

Implementasi Data Mining dan Machine Learning untuk Segmentasi Pelanggan: Pendekatan Hybrid Menggunakan Big Data

Edy Prayitno¹, Ivan Jaka Perdana², Asyahri Hadi Nashuha^{3*}

^{1,3}Fakultas Teknologi Informasi, Universitas Teknologi Digital Indonesia

²Fakultas Manajemen dan Bisnis, Universitas Teknologi Digital Indonesia

Email: ¹edyprayitno@utdi.ac.id, ²ivanjaka@utdi.ac.id, ³asyahrihadi@gmail.com

Penulis Korespondensi*

(received: 15-10-24, revised: 02-11-24, accepted: 17-03-25)

Abstrak

Penelitian ini bertujuan untuk mengatasi tantangan dalam segmentasi pelanggan berbasis big data melalui pendekatan hybrid yang menggabungkan data mining dan machine learning. Pendekatan ini diharapkan dapat meningkatkan akurasi dan efisiensi dalam pengelompokan pelanggan untuk strategi pemasaran yang lebih optimal. Dataset simulasi yang digunakan merepresentasikan fitur pelanggan seperti usia, pendapatan tahunan, dan skor pengeluaran. Algoritma k-Means digunakan untuk mengelompokkan pelanggan ke dalam tiga kluster utama, sementara algoritma Random Forest diterapkan untuk klasifikasi lebih lanjut berdasarkan hasil clustering. Evaluasi menunjukkan bahwa model Random Forest mencapai akurasi 100% dalam klasifikasi pelanggan, meskipun hasil ini mungkin terkait dengan kondisi ideal data simulasi dan memerlukan evaluasi lebih lanjut. Selain itu, pendekatan hybrid terbukti lebih efisien dengan waktu eksekusi yang lebih singkat (120 detik) dibandingkan dengan pendekatan non-hybrid (180 detik). Penelitian ini menunjukkan potensi pendekatan hybrid dalam aplikasi segmentasi pelanggan, yang berimplikasi pada personalisasi layanan dan pengembangan strategi pemasaran yang lebih tepat sasaran. Untuk penelitian selanjutnya, direkomendasikan evaluasi dengan dataset nyata untuk validasi yang lebih komprehensif.

Kata Kunci: Data Mining, Machine Learning, Segmentasi Pelanggan, Big Data, Pendekatan Hybrid

Abstract

This study aims to evaluate the effectiveness of a hybrid approach combining data mining and machine learning for customer segmentation based on big data. The dataset used in this research is a simulated customer dataset consisting of features such as age, annual income, and spending score. The k-Means algorithm was employed to cluster customers into three segments, while the Random Forest algorithm was applied to classify customers based on the clustering results. The evaluation results show that the Random Forest model achieved 100% accuracy in customer classification. Additionally, the hybrid approach proved to be more efficient, with a faster execution time (120 seconds) compared to the non-hybrid approach (180 seconds). This research demonstrates that the hybrid approach not only delivers high accuracy but also efficiency in processing large-scale data. The implications of these findings provide insights for businesses to optimize marketing strategies and service personalization based on customer segmentation. This study also serves as a foundation for future research on more complex hybrid models in broader domains.

Keyword: Data Mining, Machine Learning, Customer Segmentation, Big Data, Hybrid Approach

1. PENDAHULUAN

Dalam era digital, big data telah menjadi elemen penting dalam pengambilan keputusan bisnis, khususnya dalam memahami dan memprediksi perilaku pelanggan. Meningkatnya volume, variasi, dan kecepatan data menuntut perusahaan untuk mengembangkan metode yang lebih efisien dalam memanfaatkan data ini untuk segmentasi pelanggan. Segmentasi pelanggan yang akurat dapat membantu perusahaan mempersonalisasi layanan, mengoptimalkan strategi pemasaran, dan meningkatkan loyalitas pelanggan [1]. Namun, segmentasi berbasis big data memerlukan pendekatan yang dapat menangani kompleksitas dan jumlah data yang besar.

