

Pemanfaatan Metode *Naïve Bayes* untuk Analisis Sentimen Terkait Insiden Server Down dan Peretasan pada Pusat Data Nasional

¹Oktaviani, ²Endah Kurniasari, ³Rosdiana, ⁴Faturachman Kurniawan Putra
^{1,2,3,4}Fakultas Ilmu Komputer dan Teknologi Informasi,
Universitas Gunadarma, Jakarta, Indonesia
¹oktaviani@staff.gunadarma.ac.id, ²endah_ks@staff.gunadarma.ac.id,
³rosdiana@staff.gunadarma.ac.id, ⁴faturachman1102@gmail.com

Submit : 13 Jun 2025 | Diterima : 21 Jun 2025 | Terbit : 22 Jun 2025

ABSTRAK

Penelitian ini bertujuan untuk menganalisis sentimen pengguna platform X terhadap insiden gangguan layanan dan peretasan yang menimpa Pusat Data Nasional (PDN). Metode *Naïve Bayes* digunakan untuk mengklasifikasikan 1004 tweet berbahasa Indonesia yang dikumpulkan selama bulan Juni hingga Juli 2024. Proses penelitian mencakup tahapan *crawling* data, *preprocessing*, klasifikasi sentimen, evaluasi model, hingga visualisasi hasil. Dataset dibagi menjadi 80% data latih dan 20% data uji. Hasil menunjukkan bahwa model mencapai akurasi sebesar 83% dengan *f1-score* untuk sentimen positif sebesar 0,83 dan negatif sebesar 0,87, sementara performa untuk sentimen netral masih rendah. Temuan ini memperlihatkan efektivitas metode *Naïve Bayes* dalam memetakan opini publik terhadap isu keamanan data nasional.

Kata kunci : *Naïve Bayes*, Analisis Sentimen, Pusat Data Nasional, Media Sosial

PENDAHULUAN

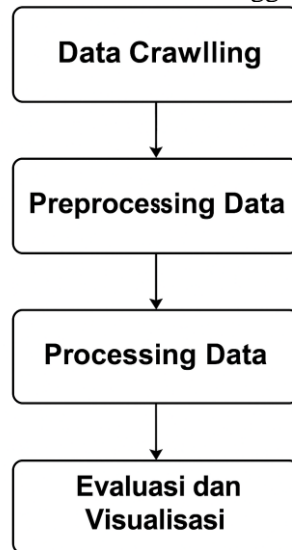
Insiden peretasan dan gangguan layanan yang menimpa Pusat Data Nasional (PDN) pada pertengahan tahun 2024 menjadi perhatian publik yang luas. Sebagai wadah penyimpanan dan pengelolaan data instansi pemerintah, gangguan terhadap PDN berdampak signifikan terhadap kepercayaan masyarakat terhadap keamanan digital nasional. Reaksi masyarakat terhadap insiden tersebut banyak terekam melalui platform media sosial X (sebelumnya Twitter), yang menjadi ruang ekspresi publik secara *real-time*. Tingginya intensitas penggunaan platform X dari tahun ke tahun, terutama sepanjang 2020 hingga 2024, menjadikan opini publik yang tersebar di dalamnya semakin penting untuk ditelusuri dan dianalisis (Annur, 2024).

Insiden yang menimpa PDN tidak hanya menyoroti permasalahan teknis, tetapi juga memicu kekhawatiran serius terkait aspek keamanan siber dan perlindungan data pribadi masyarakat. Kominfo (2024) menegaskan bahwa dampak dari peristiwa tersebut menyentuh berbagai layanan pemerintah yang bergantung pada integritas sistem data nasional. Fenomena maraknya opini masyarakat di media sosial pasca kejadian ini menunjukkan bahwa analisis sentimen menjadi salah satu langkah strategis untuk memahami kecenderungan respons publik. Temuan dari analisis tersebut dapat menjadi dasar pertimbangan dalam menyusun kebijakan yang lebih tepat terkait penguatan sistem keamanan digital nasional.

