

KLASIFIKASI TINGKAT PRESTASI SISWA BERDASARKAN DATA NILAI AKADEMIK DENGAN ALGORITMA C4.5

**Relita Buaton^{1*}, Husnul Khair²,
Imeldawaty Gultom³, Siti Nur Azizah⁴,
Novita Anggraini⁵**

Program Studi Teknik Informatika¹,
Program Studi Teknik Informatika²,
Program Studi Komputerisasi Akuntansi³,
Program Studi Sistem Informasi⁴,
Program Studi Sistem Informasi⁵
STMIK Kaputama¹, STMIK Kaputama², STMIK Kaputama³,
STMIK Kaputama⁴, STMIK Kaputama⁵

*Correspondent Author: relitabuaton@kaputama.ac.id,

Author Email: Husnul.khair@gmail.com²,
imeldagultom81@gmail.com³, azizahasm31@gmail.com⁴,
novi200327@gmail.com⁵

Received: March 25, 2026. **Revised:** May 4, 2026. **Accepted:** May 12 2026. **Issue Period:** Vol.10 No.2 (2026), Pp. 628-635

Abstrak: Data akademik sangat dibutuhkan sebagai bahan evaluasi, pengolahan data akademik siswa secara sistematis diperlukan untuk mendukung identifikasi tingkat prestasi secara lebih objektif. Dalam penelitian ini masalah yang ditemukan adalah belum optimalnya pemanfaatan data nilai akademik sebagai dasar klasifikasi prestasi siswa dengan menerapkan algoritma C4.5 dengan data set sebanyak 507 data dan 11 atribut. Metode yang digunakan meliputi preprocessing data, perhitungan entropy, information gain, split information, gain ratio, pembentukan pohon keputusan, serta evaluasi model dengan pembagian data 80:20 secara stratified dan validasi silang 5-fold. Hasil penelitian menunjukkan bahwa atribut IPA menjadi pemisah awal terbaik dengan nilai entropy 1,5108, information gain 0,0742, dan gain ratio 0,0524. Model yang dihasilkan memperoleh akurasi pengujian sebesar 56,86% dan rerata akurasi validasi silang sebesar 55,42%.

Kata kunci: Data Akademik; Algoritma C4.5; Pohon Keputusan; Klasifikasi Prestasi Siswa.

Abstract: Academic data are essential as evaluation materials, and the systematic processing of student academic data is needed to support a more objective identification of achievement levels. The problem addressed in this study is the suboptimal use of academic score data as a basis for classifying student achievement by applying the C4.5 algorithm to a dataset of 507 records and 11 attributes. The methods employed include data preprocessing, the calculation of entropy,



DOI: 10.52362/jisamar.v10i2.2410

Ciptaan disebarluaskan di bawah [Lisensi Creative Commons Atribusi 4.0 Internasional](https://creativecommons.org/licenses/by/4.0/).

information gain, split information, and gain ratio, decision tree construction, and model evaluation using an 80:20 stratified data split and 5-fold cross-validation. The results show that Science is the best initial splitting attribute, with an entropy value of 1.5108, information gain of 0.0742, and gain ratio of 0.0524. The resulting model achieved a testing accuracy of 56.86% and an average cross-validation accuracy of 55.42%.

Keywords: *Academic Data; C4.5 Algorithm; Decision Tree; Student Achievement Classification.*

I. PENDAHULUAN

Pemanfaatan data akademik dalam dunia pendidikan tidak lagi cukup hanya untuk keperluan pencatatan dan pelaporan nilai, tetapi juga perlu diarahkan sebagai dasar evaluasi dan pengambilan keputusan akademik. Data nilai siswa yang tersimpan di sekolah sesungguhnya memiliki potensi untuk digunakan dalam mengidentifikasi tingkat prestasi, memetakan capaian belajar, dan membantu sekolah dalam menentukan langkah pembinaan secara lebih dini. Namun, pada praktiknya, pemanfaatan data akademik di banyak sekolah masih belum optimal karena sebagian besar data hanya digunakan sebagai arsip administratif dan belum diolah menjadi informasi yang mendukung klasifikasi prestasi siswa secara sistematis.

