



# Implementasi Algoritma C4.5 Untuk Klasifikasi Pengenalan Warna Dasar di Taman Kanak-Kanak

Anandaya Difi Dzulardi Kalimanti, Vilianty Rafida, Aisyah Fajriantini\*

Program Studi Teknik Informatika, STMIK Widya Cipta Dharma, Kota Samarinda, Indonesia

Email: <sup>1</sup>2243016@wicida.ac.id, <sup>2</sup>vilianty@wicida.ac.id, <sup>3</sup>\*aisyah@wicida.ac.id

Email Penulis Korespondensi: aisyah@wicida.ac.id

**Abstrak**—Ketimpangan kemampuan anak usia dini dalam mengenali warna dasar di TK Negeri 01 Barong Tongkok menunjukkan perlunya sistem evaluasi yang tersusun secara sistematis agar penilaian dapat dilakukan secara objektif. Pengelompokan kemampuan tersebut dilakukan dengan memanfaatkan algoritma C4.5 dalam kerangka penelitian kuantitatif berdesain eksperimen. Data diperoleh melalui observasi dan tes pengenalan warna yang melibatkan 35 anak sebagai responden, kemudian diolah berdasarkan atribut yang telah ditentukan untuk membangun model klasifikasi. Hasil analisis mengelompokkan kemampuan anak ke dalam tiga kategori, yaitu cukup mengenal (rendah), mengenal (sedang), dan sangat mengenal (tinggi). Hasil penelitian menunjukkan bahwa algoritma C4.5 sangat efektif dalam memetakan kemampuan anak dengan capaian rata-rata akurasi sebesar 85,71% melalui pengujian *5-Fold Cross Validation*. Lebih lanjut, pohon keputusan yang dihasilkan menyediakan struktur yang intuitif dan transparan, sehingga memudahkan tenaga pendidik dalam menginterpretasikan hasil evaluasi serta memahami variabel dominan yang menentukan tingkat keberhasilan belajar siswa dibandingkan dengan penggunaan model *black-box* lainnya. Kontribusi utama penelitian ini terletak pada penyediaan model evaluasi berbasis data yang menghasilkan aturan keputusan (*if-then rules*) yang terukur secara empiris, sekaligus menjadi jembatan metodologis untuk menciptakan strategi pembelajaran diferensiasi di tingkat PAUD. Dengan demikian, implementasi algoritma C4.5 mewakili alternatif yang strategis, efisien, dan dapat dipertanggungjawabkan secara ilmiah untuk meningkatkan efektivitas pedagogis serta pemantauan kognitif dalam pendidikan anak usia dini.

**Kata Kunci:** Klasifikasi; C4.5; Pengenalan Warna; Anak Usia Dini; Decision Tree

**Abstract**—Differences in early childhood ability to recognize basic colors at TK Negeri 01 Barong Tongkok indicate the need for a structured evaluation system to ensure objective assessment. The classification of these abilities is carried out by applying the C4.5 algorithm within a quantitative experimental framework. Data are collected through observations and color recognition tests involving 35 children as respondents, then processed using predefined attributes to construct a classification model. The analysis results group children's abilities into three categories: *Sangat Mengenal* (High), *Mengenal* (Moderate), and *Cukup Mengenal* (Low). The experimental results indicate that the C4.5 algorithm is highly effective and stable, achieving an average classification accuracy of 85.71% through 5-Fold Cross-Validation. Furthermore, the resulting decision tree provides an intuitive and transparent structure that assists educators in interpreting evaluation outcomes and understanding the dominant variables that determine student learning success more clearly than black-box models. The primary contribution of this study lies in the provision of a data-driven evaluation model that generates empirically measurable decision rules (if-then rules), while simultaneously serving as a methodological bridge to create differentiated learning strategies at the early childhood education (PAUD) level. Consequently, the implementation of the C4.5 algorithm represents a strategic, efficient, and scientifically accountable alternative for enhancing pedagogical effectiveness and cognitive monitoring in early childhood education.

