

Optimasi Algoritma Decision Tree Dengan Seleksi Fitur Dalam Klasifikasi Prestasi Akademik Siswa Sekolah

I Made Ryan Prana Dhita^{a1}, Gst. Ayu Vida Mastrika Giri^{a2}, I Putu Gede Hendra Suputra^{a3},
Anak Agung Istri Ngurah Eka Karyawati^{a4}

^aProgram Studi Informatika, Fakultas Matematika dan Ilmu Pengetahuan Alam,
Universitas Udayana

Jalan Raya Kampus Unud, Jimbaran, Bali, 80361, Indonesia

¹ryanprana555@gmail.com

²vida@unud.ac.id

³hendra.suputra@unud.ac.id

⁴eka.karyawati@unud.ac.id

Abstract

Quality education is an important foundation for the progress of the nation. In an effort to improve the quality of education, student academic achievement is a significant indicator. With the development of technology, the application of machine learning algorithms, such as Decision Tree, allows for more accurate prediction of student academic achievement. This study aims to optimize the Decision Tree C4.5 algorithm through a combined feature selection between Gain Ratio and Chi-Square to improve the performance of student academic achievement classification. The research data were collected from students of SMA N 1 Blahbatuh and went through a preprocessing process, feature selection, and evaluation using accuracy, precision, recall, and F1-Score metrics. The results showed that the combined feature selection method succeeded in improving the performance of the C4.5 algorithm with an accuracy of 82.2%, much higher than the model without feature selection (54.2%). The implementation of a web-based system was also developed to support practical predictions. Thus, the results of this study contribute to the development of educational data analysis methods to improve the quality of education in the future.

Keywords: Academic Achievement, Decision Tree C4.5, Gain Ratio, Chi-Square, Feature Selection, Machine Learning

1. Pendahuluan

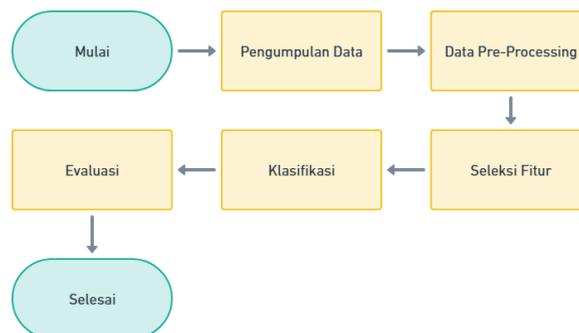
Pendidikan merupakan salah satu hal penting dalam kehidupan manusia. Kualitas pendidikan yang baik akan memberikan pengaruh positif terhadap perkembangan individu dan kemajuan bangsa, namun sebaliknya kualitas Pendidikan yang buruk akan membuat bangsa atau negara tersebut mengalami ketertinggalan[1], rendahnya kualitas Pendidikan di Indonesia membuat pemerintah berupaya lebih untuk meningkatkan mutu pendidikan dengan program serta kurikulum - kurikulum baru yang dibuat. Penerapan kurikulum ini diharapkan dapat menjadi kesempatan yang bagus untuk Indonesia dalam meningkatkan kualitas pendidikannya dan meningkatkan daya saing agar setara dengan negara-negara lain[2]. Dalam era digital saat ini, banyak sekolah dan perguruan tinggi yang memiliki data prestasi akademik siswa dalam bentuk dataset. Data ini dapat dimanfaatkan untuk melakukan analisis dan prediksi prestasi akademik siswa menggunakan algoritma machine learning. Salah satu algoritma machine learning yang dapat digunakan adalah decision tree. Algoritma pembelajaran mesin decision tree merupakan metode klasifikasi dan prediksi yang sangat kuat dan terkenal. Metode pohon keputusan mengubah fakta yang sangat besar menjadi pohon keputusan yang merepresentasikan aturan[3] serta decision tree adalah sebuah struktur yang dapat digunakan untuk membagi kumpulan data yang besar menjadi himpunan himpunan record yang lebih kecil dengan menerapkan serangkaian aturan keputusan. Dengan masing-masing rangkaian pembagian, anggota himpunan hasil menjadi mirip satu dengan yang lain[4]. Namun untuk mengoptimalkan prestasi metode seperti seleksi fitur perlu diterapkan sehingga dapat menghasilkan prediksi yang lebih baik dan akurat.