Dalam konteks penelitian ini, metode *Naïve Bayes* dipilih karena memiliki keunggulan dalam proses klasifikasi data teks yang besar, cepat, dan efisien. Selain itu, algoritma ini mampu bekerja dengan baik meskipun data yang dianalisis memiliki distribusi yang tidak seimbang antar kelas. Dengan pendekatan ini, opini masyarakat dapat dipetakan menjadi kategori sentimen positif, negatif, maupun netral, sehingga memberikan gambaran yang lebih jelas mengenai bagaimana masyarakat menanggapi insiden keamanan siber seperti yang dialami oleh PDN (Khotimah & Utami, 2022).

METODOLOGI PENELITIAN

Metode yang digunakan dalam penelitian ini adalah pendekatan kuantitatif dengan memanfaatkan algoritma Naïve Bayes untuk klasifikasi sentimen. Data diambil dari 1004 tweet yang mengandung kata kunci "Pusat Data Nasional" menggunakan alat tweet-harvest.



Gambar 1. Alur Penelitian

Pada gambar 1 dijelaskan bahwa pengumpulan data dalam penelitian ini dilakukan melalui proses *crawling* pada media sosial X, dengan periode pengambilan data selama bulan Juni hingga Juli 2024. Data yang terkumpul kemudian diproses melalui beberapa tahapan *preprocessing* untuk memastikan kualitasnya sebelum dilakukan analisis. Tahapan *preprocessing* meliputi pembersihan data dari karakter yang tidak relevan (*cleaning*), normalisasi kata agar sesuai dengan bentuk baku, konversi huruf menjadi huruf kecil (*case folding*), penghapusan kata-kata tidak penting (*stopword removal*), pemisahan kalimat menjadi kata-kata terpisah (tokenisasi), pengembalian kata ke bentuk dasarnya (*stemming*), serta penerjemahan teks ke dalam bahasa Inggris untuk meningkatkan kompatibilitas dengan pustaka analisis yang digunakan.

Selanjutnya, data yang telah diproses dibagi menjadi dua bagian, yaitu 80% untuk proses pelatihan (*training*) dan 20% untuk pengujian (*testing*). Model *Naïve Bayes* digunakan dalam proses klasifikasi untuk menentukan kategori sentimen dari setiap data, baik sebagai sentimen positif, negatif, maupun netral. Setelah proses klasifikasi selesai, tahap berikutnya adalah evaluasi kinerja model. Kinerja model dievaluasi dengan menggunakan beberapa metrik evaluasi, antara lain akurasi, *precision*, *recall*, dan *f1-score*. Perhitungan metrik evaluasi menggunakan persamaan berikut:

Rumus *Precision*:

$$Precision = \frac{TP}{TP+FP} \dots\dots\dots 1$$

Rumus *Recall*:

$$Recall = \frac{TP}{TP+FN} \dots\dots\dots 2$$

Rumus *F1-Score*:

$$F1-Score = 2 \times \frac{(Precision \times Recall)}{Precision+Recall} \dots\dots\dots 3$$

Keterangan:

- True Positive (TP)** merupakan kondisi ketika model berhasil melakukan prediksi dengan tepat, di mana hasil prediksi sesuai dengan kelas sebenarnya.
- True Negative (TN)** adalah kondisi ketika model memprediksi sebuah data sebagai bukan bagian dari kelas tertentu, dan hasil tersebut sesuai dengan kondisi sebenarnya.
- False Positive (FP)** menggambarkan situasi ketika model memberikan prediksi positif, padahal kenyataannya data tersebut bukan termasuk dalam kelas yang dimaksud.
- False Negative (FN)** terjadi ketika model memprediksi sebuah data sebagai negatif, sedangkan pada kenyataannya data tersebut justru termasuk dalam kelas yang seharusnya diprediksi positif.