**Keywords:** Classification; C4.5; Color Recognition; Early Childhood; Decision Tree

## 1. PENDAHULUAN

Kemampuan mengenali warna dasar menjadi bagian penting dalam perkembangan kognitif anak usia dini dan berperan sebagai landasan dalam proses pembelajaran. Penguasaan kemampuan ini membantu anak memahami lingkungan sekitarnya sekaligus mendukung perkembangan berpikir, komunikasi, serta kesiapan dalam mengikuti pembelajaran pada tahap selanjutnya [1]. Dari sudut pandang kognitif, pengenalan warna juga berkaitan dengan kemampuan anak dalam melakukan identifikasi, klasifikasi, dan pengelompokan objek yang dijumpai dalam kehidupan sehari-hari [2]. Secara neurologis, proses ini melibatkan koordinasi antara persepsi visual dan memori jangka pendek, di mana anak belajar mengasosiasikan label verbal dengan stimulasi kromatik yang diterima indra penglihatan. Namun demikian, praktik pembelajaran di taman kanak-kanak masih menunjukkan bahwa sebagian anak belum mampu membedakan warna dasar seperti merah, kuning, dan biru secara optimal. Kondisi ini dipengaruhi oleh berbagai faktor, di antaranya kurangnya variasi metode pembelajaran, keterbatasan media pembelajaran interaktif, serta perbedaan tingkat perkembangan kognitif antar anak [3]. Hambatan dalam penguasaan konsep warna ini jika tidak terdeteksi sejak dini dapat berdampak pada kepercayaan diri anak dalam berinteraksi dengan materi pembelajaran yang lebih kompleks di masa depan. Situasi tersebut menegaskan perlunya pendekatan yang mampu mengidentifikasi kemampuan anak secara lebih sistematis dan akurat sebagai dasar evaluasi pembelajaran.

Sejumlah penelitian terdahulu menunjukkan bahwa penggunaan media pembelajaran visual dan interaktif mampu meningkatkan kemampuan anak dalam mengenali warna secara signifikan [4]. Pendekatan pembelajaran yang melibatkan aktivitas langsung dan pengalaman nyata juga terbukti efektif dalam memperkuat pemahaman anak terhadap konsep warna [5]. Selain itu, pemanfaatan teknologi digital dalam evaluasi pembelajaran dinilai mampu meningkatkan efektivitas sekaligus keterlibatan peserta didik dibandingkan metode konvensional, sehingga hasil pembelajaran dapat diukur secara lebih objektif dan terstruktur [6].

Evaluasi berbasis teknologi memberikan keunggulan dalam hal standarisasi penilaian, sehingga subjektivitas pendidik dapat diminimalisir melalui parameter ukur yang konsisten. Seiring perkembangan teknologi, penerapan metode



berbasis data mining dalam pendidikan semakin meluas karena memungkinkan analisis data dilakukan secara sistematis untuk menemukan pola tersembunyi [7]. Salah satu metode yang sering digunakan adalah algoritma Decision Tree, yang dikenal mampu menghasilkan model klasifikasi dengan tingkat akurasi yang baik serta mudah dipahami [8]. Algoritma ini bekerja dengan membentuk struktur pohon keputusan berdasarkan atribut yang digunakan, sehingga data dapat dikelompokkan ke dalam kategori tertentu secara terstruktur. Dalam pendidikan anak usia dini, pendekatan ini dapat dimanfaatkan untuk membantu guru mengelompokkan kemampuan anak berdasarkan hasil evaluasi pembelajaran agar strategi yang diterapkan lebih tepat sasaran. Selain itu, teknik klasifikasi berbasis machine learning telah banyak digunakan dalam mendukung pengambilan keputusan di bidang pendidikan [9]. Penggunaan model komputasi ini terbukti lebih unggul dalam mengolah variabel yang kompleks dibandingkan dengan metode statistik manual yang seringkali terbatas pada analisis deskriptif sederhana.

Salah satu algoritma yang umum digunakan adalah C4.5, yang merupakan pengembangan dari Decision Tree dengan memanfaatkan konsep entropy dan gain ratio dalam menentukan atribut terbaik [10]. Melalui penghitungan Gain Ratio, algoritma C4.5 mampu mengoreksi kelemahan pendahulunya, yaitu ID3, dalam hal bias terhadap atribut yang memiliki banyak nilai unik, sehingga menghasilkan pohon keputusan yang lebih efisien dan akurat. Keunggulan algoritma ini terletak pada kemampuannya menghasilkan model yang mudah diinterpretasikan serta mampu menangani data numerik dan kategorikal. Karakteristik pohon keputusan yang menyerupai alur berpikir manusia menjadikannya instrumen yang sangat relevan untuk diaplikasikan dalam lingkungan pendidikan yang memerlukan transparansi logika. Berbagai penelitian menunjukkan bahwa C4.5 efektif dalam berbagai kasus klasifikasi dengan tingkat akurasi yang baik [11]. Di sisi lain, penerapan metode klasifikasi dalam pendidikan terbukti mampu meningkatkan ketepatan pengambilan keputusan dan memberikan rekomendasi pembelajaran yang lebih optimal [12].