Dalam klasifikasi dengan decision tree metode seleksi fitur dengan gain ratio dapat digunakan, metode ini digunakan untuk memilih atribut terbaik yang akan digunakan untuk membagi data di setiap

langkah dalam proses pembuatan pohon keputusan. Gain ratio mengatasi beberapa kelemahan dari metode "information gain" lainnya, seperti information gain yang cenderung memilih atribut dengan banyak nilai (atribut dengan banyak nilai cenderung memiliki information gain yang lebih tinggi), penerapan fitur seleksi gain ratio sangat bermanfaat dalam klasifikasi. Hasil penelitian tersebut menyimpulkan bahwa penggunaan fitur seleksi gain ratio dapat secara drastis mempercepat waktu pembangunan model klasifikasi dan juga hal ini dapat membantu menghindari bias terhadap atribut dengan banyak nilai dan menghasilkan pohon keputusan yang lebih seimbang[5].

2. Metode Penelitian

Penelitian dimulai dari tahapan pengumpulan dataset lalu pemilihan algoritma yaitu algoritma decision tree, melakukan normalisasi data serta data encoding, seleksi fitur, klasifikasi lalu evaluasi hasil dari data yang telah diproses, dapat dilihat pada gambar 1 merupakan tahapan penelitian merupakan ilustrasi langkah- langkah metode yang akan dikerjakan.



Gambar 1. Metode Penelitian

2.1. Pengumpulan Data

pengumpulan data dilakukan dengan cara melakukan wawancara kepada kepala sekolah dari instansi terkait serta menyebar kuesioner kepada siswa sekolah (<https://forms.gle/r9Gqs8DdVB1iJqFT7>) dan nantinya data yang telah diperoleh melalui wawancara serta kuesioner yang dibagikan akan digunakan sebagai dataset dari penelitian ini.

2.2. Data Pre- Processing

2.2.1. Pembersihan Data

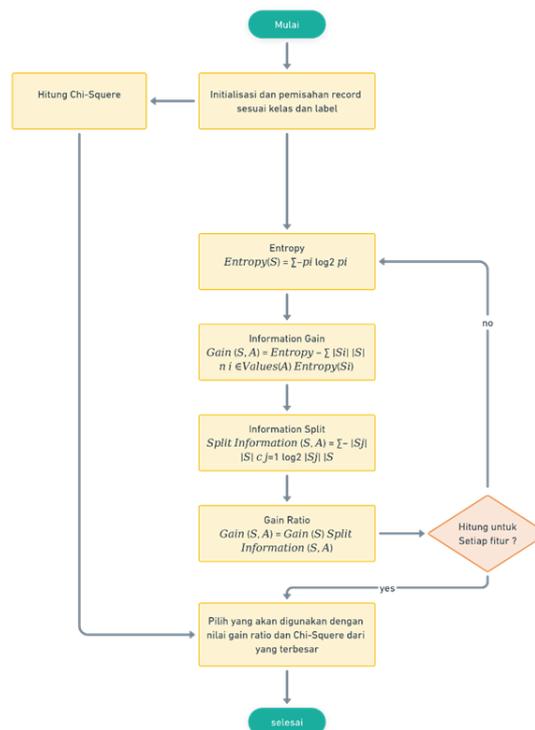
Data *cleaning* merupakan proses yang bertujuan untuk menganalisis kualitas data dengan melakukan tindakan seperti menghapus data duplikat, memperbaiki atau menghapus data yang hilang (*missing value*), serta menghilangkan *outlier*. Proses ini penting untuk memastikan data yang digunakan dalam analisis atau pemodelan memiliki integritas yang baik dan dapat diandalkan. Dengan melakukan data *cleaning*, kita dapat memastikan bahwa data yang digunakan dalam analisis tidak tercemar oleh duplikasi, informasi yang hilang, atau nilai yang tidak wajar yang dapat mempengaruhi hasil dan kesimpulan yang diambil dari data tersebut.