Berbagai penelitian menunjukkan bahwa data akademik dapat dimanfaatkan untuk memprediksi dan mengklasifikasikan performa belajar siswa melalui pendekatan educational data mining [1][2]. Kajian lain juga menegaskan bahwa nilai antarmata pelajaran dapat memberikan pola yang relevan untuk membedakan capaian akademik siswa [3]. Selain itu, penelitian terdahulu menunjukkan bahwa penerapan algoritma klasifikasi pada data pendidikan perlu disertai evaluasi yang terukur melalui accuracy, precision, recall, dan F1-score agar kualitas model dapat dinilai secara objektif [4][12]. Temuan-temuan tersebut memperlihatkan bahwa pengolahan data pendidikan tidak hanya berorientasi pada hasil prediksi, tetapi juga pada keterjelasan model agar dapat dimanfaatkan dalam konteks pembelajaran.

Meskipun demikian, sebagian besar penelitian sebelumnya masih banyak menggunakan data dari perguruan tinggi, learning management system, atau data yang menggabungkan variabel demografis dan perilaku belajar, sedangkan pemanfaatan data nilai rapor di tingkat sekolah menengah pertama masih relatif terbatas [5][6]. Padahal, data nilai rapor merupakan sumber data yang paling tersedia dan paling dekat dengan kebutuhan sekolah untuk melakukan pemetaan prestasi siswa. Kondisi ini menunjukkan adanya kebutuhan terhadap model klasifikasi yang mampu mengolah data akademik internal sekolah secara sederhana, terukur, dan mudah dipahami.

Salah satu metode yang dapat digunakan untuk memenuhi kebutuhan tersebut adalah algoritma C4.5. Algoritma ini dipilih karena mampu membentuk pohon keputusan yang bersifat interpretatif, sehingga hasil klasifikasi tidak hanya menunjukkan kategori prestasi siswa, tetapi juga memberikan aturan keputusan yang dapat dijelaskan kembali kepada guru dan pihak sekolah [7][8]. Dengan dukungan implementasi berbasis Python, proses pengolahan data, pembentukan model, dan evaluasi performa dapat dilakukan secara lebih sistematis dan efisien. Berdasarkan kondisi tersebut, penelitian ini menerapkan algoritma C4.5 berbasis Python untuk mengklasifikasikan tingkat prestasi siswa SMP Swasta Teladan berdasarkan data nilai akademik. Penelitian ini menggunakan atribut nilai mata pelajaran sebagai variabel prediktor dan membagi tingkat prestasi ke dalam kategori Rendah, Sedang, dan Tinggi. Dengan demikian, penelitian ini diharapkan dapat menghasilkan model klasifikasi yang terukur, interpretatif, dan relevan sebagai dasar pendukung evaluasi akademik di sekolah.

II. METODE DAN MATERI

Metodologi penelitian ini disusun untuk mendukung proses klasifikasi tingkat prestasi siswa secara sistematis, mulai dari pengumpulan data, persiapan data, penerapan teknik data mining, pembentukan model klasifikasi menggunakan algoritma C4.5, implementasi dengan Python, hingga pengujian dan evaluasi performa model. Dalam penelitian ini, pendekatan yang digunakan adalah data mining, khususnya educational data mining, karena penelitian berfokus pada penggalian pola dari data akademik siswa untuk menghasilkan informasi yang dapat mendukung pengambilan keputusan di bidang pendidikan. Tahapan inti dalam educational data mining



umumnya meliputi persiapan data, pemodelan, dan evaluasi [8], sedangkan klasifikasi merupakan salah satu pendekatan utama untuk mengubah data pendidikan menjadi dasar pengambilan keputusan yang lebih terarah [11].

