

Klasifikasi Hate Speech dan Emosi Dalam Teks Berbahasa Indonesia Pada Pengguna Twitter Menggunakan Metode Naïve Bayes Classifier

Chandra Hary Pratama, Yulian Findawati*

Universitas Muhammadiyah Sidoarjo

Abstrak: Ujaran kebencian merupakan salah satu bentuk ekspresi yang menghasut, menyebarkan, membenarkan, atau mendorong kebencian, diskriminasi serta kekerasan atas individu dan kelompok sebab berbagai alasan. Hate speech biasanya ditemukan pada sosial media yang terhubung dengan internet, salah satunya pada penelitian ini melalui sosial media twitter dengan menggunakan metode Naïve Bayes Classifier. Dataset yang digunakan pada penelitian ini berjumlah 1800 data berlabel bukan ujaran kebencian dan 2250 data berlabel ujaran kebencian dengan perbandingan 60% data latih dan 40% data uji. Hasil evaluasi data uji dengan confusion matrix diperoleh pengukuran matrix mean accuracy for hate speech classification 0,89 dan matrix mean accuracy for emotion classification 0,59. Berdasarkan hasil yang didapat tersebut dapat diambil kesimpulan bahwa untuk melakukan klasifikasi hate speech dan emosi pada Twitter menggunakan Naïve Bayes hasil paling bagus dengan Confusion Matrix tanpa melakukan seleksi fitur Information Gain.

Kata Kunci: Klasifikasi, Ujaran Kebencian, Emosi, Naïve Bayes, Twitter

DOI:

<https://doi.org/10.47134/ijat.v1i3.3105>

*Correspondence: Yulian Findawati

Email: yulianfindawati@umsida.ac.id

Received: 11-07-20224

Accepted: 14-07-2024

Published: 15-07-2024



Copyright: © 2024 by the authors. Submitted for open access publication under the terms and conditions of the Creative Commons Attribution (BY SA) license (<http://creativecommons.org/licenses/by/4.0/>).

Abstract: Hate speech is a form of expression that incites, spreads, justifies, or encourages hatred, discrimination and violence against individuals and groups for various reasons. Hate speech is usually found on social media connected to the internet, one of which is in this study through social media twitter using the Naïve Bayes Classifier method. The dataset used in this study amounted to 1800 data labeled not hate speech and 2250 data labeled hate speech with a comparison of 60% training data and 40% test data. The results of the evaluation of test data with confusion matrix obtained measurements of matrix mean accuracy for hate speech classification 0.89 and matrix mean accuracy for emotion classification 0.59. Based on the results obtained, it can be concluded that to classify hate speech and emotions on Twitter using Naïve Bayes, the best results with the Confusion Matrix without selecting the Information Gain feature.

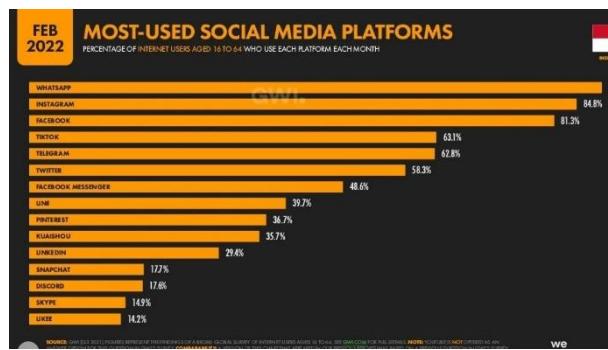
Keywords: Clasification, Hate Speech. Emotion, Naïve Bayes, Twitter

Pendahuluan

Hate speech atau ujaran kebencian adalah suatu bentuk ekspresi yang menghasut, menyebarkan, membenarkan, atau mendorong kebencian, diskriminasi serta kekerasan atas individu dan kelompok sebab berbagai alasan (Liu & Sari, 2019). Hate speech atau ujaran

kebencian tidak jarang kita jumpai pada kehidupan sehari-hari. Ujaran kebencian sangat sering digunakan pada status, komentar, atau postingan pada media sosial(Abro et al., 2020).