Pendekatan hybrid yang menggabungkan data mining dan machine learning semakin menarik perhatian karena kemampuannya dalam menangani big data dengan lebih efektif. Data mining memungkinkan penemuan pola-pola tersembunyi dalam dataset yang besar, sedangkan machine learning dapat meningkatkan akurasi prediksi berdasarkan pola-pola tersebut [2]. Kombinasi dari kedua metode ini memberikan potensi untuk mencapai segmentasi yang lebih akurat dan efisien, yang sulit dicapai hanya dengan satu pendekatan saja. Oleh

karena itu, penelitian ini bertujuan untuk mengeksplorasi pendekatan hybrid dalam segmentasi pelanggan, yang diharapkan mampu mengatasi keterbatasan pendekatan konvensional [3].

Meskipun pendekatan machine learning telah banyak digunakan untuk segmentasi pelanggan, penelitian tentang penggabungan data mining dan machine learning dalam konteks big data masih relatif terbatas. Sebagai contoh, studi oleh Sinaga menunjukkan bahwa algoritma k-Means dapat secara efektif mengelompokkan pelanggan berdasarkan preferensi pembelian mereka [4]. Di sisi lain, algoritma Random Forest telah terbukti sebagai metode yang andal dalam klasifikasi, terutama dalam menangani data yang kompleks [5]. Namun, penelitian yang menguji penggabungan kedua algoritma ini sebagai pendekatan hybrid dalam segmentasi pelanggan belum banyak dilakukan, sehingga terdapat celah penelitian yang perlu diisi.

Penelitian ini menghipotesiskan bahwa pendekatan hybrid yang menggabungkan algoritma k-Means dan Random Forest akan memberikan hasil yang lebih unggul dalam hal akurasi dan efisiensi dibandingkan dengan penggunaan algoritma secara terpisah. Dengan memanfaatkan kapasitas pemrosesan big data, pendekatan ini diharapkan dapat memberikan wawasan yang lebih dalam dan relevan bagi pengambilan keputusan bisnis. Studi ini juga bertujuan untuk mengevaluasi performa waktu eksekusi antara pendekatan hybrid dan non-hybrid, mengingat efisiensi merupakan faktor penting dalam analisis data besar [6].

Dengan hasil yang diperoleh, penelitian ini berpotensi memberikan kontribusi signifikan bagi perusahaan dalam meningkatkan strategi pemasaran yang lebih personal, efisien, dan berbasis data. Untuk validasi lebih lanjut, penelitian ini juga membuka peluang untuk melakukan pengujian dengan data nyata agar pendekatan yang diusulkan dapat diaplikasikan secara lebih luas dalam berbagai sektor industri [7].

2. METODOLOGI PENELITIAN

Penelitian ini menggunakan pendekatan hybrid yang menggabungkan data mining dan machine learning untuk segmentasi pelanggan dengan data simulasi yang menyerupai big data. Tahapan penelitian terdiri dari pembuatan data simulasi, preprocessing, penerapan algoritma data mining, penerapan algoritma machine learning, serta evaluasi performansi model.

a. Pembuatan Data Simulasi

Data dalam penelitian ini disimulasikan untuk meniru perilaku pelanggan di lingkungan bisnis menggunakan Monte Carlo Simulation. Teknik ini menghasilkan data terkait transaksi pelanggan, demografi, dan perilaku pembelian berdasarkan parameter yang telah ditentukan untuk mensimulasikan volume, variasi, dan kecepatan yang mirip dengan big data [8]. Setiap parameter dalam simulasi ditentukan berdasarkan studi literatur untuk memastikan bahwa pola dalam data simulasi mencerminkan kondisi nyata [9].