2.1. Penelitian Terdahulu

Penelitian mengenai prediksi performa akademik siswa telah berkembang pesat dalam ranah educational data mining. Data akademik siswa dapat dimanfaatkan untuk memprediksi performa belajar menggunakan algoritma machine learning, sehingga data pendidikan tidak hanya berfungsi sebagai arsip administratif, tetapi juga sebagai sumber informasi prediktif yang bernilai [6]. Temuan ini menunjukkan bahwa data akademik yang sederhana pun dapat diolah menjadi model yang mampu mendukung evaluasi dan pengambilan keputusan akademik.

Perkembangan penelitian selanjutnya memperlihatkan bahwa analisis prestasi siswa tidak hanya berfokus pada ketepatan prediksi, tetapi juga pada kemampuan menemukan pola hubungan antarnilai mata pelajaran. Data nilai antarmata pelajaran dapat mengungkap pola capaian akademik yang relevan untuk mengelompokkan performa siswa [5]. Dengan demikian, nilai rapor dapat diposisikan bukan hanya sebagai hasil evaluasi belajar, melainkan juga sebagai data yang dapat ditambang untuk memperoleh pengetahuan baru.

Di sisi lain, model klasifikasi dalam prediksi performa siswa perlu diuji menggunakan ukuran evaluasi yang jelas, seperti accuracy, precision, recall, dan F1-score, agar kualitas model dapat diukur secara objektif [3]. Hal ini penting karena model yang dihasilkan tidak cukup hanya mampu mengklasifikasikan data, tetapi juga harus dapat dibuktikan tingkat kinerjanya secara kuantitatif.

Dalam konteks pendidikan, kebutuhan terhadap model yang interpretatif tetap menjadi perhatian utama. Hasil prediksi performa siswa akan lebih bermanfaat apabila mampu diterjemahkan menjadi dasar intervensi pembelajaran [1]. Sejalan dengan itu, keputusan berbasis data dalam pendidikan akan menjadi lebih kuat apabila model dibangun dari data yang bersih, tervalidasi, dan dievaluasi secara terukur [2].

Meskipun berbagai penelitian terdahulu telah membahas prediksi performa siswa, sebagian besar penelitian menggunakan data perguruan tinggi, data LMS, atau kombinasi data akademik dengan variabel perilaku dan demografis. Sementara itu, pada tingkat sekolah menengah pertama, kebutuhan yang lebih praktis justru terletak pada pemanfaatan data nilai rapor yang telah tersedia dalam administrasi sekolah. Oleh karena itu, penelitian ini mengambil fokus yang lebih spesifik, yaitu menerapkan teknik data mining pada data nilai akademik siswa SMP Swasta Teladan dengan menggunakan algoritma C4.5 berbasis Python untuk mengklasifikasikan tingkat prestasi siswa ke dalam kategori Rendah, Sedang, dan Tinggi. Dengan demikian, penelitian ini tidak hanya berorientasi pada hasil klasifikasi, tetapi juga pada keterbacaan aturan keputusan yang dihasilkan.

2.2. Metode yang Digunakan

Metode yang digunakan dalam penelitian ini adalah metode klasifikasi dalam data mining. Klasifikasi dipilih karena tujuan penelitian adalah mengelompokkan siswa ke dalam kelas target tertentu, yaitu tingkat prestasi Rendah, Sedang, dan Tinggi. Dalam penelitian ini, algoritma yang digunakan adalah C4.5, yaitu algoritma pembentukan pohon keputusan yang memilih atribut terbaik berdasarkan nilai gain ratio.

Pemilihan algoritma C4.5 didasarkan pada kemampuannya dalam menghasilkan model yang mudah dipahami dan dapat diterjemahkan ke dalam bentuk aturan keputusan [13]. Hal ini penting dalam konteks pendidikan, karena hasil klasifikasi harus dapat dijelaskan kembali kepada guru atau pihak sekolah. Relevansi penggunaan klasifikasi pada data pendidikan juga didukung oleh penelitian yang menunjukkan bahwa model klasifikasi mampu mengidentifikasi pola performa akademik siswa dan mendukung intervensi pembelajaran yang lebih terarah [5][10].

