

Optimasi Hyperparameter Gaussian Naive Bayes Untuk Prediksi Risiko Stroke Pada Data Tidak Seimbang

Khoirun Nida*, Ridwan Mahendra, Erliyan Redi Susanto

Fakultas Teknik dan Ilmu Komputer, Program Studi Sistem Informasi, Universitas Teknokrat Indonesia, Bandar Lampung, Indonesia

Email: ¹*khoirun_nida@teknokrat.ac.id, ²ridwanmahendra@teknokrat.ac.id, ³erliyan.redy@teknokrat.ac.id

Email Penulis Korespondensi: khoirun_nida@teknokrat.ac.id

Submitted: 09/10/2025; Accepted: 06/12/2025; Published: 08/12/2025

Abstrak—Stroke merupakan salah satu penyakit serius dengan dampak global yang menuntut adanya deteksi dini berakurasi tinggi. Kesulitan signifikan dalam perancangan model prediktif berbasis *machine learning* muncul akibat kondisi data yang tidak proporsional (*imbalanced dataset*). Ini terjadi karena jumlah kasus stroke (kelas minoritas) sangat minim jika dibandingkan dengan kasus non-stroke. Situasi data yang timpang ini sering mendorong model menjadi bias dan berpotensi besar menghasilkan tingkat *false negative* yang tinggi, suatu hal yang sangat berisiko di lingkungan klinis. Penelitian ini berfokus pada peningkatan sensitivitas model Gaussian Naive Bayes (GNB) melalui optimasi *hyperparameter* dan penyesuaian *threshold* klasifikasi. Proses penelitian mencakup tahapan *preprocessing* data, pembagian *dataset* secara *stratified* (70% pelatihan dan 30% pengujian), *feature scaling*, optimasi parameter *var_smoothing* menggunakan GridSearchCV, serta penyesuaian *threshold* untuk memaksimalkan nilai *Recall*. Hasil percobaan menunjukkan bahwa model GNB standar hanya memperoleh nilai *Recall* sebesar 0,4400. Namun, setelah dilakukan optimasi *var_smoothing* ($1,00 \times 10^{-10}$) dan penyesuaian *threshold* menjadi 0,0100, nilai *Recall* meningkat signifikan hingga 0,8000. Peningkatan ini memang diikuti penurunan *Accuracy* (0,5988) dan *Precision* (0,0909). Peningkatan *Recall* yang tinggi (0,8000) menunjukkan model lebih baik untuk skrining massal (fase deteksi dini), meskipun harus diimbangi dengan proses diagnostik lanjutan karena *precision* yang rendah. Nilai *Recall* yang tinggi ini menegaskan keberhasilan model dalam meminimalkan *False Negative*, yang merupakan prioritas utama dalam kasus prediksi risiko stroke.

Kata Kunci: Gaussian Naive Bayes; Optimasi Hyperparameter; Data Tidak Seimbang; Prediksi Stroke; Recall

Abstract—Stroke is a serious disease with global impact that requires high-accuracy early detection. Significant difficulties in designing machine learning-based predictive models arise due to disproportionate data conditions (imbalanced datasets). This occurs because the number of stroke cases (minority class) is very small compared to non-stroke cases. This imbalanced data situation often causes models to become biased and potentially produce high false negative rates, which is very risky in a clinical setting. This study focuses on improving the sensitivity of the Gaussian Naive Bayes (GNB) model through hyperparameter optimization and classification threshold adjustment. The research process included data preprocessing, stratified dataset division (70% training and 30% testing), feature scaling, var_smoothing parameter optimization using GridSearchCV, and threshold adjustment to maximize the Recall value. The results showed that the standard GNB model only achieved a Recall value of 0.4400. However, after var_smoothing optimization (1.00×10^{-10}) and threshold adjustment to 0.0100, the Recall value increased significantly to 0.8000. This increase was accompanied by a decrease in Accuracy (0.5988) and Precision (0.0909). This improvement was accompanied by a decrease in Accuracy (0.5988) and Precision (0.0909). The high Recall (0.8000) indicates that the model is better for mass screening (early detection phase), although it must be balanced with further diagnostic processes due to low precision. This high Recall value confirms the model's success in minimizing False Negatives, which is a top priority in stroke risk prediction cases.

Keywords: Gaussian Naive Bayes; Hyperparameter Optimization; Imbalanced Data; Stroke Prediction; Recall

1. PENDAHULUAN

Stroke termasuk di antara faktor utama penyebab mortalitas secara global, signifikan dinegara berkembang ataupun maju [1]. Menurut *World Health Organization* (WHO), diperkirakan bahwa lebih dari enam juta orang meninggal akibat stroke dan menjadi penyebab utama kecacatan jangka panjang [2]. Di Indonesia, prevalensi stroke terus meningkat dan pada tahun 2013 mencapai 10,9 per 1.000 penduduk [3]. Peningkatan asus ini erat kaitannya dengan gaya hidup modern, penuaan populasi, serta meningkatnya prevalensi faktor risiko seperti hipertensi, diabetes, obesitas, dan kurangnya aktivitas fisik [3]. Oleh karena itu, deteksi dini risiko stroke menjadi krusial untuk mencegah komplikasi serius dan mengurangi beban sistem kesehatan.

Meski teknologi medis terus berkembang, prediksi dini stroke masih menghadapi tantangan besar. Kendala utama dalam prediksi dini adalah mutu data pasien. Data sering ditemukan tidak lengkap (mengandung *missing values*), memiliki inkonsistensi format, dan mengalami ketimpangan kelas (*class imbalance*) sebab data pasien stroke jauh lebih langka daripada pasien non-stroke[4]. Misalnya, dalam dataset yang digunakan pada penelitian ini, hanya sekitar 4,9% dari total 5.110 rekam medis yang mencatat kejadian stroke, sedangkan 95,1% merupakan pasien tanpa stroke. Fenomena data yang tidak seimbang berpotensi menyebabkan model *machine learning* menjadi bias. Model mungkin akan cenderung mengklasifikasikan semua sampel sebagai kelas mayoritas (non-stroke), sehingga kasus stroke yang sesungguhnya terabaikan. Kesalahan prediksi seperti ini memiliki risiko besar di lingkungan klinis, sebab pasien berisiko tinggi stroke bisa tidak mendapatkan intervensi dini yang diperlukan, yang pada akhirnya dapat mengakibatkan kecacatan permanen atau bahkan kematian[5].

Isu ketidakseimbangan kelas adalah permasalahan vital yang harus diatasi dalam implementasi machine learning disektor kesehatan [6]. dalam kasus prediksi stroke, jumlah pasien yang mengalami stroke jauh lebih sedikit

dibandingkan pasien sehat. Situasi ini membuat model cenderung bias ke kelas mayoritas, sehingga kemampuan mendeteksi kasus stroke menjadi kurang optimal. Jika tidak diantisipasi sejak tahap pra-pemrosesan data, model dapat menghasilkan banyak *false negative*, yang dalam konteks klinis berpotensi fatal.