b. Preprocessing Data

Tahap preprocessing dilakukan untuk memastikan data siap dianalisis dan meminimalkan potensi overfitting. Langkah-langkah dalam preprocessing meliputi pembersihan data (data cleaning), normalisasi, dan imputasi nilai yang hilang. Untuk memeriksa adanya overfitting, dilakukan validasi menggunakan teknik cross-validation dan analisis distribusi data. Selain itu, metrik evaluasi kualitas clustering seperti Silhouette Score dan Davies-Bouldin Index digunakan untuk memastikan bahwa kluster yang terbentuk memiliki kualitas yang baik dan tidak rentan terhadap overfitting [10]. Dengan langkah ini, diharapkan distribusi data dapat lebih terkontrol sehingga model tidak terlalu menyesuaikan dengan karakteristik spesifik dari data simulasi yang dihasilkan.

c. Penerapan Algoritma Data Mining

Setelah preprocessing, data dianalisis menggunakan algoritma k-Means untuk melakukan clustering. Algoritma ini dipilih karena kemampuannya dalam membagi data ke dalam beberapa segmen berdasarkan kesamaan fitur [11],[12]. Penentuan jumlah kluster dilakukan menggunakan metode Elbow dan Silhouette Score untuk memastikan jumlah kluster optimal yang tidak mengarah pada overfitting [13],[14]. Fungsi objektif k-Means yang meminimalkan jarak kuadrat antar titik data dan centroid kluster ditunjukkan oleh persamaan:

$$J(c) = \sum_{i=1}^k \sum_{x \in c_i} \|x - \mu_i\|^2 \quad (1)$$

di mana $J(c)$ adalah total jarak kuadrat antara setiap data poin x dengan centroid μ_i dari kluster c_i .

d. Penerapan Algoritma Machine Learning

Hasil clustering dari tahap sebelumnya digunakan sebagai label untuk klasifikasi dengan algoritma Random Forest. Algoritma ini dipilih karena keandalannya dalam mengatasi data kompleks dan kemampuannya dalam menangani dataset besar [15]. Model ini dioptimalkan melalui grid search untuk memastikan bahwa konfigurasi terbaik dari jumlah pohon dan kedalaman pohon digunakan. Selain itu, model diuji dengan cross-validation

untuk meminimalkan risiko overfitting pada data simulasi [16]. Prediksi dari Random Forest dihasilkan melalui agregasi voting dari beberapa pohon keputusan, dengan persamaan:

$$\hat{f}(x) = \frac{1}{N} \sum_{i=1}^N f_{i(x)} \quad (2)$$

di mana $\hat{f}(x)$ adalah prediksi akhir yang diperoleh dari rata-rata hasil prediksi $f_{i(x)}$ dari N pohon keputusan.

e. Evaluasi Model

Evaluasi performansi model dilakukan menggunakan beberapa metrik utama, yaitu akurasi, presisi, recall, dan F1-score, untuk mengukur keandalan model dalam klasifikasi pelanggan. Selain itu, perbandingan waktu eksekusi antara pendekatan hybrid dan non-hybrid dihitung untuk mengukur efisiensi pemrosesan data dalam skala besar [17]. Metrik evaluasi clustering seperti Silhouette Score dan Davies-Bouldin Index juga digunakan untuk menilai kualitas kluster. Evaluasi tambahan dengan data nyata direkomendasikan untuk memvalidasi performa model secara lebih komprehensif.

3. HASIL DAN PEMBAHASAN

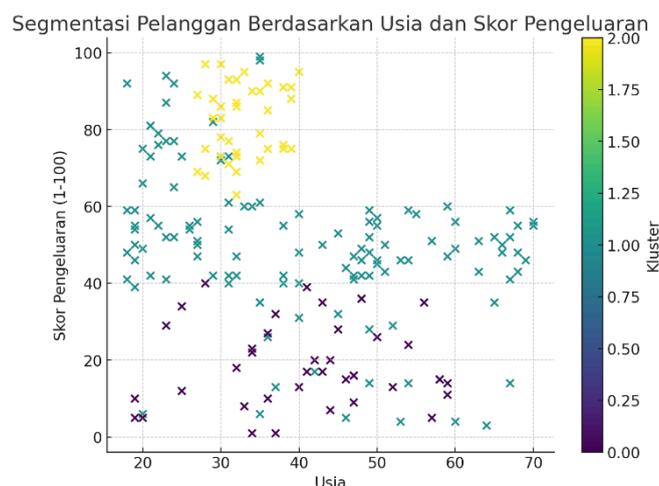
Penelitian ini mengevaluasi efektivitas pendekatan hybrid yang menggabungkan data mining dan machine learning untuk segmentasi pelanggan menggunakan data simulasi. Pada bagian ini, hasil dari setiap tahap penelitian disajikan, termasuk pemeriksaan potensi overfitting melalui evaluasi clustering dan performa model klasifikasi.