HASIL DAN PEMBAHASAN

Crawling Data

Proses pengumpulan data diawali dengan menentukan kata kunci "Pusat Data Nasional" sebagai acuan untuk menelusuri tweet berbahasa Indonesia dari pengguna platform X. Pengambilan data dilakukan dengan memanfaatkan alat bantu *tweet-harvest*, sehingga berhasil terkumpul sebanyak 1004 tweet dalam rentang waktu Juni hingga Juli 2024.



full_text	id_str	image_url	in_reply_to_screen_name	lang	loc	qot	reply_count	retweet_count
...jLuwajō Di...	https://pbs.twimg.com/media/GRXlt...	Null		in	in	Central java Indonesia	757	859
SANA PDGa...	https://pbs.twimg.com/ext_tw_video...	Null		in	in		8	18
faisal PDN	Null	Null		in	in	Indonesia	0	0
luan terakah...	https://twitter.com/esspersonBar/sta...	Null		in	in	Indonesia	0	0
isaka PDN nih...	https://pbs.twimg.com/media/GRXlt...	Null		in	in	N. Minde(public of cis)	132	150
rusak (109#)...	https://twitter.com/a#kamana/status...	Null		in	in	RICORCHESETT	132	894

Gambar 2. Hasil *Crawling* Data

Preprocessing Data

Proses *preprocessing* data dilakukan untuk menyiapkan data mentah agar sesuai dan siap diproses pada tahapan analisis berikutnya. Data yang diperoleh dari media sosial X kemudian diekstraksi untuk menghapus elemen-elemen yang tidak relevan, sehingga dapat meningkatkan akurasi hasil analisis sentimen. Hasil akhir dari tahapan *preprocessing* ini disajikan secara rinci pada Tabel 1.

Tabel 1. Hasil Tahapan *Preprocessing* Data

Proses	Hasil
<i>Cleaning</i>	pasrah data pdn yang diretas tak bisa dipulihkan beginilah akibatnya kalau yang dijadikan menteri bukan orang yang kompeten dalam bidangnya
Normalisasi	pasrah data pusat data nasional yang diretas tak bisa dipulihkan beginilah akibatnya kalau yang dijadikan menteri bukan orang yang kompeten dalam bidangnya
<i>Case Folding</i>	pasrah data pusat data nasional yang diretas tak bisa dipulihkan beginilah akibatnya kalau yang dijadikan menteri bukan orang yang kompeten dalam bidangnya
<i>Stopword</i>	pasrah data pusat data nasional diretas tak bisa pulihkan beginilah akibatnya kalau dijadikan menteri bukan orang kompeten bidangnya

Proses	Hasil
Tokenisasi	[“pasrah”, “data”, “pusat”, “data”, “nasional”, “diretas”, “tak”, “bisa”, “pulihkan”, “beginilah”, “akibatnya”, “kalau”, “dijadikan”, “menteri”, “bukan”, “orang”, “kompeten”, “bidangnya”]
Stemming	pasrah data pusat data nasional diretas tak bisa pulih begini akibat kalau jadi menteri bukan orang kompeten bidang
Translasi	I'm resigned to the fact that the hacked national data center cannot be recovered; this is what happens when someone incompetent is appointed as minister

Processing Data

Dalam penelitian ini, sebanyak 939 data yang telah melewati proses *preprocessing* kemudian dibagi dengan rasio 80:20, yang terdiri dari 751 data untuk pelatihan (*training*) dan 188 data untuk pengujian (*testing*). Berdasarkan hasil pengujian, diperoleh 112 data *true negative*, 43 *true positive*, 22 *false negative*, 6 *false positive*, dan 5 *false neutral*.

index	Unnamed: 0	full_text	tweet_english	klasifikasi	klasifikasi_bayes
30	326	kalau kasus pusat data nasional bobol bikin tobat massal buat deploying xampp prod emang udah ada obat masalah skill issue mau upgrade data center ampe nilai triliun mentok users bisa pake barang canggihcanggih	If the case of a national data center being breached, make a mass repentance for deploying XAMPP PROD, there is already a cure for the skill issue, want to upgrade the data center until the value of trillions is stuck, users can use sophisticated things, sophisticated	Positif	Positif
35	844	padahal ruang lingkup pusat data nasional tahun 2024 harus wajib susun sop kait failover minimal downtime bisa panjang urus kalau nyata salah di usaha menang	even though the scope of the national data center in 2024 must be to prepare failover hooks so that minimal downtime can take a long time to deal with if something goes wrong in trying to win	Positif	Positif

Gambar 3. Data *True Positive*

Data pada gambar 3 menunjukkan bahwa model berhasil mendeteksi opini publik yang bernada membangun, solutif, dan optimistis terhadap perbaikan sistem PDN. Ini mengindikasikan bahwa model Naïve Bayes mampu mengenali sentimen positif dengan baik.