2.3. Data Penelitian

Penelitian ini menggunakan data nilai akademik siswa SMP Swasta Teladan. Data awal yang dihimpun berjumlah 529 record, kemudian setelah dilakukan proses seleksi dan pembersihan data, diperoleh 507 record valid yang digunakan dalam pemodelan. Setiap record merepresentasikan satu siswa dengan 11 atribut nilai mata pelajaran, yaitu Pendidikan Agama, Pendidikan Pancasila, Bahasa Indonesia, Matematika, IPA, IPS, Bahasa Inggris, PJOK, TIK, SBK, dan Prakarya.

Variabel target dalam penelitian ini adalah Tingkat Prestasi, yang terdiri atas tiga kelas, yaitu Rendah, Sedang, dan Tinggi. Seluruh atribut nilai mata pelajaran digunakan sebagai variabel prediktor. Penggunaan atribut



akademik sebagai dasar klasifikasi sejalan dengan penelitian yang menunjukkan bahwa nilai akademik merupakan indikator penting dalam membangun model prediksi performa siswa [3][10].

<i>Komponen</i>	<i>Nilai</i>
Jumlah data mentah	529
Data akhir setelah cleaning	507
Data terhapus	22
Jumlah atribut prediktor	11
Jumlah kelas target	3

<i>Kelas</i>	<i>Jumlah</i>
Rendah	243
Sedang	112
Tinggi	152
Kelas	243

2.4. Preprocessing Data

Sebelum dilakukan pemodelan, data terlebih dahulu melalui tahap preprocessing agar siap digunakan dalam proses klasifikasi. Tahap ini penting karena kualitas data awal sangat memengaruhi kualitas model yang dihasilkan. Proses preprocessing dalam penelitian ini meliputi seleksi atribut yang relevan, penggabungan sheet yang memiliki struktur serupa, penyeragaman nama kolom, konversi nilai ke format numerik, penghapusan label target yang tidak valid, serta penghapusan baris yang masih memiliki nilai kosong pada atribut utama.

Selain itu, pada tahap analisis manual algoritma C4.5, beberapa atribut numerik juga dikelompokkan ke dalam kategori interval agar proses perhitungan entropy, information gain, split information, dan gain ratio dapat dijelaskan secara sistematis. Langkah preprocessing ini diperlukan karena data mentah rentan mengandung missing value, noise, dan inkonsistensi, sehingga perlu dibersihkan dan ditransformasikan sebelum digunakan dalam proses klasifikasi [9].

2.5. Algoritma C4.5

Algoritma C4.5 digunakan untuk membentuk pohon keputusan berdasarkan atribut dengan nilai gain ratio tertinggi. Proses klasifikasi diawali dengan menghitung entropy total data, kemudian dilanjutkan dengan menghitung entropy setiap partisi atribut, information gain, split information, dan gain ratio. Atribut dengan gain ratio tertinggi dipilih sebagai node akar, lalu proses yang sama diulang pada setiap cabang hingga seluruh data berada pada kelas yang homogen atau tidak ada atribut yang lebih baik untuk dipilih.

Melalui mekanisme tersebut, algoritma C4.5 mampu menghasilkan aturan keputusan yang mudah diinterpretasikan untuk mengelompokkan siswa ke dalam kategori prestasi Rendah, Sedang, atau Tinggi. Penggunaan pendekatan klasifikasi seperti ini relevan dengan penelitian yang menunjukkan bahwa model klasifikasi dapat mengidentifikasi pola performa akademik siswa secara efektif [5][10].