2.2.2. Normalisasi Data

Normalisasi data melibatkan mengubah skala data agar sesuai dengan rentang atau standar yang ditetapkan. Normalisasi dapat melibatkan pemetaan data ke rentang yang spesifik atau transformasi data untuk menghilangkan bias atau asumsi tertentu. Metode *label encode* digunakan untuk normalisasi data pada penelitian ini dengan mentransformasikan setiap nilai dalam rentang data ke rentang baru.

2.3. Seleksi Fitur

Setelah dilakukan pengolahan data awal, pada tahapan ini akan dilakukan seleksi fitur dengan menggunakan *gain ratio* serta *chi-square*. Kombinasi kedua teknik ini bisa sangat bermanfaat dalam menentukan fitur-fitur yang penting untuk memprediksi kelas atau label yang benar dalam decision tree. Kombinasi *information Gain* sebagai kriteria awal untuk memilih fitur terbaik dan kemudian menggunakan *chi-square* untuk mengkonfirmasi hubungan statistik antara fitur-fitur tersebut dengan label kelas. Langkah-langkah metode *gain ratio* akan dijelaskan pada gambar 2



Gambar 2. flowchart Seleksi fitur

Penjelasan dari alur diagram alir seleksi fitur dengan *gain ratio* pada gambar 2 adalah sebagai berikut :

- **Hitung *Chi-Square*** : hitung dengan rumus chi-square berdasarkan data yang telah ditentukan

$$X^2 = \sum \cdot \frac{(F_0 - F_h)^2}{F_h} \quad (1)$$

- **Hitung Entropi Data Asli** : Entropi mengukur tingkat ketidakpastian atau keacakan dalam data. Hitung entropi data asli sebelum membagi data menjadi subgrup berdasarkan atribut apa pun. Entropi digunakan untuk mengukur seberapa banyak informasi yang ada dalam data sebelum membangun pohon keputusan.

$$Entropy(S) = \sum_i^c - pi \log \log_2 pi \quad (2)$$

- **Hitung *information gain*** : *Information gain* adalah ukuran yang digunakan untuk mengukur seberapa banyak informasi baru yang dihasilkan ketika data dibagi berdasarkan atribut tertentu. Semakin tinggi *information gain*, semakin banyak informasi baru yang diperoleh dari pemisahan data dengan atribut tersebut. *Information gain* dihitung dengan membandingkan entropi data sebelum dan setelah pemisahan. Entropi data mengukur tingkat ketidakpastian atau keacakan dalam data. Ketika data homogen (semua contoh data memiliki label yang sama), entropi akan rendah (0). Sebaliknya, jika data heterogen (contoh data memiliki label yang berbeda-beda), entropi akan tinggi (nilai lebih dari 0). Tujuan dari *information gain* adalah untuk mencari atribut yang menghasilkan pemisahan data yang homogen, sehingga mengurangi ketidakpastian dan meningkatkan akurasi model.

$$Gain(S, A) = Entropy - \sum_{i \in Values(A)} \frac{|S_i|}{|S|} Entropy(S_i) \quad (3)$$

- **Hitung *split information*** : Split information mengukur tingkat variasi atau pemisahan yang dihasilkan oleh atribut tertentu. Semakin tinggi split information, semakin tinggi tingkat variasi atau jumlah cabang yang dihasilkan oleh atribut tersebut. Split information dihitung dengan memperhitungkan jumlah cabang yang mungkin dihasilkan oleh atribut.

$$Split\ Information(S, A) = \sum_{j=1}^c - \frac{|S_j|}{|S|} \log_2 \frac{|S_j|}{|S|} \quad (4)$$

- **Hitung *gain ratio*** : *Gain ratio* adalah ukuran yang menggabungkan gain informasi dan tingkat pemisahan atribut. Gain informasi mengukur seberapa banyak informasi baru yang dihasilkan dengan membagi data berdasarkan atribut tertentu. Gain ratio menghindari atribut yang memiliki banyak nilai (banyak cabang) yang cenderung menghasilkan pohon yang terlalu kompleks.