Kebaruan Kontribusi utama penelitian ini tidak terletak pada kerumitan model yang digunakan, melainkan pada penerapan kerangka evaluasi yang lebih adil dan relevan secara medis dalam menangani data yang tidak seimbang. Kebaruan penelitian ini terletak pada pendekatan sistematis yang mengintegrasikan preprocessing data dengan optimasi hyperparameter (*var_smoothing*) pada Gaussian Naive Bayes [7]. Strategi ini dirancang khusus untuk meningkatkan sensitivitas (Recall) terhadap kelas minoritas, yaitu kasus stroke, yang memiliki signifikansi klinis jauh lebih penting dibandingkan hanya mengejar nilai akurasi semata.

Studi terdahulu dalam prediksi risiko stroke menggunakan metode pembelajaran mesin telah mencoba berbagai pendekatan canggih, seperti *Support Vector Machine* (SVM) atau berbagai varian *Ensemble Learning* (misalnya, *Random Forest* atau *Voting System*)[8][9][10]. Meskipun model-model yang lebih kompleks ini sering kali menghasilkan nilai *Accuracy* yang tinggi, mereka sering kali mengabaikan isu *Recall* pada kelas minoritas (kasus stroke) karena bias data. Hal ini menyebabkan tingginya *False Negative* yang tidak dapat diterima dalam konteks klinis. Dengan menargetkan optimasi pada model Gaussian Naive Bayes (GNB) yang secara komputasi lebih sederhana dan ringan, penelitian ini bertujuan untuk menguji apakah optimasi *hyperparameter* dan penyesuaian *threshold* yang tepat dapat menandingi performa model yang lebih kompleks, khususnya dalam hal sensitivitas klinis (tingkat *Recall*).

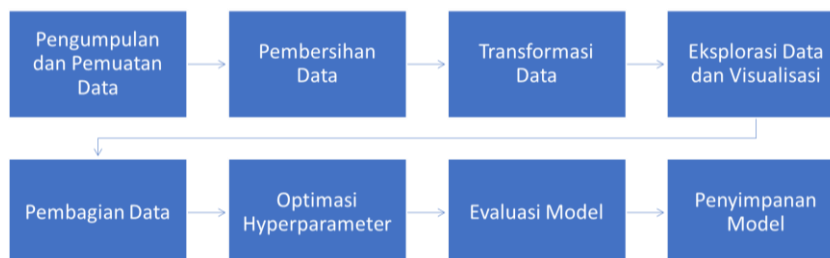
Berdasarkan tinjauan literatur mengenai deteksi risiko stroke menggunakan *Machine Learning*, teridentifikasi empat kesenjangan penelitian (*research gaps*) yang signifikan. Pertama, meskipun berbagai model klasifikasi telah digunakan untuk memprediksi stroke, mayoritas studi pada data *imbalanced* [8] gagal mengatasi bias model yang mengarah pada *False Negative* yang tinggi secara klinis. Pendekatan yang hanya berfokus pada *Accuracy* sering mengabaikan risiko ini, padahal dalam konteks klinis stroke, *False Negative* dapat berakibat fatal [11]. Kedua, model seperti Gaussian Naive Bayes (GNB) telah banyak diterapkan dalam klasifikasi data, namun kinerja GNB terbukti sangat sensitif dan rentan mengalami penurunan performa pada kasus data yang tidak seimbang. Oleh karena itu, optimasi *hyperparameter* (*tuning*) yang spesifik untuk GNB (*var_smoothing*) adalah langkah krusial [9], tetapi studi yang mengeksplorasi optimasi GNB secara mendalam pada data stroke masih sangat terbatas. Ketiga, solusi konvensional untuk ketidakseimbangan kelas seringkali melibatkan teknik *resampling* seperti SMOTE [8]. Namun, efektivitas penyesuaian batas keputusan (*threshold adjustment*) sebagai strategi pasca-pemodelan yang lebih langsung, transparan, dan mampu menggeser *trade-off* antara *Precision* dan *Recall* belum teruji sepenuhnya pada model GNB untuk kasus stroke. Keempat, masih jarang penelitian yang secara khusus menempatkan *Recall* sebagai prioritas utama, padahal metrik ini sangat krusial untuk keselamatan pasien. Dalam konteks klinis, perlu ada alasan yang jelas mengapa *Precision* dapat dikorbankan, asalkan jumlah *False Negative* ditekan seminimal mungkin demi kepentingan deteksi dini. Kesenjangan tersebut menegaskan perlunya pengembangan model prediksi stroke yang tidak hanya menekankan akurasi, tetapi juga mampu lebih peka terhadap kasus minoritas. Berdasarkan hal itu, penelitian ini difokuskan pada upaya mengoptimalkan kinerja Gaussian Naive Bayes (GNB) pada data tidak seimbang. Strategi yang digunakan mencakup optimasi *hyperparameter* (*var_smoothing*) dan penyesuaian ambang keputusan (*threshold adjustment*) untuk meningkatkan *Recall* sekaligus menekan risiko *False Negative*.

Kesenjangan-kesenjangan tersebut menegaskan perlunya pengembangan model prediksi stroke yang tidak hanya menekankan akurasi, tetapi juga sensitif terhadap kasus minoritas. Oleh karena itu, penelitian ini bertujuan mengoptimalkan kinerja Gaussian Naive Bayes (GNB) pada data tidak seimbang melalui penggabungan optimasi *hyperparameter* (*var_smoothing*) dan penyesuaian batas keputusan (*decision threshold*) demi mencapai nilai *Recall* tertinggi dan meminimalkan risiko *False Negative*.

2. METODOLOGI PENELITIAN

2.1 Tahapan Penelitian

Penelitian ini dirancang melalui serangkaian proses sistematis yang dimulai dari pengumpulan data hingga evaluasi model akhir. Seluruh alur tahapan penelitian secara komprehensif disajikan dalam Gambar 1, diikuti oleh penjelasan singkat dari setiap tahapannya.



Gambar 1. Tahapan penelitian

Proses penelitian mengikuti enam langkah utama berikut

- Pengumpulan dan Preprocessing Data, pembersihan *missing values* dan transformasi data (Encoding) dilakukan agar data siap untuk pemodelan.
- Eksplorasi Data, dilakukan analisis untuk memahami karakteristik data, termasuk identifikasi *class imbalance* yang signifikan.
- Pembagian Dataset, data dipisahkan secara Stratified sampling (80% training, 20% testing) untuk menjaga rasio kelas.
- Optimasi Hyperparameter, parameter *var_smoothing* GNB dioptimasi menggunakan GridSearchCV demi stabilitas perhitungan probabilitas.
- Strategi Threshold Adjustment, diterapkan penyesuaian Decision Threshold sebagai strategi utama penanganan *class imbalance* dan peningkatan recall.
- Pelatihan dan Evaluasi Model, model terbaik dievaluasi secara komprehensif menggunakan metrik sensitivitas klinis pada data pengujian.