Entropy dihitung menggunakan Persamaan (1).

$$Entropy(S) = - \sum_{i=1}^n p_i \log_2 p_i$$

Information Gain dihitung menggunakan Persamaan (2).

$$Gain(S, A) = Entropy(S) - \sum_{i=1}^n \frac{|S_i|}{|S|} Entropy(S_i)$$

Split Information dihitung menggunakan Persamaan (3).

$$SplitInfo(S, A) = - \sum_{i=1}^n \frac{|S_i|}{|S|} \log_2 \left(\frac{|S_i|}{|S|} \right)$$



Gain Ratio dihitung menggunakan Persamaan (4).

$$GainRatio(S, A) = \frac{Gain(S, A)}{SplitInfo(S, A)}$$

2.6. Pengujian dan Evaluasi Model

Metode pengujian model dalam penelitian ini dilakukan menggunakan train-test split sebesar 80:20 secara stratified, agar distribusi kelas target tetap terjaga pada data latih dan data uji. Selain itu, untuk memperoleh gambaran performa model yang lebih stabil, penelitian ini juga menggunakan 5-fold cross-validation. Melalui validasi silang ini, model diuji pada beberapa pembagian data yang berbeda sehingga hasil evaluasi menjadi lebih representatif.

Evaluasi model dilakukan menggunakan confusion matrix, accuracy, precision, recall, dan F1-score. Accuracy digunakan untuk melihat proporsi prediksi benar terhadap seluruh data uji, sedangkan precision, recall, dan F1-score digunakan untuk menilai kemampuan model dalam mengklasifikasikan masing-masing kelas prestasi. Penggunaan kombinasi metrik ini penting agar evaluasi model tidak hanya berfokus pada akurasi total, tetapi juga pada kualitas prediksi setiap kelas, terutama karena kelas Sedang cenderung lebih sulit dipisahkan dibandingkan kelas lainnya. Pemilihan metrik evaluasi tersebut sejalan dengan penelitian yang juga menggunakan confusion matrix, accuracy, precision, recall, dan F1-score untuk menilai kualitas model prediksi performa akademik [3][7].

2.7. Tahapan Penelitian

Tahapan penelitian terdiri atas: (1) pengumpulan data nilai akademik siswa, (2) seleksi atribut yang relevan, (3) pembersihan dan transformasi data, (4) pembagian data menjadi data latih dan data uji, (5) perhitungan entropy, gain, split information, dan gain ratio, (6) pembentukan pohon keputusan menggunakan algoritma C4.5, (7) implementasi model menggunakan Python, (8) evaluasi model dengan confusion matrix, accuracy, precision, recall, dan F1-score, serta (9) penarikan kesimpulan berdasarkan hasil pengujian model. Urutan ini sesuai dengan alur umum educational data mining yang dirangkum oleh [8][11], yaitu dimulai dari persiapan data, pemodelan, hingga evaluasi untuk menghasilkan keputusan berbasis data.

III. PEMBAHASA DAN HASIL

3.1. Deskripsi Data

Data penelitian ini terdiri atas 507 data valid dengan 11 atribut nilai mata pelajaran dan satu variabel target, yaitu Tingkat Prestasi. Distribusi data menunjukkan 243 siswa berada pada kategori Rendah, 112 siswa pada kategori Sedang, dan 152 siswa pada kategori Tinggi. Sebaran ini menunjukkan bahwa kelas Rendah lebih dominan, sedangkan kelas Sedang memiliki jumlah paling sedikit. Selain itu, rata-rata nilai tertinggi terdapat pada SBK sebesar 79,60, sedangkan rata-rata nilai terendah terdapat pada Bahasa Inggris sebesar 76,39, sehingga variasi nilai antarmata pelajaran menjadi dasar penting dalam pembentukan model klasifikasi.