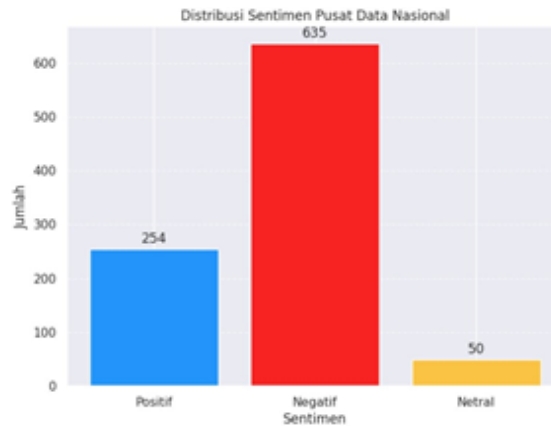
index	Unnamed: 0	full_text	tweet_english	klasifikasi	klasifikasi_bayes
28	70	data pusat data nasional kena encrypted pak arti data lu bocor jebol doboll encrypted dkk banyak orang pinter indonesia orang kayak gin jadi menkominfo	National data center data has been encrypted, sir, what does your data leak mean? Doboll encrypted, etc. There are lots of smart Indonesian people, people like Gin become Minister of Communication and Information.	Negatif	Negatif
31	581	wow lemah sistem it indonesia serang siber server pusat data nasional kominfo down antrre panjang bandara soetta 40 auto gate mati	wow weak Indonesian IT system cyber attack national data center server Kominfo down long queue at Soetta 40 airport auto gate down	Negatif	Negatif

Gambar 4. Data *True Negative*

Data pada gambar 4 menjelaskan, bahwa isi tweet berisi kritik terhadap lemahnya sistem keamanan IT nasional yang menyebabkan dampak sistemik, seperti gangguan layanan di bandara. Nada negatif disebabkan oleh ketidakpuasan terhadap buruknya manajemen IT nasional.

Hasil dan Evaluasi

Tahapan evaluasi dan visualisasi dilakukan untuk mengetahui sejauh mana tingkat keberhasilan metode klasifikasi dalam mengelompokkan data secara tepat. Evaluasi dilakukan dengan menghitung tingkat akurasi berdasarkan hasil prediksi yang dihasilkan oleh model klasifikasi. Selanjutnya, hasil evaluasi ini divisualisasikan dalam bentuk diagram batang dan *word cloud*, untuk memberikan gambaran yang lebih jelas mengenai distribusi sentimen yang terkandung dalam data.



Gambar 5. Diagram Batang

Gambar 5 menyajikan diagram batang yang merepresentasikan jumlah total data sebanyak 939 tweet. Berdasarkan hasil analisis, sentimen dengan kategori **positif** berjumlah 254 data, sedangkan kategori **negatif** mendominasi dengan jumlah sebanyak 635 data. Sementara itu, untuk kategori **netral** hanya ditemukan sebanyak 50 data. Visualisasi ini memperlihatkan bahwa opini negatif terkait insiden Pusat Data Nasional mendominasi percakapan publik di media sosial.