2.2 Preprocessing Data dan Feature Scaling

Pembersihan data dimulai dengan normalisasi nilai hilang. Placeholder seperti N/A, *Unknown*, atau string kosong dikonversi menjadi NaN. Kolom numerik seperti BMI dan kadar glukosa diubah ke tipe numerik. Untuk mencegah kesalahan, digunakan parameter `errors='coerce'` sehingga nilai tidak valid otomatis diganti NaN [4].

Nilai hilang pada fitur numerik (BMI dan glukosa) diisi menggunakan median dengan metode *SimpleImputer*. Sedangkan pada fitur kategorikal (gender, status pernikahan, jenis pekerjaan, tempat tinggal, dan status merokok), nilai kosong dilengkapi dengan modus atau nilai yang paling sering muncul. Kolom id dihapus karena tidak berkontribusi terhadap prediksi.

Tabel 2. Data yang sudah dibersihkan

Gender	Age	Hypertension	heart_disease	ever_married	work_type	Residence_type	avg_glucose_level	Bmi	Smoking
Male	67	0	1	Yes	Private	Urban	228.69	36.60	formerly smoked
Famale	61	0	0	Yes	Self-employed	Rural	202.21	33.45	Never smoked
Male	80	0	1	Yes	Private	Rural	105.92	32.50	Never smoked
Famale	49	0	0	Yes	Private	Urban	171.23	34.40	smokes
Famale	79	1	0	Yes	Self-employed	Rural	174.12	24.00	Never smoked

Seluruh fitur kategorikal diubah ke dalam bentuk numerik dengan metode Label Encoding, yaitu proses yang mengonversi setiap kategori menjadi bilangan bulat. Langkah ini penting karena algoritma machine learning hanya dapat bekerja dengan data numerik. Untuk menjaga konsistensi hasil dan memastikan model tetap dapat digunakan pada data baru, setiap LabelEncoder disimpan menggunakan `joblib`.

2.3 Eksplorasi Data dan Visualisasi

Analisis eksploratif dilakukan untuk memahami karakteristik data. Deskripsi statistik disajikan untuk semua fitur, baik numerik maupun kategorikal. Visualisasi mencakup:

- Heatmap korelasi untuk melihat hubungan antar fitur,
- Count plot distribusi kelas stroke,
- Histogram distribusi usia,
- Boxplot BMI dan kadar glukosa,
- Scatter plot usi dan kadar glukosa berdasarkan status stroke.

Visualisasi membantu mengidentifikasi pola, outlier, dan distribusi kelas, termasuk ketidakseimbangan yang signifikan antara pasien dengan dan tanpa stroke [12].

2.4 Pembagian Data

Pembagian data dilakukan menggunakan Stratified Sampling (`stratify=y`) untuk menjamin bahwa proporsi kelas minoritas (kasus *stroke*) didistribusikan secara merata di *data training* (80%) dan *data testing* (20%), sehingga hasil pengujian menjadi lebih valid.

2.5 Penanganan Ketidakseimbangan Data (Class Imbalance)

Data risiko stroke yang digunakan dalam penelitian ini merupakan *Imbalanced Dataset*, suatu kondisi yang secara inheren mendorong model untuk memprioritaskan *Accuracy* keseluruhan dan menghasilkan tingkat kesalahan *False Negative* yang tinggi suatu risiko besar dalam konteks diagnosis klinis. Untuk mengatasi masalah ketidakseimbangan

data ini, penelitian ini mengadopsi pendekatan Cost-Sensitive Classification yang diimplementasikan melalui teknik Penyesuaian Ambang Batas Klasifikasi (*Decision Threshold Adjustment*). Kami tidak menggunakan teknik penyeimbangan data berbasis *sampling* (seperti SMOTE, *Oversampling*, atau *Undersampling*). Strategi penanganan *Class Imbalance* difokuskan pada dua tahap:

- Optimasi *Hyperparameter* GNB (*var_smoothing*). Proses optimasi ini menggunakan GridSearchCV (dijelaskan di sub-bab 2.8) untuk menemukan pengaturan *hyperparameter* internal terbaik yang secara probabilistik mendukung prediksi kelas minoritas (stroke).
- Teknik Utama *Decision Threshold Adjustment*, Nilai ambang batas keputusan model diubah secara sistematis dari nilai standar 0,5 menjadi nilai optimal 0,0100. Perubahan *threshold* ini secara efektif memberikan biaya yang lebih tinggi pada kesalahan prediksi *False Negative* (skor yang paling dihindari), sehingga memaksa model untuk mengklasifikasikan kasus yang meragukan ke dalam kelas stroke. Teknik ini adalah strategi utama kami untuk mencapai peningkatan signifikan pada nilai *Recall*[13].

2.6 Optimasi Hyperparameter GNB Menggunakan GridSearchCV

Model Gaussian Naive Bayes (GNB) memiliki *hyperparameter* kunci yang dikenal sebagai *var_smoothing*. Parameter ini berfungsi menambahkan nilai ke varians data untuk mengatasi kasus *varians nol*, suatu penyesuaian yang sangat penting dalam konteks data tidak seimbang karena dapat memengaruhi distribusi probabilitas kelas minoritas. Untuk menemukan nilai optimal dari *var_smoothing* secara sistematis, penelitian ini menggunakan metode GridSearchCV sebuah teknik pencarian grid yang formal dan komprehensif. GridSearchCV diimplementasikan dengan mengeksplorasi rentang nilai *var_smoothing* dari 1×10^{-9} hingga 1×10^{-11} . Karena fokus utama penelitian adalah pada minimasi *False Negative*, metrik evaluasi yang digunakan untuk memandu GridSearchCV adalah *Recall* (Sensitivitas). Hasil optimasi menunjukkan bahwa nilai *var_smoothing* optimal yang berhasil memaksimalkan *Recall* pada data pelatihan adalah $1,00 \times 10^{-10}$, dan nilai ini kemudian digunakan untuk melatih model GNB sebelum diterapkan *Decision Threshold Adjustment*.