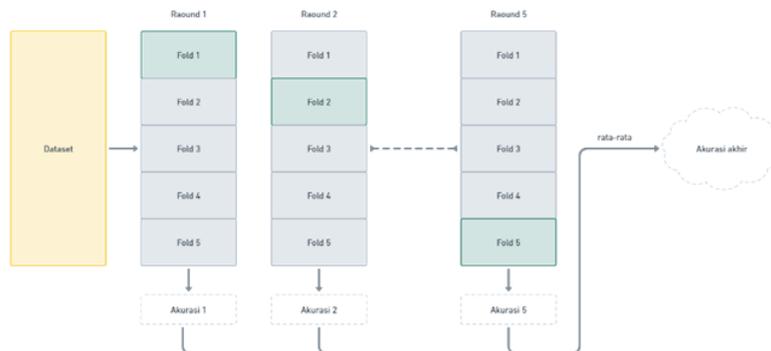
$$Gain(S, A) = \frac{Gain(S)}{Split\ Information(S, A)} \quad (5)$$

Dengan menggabungkan *gain ratio* dan *chi-square* yang terlihat pada persamaan 5, kita dapat memilih atribut terbaik untuk digunakan dalam membangun pohon keputusan dengan metode seleksi fitur gabungan tersebut.

2.4. Klasifikasi

2.4.1. K-Fold Cross Validation

Proses klasifikasi dilakukan melalui penggunaan cross-validation, di mana data dibagi secara bergantian menjadi data pelatihan (misalnya K = 5) dengan ukuran yang sama. Setiap subset disebut sebagai "fold".



Gambar 3. Model *K- Fold Cross Validation*

Seperti yang terlihat pada gambar 3 setiap iterasi dalam proses klasifikasi satu fold digunakan sebagai data validasi, sementara sisanya digunakan sebagai data pelatihan misalnya, jika klasifikasi berada pada iterasi ke-3, fold ke-3 digunakan sebagai data validasi, dan sisa fold (1, 2, 4, dan 5) digunakan sebagai data pelatihan [6].

2.4.2. Split Data

Pada setiap iterasi k dalam K-Fold Cross Validation secara otomatis terdapat pembagian data atau split data dimana 80% akan menjadi data latih dan 20% akan menjadi data uji. Sehingga, setelah dihubungkan dengan pembagian pada *K-Fold Cross Validation*, maka 1 fold akan menjadi data uji dan 4 fold lainnya akan menjadi data latih.

2.4.3. Algoritma C4.5

Setelah dilakukan seleksi fitur, atribut terpilih akan digunakan dalam proses klasifikasi menggunakan algoritma C4.5 untuk membuat pohon keputusan dari set data training yang digunakan untuk klasifikasi kelas dari data testing.

2.5. Evaluasi Dan Analisis

Tahap evaluasi terhadap metode dalam penelitian ini dilakukan menggunakan metode *confusion matrix*. *Confusion matrix* merupakan sebuah metode evaluasi yang umum digunakan untuk mengukur prestasi algoritma *Decision tree* dalam melakukan klasifikasi. *Confusion matrix* memberikan gambaran tentang seberapa baik model *Decision tree* dapat mengklasifikasikan data ke dalam kelas yang benar[7].

3. Hasil dan Pembahasan

3.1. Pengumpulan Data

Pada tahap ini adalah bagian pengumpulan data yang diperlukan untuk menunjang penelitian ini, Pengumpulan data dilakukan di SMA N 1 BLAHBATUH dengan membagikan form maka diperoleh jumlah data dengan tingkatan kelas seperti tabel 1.

Tabel 1. Jumlah Berdasarkan Kelas Siswa

No	Kelas	Jumlah
1	Cukup	64
2	Baik	46
3	Kurang	5
4	Sangat Baik	3
Total		118

Berdasarkan data pada tabel 1 maka data yang di dapat sejumlah 118 dengan jumlah kelas kurang 5, cukup 64, baik 46, sangat baik 3 untuk chart data lebih lengkapnya dapat dilihat pada lampiran.