Linschoten (Al Baqi, 2015) menjelaskan bahwa emosi orang dibagi menjadi tiga bagian menurut kategorinya, yaitu suasana hati, suasana perasaan, dan emosi(Mozafari et al., 2020). Dalam arti luas emosi adalah salah satu bagian dari perasaan. Emosi hadir dari perasaan yang bergolak, sehingga yang terlibat dapat mengalami perubahan dalam situasi emosi tersendiri. Pada penelitian (Martins et al., 2018) menunjukkan bahwa emosi tertentu seperti kemarahan dan kebencian lebih berkorelasi dengan ujaran kebencian di Twitter. Klasifikasi emosi terdiri dari “*Anger*”, “*Anticipation*”, “*Disgust*”, “*Fear*”, “*Joy*”, “*Sadness*”, “*Surprise*” dan “*Trust*”.



Gambar 1. We Are Social Hootsuite (2022)

Survei *We Are Social* menyebutkan dalam tinjauan media sosial penduduk Indonesia yang aktif bermain media sosial mencapai 4,62 miliar orang pada tahun 2022, dan jumlah perangkat mobile yang terhubung mencapai 8,28 miliar(Perifanos & Goutsos, 2021). Media sosial merupakan sebuah wadah yang sering disalahgunakan sebagai tempat ujaran kebencian dan berekspresi. Twitter merupakan salah satu media sosial yang paling banyak digunakan. Pada survei *We Are Social*, Twitter menempati peringkat ke enam dengan persentase 58,3%. Pengguna twitter seringkali melayangkan komentar, status, bahkan postingan yang mengandung *hate speech* dan ekspresi emosi(Chiril et al., 2022). Pengguna diberikan suatu kebebasan untuk menyalurkan ekspresi dan perasaan emosi di twitter (Ahmad Gozali & Alfan Rosid, 2020).

Penelitian mengenai *hate speech* dan emosi ini sebelumnya telah dilakukan dengan pembahasan mengenai deteksi *hate speech* bahasa Indonesia pada twitter (Hakiem et al., 2019). Serta Deteksi sentimen di twitter menggunakan metode *Naive Bayes* (Fanesya et al., 2019) Hasil penelitian ini menunjukkan klasifikasi ujaran kebencian dan deteksi emosi pada akun twitter yang dimana dalam penelitian ini metode yang digunakan adalah *naïve bayes* berbasis n-gram serta seleksi fitur *Information Gain*(Kovács et al., 2021).

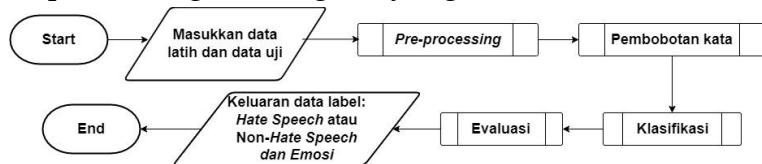
Penelitian telah yang dilakukan oleh (Afif et al., 2021) memaparkan suatu koneksi negatif antara kematangan batin dan aksi ujaran kebencian (*Hate Speech*), semakin tinggi

kematangan batin seseorang, semakin sedikit pula aksi ujaran kebencian (*Hate Speech*) dan juga sebaliknya. Maka dari itu pada penelitian ini akan menggabungkan kedua klasifikasi yaitu klasifikasi *hate speech* dan emosi dengan menggunakan metode *naïve bayes*(Chhabra & Vishwakarma, 2023).

Metode Penelitian

Metodologi penelitian adalah suatu tahapan atau langkah-langkah yang akan dilakukan untuk memecahkan suatu permasalahan guna menemukan solusi yang tepat pada sebuah penelitian. Pada penelitian ini, penulis membahas studi kasus klasifikasi *hate speech* dan emosi terhadap pengguna twitter menggunakan *Naive Bayes Classifier* (Ahmad Wildan Attabi et al., 2018).

Penelitian ini dilandaskan dengan mengimplementasikan naive bayes classifier untuk mengklasifikasikan hate speech dan emosi berbahasa Indonesia terhadap pengguna twitter. Berikut merupakan langkah-langkah yang dilakukan(Florio et al., 2020):



Gambar 2. Diagram Alir Penelitian

Berdasarkan diagram alir pada Gambar 1dapat dilihat proses dalam tahapan penelitian yang akan menjadi acuan dalam penggerjaan penelitian ini(Anderson & Barnes, 2022). Berikut adalah penjelasan dari flowchart Naïve Bayes Classifier(Alkomah & Ma, 2022):

1. Dataset

Dataset pada penelitian ini diambil dari sosial media twitter. Kemudian dilakukan identifikasi kata dasar, identifikasi kata-kata yang sering muncul atau stopword serta dilakukan identifikasi kategori.