2.7 Evaluasi Model

Evaluasi dilakukan setelah model dilatih dengan data uji. Karena dataset stroke tidak seimbang, akurasi saja tidak cukup untuk menilai kinerja model. Sebagai contoh, dengan proporsi pasien stroke hanya 4,9%, model yang selalu memprediksi “tidak stroke” bisa mencapai akurasi 95% tetapi gagal mendeteksi kasus sebenarnya. Dalam konteks klinis, kondisi ini berisiko tinggi karena pasien dengan potensi stroke dapat terlewat [14]. Oleh karena itu, evaluasi model tidak cukup hanya mengandalkan akurasi. Fokus utama adalah memastikan sistem mampu mengenali pasien yang benar-benar berisiko. Karena itu, penilaian kinerja lebih ditekankan pada recall (kemampuan mendeteksi seluruh kasus stroke nyata), precision (ketepatan prediksi positif), serta distribusi hasil melalui matriks konfusi [15]. Pendekatan ini sejalan dengan prinsip evaluasi pada ketidakseimbangan data, dengan kondisi *false negative* penderita stroke tidak terdeteksi memiliki dampak klinis yang jauh lebih serius dibandingkan *false positive* [11]. Dengan demikian, metrik-metrik tersebut diposisikan bukan sekadar angka akhir, melainkan sebagai indikator kualitas deteksi dini yang relevan secara medis.

2.8 Gaussian Naive Bayes

Varian *Gaussian Naive Bayes* (GNB) adalah metode hasil pengelompokan diadaptasi dari algoritma *Naive Bayes* dan berakar pada Teorema Bayes. Model ini berfungsi untuk menghasilkan keputusan klasifikasi dengan mengutamakan probabilitas posterior tertinggi [16]. Sistem prediktif yang dihasilkan berfungsi sebagai upaya memprediksi kemungkinan data masuk ke dalam kelas tertentu[17]. Setiap fitur dalam data direpresentasikan sebagai vektor, sedangkan label kelas ditentukan melalui pendekatan probabilistik. Pada tahap pelatihan, GaussianNB mengasumsikan setiap variabel dianggap tidak memiliki keterkaitan dengan variabel lain. Namun, model ini memiliki keterbatasan, terutama ketika jumlah fitur semakin banyak, karena dapat menimbulkan bias pada tabel probabilitas [18]. Dalam praktiknya, GaussianNB menyimpan model terbaik yang telah melewati proses validasi.

- Menentukan probabilitas prior ($P(C)$):

Untuk mengestimasi probabilitas *prior* ($P(C)$), digunakan perhitungan rasio antara kardinalitas kelas C dan total *instance* (sampel) pada datasetl.

- Mengestimasi nilai likelihood:

$$P(X|C) = \left(\frac{1}{\sqrt{2\pi\sigma^2}}\right) * e^{-\frac{(X-\mu)^2}{2\sigma^2}} \quad (1)$$

Dalam analisis data, variabel X merepresentasikan nilai observasi dari suatu fitur yang sedang dikaji. Nilai ini dibandingkan dengan μ (mu), yaitu rata-rata fitur tersebut untuk seluruh sampel dalam kelas target tertentu. Sementara itu, σ^2 (sigma kuadrat) menggambarkan tingkat keragaman atau sebaran data fitur di dalam kelas target (C). Ketiga komponen ini X , μ , dan σ^2 berperan penting dalam menjelaskan pola distribusi data pada setiap kelas. Pemahaman terhadap hubungan ketiganya menjadi dasar dalam penerapan berbagai metode statistik dan algoritma pembelajaran mesin, seperti *Naive Bayes classifier*, yang bergantung pada distribusi probabilistik untuk melakukan prediksi secara akurat.



- c. Hitung probabilitas posterior dengan rumus Teorema Bayes

$$P(C|X) = (P(X|C) * P(C))/P(X) \tag{2}$$

Dalam model probabilistik seperti *Naive Bayes classifier*, $P(C|X)$ merepresentasikan *probabilitas posterior*, yaitu kemungkinan suatu data termasuk ke dalam kelas C setelah mempertimbangkan nilai fitur X. Selanjutnya, $P(X|C)$ disebut *probabilitas likelihood*, yang menggambarkan seberapa besar kemungkinan fitur X muncul jika data berasal dari kelas C. Kemudian, $P(C)$ merupakan *probabilitas awal (prior probability)* dari kelas target C sebelum adanya pengamatan terhadap data, sedangkan $P(X)$ disebut *probabilitas bukti (evidence probability)*, yaitu probabilitas terjadinya fitur X secara keseluruhan pada dataset tanpa mempertimbangkan kelas tertentu. Keempat komponen ini saling berkaitan dan menjadi dasar dalam menghitung probabilitas klasifikasi menggunakan Teorema Bayes

- d. Label kelas final ditentukan berdasarkan nilai probabilitas posterior tertinggi

2.8 Strategi Penyesuaian Decision Threshold

Setelah model GNB mencapai `var_smoothing` optimal, dilakukan penyesuaian batas keputusan (*decision threshold*). Strategi pasca-pemodelan ini secara eksplisit dirancang untuk mengatasi Kesenjangan Penelitian mengenai fokus berlebihan pada teknik resampling konvensional. Model Gaussian Naive Bayes pada dasarnya menghasilkan skor probabilitas $P(\text{Stroke}=1|\text{Data})$ yaitu peluang bahwa sampel data tersebut termasuk dalam kelas stroke [10]. Secara default, label kelas ditentukan oleh aturan: Jika $P(\text{Stroke}=1) \geq 0.5$ maka Kelas = 'Stroke'. Namun, pada data yang tidak seimbang (hanya 4.9% kasus stroke), threshold standar 0.5 terlalu tinggi, menyebabkan model cenderung konservatif dan menghasilkan tingkat False Negative yang tinggi. Oleh karena itu, penelitian ini mengeksplorasi nilai threshold yang lebih rendah (misalnya 0.40, 0.30, 0.20, dst.) untuk menggeser trade-off antara Precision dan Recall. Penurunan threshold secara efektif membuat model lebih sensitif, meningkatkan peluang klasifikasi sebagai stroke untuk kasus yang sebelumnya dianggap non-stroke. Kriteria Pemilihan Threshold Optimal: Prioritas Klinis, Threshold optimal dipilih berdasarkan nilai yang memberikan Recall tertinggi pada kelas stroke. Hal ini sejalan dengan prinsip Cost-Sensitive Learning yang menyatakan bahwa biaya klinis False Negative (pasien berisiko terlewatkan) jauh lebih tinggi daripada False Positive (pasien sehat menjalani tes lanjutan yang tidak perlu). Pembatasan Precision, Meskipun Recall diprioritaskan, threshold tidak boleh diturunkan terlalu ekstrem. Threshold optimal harus tetap mempertahankan nilai Precision di atas batas minimum yang dapat ditoleransi (misalnya, Precision harus tetap di atas 5%) untuk menjaga relevansi prediksi [19]. Strategi ini secara langsung memprioritaskan keselamatan pasien dengan menekan False Negative seminimal mungkin, sekaligus menunjukkan alternatif metodologis yang efektif selain resampling.