3.2. Preprocessing Data

Data yang telah dimuat dan digunakan pada penelitian sejumlah 118 yang kemudian dilakukan *cleaning dataset* yang meliputi penghapusan data duplikat dan penghapusan data yang mengandung nilai kosong serta penghapusan data yang tidak ingin digunakan. mengisi nilai kosong (*NaN*) di dataset dengan meneruskan nilai dari baris sebelumnya di kolom yang sama tujuannya adalah memastikan tidak ada nilai kosong dalam dataset, sehingga data tetap konsisten untuk dianalisis lebih lanjut.

3.3. Seleksi Fitur

Dalam penelitian ini penulis akan melakukan pengembangan metode dengan melakukan seleksi fitur seleksi fitur ini dilakukan dengan tujuan untuk menghilangkan fitur-fitur dengan tingkat kebergunaan yang rendah dalam model sehingga dapat mengurangi beban model ketika melakukan klasifikasi hasil seleksi fitur yang hasilnya dapat terlihat pada tabel 2 dan tabel 3.

Tabel 2. Fitur Data Hasil Chi-Square dan *Gain Ration*

Fitur	Skor Chi-Square	Skor Gain Ration
Kenyamanan dalam lingkungan keluarga	32.720455	0.279771
Berapa jarak rumah ke sekolah	20.036781	0.061604
Jam tidur siswa	13.352250	0.075025
sering atau sedang memiliki permasalahan dengan orang tua/keluarga/teman kelas atau lingkungan sekitar	12.009175	0.087716
Waktu yang dihabiskan untuk bermain	11.311554	0.035653
Jenis kelamin	9.355097	0.065582
Kepercayaan diri siswa di sekolah	9.098040	0.061571
Interaktivitas tenaga pengajar di sekolah	8.774870	0.104061
Minat belajar siswa di sekolah	8.503138	0.232432
Masalah ekonomi keluarga	2.255357	0.014231
Apakah sedang mempunyai pasangan/pacar	1.748603	0.033164
Kenyamanan dalam lingkungan sekolah	1.634160	0.023850
Partisipasi dalam organisasi	1.559116	0.016759
Kenyamanan dalam lingkungan kelas/teman kelas	1.118090	0.016254

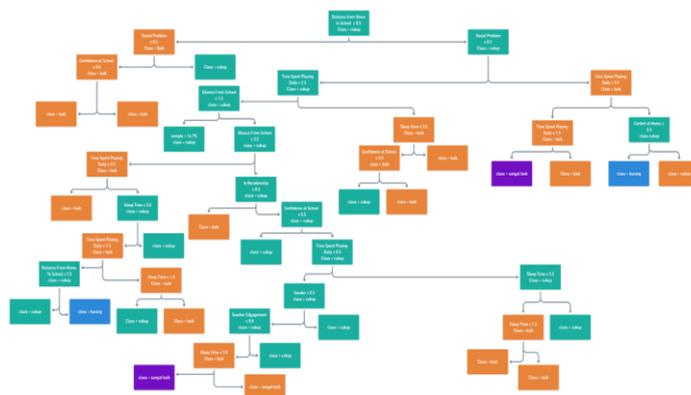
Pada tabel 2 menampilkan skor seleksi fitur *gain ratio* dan chi-square yang telah dimana semakin besar skor maka semakin relevan relasi fitur dengan hasil prediksi nantinya setelah mendapatkan data hasil dari metode tersebut maka akan diurutkan dengan mengambil fitur terbaik di setiap metode yang telah yang telah diujikan maka didapatkan 10 fitur terbaik yang relevan pada masing-masing seleksi fitur tersebut yang dapat dilihat pada tabel 3.