2. Text Pre-processing

Pre-processing merupakan sebuah proses awal pengklasifikasikan dokumen dengan tujuan menyiapkan data agar data tersebut mempunyai struktur. Text pre-processing akan menghasilkan nilai yang digunakan sebagai data untuk dilakukan proses selanjutnya. Pre-processing dibagi menjadi beberapa proses antara lain case folding, tokenizing, filtering, stemming, dan penghitungan bobot kata (Socrates et al., 2016).

3. Pembobotan Kata

Pembobotan kata (Term Weighting) yaitu suatu mekanisme pemberian nilai pada setiap kata berdasarkan indeks (Hidayat & Rosid, 2020).

4. Klasifikasi

Tujuan dari proses klasifikasi ini yaitu untuk memperoleh hasil dari data uji untuk mendapatkan luaran berupa label hate speech atau non-hate speech dan emosi yang terkandung pada data yang telah dimasukkan. Semua proses atau langkah selesai ketika diperoleh hasil luaran berupa label kelas tersebut.

5. Evaluasi

Tahapan akhir yang dilakukan adalah proses evaluasi, proses evaluasi yang bertujuan menguji hasil dari klasifikasi dengan cara pengukuran nilai kebenaran dari sistem tersebut.

Hasil dan Pembahasan

A. Dataset

Data dari media sosial Twitter yang berhasil dikumpulkan sebanyak 3.972 tweet. Memiliki 15 atribut tantara lain URL, Tanggal, Tweet, ID, Username, Likes, Quotes dan sebagainya yang disimpan dalam format .csv (Deolika et al., 2019). Data selanjutnya diolah ke tahapan preprocessing untuk meningkatkan struktur dari data. Tabel 1 menunjukkan hasil crawling tweet(Gould, 2019).

Tabel 1. Sampel tweet hasil *Crawling*

NO	Tweet
1	@arunariftan makin gila lihat LBP mewarnai Indonesia tercinta, Barisan Saku Hati Kadrun2 sich ya mampus saja ha ha ha racain.
2	Ha ha ha sigundul penguasa ancol karena selama ini taunya hanya jilat2 gabenor mulai dibuka jeroannya, KPK Kejaksaan Agung Mabes Polri tolong segera turun/selidiki MERDEKA. https://t.co/Zzp3e2IZ67
....
3972	Buat KADRUN2 nih

B. Pre-processing

Tahapan preprocessing data dilakukan untuk meningkatkan performa masing-masing algoritma klasifikasi dalam melakukan prediksi sehingga didapatkan data yang lebih presisi. Tahapan ini meliputi cleaning data, Setelah dilakukan proses preprocessing data jumlah data yang dihasilkan sejumlah 3.972 tweet(Carlson, 2021).

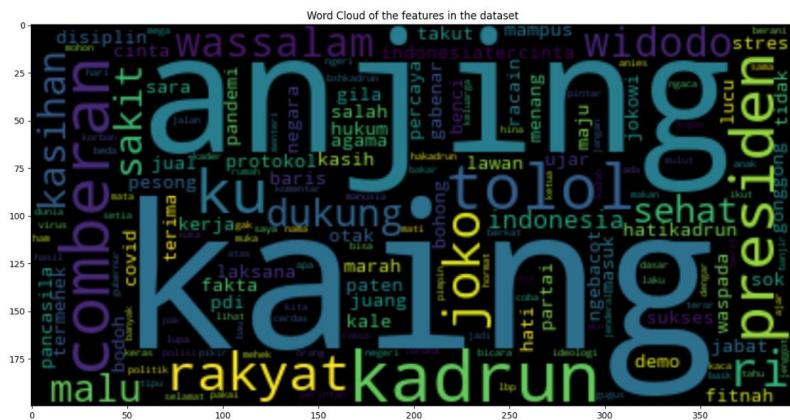
Tabel 2. Proses *Preprocessing* Data

Proses	Tweet
Cleaning data	@arunariftan makin gila lihat LBP mewarnai Indonesia tercinta, Barisan Saku Hati Kadrun2 sich ya mampus saja ha ha ha racain.