Meskipun demikian, penerapan algoritma klasifikasi untuk mengidentifikasi kemampuan pengenalan warna dasar pada anak usia dini masih terbatas, karena sebagian besar penelitian lebih berfokus pada jenjang pendidikan yang lebih tinggi atau aspek pembelajaran yang bersifat umum. Selain itu, integrasi hasil klasifikasi sebagai dasar evaluasi pembelajaran di tingkat taman kanak-kanak masih jarang dilakukan, sehingga menunjukkan adanya kesenjangan penelitian yang perlu dikembangkan lebih lanjut. Permasalahan lain yang sering muncul adalah ketidakseimbangan jumlah data pada setiap kategori (*imbalanced data*), yang dapat memengaruhi kinerja model dan menurunkan tingkat akurasi klasifikasi [13]. Kekosongan penelitian (*research gap*) ini terlihat jelas dari analisis literatur sejenis; misalnya, penelitian oleh Anwar dkk. [10] serta Alfayyadh & Assegaff [19] yang berfokus pada prediksi kelulusan dan performa akademik mahasiswa, penelitian Nugraha dkk. [11] yang terbatas pada penentuan penerima beasiswa di tingkat sekolah menengah, serta penelitian Edy & Lasut [12] yang menyorot pada sistem rekomendasi penjurusan berbasis web.

Kebanyakan sistem pakar yang ada saat ini masih menyorot pada diagnosa medis atau performa akademik mahasiswa, sementara instrumentasi untuk memetakan perkembangan kognitif dasar anak masih sangat minim. Oleh karena itu, penelitian ini tidak hanya berfokus pada klasifikasi semata, tetapi juga memperhatikan validasi model melalui teknik cross-validation untuk memastikan reliabilitas hasil. Penggunaan 5-Fold Cross Validation dalam penelitian ini dimaksudkan untuk membagi data secara proporsional sehingga setiap data memiliki kesempatan yang sama untuk menjadi data latih dan data uji, yang pada akhirnya meminimalkan risiko overfitting. Berdasarkan kondisi tersebut, penelitian ini diarahkan pada klasifikasi kemampuan pengenalan warna dasar pada anak usia dini di TK Negeri 01 Barong Tongkok dengan memanfaatkan algoritma Decision Tree.

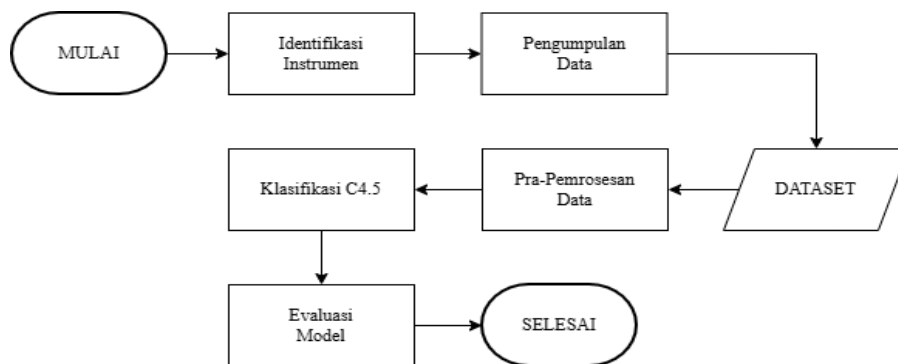
Data yang digunakan berasal dari hasil observasi dan tes pengenalan warna, dengan tujuan menghasilkan model klasifikasi yang mampu mengelompokkan kemampuan anak secara akurat dan mudah dipahami. Hasil penelitian diharapkan dapat membantu guru dalam melakukan evaluasi pembelajaran serta menyusun strategi pembelajaran yang lebih efektif sesuai dengan kebutuhan masing-masing anak, sekaligus memberikan kontribusi terhadap pengembangan pendekatan pembelajaran berbasis data pada pendidikan anak usia dini. Secara akademis, kontribusi penelitian ini terletak pada pemodelan pohon keputusan yang dapat dijadikan acuan baku dalam pemetaan kompetensi anak berdasarkan distribusi statistik normal.