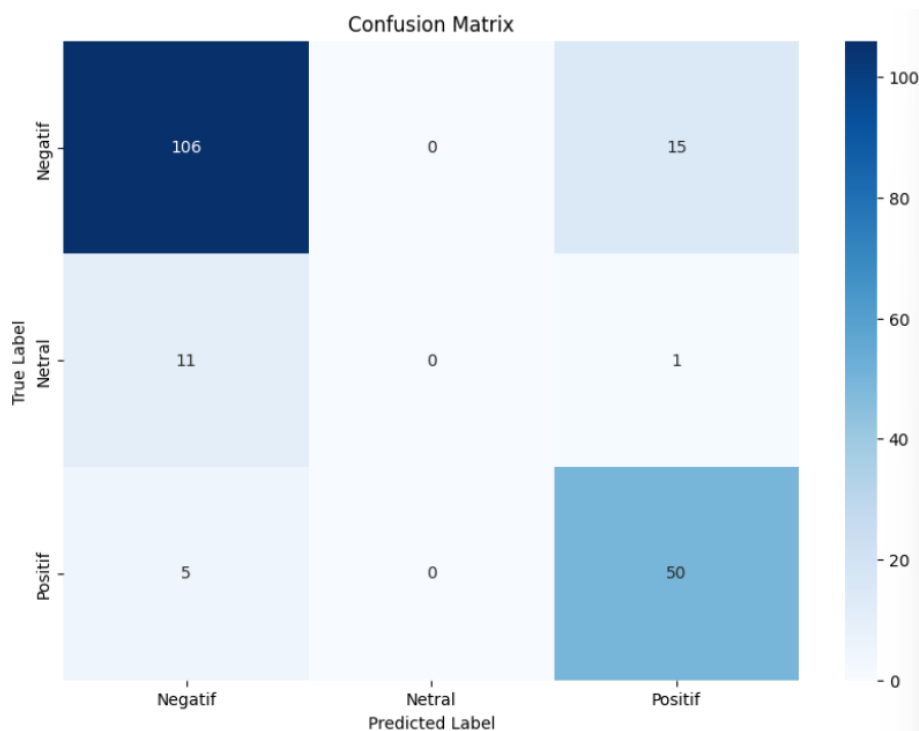
Gambar 6. Wordcloud

Gambar 6 menampilkan visualisasi *word cloud* yang memetakan kata-kata paling sering muncul dalam data yang dianalisis terkait insiden Pusat Data Nasional (PDN). Beberapa kata yang tampak menonjol antara lain “pusat”, “data”, “nasional”, “serang”, “system”, “PDN”, “Kominfo”, “ransomware”, “pemerintah”, dan “retas”. Berdasarkan visualisasi ini dapat disimpulkan bahwa topik pembicaraan publik mayoritas berfokus pada isu keamanan siber yang melibatkan Pusat Data Nasional. Dominasi kata seperti “ransomware”, “retas”, dan “serang” mengindikasikan adanya kekhawatiran terhadap ancaman serangan siber, khususnya yang berdampak pada sistem data nasional. Selain itu, perhatian publik juga banyak tertuju pada peran Kementerian Komunikasi dan Informatika (Kominfo), serta upaya penanganan yang dilakukan oleh Badan Siber dan Sandi Negara (BSSN). Secara keseluruhan, *word cloud* ini mencerminkan keresahan masyarakat terhadap isu keamanan data serta dampaknya terhadap tata kelola pemerintahan.

Classification report :				
	precision	recall	f1-score	support
Negatif	0.87	0.88	0.87	121
Netral	0.00	0.00	0.00	12
Positif	0.76	0.91	0.83	55
accuracy			0.83	188
macro avg	0.54	0.60	0.57	188
weighted avg	0.78	0.83	0.80	188