Table III. Lima Atribut Dengan Gain Ratio Tertinggi

<i>Atribut</i>	<i>Gain Ratio</i>
IPA	0.0524
SBK	0.0477
IPS	0.0413
Bahasa Inggris	0.0399
Bahasa Indonesia	0.0397

3.2. Implementasi Python

Implementasi model dilakukan menggunakan Python karena mampu mendukung proses pengolahan data, pembentukan model, dan evaluasi hasil secara terintegrasi. Pada tahap ini, data nilai akademik digunakan sebagai variabel prediktor, sedangkan Tingkat Prestasi dijadikan variabel target. Selanjutnya, data dibagi menjadi data



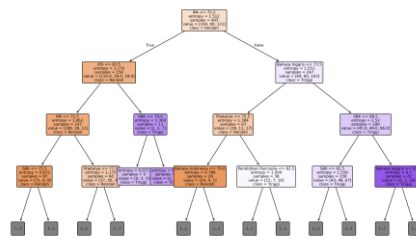
DOI: 10.52362/jisamar.v10i2.2410

Ciptaan disebarluaskan di bawah [Lisensi Creative Commons Atribusi 4.0 Internasional](https://creativecommons.org/licenses/by/4.0/).

latih dan data uji dengan proporsi 80:20 secara stratified agar distribusi kelas tetap terjaga. Model kemudian dibentuk menggunakan *DecisionTreeClassifier* dengan *criterion=entropy*, dilatih pada data latih, lalu digunakan untuk memprediksi data uji. Hasil prediksi tersebut selanjutnya dievaluasi menggunakan *accuracy*, *confusion matrix*, *precision*, *recall*, dan *F1-score*, sehingga kinerja model dapat diukur secara menyeluruh [3][7].

3.3. Hasil Pemodelan C4.5 Berbasis Python

Model menghasilkan akurasi data uji sebesar 53.92% dengan rerata akurasi validasi silang 58.58% dan simpangan baku 2.88%. Nilai ini menunjukkan bahwa model telah mampu menangkap pola dasar prestasi siswa, namun pemisahan antar kelas yang berdekatan, terutama kelas Sedang, masih menimbulkan salah klasifikasi. Oleh karena itu, pembacaan hasil tidak cukup hanya dari akurasi, tetapi juga perlu melihat *precision*, *recall*, dan *F1-score* per kelas agar kualitas klasifikasi dapat dipahami secara lebih utuh.

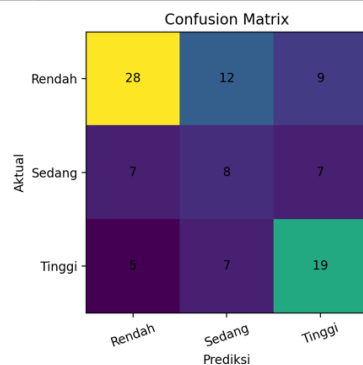


Gambar 1. Pohon keputusan hasil pemodelan Python

Gambar 1 memperlihatkan bahwa percabangan awal dimulai dari atribut IPA. Setelah itu model banyak menggunakan Bahasa Inggris, SBK, dan IPS sebagai pemisah lanjutan. Pola ini menunjukkan bahwa kombinasi penguasaan mata pelajaran sains, bahasa, dan seni-sosial menjadi indikator yang cukup kuat dalam membedakan tingkat prestasi siswa pada data SMP Swasta Teladan.

Table 1 Kinerja Model Per Kelas

Kelas	Precision	Recall	F1-Score	Support
Rendah	0.7000	0.5714	0.6292	49
Sedang	0.2963	0.3636	0.3265	22
Tinggi	0.5429	0.6129	0.5758	31
Macro Avg	0.5131	0.5160	0.5105	102



Gambar 2. Confusion matrix model klasifikasi

3.4. Evaluasi Kinerja Model

Berdasarkan confusion matrix, jumlah prediksi benar berada pada diagonal utama, yaitu 28 data untuk kelas Rendah, 8 data untuk kelas Sedang, dan 19 data untuk kelas Tinggi. Dengan total data uji sebanyak 102 siswa, *accuracy* model dihitung seperti pada Persamaan (7).



$$\begin{aligned} \text{Accuracy} &= \frac{TP_{\text{Rendah}} + TP_{\text{Sedang}} + TP_{\text{Tinggi}}}{N} \\ &= \frac{28 + 8 + 19}{102} = 0.5392 = 53.92\% \end{aligned}$$