Permasalahan dan kesenjangan yang telah diidentifikasi mendorong dilakukannya klasifikasi kemampuan pengenalan warna dasar pada anak usia dini di TK Negeri 01 Barong Tongkok melalui pemanfaatan algoritma Decision Tree. Data yang diperoleh dari observasi dan tes pengenalan warna diolah untuk membentuk model klasifikasi yang mampu mengelompokkan kemampuan anak secara akurat serta mudah diinterpretasikan. Model tersebut diharapkan dapat mendukung guru dalam melakukan evaluasi pembelajaran sekaligus merancang strategi pembelajaran yang lebih efektif sesuai dengan kebutuhan masing-masing anak. Dengan demikian, kajian ini diharapkan memberikan kontribusi dalam pengembangan pendekatan pembelajaran berbasis data pada pendidikan anak usia dini.

## 2. METODOLOGI PENELITIAN

### 2.1 Tahapan Penelitian

Tahapan penelitian ini disusun secara sistematis untuk memastikan akurasi klasifikasi tingkat kemampuan pengenalan warna pada anak usia dini. Alur penelitian dimulai dari identifikasi instrumen, pengumpulan data melalui kuesioner interaktif, pra-pemrosesan data, klasifikasi dengan algoritma C4.5, hingga tahap evaluasi menggunakan *K-Fold Cross Validation*.



Gambar 1. Flowchart Tahapan Penelitian

Gambar 1 menunjukkan tahapan penelitian dimulai dengan identifikasi instrumen melalui penyusunan 10 pertanyaan yang mencakup aspek pengenalan warna dasar dan sekunder, kemudian dilanjutkan dengan pengumpulan data dari 35 siswa di TK Negeri 01 Barong Tongkok menggunakan kuesioner interaktif. Data yang diperoleh selanjutnya diproses pada tahap pra-pemrosesan untuk dilakukan pembersihan serta pemberian label kategori kemampuan (Sangat Mengenal, Mengenal, Cukup Mengenal) dengan mengacu pada perhitungan statistik Distribusi Normal. Setelah melalui tahap tersebut, data dianalisis menggunakan algoritma C4.5, yang bekerja dengan menghitung nilai Entropy dan Gain Ratio guna membentuk struktur pohon keputusan yang sistematis. Pada tahap akhir, model klasifikasi yang dihasilkan dievaluasi melalui metode *5-Fold Cross Validation* untuk menguji konsistensi dan kestabilan tingkat akurasi secara objektif.

### 2.2 Pengumpulan Data dan Subjek Penelitian

Sebanyak 35 siswa aktif di TK Negeri 01 Barong Tongkok dilibatkan sebagai responden dalam penelitian ini. Data utama diperoleh melalui kegiatan observasi serta pelaksanaan tes pengenalan warna yang memuat 10 pertanyaan terkait warna dasar dan sekunder. Variabel yang dianalisis mencakup umur, jenis kelamin, dan skor pengenalan warna. Umur dikategorikan sebagai data numerik dengan rentang 4–6 tahun, sementara jenis kelamin merupakan data kategorial yang terdiri atas laki-laki dan perempuan. Skor warna dihitung berdasarkan perolehan poin dari setiap butir soal dengan total 10 pertanyaan, kemudian digunakan untuk menentukan kategori kemampuan, yaitu Cukup Mengenal, Mengenal, dan Sangat Mengenal. Keseluruhan variabel tersebut menjadi landasan dalam proses analisis guna membangun model klasifikasi yang merepresentasikan tingkat kemampuan pengenalan warna anak secara sistematis. Yang ditunjukkan pada tabel 1 sebagai berikut.