3. HASIL DAN PEMBAHASAN

3.1 Pengumpulan dan Pemuatan Data

Dataset yang dimanfaatkan dalam studi ini berjudul "healthcare-dataset-stroke-data" yang tersedia di Kaggle dapat diakses pada <https://www.kaggle.com/datasets/fedesoriano/stroke-prediction-dataset>. Dataset ini berisi 5.110 rekam medis pasien dengan 11 atribut, namun kolom id dihapus karena tidak memiliki nilai prediktif. Dataset risiko stroke ini memanfaatkan sepuluh variabel berbeda. Variabel-variabel tersebut terdiri dari faktor klinis (riwayat hipertensi, kondisi penyakit jantung, level glukosa rata-rata, nilai BMI, serta riwayat merokok), disertai dengan data demografi (gender, status menikah, kategori profesi, dan tipe hunian) [20]. Dataset diimpor dengan deteksi otomatis *delimiter*, lalu melalui tahap pembersihan dengan menghapus nilai tidak valid seperti N/A dan *Unknown*, sehingga siap digunakan untuk pemodelan [21].

Tabel.1 Hasil Sample Pengumpulan data

Id	Gender	Age	Hypertension	heart_disease	ever_married	Kategori profesi	Jenis domisili	avg_glucose_level	Bmi	Smoking
9046	Male	67	0	1	Yes	Private	Urban	228.69	36.6	formerly smoked
51676	Female	61	0	0	Yes	Self-employed	Rural	202.21	N/A	Never smoked
31112	Male	80	0	1	Yes	Private	Rural	105.92	32.5	Never smoked
60182	Female	49	0	0	Yes	Private	Urban	171.23	34.4	smokes
1665	Female	79	1	0	Yes	Self-employed	Rural	174.12	24	Never smoked

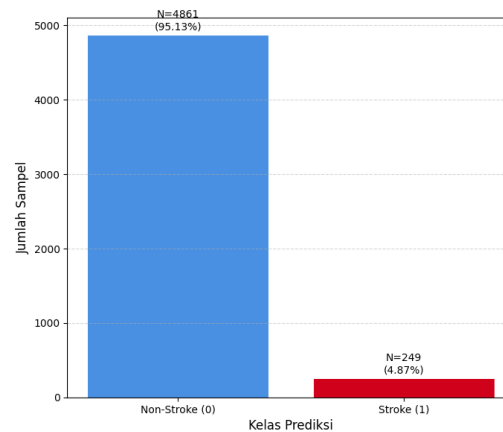
3.2 Distribusi Kelas dan Karakteristik Dataset

Sebelum melakukan pemodelan, dilakukan analisis terhadap distribusi kelas pada variabel target stroke. Hasilnya menunjukkan bahwa dataset bersifat tidak seimbang (imbalanced).

Tabel 3. Jumlah Kelas Stroke

Stroke	0	No	4861
	1	Yes	249

Ketidakeimbangan kelas ini dapat divisualisasikan lebih lanjut pada Gambar 2.

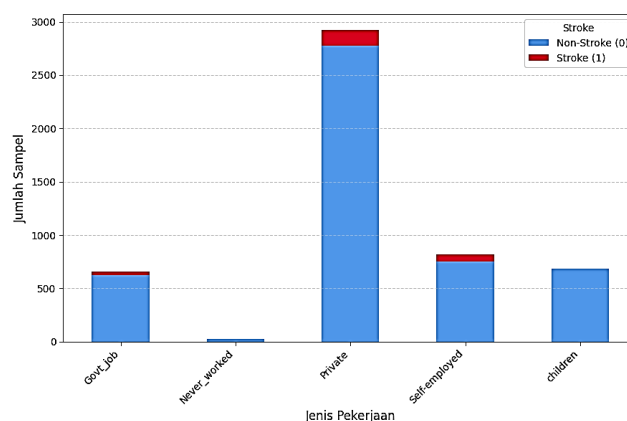


Gambar 2. Distribusi Kelas Target

Analisis awal dilakukan pada distribusi variabel target stroke. Hasil menunjukkan bahwa dataset tidak seimbang: 4.861 pasien (95,1%) tercatat tidak mengalami stroke, sedangkan hanya 249 pasien (4,9%) yang mengalami stroke. Kondisi ini mencerminkan realitas bahwa stroke merupakan kasus dengan prevalensi rendah, namun penting secara klinis. Tantangan utama dari distribusi ini adalah potensi bias model terhadap kelas mayoritas, sehingga meningkatkan risiko *false negative* kasus stroke yang terlewat dalam prediksi.

Visualisasi data memberikan pemahaman awal mengenai karakteristik dataset dan keterkaitan antar variabel dalam konteks prediksi stroke. Distribusi kelas menunjukkan ketidakeimbangan yang cukup besar, di mana terdapat 4.861 pasien tidak mengalami stroke (95,1%) dan hanya 249 pasien yang mengalami stroke (4,9%). Kondisi ini merefleksikan realitas, bahwa kasus stroke relatif jarang terjadi. Dari sisi usia, pola distribusi menunjukkan bentuk bimodal dengan konsentrasi pada kelompok usia muda (0–30 tahun) dan lansia (50–80 tahun). Risiko stroke terlihat meningkat seiring bertambahnya usia.

Hasil heatmap korelasi mengindikasikan adanya hubungan positif moderat antara stroke dengan usia, hipertensi, dan penyakit jantung [22]. Namun, tidak ada variabel yang menunjukkan korelasi kuat ($>0,5$), sehingga prediksi stroke memerlukan integrasi dari berbagai faktor. Analisis boxplot pada BMI dan kadar glukosa rata-rata memperlihatkan bahwa pasien stroke cenderung memiliki glukosa lebih tinggi [23]. Meski demikian, distribusi BMI antara kelompok stroke dan non-stroke sangat tumpang tindih, sehingga BMI tidak dapat dijadikan indikator tunggal. Scatter plot hubungan usia dan kadar glukosa menunjukkan adanya peningkatan glukosa seiring pertambahan usia. Namun, batas yang jelas antara pasien stroke dan non-stroke tetap sulit diidentifikasi, sehingga diperlukan model probabilistik yang mampu menangkap kompleksitas interaksi antarfitur. Visualisasi data kategorikal pada Gambar 2 memperlihatkan perbandingan jumlah kasus stroke berdasarkan kategori jenis pekerjaan.



Gambar 3. Proporsi kejadian stroke berdasarkan jenis pekerjaan

Secara keseluruhan, visualisasi ini berfungsi sebagai landasan penting dalam memahami struktur data sebelum pemodelan.

