibandingkan K-Nearest Neighbor pada komentar Instagram. Perbedaan nilai accuracy pada penelitian ini dibandingkan penelitian sebelumnya diduga disebabkan oleh karakteristik dataset yang lebih kompleks. Komentar terkait demonstrasi politik umumnya mengandung bahasa informal, singkatan, slang, campuran bahasa, serta konteks yang ambigu sehingga proses klasifikasi menjadi lebih sulit dibandingkan data sentimen umum atau cyberbullying. Hasil penelitian juga menunjukkan bahwa SVM memperoleh performa lebih rendah dibandingkan Naïve Bayes. Kondisi ini berbeda dengan beberapa penelitian yang menyatakan bahwa SVM memiliki performa tinggi pada data teks berdimensi besar. Perbedaan tersebut diduga dipengaruhi oleh karakteristik komentar Instagram yang tidak terstruktur serta penggunaan parameter default tanpa optimasi hyperparameter. Sementara itu, Decision Tree menghasilkan accuracy tertinggi namun memiliki recall kelas provokatif yang rendah. Hal ini menunjukkan bahwa model cenderung menghasilkan false negative yang cukup tinggi sehingga kurang optimal untuk mendeteksi seluruh komentar provokatif. Oleh karena itu, meskipun accuracy Decision Tree lebih tinggi, algoritma Naïve Bayes dinilai lebih sesuai digunakan pada penelitian ini karena mampu memberikan performa yang lebih seimbang dalam mendeteksi kedua kelas komentar.

4. KESIMPULAN

Berdasarkan hasil evaluasi menggunakan confusion matrix, model Naïve Bayes memperoleh accuracy sebesar 72,15% dengan nilai precision dan recall yang relatif seimbang pada kedua kelas. Hasil tersebut menunjukkan bahwa algoritma Naïve Bayes memiliki performa yang cukup baik dalam mengklasifikasikan komentar provokatif dan non-provokatif pada media sosial Instagram, meskipun masih terdapat beberapa kesalahan klasifikasi akibat penggunaan bahasa informal, singkatan, slang, typo, serta konteks kalimat yang ambigu. Berdasarkan hasil komparasi, algoritma Naïve Bayes menunjukkan performa yang lebih seimbang dibandingkan SVM dan Decision Tree dalam mendeteksi kedua kelas komentar. Temuan penelitian ini menunjukkan bahwa pendekatan text mining dapat dimanfaatkan sebagai salah satu pendukung moderasi konten pada media sosial untuk membantu mengidentifikasi komentar provokatif secara

otomatis di ruang digital. Selain itu, penelitian ini memberikan kontribusi berupa analisis komparatif performa algoritma Naïve Bayes, SVM, dan Decision Tree pada klasifikasi komentar provokatif dalam konteks demonstrasi politik di Instagram yang memiliki karakteristik bahasa tidak terstruktur dan dinamis. Kendati demikian, penelitian ini masih memiliki keterbatasan pada jumlah dataset yang digunakan serta belum dilakukannya optimasi parameter pada masing-masing algoritma klasifikasi. Oleh sebab itu, penelitian selanjutnya diharapkan dapat menggunakan dataset yang lebih besar, menerapkan penanganan bahasa informal secara lebih mendalam, memperbaiki proses preprocessing terhadap slang dan typo, serta membandingkan performa metode machine learning dengan pendekatan berbasis Deep Learning atau Transformer guna meningkatkan hasil klasifikasi.