Gambar 7. Hasil Evaluasi

Gambar 7 memperlihatkan hasil evaluasi model berdasarkan *classification report*, yang menggambarkan performa model dalam menganalisis data. Untuk kategori **Negatif**, model menunjukkan *precision* sebesar 0,87, *recall* sebesar 0,88, dan *f1-score* sebesar 0,87, dengan *support* sebanyak 121 data. Hal ini mengindikasikan bahwa model mampu mendeteksi mayoritas data bernuansa negatif dengan baik. Sebaliknya, untuk kategori **Netral**, model menunjukkan performa yang sangat rendah, dengan *precision*, *recall*, dan *f1-score* semuanya bernilai 0, dan hanya memiliki *support* sebanyak 12 data. Kinerja model untuk kategori **Positif** juga belum optimal, dengan *precision* sebesar 0,76 dan *recall* sebesar 0,91, menghasilkan *f1-score* sebesar 0,83, dengan *support* sebanyak 55 data. Secara keseluruhan, akurasi model mencapai 83%, yang berarti 83% dari prediksi yang dihasilkan sesuai dengan kelas sebenarnya. Akan tetapi, rata-rata makro untuk *precision*, *recall*, dan *f1-score* masing-masing sebesar 0,54; 0,60; dan 0,57 menunjukkan adanya ketidakseimbangan performa antar kelas. Sementara itu, rata-rata tertimbang menunjukkan *precision* sebesar 0,78, *recall* sebesar 0,83, dan *f1-score* sebesar 0,80, mencerminkan performa umum model dengan mempertimbangkan distribusi kelas. Meskipun akurasi cukup tinggi, hasil ini menunjukkan bahwa model memiliki kelemahan signifikan dalam mendeteksi kategori **Netral**.



Gambar 8. Heat Map Confussion Matrix

Gambar 8 memperlihatkan hasil pengujian *confusion matrix* yang menggambarkan distribusi prediksi model terhadap masing-masing kategori sentimen. Berdasarkan hasil penelitian, diperoleh 50 data *true positive*, 5 *false positive*, 0 *true neutral*, 11 *false neutral*, 106 *true negative*, dan 15 *false negative*. Kategori *true positive* menunjukkan jumlah kalimat positif yang berhasil diprediksi positif oleh model. Sementara itu, *false positive* merupakan kalimat positif yang justru diprediksi negatif atau netral oleh model. Kategori *true neutral* mengacu pada kalimat netral yang berhasil dikenali sebagai netral, sedangkan *false neutral* merupakan kalimat netral yang salah diklasifikasikan sebagai positif atau negatif. *True negative* mencerminkan kalimat negatif yang berhasil diidentifikasi sebagai negatif, sedangkan *false negative* adalah kalimat negatif yang salah diprediksi sebagai positif atau netral oleh model.

KESIMPULAN

Penelitian ini telah berhasil menyelesaikan proses analisis sentimen terhadap tanggapan pengguna X terkait insiden gangguan layanan (*server down*) dan peretasan pada Pusat Data Nasional dengan memanfaatkan metode *Naïve Bayes*. Dari total 939 tweet yang dianalisis, data dibagi dengan rasio 80% untuk pelatihan dan 20% untuk pengujian, masing-masing berjumlah 751 dan 188 data. Berdasarkan hasil pengujian, model mencapai tingkat akurasi sebesar 83%. Sentimen positif menunjukkan *precision* sebesar 0,76, *recall* sebesar 0,91, dan *f1-score* sebesar 0,83 dengan *support* sebanyak 55 data. Sebaliknya, pada kategori sentimen netral, model gagal memberikan prediksi yang akurat, dengan *precision*, *recall*, dan *f1-score* bernilai nol, serta *support* hanya sebanyak 12 data. Untuk sentimen negatif, model menunjukkan performa yang baik dengan *precision* sebesar 0,87, *recall* sebesar 0,88, dan *f1-score* sebesar 0,87, didukung oleh *support* sebanyak 121 data. Berdasarkan temuan tersebut, dapat disimpulkan bahwa metode *Naïve Bayes* mampu memberikan performa yang cukup baik dalam melakukan analisis sentimen terhadap isu keamanan siber yang terjadi pada Pusat Data Nasional, khususnya untuk kategori positif dan negatif.