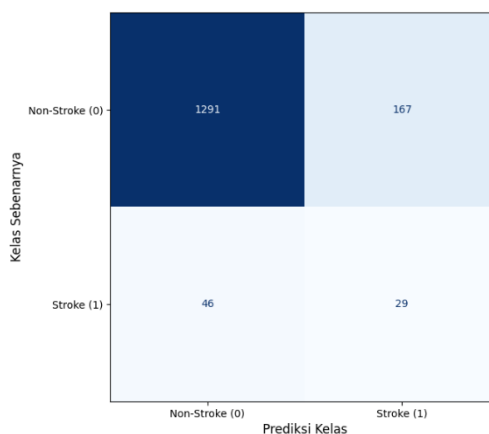
3.3 Komparasi Model dan Analisis Kinerja Kritis

Laporan hasil klasifikasi memberikan gambaran rinci tentang performa model per kelas pada Tabel 5.

Tabel 5. Tabel Classification Report

	Precision	Recall	F1-Score	Support
0	0.98	0.59	0.74	972
1	0.09	0.80	0.16	50
Accuracy			0.60	1022
Macro Avg	0.54	0.69	0.45	1022
Weighted avg	0.94	0.60	0.71	1022

Visualisasi hasil prediksi secara menyeluruh, yang mencakup jumlah *True Positive* dan *False Negative* pada masing-masing kelas, dapat dilihat pada Gambar 4.



Gambar 4. Matriks Kebingungan (*Confusion Matrix*) Model GNB

Evaluasi model ditampilkan pada *classification report* (Tabel 5). Hasil menunjukkan bahwa model yang dioptimasi berhasil mencapai *Recall* sebesar 0,80 pada kelas minoritas (stroke). Artinya, 80% pasien stroke nyata berhasil terdeteksi. Capaian ini jauh lebih baik dibandingkan model dasar, yang hanya mendeteksi 44% kasus stroke. Namun, trade-off terjadi pada *Precision* dan *Accuracy*. Meskipun nilai *Recall* meningkat, terdapat *trade-off* pada metrik *Precision* dan *Accuracy*. Nilai *Precision* untuk kelas stroke menurun signifikan menjadi 0,09. Artinya, hanya 9% dari seluruh prediksi positif benar-benar merupakan pasien stroke, sementara sisanya termasuk ke dalam *false positive*. Jumlah kasus *false positive* ini dapat dilihat pada sel kanan atas di Gambar 4. Selain itu, *Accuracy* keseluruhan model juga menurun menjadi 0.60. Meskipun *Precision* dan *Accuracy* rendah, hasil ini tetap sejalan dengan tujuan penelitian. Keselamatan pasien menjadi prioritas utama, sehingga model harus mampu mendeteksi risiko stroke sebanyak mungkin. Trade-off ini dapat diterima karena dampak dari *False Negative* pada penyakit stroke jauh lebih serius dibandingkan *False Positive*. Penurunan *Accuracy* keseluruhan (menjadi 0.5988) dan *Precision* yang sangat rendah (0.0909) merupakan konsekuensi yang diterima (*acceptable trade-off*) demi kepentingan klinis. Rendahnya *Precision* berarti terdapat banyak *False Positive* (FP) (pasien sehat diprediksi *stroke*). Namun, dalam penanganan *stroke*, konsekuensi dari *False Negative* (FN) yaitu pasien *stroke* terlewat adalah kematian atau kecacatan permanen [24], yang memiliki biaya sosial dan kemanusiaan jauh lebih tinggi daripada biaya *False Positive* (pemeriksaan lanjutan yang tidak perlu). Temuan ini sejalan dengan prinsip *Cost-Sensitive Learning*, di mana biaya kegagalan untuk mendeteksi (FN) diberi bobot yang jauh lebih tinggi daripada biaya prediksi salah (FP). Dengan kata lain, model ini dirancang untuk menjadi alarm yang lebih sensitif di lingkungan klinis.

Setelah dilakukan *hyperparameter tuning* dengan *GridSearchCV*, diperoleh nilai optimal untuk *var_smoothing* yang kemudian digunakan pada pelatihan model akhir. Proses optimasi ini bertujuan agar model menjadi lebih stabil sekaligus lebih sensitif dalam memprediksi kelas minoritas. Untuk memberikan gambaran yang lebih komprehensif, sangat disarankan menambahkan tabel atau pembahasan yang membandingkan performa *Gaussian Naive Bayes* (GNB) Default dengan GNB yang telah dioptimasi [25]. Hasil dari evaluasi tersebut ditampilkan pada Tabel 4 sebagai rangkuman kinerja model dalam keseluruhan proses pemodelan.

Tabel 4. Komparasi (Optimasi Hyperparameter & Threshold)

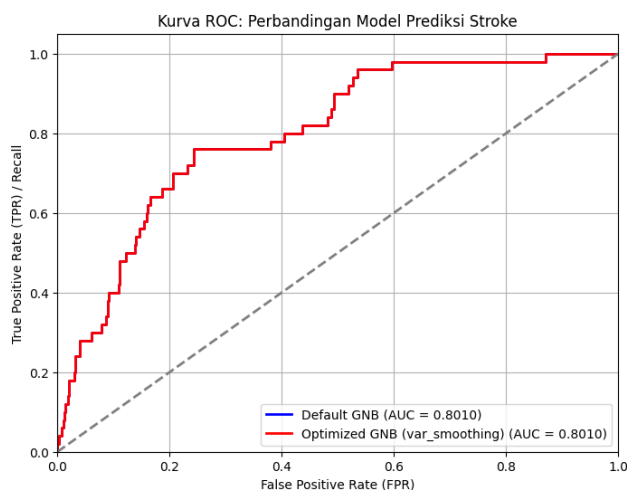
Model	Accuracy	Precision	Recall	F1-Score	Roc-Auc	Threshold
Default GNB	0.8669	0.1692	0.4400	0.2444	0.8010	0.5000
Optimized GNB + Threshold	0.5988	0.0909	0.8000	0.1633	0.8010	0.0100

Hasil komparasi pada Tabel 4 menunjukkan bahwa model Gaussian Naive Bayes telah melalui dua tahap optimasi penting dalam menangani ketidakseimbangan data. Tahap pertama dilakukan dengan optimasi *hyperparameter var_smoothing* menggunakan GridSearchCV, di mana nilai terbaik yang diperoleh adalah 1.00×10^{-10} . Optimasi ini berfungsi menjaga stabilitas perhitungan probabilitas pada GNB, yang menjadi dasar penting sebelum proses penyesuaian batas keputusan. Meskipun optimasi pada parameter *var_smoothing* tidak memberikan perubahan signifikan terhadap hasil prediksi biner ditunjukkan oleh nilai AUC-ROC yang tetap stabil pada 0.8010 fokus utama penelitian ini diarahkan pada optimasi *threshold*. Pada model dasar Gaussian Naive Bayes (GNB) dengan *threshold* 0.5000, nilai Recall hanya mencapai 0.4400. Hal ini berarti model melewatkan sekitar 56% kasus stroke nyata (*false negative*), yang tentu sangat berisiko dalam konteks klinis.