Tabel 1. Jenis jenis database

Nama Variabel	Jenis Data	Deskripsi
Umur	Numerik	Usia Responden(4-6tahun)
Jenis Kelamin	kategorial	Laki-laki dan Perempuan
Skor Warna	Numerik	Nilai poin per soal(10)
Kategori	Label	Cukup mengenal, Mengenal, Sangat Mengenal

Pengaturan variabel-variabel tersebut menjadi dasar penting dalam proses pengolahan data pada tahap analisis berikutnya, terutama dalam penerapan algoritma klasifikasi. Setiap atribut yang telah ditetapkan memungkinkan data dari responden diolah secara lebih sistematis sehingga keterkaitan antarvariabel dapat diamati dengan lebih jelas. Dengan demikian, hasil pengukuran tidak hanya merepresentasikan kondisi masing-masing anak secara individu, tetapi juga memperkuat pembentukan model yang mampu mengelompokkan tingkat kemampuan secara objektif, terstruktur, dan konsisten.

### 2.3 Distribusi Normal untuk Pelabelan

Proses pelabelan data dalam penelitian ini tidak dilakukan secara subjektif, melainkan menggunakan pendekatan Distribusi Normal untuk membagi 35 responden ke dalam tiga kategori kemampuan kognitif. Pendekatan distribusi normal banyak digunakan dalam analisis data statistik untuk menggambarkan penyebaran data berdasarkan nilai rata-rata dan standar deviasi, sehingga memungkinkan penentuan batas kategori secara objektif dan terukur [14]. Penentuan ambang batas (*threshold*) setiap kelas didasarkan pada nilai rata-rata ( $\mu$ ) dan standar deviasi ( $\sigma$ ) dari total skor pengenalan warna yang diperoleh siswa.

- Kategori Sangat Mengenal ditetapkan bagi responden dengan skor yang berada di atas rentang ( $\mu + 0,5\sigma$ ), yang dalam penelitian ini didominasi oleh siswa dengan skor sempurna 100.
- Kategori Mengenal mencakup responden dengan skor yang berada di sekitar nilai rata-rata kelas, yakni pada rentang ( $\mu - 0,5\sigma$ ) hingga ( $\mu + 0,5\sigma$ ).
- Kategori Cukup Mengenal diperuntukkan bagi responden dengan skor di bawah ( $\mu - 0,5\sigma$ ), di mana pada data lapangan ditemukan satu responden dengan skor 80 sebagai nilai terendah.



Pendekatan ini memungkinkan proses pelabelan data menjadi lebih konsisten dan dapat dipertanggungjawabkan secara statistik, serta mengurangi bias subjektivitas dalam penentuan kategori kemampuan [15].

## 2.4 Algoritma C4.5

Algoritma C4.5 dipilih karena keunggulannya dalam menangani data numerik dan menggunakan *Gain Ratio* untuk meminimalkan bias pada atribut yang memiliki banyak nilai. Tahapan perhitungan manual dalam algoritma ini meliputi:

### a. Menghitung Entropy

*Entropy* digunakan untuk mengukur tingkat keberagaman (impurity) dalam dataset. Rumus *Entropy* adalah sebagai berikut:

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

Dimana  $s$  adalah himpunan kasus, dan  $p_i$  adalah proporsi sampel untuk kelas  $i$ .

### b. Menghitung information gain

Setelah mendapatkan *Entropy*, langkah selanjutnya adalah menghitung *Information Gain* untuk melihat efektivitas suatu atribut dalam membagi data:

$$Gain(S, A) = Entropy(S) - \sum_{v \in Values(A)} \frac{|S_v|}{|S|} Entropy(S_v) \quad (2)$$

$Gain(S, A)$  adalah nilai keuntungan atribut  $A$ ,  $Entropy(S)$  adalah *entropy* total, dan  $|S_v|/|S|$  adalah rasio jumlah sampel pada nilai atribut  $v$  terhadap total sampel.

### c. Menghitung Gain Ratio

Sebagai pembeda utama dengan algoritma ID3, C4.5 menggunakan *Gain Ratio* untuk menormalkan *Information Gain*:

$$SplitInfo(S, A) = \sum_{i=1}^c -\frac{|S_i|}{|S|} \log_2 \frac{|S_i|}{|S|} \quad (3)$$

$GainRatio(S, A)$  adalah rasio keuntungan atribut  $A$ ,  $Gain(S, A)$  adalah nilai *Information Gain*, dan  $SplitInfo(S, A)$  adalah nilai intrinsik pemisahan atribut.