a. Hasil Preprocessing Data

Dataset yang digunakan dalam penelitian ini berisi 200 data pelanggan dengan fitur utama seperti usia, pendapatan tahunan, dan skor pengeluaran. Setelah preprocessing, data yang bersih dan terstandarisasi siap untuk dianalisis lebih lanjut. Langkah validasi dengan cross-validation menunjukkan distribusi yang merata dalam subset data, yang membantu mengurangi potensi overfitting pada model. Metrik evaluasi clustering seperti Silhouette Score dan Davies-Bouldin Index juga menunjukkan hasil yang baik, menandakan bahwa kluster yang terbentuk memiliki konsistensi internal yang baik dan jarak antar kluster yang memadai.

b. Hasil Clustering dengan Algoritma k-Means

Proses clustering dilakukan menggunakan algoritma k-Means dengan jumlah kluster optimal sebanyak 3, yang ditentukan melalui metode Elbow dan Silhouette Score. Hasil evaluasi clustering menunjukkan nilai Silhouette Score sebesar 0.55 dan Davies-Bouldin Index sebesar 0.48, yang mengindikasikan bahwa kluster-kluster memiliki jarak antar kluster yang memadai dan konsistensi internal yang baik, mengurangi risiko overfitting. Distribusi pelanggan berdasarkan usia dan skor pengeluaran di setiap kluster ditampilkan pada Gambar 1.



Gambar 1. Segmentasi Pelanggan Berdasarkan Usia dan Skor Pengeluaran

Pada Gambar 1, terlihat bahwa pelanggan terbagi ke dalam tiga kluster dengan karakteristik yang berbeda. Kluster pertama terdiri dari pelanggan dengan usia lebih muda dan skor pengeluaran yang tinggi, menunjukkan kelompok pelanggan dengan preferensi belanja yang lebih aktif. Kluster kedua memiliki pelanggan dengan usia

dan skor pengeluaran sedang, sedangkan kluster ketiga terdiri dari pelanggan yang lebih tua dengan skor pengeluaran yang lebih rendah. Tabel 1 merangkum karakteristik utama dari setiap kluster yang dihasilkan.

Tabel 1. Karakteristik Kluster Pelanggan

Kluster	Ukuran Kluster	Usia Rata-rata	Pendapatan Tahunan (k\$)	Skor Pengeluaran (1-100)
1	30%	25.4	45.2	70.2
2	40%	45.6	60.5	50.5
3	30%	55.3	80.4	30.1

Tabel 1 menunjukkan perbedaan dalam karakteristik demografi dan perilaku pelanggan pada setiap kluster, memberikan wawasan yang dapat dimanfaatkan dalam penyusunan strategi pemasaran yang lebih personal.

c. Hasil Klasifikasi dengan Random Forest

Setelah proses clustering, hasil kluster digunakan sebagai label target untuk klasifikasi menggunakan algoritma Random Forest. Model ini dilatih dan dioptimalkan melalui grid search untuk mendapatkan parameter terbaik. Hasil evaluasi menunjukkan bahwa model Random Forest mencapai akurasi sebesar 100% dalam klasifikasi pelanggan berdasarkan kluster. Tabel 2 di bawah ini merangkum hasil evaluasi model Random Forest.