Fokus utama penelitian adalah pada penyesuaian *threshold*. Pada konfigurasi standar (*threshold* 0,5000), *Recall* hanya mencapai 0,4400 sehingga 56% kasus stroke terlewat. Setelah *threshold* diturunkan menjadi 0,0100, nilai *Recall* melonjak ke 0,8000. Artinya, model mampu mendeteksi 80% kasus stroke nyata, meskipun harus mengorbankan *Precision* dan *Accuracy*. Temuan ini menegaskan bahwa dalam konteks klinis, prioritas utama adalah mengurangi *false negative*, bukan sekadar mempertahankan akurasi keseluruhan [24].

3.3 Menekankan Kemampuan Diskriminatif

Untuk menilai kemampuan diskriminatif model tanpa bergantung pada *threshold*, digunakan kurva ROC-AUC (Gambar 5). Grafik menunjukkan bahwa model GNB sebelum dan sesudah optimasi *var_smoothing* memiliki kurva hampir identik dengan nilai AUC-ROC tetap 0,8010. Hal ini mengindikasikan bahwa optimasi parameter tersebut tidak banyak meningkatkan kemampuan model dalam meranking probabilitas, yang pada gilirannya mengimplikasikan bahwa model GNB telah mencapai batas performa diskriminatif intrinsiknya (*inherent limit*) dalam memisahkan kedua kelas. Nilai AUC-ROC sebesar 0,8010 sendiri dikategorikan sebagai kemampuan diskriminatif yang baik, namun keterbatasan inilah yang membenarkan keputusan penelitian untuk fokus pada strategi pasca-pemodelan. Temuan ini semakin menguatkan argumentasi bahwa dalam penanganan data klinis yang sangat tidak seimbang, modifikasi pada *output* model (seperti *Decision Threshold*) seringkali lebih efektif dan berdampak langsung dibandingkan penyesuaian parameter internal yang hanya memengaruhi distribusi probabilitas dasar. Oleh karena itu, AUC-ROC (0,8010) berfungsi sebagai batas atas kemampuan prediksi, sementara penyesuaian *threshold* berperan sebagai alat strategis untuk mencapai kinerja *Recall* yang diinginkan untuk keselamatan pasien.



Gambar 5. Perbandingan Kurva ROC antara Model GNB Default dan GNB Optimized

Kurva ROC pada kedua model terlihat hampir berimpit, dengan nilai AUC-ROC sebesar 0.8010. Temuan ini mengindikasikan bahwa optimasi *var_smoothing* tidak memberikan peningkatan berarti terhadap kemampuan model dalam meranking probabilitas.

4. KESIMPULAN

Penelitian ini berfokus pada optimalisasi model Gaussian Naive Bayes (GNB) dalam menghadapi tantangan prediksi risiko stroke yang muncul akibat distribusi data yang sangat tidak seimbang (*imbalanced dataset*), di mana kasus stroke hanya mewakili sekitar 4,9% dari total data. Proses optimalisasi dilakukan melalui dua langkah utama, yaitu penyetelan *hyperparameter* dan penyesuaian ambang keputusan (*decision threshold*) untuk meningkatkan kemampuan deteksi terhadap kelas minoritas. Walaupun penyetelan *hyperparameter* berperan penting dalam menstabilkan perhitungan probabilitas GNB terlihat dari nilai AUC-ROC yang tetap konstan pada 0,8010 peningkatan performa terbesar justru berasal dari penyesuaian *threshold*. Pada konfigurasi standar (*threshold* = 0,5000), model hanya mencapai nilai Recall 0,4400, yang berarti lebih dari separuh kasus stroke nyata tidak terdeteksi (*false negative*).



Setelah dilakukan eksplorasi terhadap batas keputusan, diperoleh threshold optimal pada nilai 0,0100, sehingga Recall meningkat signifikan menjadi 0,8000. Dengan kata lain, model berhasil mengenali sekitar 80% dari total kasus stroke yang sebenarnya ada. Peningkatan Recall ini memang diikuti penurunan pada metrik Accuracy (0,5988) dan Precision (0,0909) akibat bertambahnya false positive. Namun, kompromi ini masih dapat diterima dalam konteks klinis, mengingat pendekatan tersebut sejalan dengan prinsip Cost-Sensitive Learning. Dalam deteksi dini stroke, risiko kegagalan mendeteksi pasien berisiko (false negative) jauh lebih fatal karena dapat menyebabkan kecacatan permanen atau kematian, dibandingkan dampak dari false positive yang hanya menimbulkan pemeriksaan tambahan. Oleh karena itu, strategi penyesuaian ambang keputusan (threshold adjustment) terbukti efektif dalam menekan risiko false negative dan menghasilkan model dengan sensitivitas tinggi, menjadikannya relevan serta bernilai signifikan bagi penerapan di lingkungan klinis.