Tabel 3. Fitur Data Hasil Akhir Seleksi Fitur Gabungan

Fitur	Skor Chi-Square	Skor Gain Ration
Kenyamanan dalam lingkungan keluarga	32.720455	0.279771
Berapa jarak rumah ke sekolah	20.036781	0.061604
Jam tidur siswa	13.352250	0.075025
sering atau sedang memiliki permasalahan dengan orang tua/keluarga/teman kelas atau lingkungan sekitar	12.009175	0.087716
Waktu yang dihabiskan untuk bermain	11.311554	0.035653
Jenis kelamin	9.355097	0.065582
Kepercayaan diri siswa di sekolah	9.098040	0.061571
Interaktivitas tenaga pengajar di sekolah	8.774870	0.104061
Apakah sedang mempunyai pasangan saat ini?	0.033164	1.748603
Minat belajar siswa di sekolah	8.503138	0.232432

3.4. Klasifikasi

Implementasi model klasifikasi dengan metode *decision tree* C4.5 serta akan dilakukan juga *k-fold cross validation* untuk melihat tingkat konsistensi model dalam melakukan klasifikasi sehingga model tidak hanya terlihat bagus dalam satu kondisi saja dan menghasilkan model terbaik yang visualnya dapat dilihat pada gambar 4.

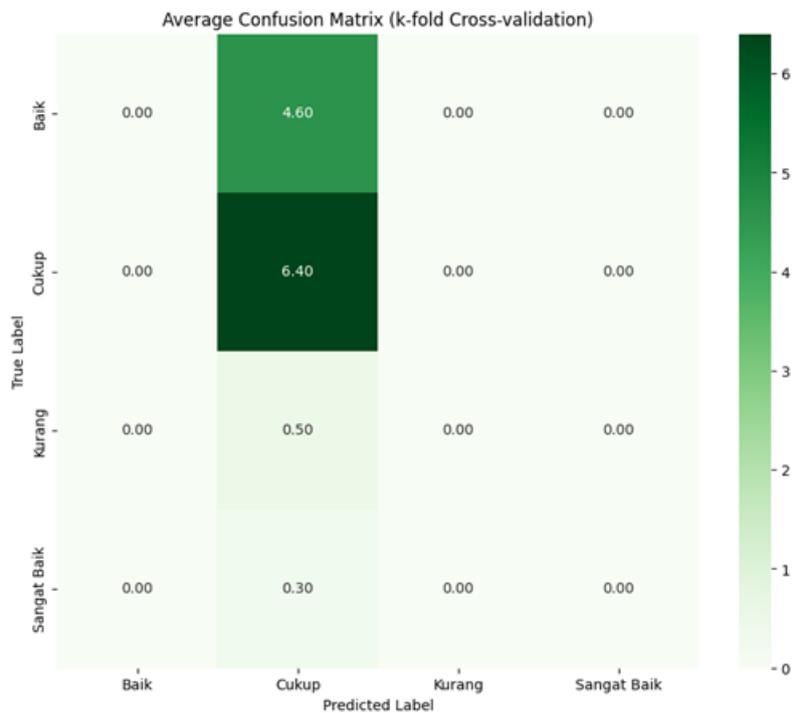


Gambar 4. Visualisasi Model Decision Tree

3.5. Evaluasi Dan Analisis

Pada Implementasi penulis melakukan evaluasi rata-rata *confusion matrix* dari akurasi terbaik k-fold cross-validation, seperti True Positive (TP), False Positive (FP), False Negative (FN), dan True Negative (TN), yang memberikan gambaran detail performa model serta membedakan antara kelas positif dan negatif. Hasil evaluasi juga divisualisasikan dalam bentuk heatmap menggunakan *seaborn*,

yang mempermudah interpretasi *confusion matrix* dengan tampilan yang informatif, menampilkan hubungan antara prediksi dan label sebenarnya untuk setiap kelas.