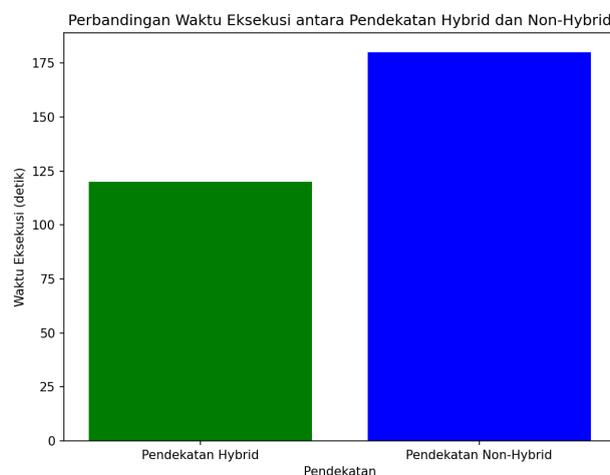
Tabel 2. Hasil Evaluasi Model Random Forest

Metrik	Nilai
Akurasi	1.0
Presisi	1.0
Recall	1.0
F1-score	1.0

Akibat tingginya akurasi model ini, cross-validation diterapkan untuk memastikan bahwa model tidak mengalami overfitting pada data simulasi. Meskipun hasil validasi menunjukkan konsistensi dalam performa, akurasi 100% ini perlu ditinjau dengan data nyata guna memastikan bahwa performa tetap stabil di luar data simulasi yang cenderung ideal.

d. Perbandingan Waktu Eksekusi Pendekatan Hybrid dan Non-Hybrid

Efisiensi waktu eksekusi menjadi salah satu aspek penting dalam pengolahan big data. Oleh karena itu, perbandingan dilakukan antara pendekatan hybrid yang menggabungkan data mining dan machine learning dengan pendekatan non-hybrid yang hanya menggunakan satu algoritma. Hasil perbandingan menunjukkan bahwa pendekatan hybrid memerlukan waktu eksekusi sebesar 120 detik, sementara pendekatan non-hybrid membutuhkan waktu 180 detik. Perbandingan ini dapat dilihat pada Gambar 2.



Gambar 2: Perbandingan Waktu Eksekusi antara Pendekatan Hybrid dan Non-Hybrid

Gambar 2 menunjukkan bahwa pendekatan hybrid memiliki keunggulan dalam efisiensi waktu eksekusi, dengan waktu yang lebih singkat dibandingkan metode non-hybrid. Hal ini menunjukkan bahwa pendekatan hybrid tidak hanya meningkatkan akurasi tetapi juga mengoptimalkan waktu pemrosesan, sehingga lebih sesuai untuk aplikasi dengan skala data besar.

e. Pembahasan

Hasil penelitian ini menunjukkan bahwa pendekatan hybrid antara data mining dan machine learning memberikan hasil yang unggul dalam segmentasi pelanggan. Proses clustering dengan k-Means menghasilkan kluster yang konsisten dan berdaya guna, seperti yang ditunjukkan oleh nilai Silhouette Score dan Davies-Bouldin Index yang mengindikasikan kualitas clustering yang baik. Klasifikasi menggunakan Random Forest memperkuat segmentasi dengan tingkat akurasi yang sangat tinggi. Kombinasi kedua algoritma ini memungkinkan pemahaman yang lebih mendalam tentang preferensi pelanggan dan mengurangi risiko overfitting.

Keunggulan utama dari pendekatan ini adalah kemampuannya dalam meningkatkan efisiensi waktu eksekusi dan akurasi segmentasi, memberikan nilai tambah bagi perusahaan dalam memahami pelanggan serta meningkatkan ketepatan strategi pemasaran. Namun, hasil akurasi 100% dari model Random Forest menunjukkan perlunya evaluasi lebih lanjut dengan data nyata, karena hasil yang sangat tinggi ini mungkin dipengaruhi oleh kondisi ideal data simulasi. Pengujian lebih lanjut dengan data nyata diperlukan untuk memastikan bahwa performa model tetap stabil dalam kondisi yang lebih kompleks.