Untuk menunjukkan perhitungan precision, recall, dan F1-score secara manual, kelas Sedang digunakan sebagai contoh karena kelas ini memiliki performa paling rendah dan menjadi sumber salah klasifikasi terbesar. Perhitungan rinci ditunjukkan pada Persamaan (8).

$$\text{Precision}_{\text{Sedang}} = \frac{TP}{TP + FP} = \frac{8}{27} = 0.2963$$

$$\text{Recall}_{\text{Sedang}} = \frac{TP}{TP + FN} = \frac{8}{22} = 0.3636$$

$$F1_{\text{Sedang}} = \frac{2 \times 0.2963 \times 0.3636}{0.2963 + 0.3636} = 0.3265$$

Hasil pada Persamaan (8) konsisten dengan Tabel 5, yaitu precision kelas Sedang sebesar 0.2963, recall 0.3636, dan F1-score 0.3265. Nilai tersebut menunjukkan bahwa model masih kesulitan membedakan kelas Sedang dari dua kelas lain yang memiliki rentang nilai berdekatan. Sebaliknya, kelas Rendah memiliki precision 0.7000 dan recall 0.5714, sedangkan kelas Tinggi memiliki precision 0.5429 dan recall 0.6129. Pola ini menunjukkan bahwa model lebih stabil mengenali dua kelas ekstrem dibandingkan kelas menengah.

3.5. Interpretasi Aturan Keputusan

Interpretasi aturan keputusan menunjukkan bahwa atribut IPA menjadi akar pohon keputusan, sehingga nilai mata pelajaran ini paling menentukan pemisahan awal tingkat prestasi siswa. Setelah itu, model memanfaatkan atribut Bahasa Inggris, SBK, dan IPS sebagai pemisah lanjutan. Hal ini menunjukkan bahwa siswa dengan capaian yang konsisten pada atribut-atribut tersebut cenderung lebih mudah diarahkan ke kategori prestasi tertentu.

Secara umum, cabang keputusan memperlihatkan bahwa siswa dengan kombinasi nilai yang lebih baik pada atribut utama cenderung masuk ke kategori Tinggi, sedangkan kombinasi nilai yang lebih rendah cenderung masuk ke kategori Rendah. Sementara itu, kategori Sedang berada pada area transisi, sehingga lebih sering mengalami tumpang tindih dengan dua kelas lainnya. Temuan ini sejalan dengan hasil evaluasi yang menunjukkan bahwa kelas Sedang memiliki precision, recall, dan F1-score paling rendah.

IV. KESIMPULAN

Berdasarkan hasil pengujian yang telah dilakukan, algoritma C4.5 mampu digunakan untuk mengklasifikasikan tingkat prestasi siswa berdasarkan data nilai akademik dengan hasil yang cukup baik. Atribut IPA menjadi pemisah awal terbaik dalam pembentukan pohon keputusan, kemudian diikuti oleh Bahasa Inggris, SBK, dan IPS sebagai atribut lanjutan yang membantu membedakan kelas Rendah, Sedang, dan Tinggi. Hasil ini menunjukkan bahwa kombinasi nilai sains, bahasa, dan mata pelajaran pendukung lainnya memiliki kontribusi penting dalam proses klasifikasi prestasi siswa.

Dari sisi evaluasi, model memperoleh akurasi data uji sebesar 53.92% dengan rerata akurasi validasi silang sebesar 58.58%. Pada data uji, model lebih baik dalam mengenali kelas Rendah dan Tinggi, sedangkan kelas Sedang masih menjadi kelas yang paling sulit dipisahkan karena memiliki nilai precision 0.2963, recall 0.3636, dan F1-score 0.3265. Dengan demikian, penelitian ini memberikan kontribusi berupa model klasifikasi yang dapat digunakan sebagai dasar evaluasi akademik awal di sekolah. Untuk penelitian berikutnya, hasil ini dapat dikembangkan melalui penambahan atribut nonakademik, penyeimbangan kelas, atau perbandingan dengan algoritma klasifikasi lain agar performa model menjadi lebih baik.