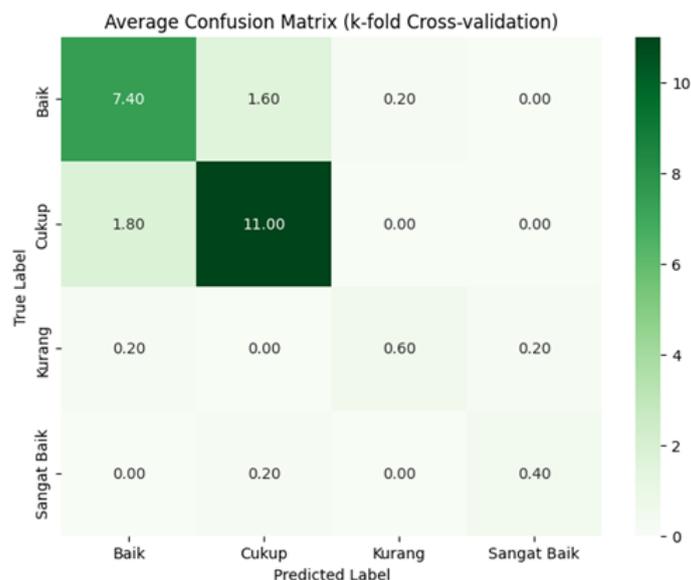
Gambar 5. Confusion Matrix Model Decision Tree c4.5

Pada gambar 5 merupakan hasil *confusion matrix* model *decision tree* tanpa implementasi seleksi fitur. Detail matriks evaluasi dapat dilihat pada tabel 4

Tabel 4. Matriks Evaluasi Model *Decision Tree* tanpa Seleksi Fitur

Matrik Evaluasi	Nilai
<i>Accuracy</i>	0.542
<i>Precision</i>	0.215
<i>Recall</i>	0.391
<i>F1-Score</i>	0.276

Sedangkan hasil *confusion matrix* model *decision tree* dengan implementasi seleksi fitur dapat dilihat pada gambar 6



Gambar 6. Confusion Matrix Model Decision Tree c4.5 dengan Implementasi Seleksi Fitur

Untuk detail matriks evaluasi model decision tree c4.5 dengan implementasi seleksi fitur dapat dilihat pada tabel 5.

Tabel 5. Matriks Evaluasi Model *Decision Tree* dengan Seleksi Fitur

Matrik Evaluasi	Nilai
<i>Accuracy</i>	0.820
<i>Precision</i>	0.799
<i>Recall</i>	0.809
<i>F1-Score</i>	0.792

Berdasarkan hasil implementasi model klasifikasi telah diperoleh sebanyak dua model klasifikasi yang dapat dibandingkan antara satu dengan yang lain untuk memilih model terbaik yang akan diimplementasikan dalam sistem. Adapun resume hasil implementasi dapat dilihat pada tabel 6.

Tabel 6. Resume Hasil Evaluasi Model

Model	TP	TN	FP	FN	AC	PR	RC	F1
A	6	29	5	4	0.542	0.215	0.391	0.276
B	8	32	1	1	0.822	0.788	0.773	0.769

Model A merupakan model awal yang memperoleh accuracy 0.542, precision 0.215, recall 0.391, dan F-1 score 0.276. Dimana berdasarkan tabel 4.10 nilai tersebut menunjukkan performa model yang tergolong kurang baik dalam melakukan klasifikasi. Kemudian pada implementasi selanjutnya dengan mengimplementasikan seleksi fitur, didapatkan sembilan fitur yang memiliki korelasi terbaik yaitu “Kenyamanan dalam lingkungan keluarga”, “Berapa jarak rumah ke sekolah”, “Jam tidur”, “sering atau sedang memiliki permasalahan dengan orang tua/keluarga/teman kelas atau lingkungan sekitar”, “Waktu yang dihabiskan untuk bermain”, “Jenis kelamin”, “Kepercayaan diri siswa di sekolah”, “Interaktivitas tenaga pengajar di sekolah”, “Minat belajar siswa di sekolah” fitur tersebut yang digunakan dalam implementasi selanjutnya yang menghasilkan Model B. Model B yang merupakan model hasil implementasi

seleksi fitur terbaik yang memperoleh accuracy 0.820, precision 0.799, recall 0.809, F-1 score 0.782 lebih baik dari model sebelumnya, Untuk detail hasil pengujian fitur lainnya dapat dilihat pada lampiran.