Penelitian ini masih memiliki potensi untuk dikembangkan lebih lanjut melalui beberapa pendekatan yang dapat diterapkan. Salah satu upaya yang dapat dilakukan adalah mengatasi permasalahan ketidakseimbangan kelas dengan menerapkan teknik *oversampling* atau *undersampling*, sehingga distribusi data antar kategori menjadi lebih proporsional. Selain itu, pengayaan fitur dengan menerapkan *n-grams* dan *word embeddings*, serta penerapan proses *preprocessing* yang lebih komprehensif, diharapkan mampu meningkatkan akurasi prediksi. Penggunaan teknik *cross-validation* juga dapat memberikan evaluasi model yang lebih optimal. Tidak kalah penting, penambahan jumlah dataset khususnya pada kelas minoritas akan membantu model dalam melakukan proses pembelajaran yang lebih baik. Dengan menerapkan langkah-langkah tersebut, diharapkan dapat menghasilkan model klasifikasi yang lebih optimal dan akurat untuk penelitian di masa mendatang.

DAFTAR PUSTAKA

- Mairita, D., & Abdullah, A. (2024). Media Sosial sebagai Media Kampanye Politik Menjelang Pemilu 2024. *JURNAL SIMBOLIKA Research and Learning in Communication Study*, 10(1), 72-81.
<https://doi.org/10.31289/symbolika.v10i1.11468>
- Sakdiah, H., Nadiyah, N., Ginting, G. T., Gea, M., Aisyah, N., Situmorang, A., ... & Ramadhan, T. (2024). Correlation of Rights & Obligations of Citizens & States in Personal Data Protection (Highlighting National Data Hacking Cases). *QISTINA: Jurnal Multidisiplin Indonesia*, 3(2), 1303-1308.
<https://doi.org/10.57235/qistina.v3i2.4049>

- Ramdhan, T. W., Florina, I. D., & Permadi, D. (2024). Analisis Framing Pemberitaan Peretasan Pusat Data Nasional (PDN) di Media Online Tempo. co. *Journal of Education Research*, 5(3), 3368-3379. <https://doi.org/10.37985/jer.v5i3.1491>
- Mono, J. R. (2024). Perlindungan Hukum Terhadap Hak Privasi Subjek Data Pribadi dalam Insiden Serangan Siber Pusat Data Nasional Sementara. *Jurnal Ilmu Hukum, Humaniora dan Politik (JIHHP)*, 5(1).
<https://doi.org/10.38035/jihhp.v5i1.3195>
- Handayani, A., & Zufria, I. (2023). Analisis sentimen terhadap bakal capres ri 2024 di twitter menggunakan algoritma svm. *J. Inf. Syst. Res*, 5(1), 53-63.
<https://doi.org/10.47065/josh.v5i1.4379>
- Pahlawan, M. R., Setyanto, A., & Arief, M. R. (2024). A Comprehensive Review of Classifier used with Imbalanced Data in Machine Learning. *Journal of Electrical Engineering and Computer (JEECOM)*, 6(1), 177-185.
<https://doi.org/10.33650/jeecom.v6i1.8510>
- Suprihanto, S., Awaludin, I., Fadhil, M., & Zulfikor, M. A. Z. (2022). Analisis kinerja resnet-50 dalam klasifikasi penyakit pada daun kopi robusta. *J. Inform*, 9(2), 116-122.
<https://doi.org/10.31294/inf.v9i1.13049>
- Atimi, R. L., & Pratama, E. E. (2022). Implementasi Model Klasifikasi Sentimen Pada Review Produk Lazada Indonesia. *Jurnal Sains Dan Informatika*, 8(1), 88-96.
<https://doi.org/10.34128/jsi.v8i1.419>
- Suryadewiansyah, M. K., & Tju, T. E. E. (2022). Naïve Bayes dan Confusion Matrix untuk Efisiensi Analisa Intrusion Detection System Alert. *Jurnal Nasional Teknologi dan Sistem Informasi*, 8(2), 81-88.
<https://doi.org/10.25077/teknosi.v8i2.2022.81-88>
- Riyanto, S., Imas, S. S., Djatna, T., & Atikah, T. D. (2023). Comparative analysis using various performance metrics in imbalanced data for multi-class text classification. *International Journal of Advanced Computer Science and Applications*, 14(6).
<https://doi.org/10.14569/ijacsa.2023.01406116>