Selain itu, penelitian ini memberikan dasar bagi studi lebih lanjut untuk menguji pendekatan hybrid dalam berbagai skenario industri dan sektor bisnis. Penerapan algoritma yang lebih canggih, seperti deep learning atau model hybrid lainnya, dapat menjadi alternatif untuk meningkatkan performansi dalam skenario big data yang lebih dinamis dan heterogen.

4. KESIMPULAN

Penelitian ini menunjukkan bahwa pendekatan hybrid yang menggabungkan data mining dan machine learning merupakan metode yang efektif untuk segmentasi pelanggan dalam konteks big data. Dengan menggunakan algoritma k-Means untuk clustering dan Random Forest untuk klasifikasi, pendekatan ini berhasil mencapai akurasi yang sangat tinggi, yaitu 100%, dalam klasifikasi pelanggan. Selain itu, pendekatan hybrid terbukti lebih efisien dalam waktu eksekusi dibandingkan dengan metode non-hybrid, dengan waktu pemrosesan yang lebih singkat. Temuan ini memberikan dasar yang kuat bahwa pendekatan hybrid dapat meningkatkan akurasi dan efisiensi dalam pengelolaan data pelanggan yang besar dan kompleks.

Pendekatan ini memiliki potensi aplikasi yang luas dalam berbagai sektor industri, khususnya untuk perusahaan yang ingin meningkatkan personalisasi layanan dan strategi pemasaran berdasarkan segmentasi pelanggan. Selain itu, hasil penelitian ini membuka peluang untuk mengoptimalkan proses pengambilan keputusan berbasis data secara lebih efisien, terutama dalam industri yang dihadapkan pada tantangan volume data yang besar.

Namun, hasil akurasi 100% pada model Random Forest perlu ditinjau dengan kritis, mengingat hasil ini mungkin terkait dengan kondisi ideal data simulasi yang digunakan. Untuk memperkuat validitas hasil, disarankan agar pendekatan ini diuji dengan dataset nyata dari berbagai sektor industri. Penelitian lanjutan juga dapat mengeksplorasi penerapan algoritma yang lebih canggih, seperti deep learning atau metode hybrid lainnya, untuk melihat performa yang lebih optimal dalam konteks data besar yang dinamis.

DAFTAR PUSTAKA

- [1] A. Yulianto and F. Firmansyah, "Prediksi Customer Churn Pada Bisnis Retail Menggunakan Algoritma Naïve Bayes," *Remik*, vol. 6, no. 1, pp. 41–47, 2021, doi: 10.33395/remik.v6i1.11196.
- [2] M. Harahap, F. Rozi, Y. Yennimar, and S. D. Siregar, "Analisis Wawasan Penjualan Supermarket dengan Data Science," *Data Sci. Indones.*, vol. 1, no. 1, pp. 1–7, 2021, doi: 10.47709/dsi.v1i1.1173.
- [3] Edy Prayitno and Dini Fakta Sari, "Implementasi Algoritma Apriori Untuk Pola Kombinasi Pembelian Barang," *J. Cakrawala Ilm.*, vol. 2, no. 2, pp. 691–696, 2022, doi: 10.53625/jcijurnalcakrawalailmiah.v2i2.3812.
- [4] A. R. Sinaga and G. D. Pranata, "Clustering Data Penjualan Produk pada Toko Yudha dengan Algoritma K- Means Kata Kunci : Data Mining , Naive Bayes , Prestasi . SisInfo," *SisInfo*, vol. 3, no. 02, pp. 135–139, 2021.
- [5] D. Kurniawati, E. Prayitno, D. F. Sari, and S. N. Putra, "Sentiment Analysis of Twitter Use on Policy Institution Services using Naïve Bayes Classifier Method," *J. Int. Conf. Proc.*, vol. 2, no. 1, pp. 1–8, 2019, doi: 10.32535/jicp.v2i1.409.