Tokenization	['makin', 'gila', 'lihat', 'lbp', 'mewarnai', 'indonesia', 'tercinta', 'barisan', 'saku', 'hati', 'kadrun', 'sich', 'ya', 'mampus', 'saja', 'ha', 'ha', 'ha', 'racain']
Stopword Removal	['makin', 'gila', 'lihat', 'lbp', 'mewarnai', 'indonesia', 'tercinta', 'barisan', 'saku', 'hati', 'kadrun', 'sich', 'ya', 'mampus', 'saja', 'ha', 'ha', 'ha', 'racain']
Stemming	makin gila lihat lbp mewarnai indonesia tercinta barisan saku hati kadrun2 sich ya mampus saja ha ha ha racain

C. Pembobotan Kata

Pembobotan kata (Term Weighting) yaitu suatu mekanisme pemberian nilai pada setiap kata berdasarkan indeks. Proses pembobotan kata ini bertujuan untuk memperoleh jumlah kemunculan kata setelah dilakukan proses pre-processing pada dataset. Pada penelitian ini metode pembobotan kata yang digunakan adalah TF-IDF (Term Frequency Invers Document Frequency)(Kumari, 2014). Metode ini digunakan untuk menentukan keterhubungan kata terhadap dokumen dengan cara memberikan bobot pada setiap kata. Pada perhitungan TF-IDF terlebih dahulu menghitung nilai TF pada setiap kata, dengan bobot satu kata adalah 1(Khan et al., 2021).

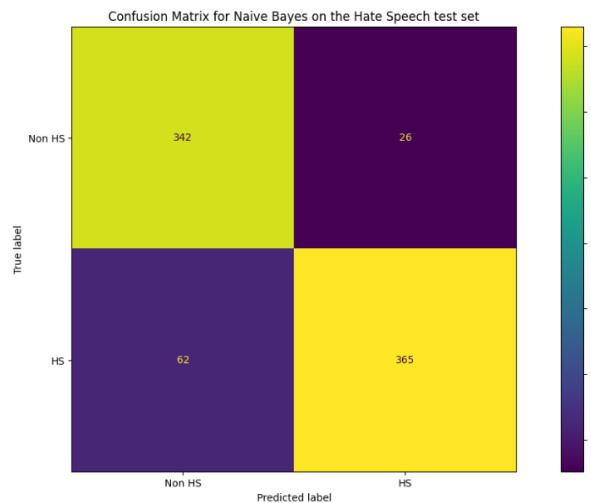


Gambar 3. Word Cloud

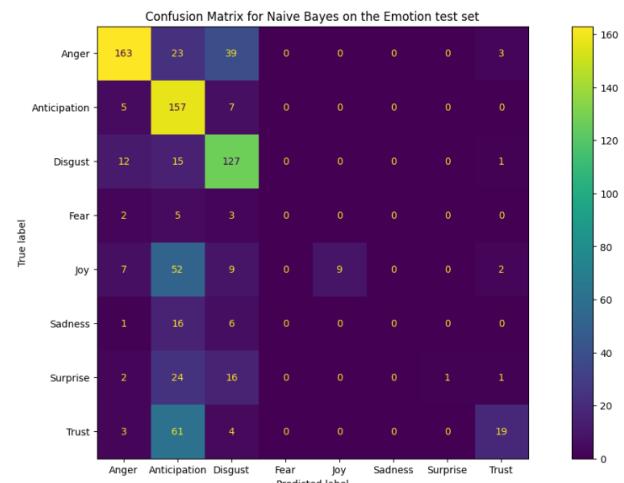
D. Klasifikasi

Setelah melalui tahapan proses pembobotan kata, langkah berikutnya adalah memasukkan data ke dalam fase pemodelan klasifikasi. Sebelumnya, data akan dibagi menjadi dua bagian, yakni data latih dan data uji, dengan rasio 60:40, di mana 60% digunakan untuk data latih dan 40% untuk data uji dalam eksperimen ini (Dwitama, 2021). Data uji berperan sebagai alat evaluasi model, sementara data latih berfungsi untuk pembangunan model dan identifikasi pola(Pettersson, 2019). Setelah pembagian