REFERENCES

- [1] S. Kemp, “Digital 2024: Indonesia,” 2024. Accessed: Sep. 16, 2025. [Online]. Available: <https://datareportal.com/reports/digital-2024-indonesia>
- [2] K. H. Yunior, A. V. Vitianingsih, S. Kacung, A. Lidya Maukar, and A. Dwi Arumsari, “Sentiment Analysis of Cyberbullying Detection on Social Networks using the Sentistrength Method,” *Sistemasi: Jurnal Sistem Informasi*, vol. 13, no. 4, 2024, doi: 10.32520/stmsi.v13i4.4226.
- [3] M. O. Ibrohim and I. Budi, “Hate speech and abusive language detection in Indonesian social media: Progress and challenges,” *Heliyon*, vol. 9, no. 8, p. e18647, 2023, doi: 10.1016/j.heliyon.2023.e18647.
- [4] Anjum and R. Katarya, “Hate speech, toxicity detection in online social media: a recent survey of state of the art and opportunities,” *Int. J. Inf. Secur.*, vol. 23, no. 1, pp. 577–608, 2024, doi: 10.1007/s10207-023-00755-2.
- [5] G. Ramos *et al.*, “A comprehensive review on automatic hate speech detection in the age of the transformer,” *Soc. Netw. Anal. Min.*, vol. 14, no. 1, p. 204, 2024, doi: 10.1007/s13278-024-01361-3.
- [6] S. Saha *et al.*, “Bengali cyberbullying detection in social media using machine learning algorithms,” in *2023 5th International Conference on Sustainable Technologies for Industry 5.0 (STI)*, IEEE, 2023, pp. 1–6. doi: 10.1109/STI59863.2023.10464740.
- [7] F. Bachaqi and N. Cahyono, “Analisis Sentimen Terhadap Cyberbullying Pada Komentar Di Instagram Menggunakan Algoritma Naïve Bayes,” *The Indonesian Journal of Computer Science*, vol. 13, no. 1, pp. 1051–1063, 2024, doi: 10.33022/ijcs.v13i1.3301.
- [8] Stephanie and H. Irsyad, “Analisis Sentimen Opini Publik Terhadap Cyberbullying Pada Komentar Instagram Menggunakan Multinomial Naïve Bayes,” *Jurnal Nasional Teknologi Komputer*, vol. 4, no. 4, pp. 44–57, Oct. 2024, doi: 10.61306/jnastek.v4i4.157.
- [9] F. Anisa Nirmala, M. Jazman, N. E. Rozanda, and F. N. Salisah, “Cyberbullying Sentiment Analysis Of Instagram Comments Using Naïve Bayes Classifier And K-Nearest Neighbor Algorithm Methods,” *Jurnal Teknik Informatika (Jutif)*, vol. 5, no. 5, pp. 1213–1219, May 2024, doi: 10.52436/1.jutif.2024.5.5.1997.
- [10] D. Warraihan, I. Permana, M. Mustakim, and R. Novita, “Analisis Sentimen Pengguna Transportasi Online Maxim Pada Instagram Menggunakan Naïve Bayes Classifier dan K-Nearest Neighbor,” *Jurnal Media Informatika Budidarma*, vol. 7, no. 3, pp. 1134–1143, 2023, doi: 10.30865/mib.v7i3.6336.
- [11] Y. Mubarak, D. Sudana, and W. Gunawan, “Hate Speech in the Comments’ Column Instagram: A Discourse Analysis,” *Journal of Languages and Language Teaching*, vol. 12, no. 1, pp. 439–450, 2024, doi: 10.33394/jollt.v12i1.9050.
- [12] A. Jalili, H. Tabrizchi, A. Mosavi, and A. R. Varkonyi-Koczy, “Enhancing Language Model Performance with a Novel Text Preprocessing Method,” *Acta Phys. Pol. A*, vol. 146, pp. 542–552, Nov. 2024, doi: 10.12693/APhysPolA.146.542.
- [13] D. A. B. Listiawan, W. A. Wicaksono, and Wiyanto, “Analisis Sentimen Publik pada Komentar Instagram Postingan Akun Tempodotco terhadap Penerapan Kurikulum AI di Indonesia Menggunakan Metode Naïve Bayes,” *JATI (Jurnal Mahasiswa Teknik Informatika)*, vol. 9, no. 5, pp. 7881–7889, 2025, doi: 10.36040/jati.v9i5.14934.
- [14] S. K. Papia, M. A. Khan, T. Habib, M. Rahman, and M. N. Islam, “DistilRoBiLSTMFuse: an efficient hybrid deep learning approach for sentiment analysis,” *PeerJ Comput. Sci.*, vol. 10, p. e2349, 2024, doi: 10.7717/peerj-cs.2349.
- [15] A. D. Saruksuk, E. Sianturi, M. I. Ahda, T. O. Panggabean, and M. W. Siregar, “Analisis Semantik di Balik Ujaran Provokatif di Media Sosial dan Implikasi Hukum dalam Kasus ‘Gus Mifta Menghina Penjual Es,’” *Jurnal Pendidikan Tambusai*, vol. 8, no. 3, pp. 50327–50332, Dec. 2024, Accessed: May 29, 2026. [Online]. Available: <https://jptam.org/index.php/jptam/article/view/23838>
- [16] R. A. Perdana, C. E. Widodo, and R. Santoso, “Sentiment Analysis of Naïve Bayes, Decision Tree, and K-Nearest Neighbor (K-NN) Algorithms for Cyberbullying Comments on Instagram Accounts,” in *2024 11th International Conference on Information Technology, Computer, and Electrical Engineering (ICITACEE)*, 2024, pp. 253–258. doi: 10.1109/ICITACEE62763.2024.10762775.
- [17] A. Anggara, Nurdin, and R. Meiyanti, “Sentiment Analysis of the MK Decision Trial of the Result of the 2024 President and Vice President General Election on Social Media X Using the Support Vector Machine Method,” *International Journal of Engineering, Science and Information Technology*, vol. 4, no. 4, pp. 125–134, 2024, doi: 10.52088/ijesty.v4i4.591.
- [18] E. Helmut, F. Fitriyani, and P. Romadiana, “Classification Comparison Performance of Supervised Machine Learning Random Forest and Decision Tree Algorithms Using Confusion Matrix,” *Jurnal Sisfokom (Sistem Informasi dan Komputer)*, vol. 13, no. 1, pp. 92–97, Feb. 2024, doi: 10.32736/sisfokom.v13i1.1985.
- [19] N. W. S. Saraswati, C. P. Yanti, I. D. M. K. Muku, and D. A. P. R. Dewi, “Evaluation Analysis of the Necessity of Stemming and Lemmatization in Text Classification,” *Matrik: Jurnal Manajemen, Teknik Informatika, dan Rekayasa Komputer*, vol. 24, no. 2, pp. 321–332, 2025, doi: 10.30812/matrik.v24i2.4833.
- [20] S. F. C. Haviana, S. Mulyono, and Badie’Ah, “The Effects of Stopwords, Stemming, and Lemmatization on Pre-trained Language Models for Text Classification: A Technical Study,” in *2023 10th International Conference on Electrical Engineering, Computer Science and Informatics (EECSI)*, 2023, pp. 521–527. doi: 10.1109/EECSI59885.2023.10295797.



- [21] E. Qais and V. M. N, "TxtPrePro: Text Data Preprocessing Using Streamlit Technique for Text Analytics Process," in *2023 International Conference on Network, Multimedia and Information Technology (NMITCON)*, 2023, pp. 1-6. doi: 10.1109/NMITCON58196.2023.10275887.
- [22] M. D. Bimantara and I. Zufria, "Text Mining Sentiment Analysis on Mobile Banking Application Reviews using TF-IDF Method with Natural Language Processing Approach," *JINAV: Journal of Information and Visualization*, vol. 5, no. 1, pp. 115-123, Jul. 2024, doi: 10.35877/454RI.jinav2772.