REFERENSI



DOI: 10.52362/jisamar.v10i2.2410

Ciptaan disebarluaskan di bawah [Lisensi Creative Commons Atribusi 4.0 Internasional](https://creativecommons.org/licenses/by/4.0/).

- [1] Angeioplastis, A., Aliprantis, J., Konstantakis, M., & Tsimpiris, A. (2025). Predicting Student Performance and Enhancing Learning Outcomes: A Data-Driven Approach Using Educational Data Mining Techniques. *Computers*, 14(3), 83.
- [2] Gul, M. N., Abbasi, W., Babar, M. Z., Aljohani, A., & Arif, M. (2025). Data Driven Decisions in Education Using a Comprehensive Machine Learning Framework for Student Performance Prediction. *Discover Computing*, 28, 153.
- [3] Khairy, D., Alharbi, N., Amasha, M. A., Areed, M. F., Alkhalaf, S., & Abougalala, R. A. (2024). Prediction of Student Exam Performance Using Data Mining Classification Algorithms. *Education and Information Technologies*.
- [4] Putro, A. W. G., & Setiadi, T. (2023). Penerapan Klasifikasi Decision Tree (C4.5) untuk Memprediksi Kelulusan Siswa Sekolah Dasar di Kecamatan Juai. *Jurnal Format*, 12(2), 151-156.
- [5] Sarker, S., Paul, M. K., Thasin, S. T. H., & Hasan, M. A. M. (2024). Analyzing Students' Academic Performance Using Educational Data Mining. *Computers and Education: Artificial Intelligence*, 7, 100263.
- [6] Yağcı, M. (2022). Educational Data Mining: Prediction of Students' Academic Performance Using Machine Learning Algorithms. *Smart Learning Environments*, 9, 11.
- [7] Dawar, I., Negi, S., Lamba, S., & Kumar, A. (2024). Enhancing Student Academic Performance Forecasting: A Comparative Analysis of Machine Learning Algorithms. *SN Computer Science*, 5, 758.
- [8] Kalita, E., Oyelere, S. S., Gaftandzhieva, S., Rajesh, K. N. V. P. S., Jagatheesaperumal, S. K., Mohamed, A., Elbarawy, Y. M., Desuky, A. S., Hussain, S., Cifci, M. A., Theodorou, P., Hilcenko, S., Hazarika, J., & Ali, T. (2025). Educational data mining: a 10-year review. *Discover Computing*, 28, 81.
- [9] Maharana, K., Mondal, S., & Nemade, B. (2022). A review: Data pre-processing and data augmentation techniques. *Global Transitions Proceedings*, 3(1), 91-99.
- [10] Nayak, P., Vaheed, S., Gupta, S., & Mohan, N. (2023). Predicting students' academic performance by mining the educational data through machine learning-based classification model. *Education and Information Technologies*, 28, 14611-14637.
- [11] Papadogiannis, I., Wallace, M., & Karountzou, G. (2024). Educational Data Mining: A Foundational Overview. *Encyclopedia*, 4(4), 1644-1664.
- [12] Ginting, S., Buaton, R., & Khadapi, M. (2025). Implementasi Algoritma K-Nearest Neighbor untuk Klasifikasi Tingkat Bullying di Kalangan Siswa SMP Negeri 1 Salapian. *Global Research and Innovation Journal*, 1(3), 887–893.
- [13] Patrisyah, A., Buaton, R., & Sitompul, J. N. (2024). Klasifikasi tingkat pemahaman siswa pada pelajaran matematika di MTSS PAB 5 Klambir Lima. *Saturnus: Jurnal Teknologi dan Sistem Informasi*, 2(4), 146–156. <https://doi.org/10.61132/saturnus.v2i4.345>.