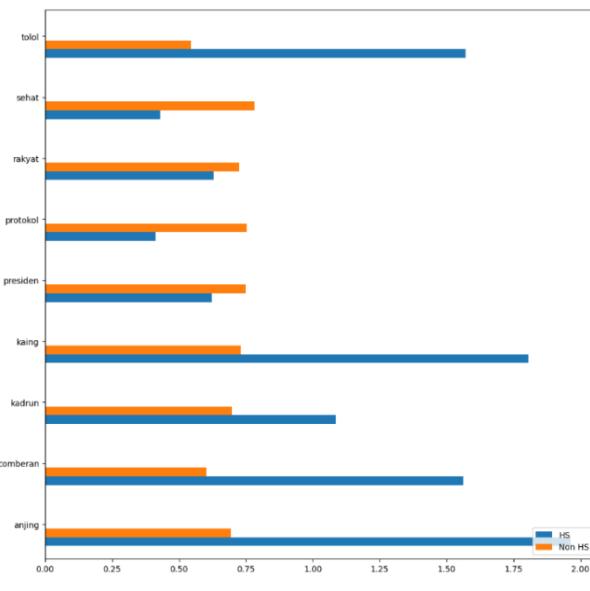
data diselesaikan, langkah selanjutnya melibatkan proses klasifikasi model. Algoritma pertama yang akan dilakukan percobaan adalah Naive Bayes (Ghassani Saskia, 2021).



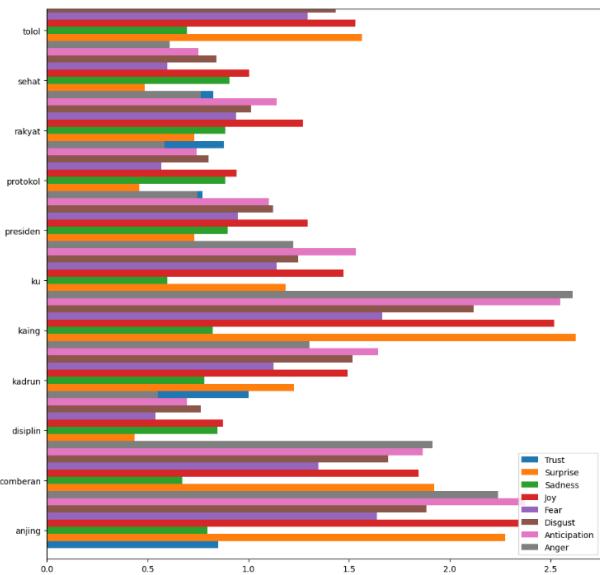
Gambar 4. Confusion Matrix
Hate Speech Set



Gambar 5. Confusion Matrix Emosi



Gambar 6. Plot Hate Speech



Gambar 7. Plot Emosi

E. Evaluasi

Tahapan akhir yang dilakukan adalah proses evaluasi, proses evaluasi yang bertujuan menguji hasil dari klasifikasi dengan cara pengukuran nilai kebenaran dari sistem tersebut (Hadna et al., 2016). Tolok ukur yang digunakan sebagai acuan dalam mengukur adalah accuracy. Pada penelitian ini ekstraksi fitur yang digunakan adalah precision dengan persamaan sebagai berikut(Döring & Mohseni, 2020):

$$\text{Accuracy} = \left(\frac{TP + TN}{TP + TN + FP + FN} \right) \times 100\%$$

(1)

Tabel 3. Mean Accuracy

Model Algoritma	Hate Speech	Emosi
<i>Mean Accuracy</i>	0.889308176100629	0.5987421383647799

F. Output

Output yang dihasilkan dari program ini adalah untuk mendeteksi hate speech dan emosi dari sebuah kalimat random. Ada dua jenis output dari program ini yaitu berupa keterangan HS/Non HS dan emosi.

Tabel 4. Output klasifikasi Hate Speech dan Emosi

Text	Hate Speech	Emosi
gin ngebacot tunjukin muka tampang ditwitter bawa islam	True	Anger
islam islam sadar ideologi indonesia pancasila	False	Anticipation
komentar dukung setia joko widodo presiden ri cinta	True	Anger
rakyat kerja keras	True	Anticipation
makin gila lihat LBP mewarnai Indonesia tercinta, Barisan		
Saku Hati Kadrun2 sich ya mampus saja ha ha ha racain.		
kasihan korban akibat cuci otak kadrun tanggung terima		
kasih polisi paspampres waspada waspada waspada		

Simpulan

Dengan Output dari program ini adalah hasil klasifikasi teks tweet berbahasa Indonesia yang telah dianalisis menggunakan algoritma Naïve Bayes Classifier. Program ini mengeluarkan dua jenis klasifikasi utama:

1. Klasifikasi Hate Speech: Program akan memberikan label pada setiap tweet apakah termasuk dalam kategori ujaran kebencian atau tidak. Ini dilakukan dengan memproses teks tweet dan membandingkannya dengan pola-pola yang telah dipelajari dari dataset pelatihan.
2. Klasifikasi Emosi: Selain mendeteksi hate speech, program juga mengidentifikasi emosi yang terkandung dalam teks tweet. Emosi yang dapat dideteksi misalnya marah, sedih, bahagia, dan lain-lain. Setiap tweet akan diberi label sesuai dengan emosi yang paling dominan berdasarkan analisis fitur linguistiknya.