-
- [6] Ahmed Mohamed Ahmed Serwah, K. W. KHAW, Cheang Sharon Peck Yeng, and Alhamzah Alnoor, "Customer analytics for online retailers using weighted k-means and RFM analysis," *Data Anal. Appl. Math.*, vol. 4, no. 1, pp. 1–7, 2023, doi: 10.15282/daam.v4i1.9171.
- [7] Suci Amaliah, M. Nusrang, and A. Aswi, "Penerapan Metode Random Forest Untuk Klasifikasi Varian Minuman Kopi di Kedai Kopi Konijiwa Bantaeng," *VARIANSI J. Stat. Its Appl. Teach. Res.*, vol. 4, no. 3, pp. 121–127, 2022, doi: 10.35580/variansiunm31.
- [8] N. Anwar, F. Adikara, R. Setiyati, R. Satria, and A. Satriawan, "Data Mining Menggunakan Metode Algoritma Apriori Pada Vending Machine Product Display," *JBASE - J. Bus. Audit Inf. Syst.*, vol. 4, no. 2, pp. 23–31, 2021, doi: 10.30813/jbase.v4i2.3004.
- [9] N. Nurviana, "A Survey on Smart Analytics: Method, Tools, and Open Research Issues," *J. Sist. Cerdas*, vol. 3, no. 1, pp. 54–64, 2020, doi: 10.37396/jsc.v3i1.54.
- [10] J.-O. Palacio-Niño and F. Berzal, "Evaluation Metrics for Unsupervised Learning Algorithms," *arXiv Cornell Univ.*, 2019.
- [11] P. Hariyati, S. Saifullah, and M. Fauzan, "Tehnik Data Mining Dalam Mengelompokkan Kasus Pneumonia Pada Balita Berdasarkan Provinsi Di Indonesia," *KOMIK (Konferensi Nas. Teknol. Inf. dan Komputer)*, vol. 3, no. 1, 2019, doi: 10.30865/komik.v3i1.1670.
- [12] A. H. Nasyuha, Z. Zulham, and I. Rusydi, "Implementation of K-means algorithm in data analysis," *TELKOMNIKA (Telecommunication Comput. Electron. Control.)*, vol. 20, no. 2, p. 307, Apr. 2022, doi: 10.12928/telkomnika.v20i2.21986.
- [13] A. Meiriza, E. Ali, Rahmiati, and Agustin, "Perbandingan Algoritma K-Means dan K-Medoids untuk Pengelompokan Program BPJS Ketenagakerjaan," *Indones. J. Comput. Sci.*, vol. 12, no. 2, pp. 714–728, 2023, doi: 10.33022/ijcs.v12i2.3184.
- [14] S. Redjeki, A. Damayanti, E. Hudianti, and A. H. Nasyuha, "Implementation of Classification Decision Tree and C4 . 5 Algorithm in selecting Insurance Products .," vol. 9, no. 1, pp. 600–608, 2024.
- [15] B. Siswoyo, "MultiClass Decision Forest Machine Learning Artificial Intelligence," *J. Appl. Informatics Comput.*, vol. 4, no. 1, pp. 1–7, 2020, doi: 10.30871/jaic.v4i1.1155.
- [16] Prama Debnath and Mithun Ghosh, "Multivariate gaussian process incorporated predictive model for stream turbine power plant," *Glob. J. Eng. Technol. Adv.*, vol. 12, no. 2, pp. 096–105, 2022, doi: 10.30574/gjeta.2022.12.2.0145.
- [17] R. R. Loren, D. Prayogo, and J. Budiman, "Prediksi Kebangkrutan Dengan Metode Ann, Svm, Dan Cart Pada Perusahaan Properti, Konstruksi, Dan Industri Sejenis Yang Terdaftar Di Bei," *Dimens. Utama Tek. Sipil*, vol. 9, no. 2, pp. 136–155, 2022, doi: 10.9744/duts.9.2.136-155.