REFERENCES

- [1] Y. Y. Saifullah, Mochammad Erwin Rachman, Ramlihan, Lilian Triana Limoa, and Nurussyariah Hamado, "Literature Review: Hubungan Hipertensi dengan Kejadian Stroke Iskemik dan Stroke Hemoragik," *Fakumi Med. J. J. Mhs. Kedokt.*, vol. 4, no. 10, pp. 695–708, Oct. 2024, doi: 10.33096/fmj.v4i10.477.
- [2] A. T. P. Ilham Darmawan, Indhit Tri Utami, "Penerapan Range Of Motion (ROM) Exercise Bola Karet Terhadap Kekuatan Otot Pasien Stroke Non Hemoragik," *J. Cendikia Muda*, vol. 4, pp. 246–254, 2024.
- [3] E. A. Handoko and Z. Muslim, "Profil Pasien Stroke Di Rs Bhayangkara Bondowoso Tahun 2024," *Syntax Idea*, vol. 7, no. 3, pp. 475–486, Mar. 2025, doi: 10.46799/syntaxidea.v7i3.12722.
- [4] A. Adi Bhirawa and U. Pradema Sanjaya, "From Data Imbalance to Precision: SMOTE-Driven Machine Learning for Early Detection of Kidney Disease," *INOVTEK Polbeng - Seri Inform.*, vol. 10, no. 1, pp. 514–525, Mar. 2025, doi: 10.35314/7jgimg64.
- [5] D. M. Makarov and A. M. Kolker, "Viscosity of deep eutectic solvents: Predictive modeling with experimental validation," *Fluid Phase Equilib.*, vol. 587, p. 114217, Jan. 2025, doi: 10.1016/j.fluid.2024.114217.
- [6] F. Y. A'la, "Optimasi Klasifikasi Sentimen Ulasan Game Berbahasa Indonesia: IndoBERT dan SMOTE untuk Menangani Ketidakseimbangan Kelas," *Edumatic J. Pendidik. Inform.*, vol. 9, no. 1, pp. 256–265, Apr. 2025, doi: 10.29408/edumatic.v9i1.29666.
- [7] S. Ernawati and I. Maulana, "Meningkatkan Klasifikasi Penyakit Diabetes Menggunakan Metode Ensemble Softvoting Dengan SMOTE-ENN dan Optimasi Bayesian," *Evolusi J. Sains dan Manaj.*, vol. 13, no. 1, pp. 71–86, Mar. 2025, doi: 10.31294/evolusi.v13i1.8267.
- [8] M. Sulistiyono, Y. Pristiyanto, S. Adi, and G. Gumelar, "Implementasi Algoritma Synthetic Minority Over-Sampling Technique untuk Menangani Ketidakseimbangan Kelas pada Dataset Klasifikasi," *Sistemasi*, vol. 10, no. 2, p. 445, 2021, doi: 10.32520/stmsi.v10i2.1303.
- [9] Joshua Agung Nurcahyo and Theopilus Bayu Sasongko, "Hyperparameter Tuning Algoritma Supervised Learning untuk Klasifikasi Keluarga Penerima Bantuan Pangan Beras," *Indones. J. Comput. Sci.*, vol. 12, no. 3, pp. 1351–1365, 2023, doi: 10.33022/ijcs.v12i3.3254.
- [10] M. Hasan *et al.*, "Enhancing stroke disease classification through machine learning models via a novel voting system by feature selection techniques," *PLoS One*, vol. 20, no. 1, p. e0312914, Jan. 2025, doi: 10.1371/journal.pone.0312914.
- [11] A. J. Appukutty, L. E. Skolarus, M. V. Springer, W. J. Meurer, and J. F. Burke, "Increasing false positive diagnoses may lead to overestimation of stroke incidence, particularly in the young: a cross-sectional study," *BMC Neurol.*, vol. 21, no. 1, pp. 1–10, 2021, doi: 10.1186/s12883-021-02172-1.
- [12] K. Apostolidis *et al.*, "Innovative Visualization Approach for Biomechanical Time Series in Stroke Diagnosis Using Explainable Machine Learning Methods: A Proof-of-Concept Study," *Information*, vol. 14, no. 10, p. 559, Oct. 2023, doi: 10.3390/info14100559.
- [13] S. S. Rambe, A. Asriyanik, and P. Prajoko, "PENERAPAN MODEL CONVOLUTIONAL NEURAL NETWORK (CNN) BERBASIS MOBILENETV2 UNTUK KLASIFIKASI TINGKAT KESEGERAN IKAN NILA," *J. Inform. dan Tek. Elektro Terap.*, vol. 13, no. 3, Jul. 2025, doi: 10.23960/jitet.v13i3.6744.
- [14] T. Saito and M. Rehmsmeier, "The Precision-Recall Plot Is More Informative than the ROC Plot When Evaluating Binary Classifiers on Imbalanced Datasets," *PLoS One*, vol. 10, no. 3, p. e0118432, Mar. 2015, doi: 10.1371/journal.pone.0118432.
- [15] Md. Khalilur Rahman, Md. Ashikur Rahman Khan, Ishtiaq Ahammad, and Joysri Rani Das, "Comparative Evaluation of Machine Learning Models for Stroke Prediction in Clinical Settings," *Cloud Comput. Data Sci.*, pp. 196–216, May 2025, doi: 10.37256/ccds.6220256976.
- [16] P. L. Reddy and M. Amanullah, "Prediction of wireless sensor network attack by using the gradient boosting classifier algorithm compared with Gaussian Naive Bayes with improved accuracy," 2025, p. 020146. doi: 10.1063/5.0258920.
- [17] K. Phung, E. Ogunshile, and M. E. Aydin, *Domain-specific implications of error-type metrics in risk-based software fault prediction*, vol. 33, no. 1. 2025. doi: 10.1007/s11219-024-09704-1.
- [18] S. K. Singhi and H. Liu, "Feature subset selection bias for classification learning," in *Proceedings of the 23rd international conference on Machine learning - ICML '06*, New York, New York, USA: ACM Press, 2006, pp. 849–856. doi: 10.1145/1143844.1143951.
- [19] J. J. Chen, C.-A. Tsai, H. Moon, H. Ahn, J. J. Young, and C.-H. Chen, "Decision threshold adjustment in class prediction," *SAR QSAR Environ. Res.*, vol. 17, no. 3, pp. 337–352, Jun. 2006, doi: 10.1080/10659360600787700.
- [20] R. S. Rohman, R. A. Saputra, and D. A. Firmansaha, "Komparasi Algoritma C4.5 Berbasis PSO Dan GA Untuk Diagnosa Penyakit Stroke," *CESS (Journal Comput. Eng. Syst. Sci.)*, vol. 5, no. 1, p. 155, 2020, doi: 10.24114/cess.v5i1.15225.
- [21] W. Garcia, "Detecting CSV file dialects by table uniformity measurement and data type inference," *Data Sci.*, vol. 7, no. 2, pp. 55–72, Nov. 2024, doi: 10.3233/DS-240062.
- [22] T. Li *et al.*, "Comprehensive bioinformatics analysis identifies LAPTM5 as a potential blood biomarker for hypertensive patients with left ventricular hypertrophy," *Aging (Albany. NY)*, vol. 14, no. 3, pp. 1508–1528, Feb. 2022, doi:



- 10.18632/aging.203894.
- [23] S. Sidiq, Alfian, and N. S. Mabur, “Pengembangan Model Prediksi Risiko Diabetes Menggunakan Pendekatan AdaBoost dan Teknik Oversampling SMOTE,” *J. Ilm. Inform. dan Ilmu Komput.*, vol. 4, no. 1, pp. 13–23, 2025.
- [24] D. McMahon, C. Micallef, and T. J. Quinn, “Review of clinical practice guidelines relating to cognitive assessment in stroke,” *Disabil. Rehabil.*, vol. 44, no. 24, pp. 7632–7640, Nov. 2022, doi: 10.1080/09638288.2021.1980122.
- [25] T. Taslim, S. Handayani, and F. Fajrizal, “Kinerja Komparatif Optimasi Algoritma Naive Bayes dalam Klasifikasi Teks untuk Uji Klinis Kanker,” *J. Eksplora Inform.*, vol. 13, no. 1, pp. 113–123, Sep. 2023, doi: 10.30864/eksplora.v13i1.994.