Output akhir yang dihasilkan berupa data terstruktur yang menunjukkan tweet, label hate speech ("hate speech = true" atau "hate speech = false"), dan label emosi ("Anger", "Anticipation", "Disgust", "Fear", "Joy", "Sadness", "Surprise" dan "Trust"). Data ini dapat digunakan untuk analisis lebih lanjut atau untuk tindakan moderasi konten di platform Twitter.

Daftar Pustaka

- Abro, S., Shaikh, S., Khand, Z. H., Zafar, A., & ... (2020). Automatic hate speech detection using machine learning: A comparative study. *International* <https://pdfs.semanticscholar.org/0445/07a2f4d0030c05434eceb0230c40f868804d.pdf>
- Afif, Much. F. A., Nurhamidah, Y., & Mashuri, M. F. (2021). Kematangan emosi dalam perilaku ujaran kebencian pada kebijakan politik. *Cognicia*, 9(1), 25–30. <https://doi.org/10.22219/cognicia.v9i1.14234>
- Ahmad Gozali, H., & Alfan Rosid, M. (2020). Classification of Student Complaints with Naive Bayes and Literary Methods Klasifikasi Keluhan Mahasiswa dengan Metode Naive Bayes dan Sastrawi. *Network, and Computer Science* |, 3(1), 22–26.
- Ahmad Wildan Attabi, Laillil Muflikhah, & Mochammad Ali Fauzi. (2018). Penerapan Analisis Sentimen untuk Menilai Suatu Produk pada Twitter Berbahasa Indonesia dengan Metode Naïve Bayes Classifier dan Information Gain. *Jurnal Pengembangan Teknologi Informasi Dan Ilmu Komputer*, 2(11), 4548–4554.
- Al Baqi, S. (2015). Ekspresi Emosi Marah. *Buletin Psikologi*, 23(1), 22. <https://doi.org/10.22146/bpsi.10574>
- Alkomah, F., & Ma, X. (2022). A literature review of textual hate speech detection methods and datasets. *Information*. <https://www.mdpi.com/2078-2489/13/6/273>
- Anderson, L., & Barnes, M. (2022). *Hate speech*. plato.stanford.edu. <https://plato.stanford.edu/ENTRIES/hate-speech/>
- Carlson, C. R. (2021). *Hate speech*. books.google.com. [https://books.google.com/books?hl=en&lr=&id=dDoiEAAAQBAJ&oi=fnd&pg=PA1&d q=hate+speech&tots=fmclhAYLhB&sig=GxQemAMypPvYukVYNRaxHUbWOMu](https://books.google.com/books?hl=en&lr=&id=dDoiEAAAQBAJ&oi=fnd&pg=PA1&dq=hate+speech&tots=fmclhAYLhB&sig=GxQemAMypPvYukVYNRaxHUbWOMu)
- Chhabra, A., & Vishwakarma, D. K. (2023). A literature survey on multimodal and multilingual automatic hate speech identification. *Multimedia Systems*. <https://doi.org/10.1007/s00530-023-01051-8>
- Chiril, P., Pamungkas, E. W., Benamara, F., Moriceau, V., & ... (2022). Emotionally informed hate speech detection: a multi-target perspective. *Cognitive* <https://doi.org/10.1007/s12559-021-09862-5>
- Deolika, A., Kusrini, K., & Luthfi, E. T. (2019). Analisis Pembobotan Kata Pada Klasifikasi Text Mining. *Jurnal Teknologi Informasi*, 3(2), 179. <https://doi.org/10.36294/jurti.v3i2.1077>
- Döring, N., & Mohseni, M. R. (2020). Gendered hate speech in YouTube and YouNow comments: Results of two content analyses. *SCM Studies in Communication and* <https://doi.org/10.5771/2192-4007-2020-1-62>