3.6. Implementasi Pada Model Sistem

Dari hasil model yang terpilih selanjutnya diimplementasikan pada sebuah sistem sederhana untuk melakukan prediksi prestasi siswa sekolah berdasarkan hasil klasifikasi dari model. Sistem dibangun dalam aplikasi website dalam pengimplementasiannya.



Gambar 8. Halaman Lading Page

4. Kesimpulan

Terdapat peningkatan performa model dengan seleksi fitur menggunakan metode kombinasi gain ratio dan chi- yang berhasil meningkatkan performa model klasifikasi. Hal ini terlihat dari perbandingan antara model awal tanpa seleksi fitur dan model dengan seleksi fitur. Model dengan seleksi fitur memperoleh rata-rata akurasi sebesar 82.2%, *precision* 78.8%, *recall* 77.3%, dan *F1-score* 76.9%, yang jauh lebih baik dibandingkan model awal dengan akurasi 54.2%, *precision* 21.5%, *recall* 39.1%, dan *F1-score* 27.6%. Seleksi fitur chi-square dan gain ratio yang telah diimplementasikan mampu mengidentifikasi sepuluh fitur terbaik yang relevan untuk prediksi prestasi akademik siswa dimana setelah dilakukan percobaan *threshold* fitur yang terpilih pada model akurasi terbaik memiliki score gain ratio di atas 0.015 dan score chi-square di atas 1.60. Pengurangan jumlah fitur ini telah berhasil meningkatkan akurasi model dengan cukup baik.

Referensi

- [1] Kurniawati, F. N. (2022). MENINJAU PERMASALAHAN RENDAHNYA KUALITAS PENDIDIKAN DI INDONESIA DAN SOLUSI. *Academy of Education Journal*, 13, 13.
- [2] Syarifuddin. (2022). Menurunnya Kualitas Pendidikan Seiring Dengan Adanya Inovasi Pendidikan. : *Inovasi Pendidikan-AKBK3602*, 1, 6.
- [3] Sari, B. N. (216). IMPLEMENTASI TEKNIK SELEKSI FITUR INFORMATION GAIN PADA ALGORITMA KLASIFIKASI MACHINE LEARNING UNTU PREDIKSI PERFORMA AKADEMIK SISWA. *Betha Nurina Sari(2302-3805)*, 6.
- [4] Vivy Junita, F. A. (2019). Klasifikasi Aktivitas Manusia menggunakan Algoritme Decision Tree C4.5 dan Information Gain untuk Seleksi Fitur. *Jurnal Pengembangan Teknologi Informasi dan Ilmu Komputer*, 3(10), 8.

- [5] WAHYUDI, M. A. (2020). PENERAPAN FITUR SELEKSI GAIN RATIO DAN DECISION TREE C5.0 UNTUK KLASIFIKASI TINGKAT SERANGAN JARINGAN. *SUSKA RIAU*, 1(1), 95.
- [6] Msy Aulia Hasanah, S. S. (2021). Implementasi CRISP-DM Model Menggunakan Metode Decision Tree dengan Algoritma CART untuk Prediksi Curah Hujan Berpotensi Banjir. *Journal of Applied Informatics and Computing (JAIC)* , 5(2), 6.
- [7] Oman Somantri, W. E. (2022). Penerapan Feature Selection Pada Algoritma Decision Tree Untuk Menentukan Pola Rekomendasi Dini Konseling . *Jurnal Sistem Komputer dan Informatika (JSON)*, 4(2), 8.
- [8] Ozdemir, S. (2016). *Principles of Data Science*. Birmingham B3 2PB, UK.: Packt Publishing Ltd.
- [9] Ripanti, E. P. (2019). Model Prediksi Awal Masa Studi Mahasiswa Menggunakan Algoritma Decision tree c4.5. *JUSTIN (Jurnal Sistem dan Teknologi Informasi)*, 7(4), 7.
- [10] Suradarma, I. M. (2018). Penerapan Optimasi Algoritma C45 dengan Naïve Bayes pada Pemilihan Internet Service Provider. *EKSPLORA INFORMATIKA*, 7(2), 11.

This page is intentionally left blank.