- Dwitama, A. P. J. (2021). Deteksi Ujaran Kebencian Pada Twitter Bahasa Indonesia Menggunakan Machine Learning: Reviu Literatur. *Jurnal Sains, Nalar, Dan Aplikasi Teknologi Informasi*, 1(1), 31–39. <https://doi.org/10.20885/snati.v1i1.5>
- Fanesya, F., Wihandika, R. C., & Indriati. (2019). Deteksi Emosi pada Twitter Menggunakan Metode Naive Bayes dan Kombinasi Fitur. *Jurnal Pengembangan Teknologi Informasi Dan Ilmu Komputer*, 3(7), 3.
- Florio, K., Basile, V., Polignano, M., Basile, P., & Patti, V. (2020). Time of your hate: The challenge of time in hate speech detection on social media. *Applied Sciences*. <https://www.mdpi.com/2076-3417/10/12/4180>
- Ghassani Saskia, T. (2021). *Klasifikasi Hate Speech Dan Abusive Language Pada Twitter Bahasa Indonesia Dengan Metode Naive Bayes Classifier*.
- Gould, J. B. (2019). *Speak no evil: The triumph of hate speech regulation*. degruyter.com. <https://doi.org/10.7208/9780226305134>
- Hadna, N. M. S., Santosa, P. I., & Winarno, W. W. (2016). Studi Literatur Tentang Perbandingan Metode Untuk Proses Analisis Sentimen Di Twitter. *Seminar Nasional Teknologi Informasi Dan Komunikasi*, 2016(March), 57–64.
- Hakiem, M., Fauzi, M. A., & Indriati. (2019). Klasifikasi Ujaran Kebencian pada Twitter Menggunakan Metode Naïve Bayes Berbasis N-Gram Dengan Seleksi Fitur Information Gain. *Jurnal Pengembangan Teknologi Informasi Dan Ilmu Komputer*, 3(3), 2443–2451.
- Hidayat, T. E., & Rosid, A. (2020). Analysis of Community Sentiments Regarding Plans to Relocate National Capital Using the Naïve Bayes Method Analisa Sentimen Masyarakat Tentang Rencana Pemindahan Ibukota Negara Dengan Metode Naïve Bayes. *Network, and Computer Science* |, 3(2), 43–49.
- Khan, M. M., Shahzad, K., & Malik, M. K. (2021). Hate speech detection in roman urdu. *ACM Transactions on Asian and Low* <https://doi.org/10.1145/3414524>
- Kovács, G., Alonso, P., & Saini, R. (2021). Challenges of hate speech detection in social media: Data scarcity, and leveraging external resources. *SN Computer Science*. <https://doi.org/10.1007/s42979-021-00457-3>
- Kumari, A. (2014). Study on Naive Bayesian Classifier and its relation to Information Gain. *International Journal on Recent and Innovation Trends in Computing and Communication*, 2(3), 601–602.
- Liu, I., & Sari, Y. A. (2019). Klasifikasi Hate Speech Berbahasa Indonesia di Twitter Menggunakan Naive Bayes dan Seleksi Fitur Information Gain dengan Normalisasi Kata. *Jurnal Pengembangan Teknologi Informasi Dan Ilmu Komputer*, 3(5), 4914–4922.
- Martins, B., Sheppes, G., Gross, J. J., & Mather, M. (2018). Age Differences in Emotion Regulation Choice: Older Adults Use Distraction Less Than Younger Adults in High-Intensity Positive Contexts. *Journals of Gerontology - Series B Psychological Sciences and Social Sciences*, 73(4), 603–611. <https://doi.org/10.1093/geronb/gbw028>
- Mozafari, M., Farahbakhsh, R., & Crespi, N. (2020). Hate speech detection and racial bias mitigation in social media based on BERT model. *PloS One*. <https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0237861>

- Perifanos, K., & Goutsos, D. (2021). Multimodal hate speech detection in greek social media. *Multimodal Technologies and Interaction*. <https://www.mdpi.com/2414-4088/5/7/34>
- Pettersson, K. (2019). "Freedom of speech requires actions": Exploring the discourse of politicians convicted of hate-speech against Muslims. *European Journal of Social Psychology*. <https://doi.org/10.1002/ejsp.2577>
- Socrates, I. G. A., Akbar, A. L., Akbar, M. S., Arifin, A. Z., & Herumurti, D. (2016). Optimasi Naive Bayes Dengan Pemilihan Fitur Dan Pembobotan Gain Ratio. *Lontar Komputer: Jurnal Ilmiah Teknologi Informasi*, 7(1), 22. <https://doi.org/10.24843/lkjiti.2016.v07.i01.p03>