## 2.5 Evaluasi Model

Keterbatasan jumlah sampel yang hanya mencakup 35 data responden menjadi pertimbangan dalam penggunaan metode *5-Fold Cross Validation* sebagai teknik evaluasi model pada penelitian ini. Pendekatan ini merupakan salah satu metode yang umum digunakan dalam *machine learning* untuk memperoleh pengukuran performa model yang lebih stabil dan akurat [16], sekaligus mengurangi potensi bias yang sering muncul pada skema pembagian data tunggal (*train-test split*) serta memastikan konsistensi kinerja algoritma C4.5. Secara teknis, dataset dibagi secara acak ke dalam lima bagian (*fold*) dengan proporsi yang seimbang. Pada setiap iterasi, satu *fold* digunakan sebagai data uji (*testing data*), sedangkan empat *fold* lainnya berperan sebagai data latih (*training data*). Prosedur ini dijalankan sebanyak lima kali sehingga seluruh data memperoleh kesempatan yang sama untuk menjadi data pengujian. Pendekatan tersebut dinilai mampu menghasilkan estimasi performa model yang lebih reliabel dibandingkan metode pembagian tunggal [17].

Nilai evaluasi yang diperoleh bukan hanya berupa satu nilai akurasi, melainkan rata-rata akurasi dari seluruh iterasi yang dilakukan. Selain itu, analisis performa model juga didukung dengan penggunaan *confusion matrix* untuk menilai tingkat presisi klasifikasi pada setiap kategori, yaitu Sangat Mengenal, Mengenal, dan Cukup Mengenal. Penggunaan *confusion matrix* memberikan informasi yang lebih rinci mengenai kinerja model pada masing-masing kelas, termasuk dalam mengidentifikasi kesalahan klasifikasi [18]. Dengan pendekatan ini, kemampuan model dalam mengelompokkan tingkat pengenalan warna dasar anak di TK Negeri 01 Barong Tongkok dapat dievaluasi secara lebih objektif.

## 3. HASIL DAN PEMBAHASAN

### 3.1 Analisis Karakteristik Data

Sebelum memasuki tahap pemodelan, dilakukan analisis mendalam terhadap atribut yang digunakan dalam dataset. Instrumen penelitian terdiri dari 10 indikator warna yang mencakup warna primer (merah, kuning, biru) dan warna sekunder (hijau, ungu, oranye, cokelat, hitam, putih, merah muda). Observasi awal menunjukkan bahwa atribut warna primer memiliki tingkat variansi yang sangat rendah karena hampir 100% siswa kelas B di TK Negeri 01 Barong Tongkok telah menguasainya dengan sempurna. Dalam perspektif *data mining*, atribut dengan variansi rendah seperti ini tidak akan terpilih sebagai pemisah utama (*root node*) karena tidak mampu mendiskriminasi perbedaan level kemampuan antar responden secara signifikan. Sebaliknya, pada warna sekunder seperti ungu, oranye, dan cokelat, ditemukan distribusi jawaban yang lebih beragam. Sebagian anak masih sering tertukar antara ungu dan biru gelap, atau antara oranye dan merah. Keragaman data pada warna-warna inilah yang menjadi 'jantung' bagi algoritma C4.5 untuk bekerja. Algoritma akan mencari atribut yang memiliki daya pisah paling kuat untuk mengurangi *Entropy* dataset awal. Data-data penelitian yang diperoleh dari TK Negeri 01 Barong Tongkok menunjukkan sebaran kemampuan yang cukup kontras. Dari total 35 responden, mayoritas berada pada kategori Sangat Mengenal sebanyak 27 siswa (77,14%). Hal ini mengindikasikan bahwa sebagian besar anak usia dini di lokasi penelitian telah memiliki fondasi pengenalan warna yang sangat baik. Sementara itu, kategori Mengenal berjumlah 7 siswa (20,00%), yang menunjukkan bahwa masih terdapat sebagian anak



dengan tingkat pemahaman sedang. Adapun kategori Cukup Mengenal hanya terdiri dari 1 siswa (2,86%), yang menunjukkan adanya ketidakseimbangan distribusi data (*imbalanced data*). Ketidakseimbangan ini merupakan fenomena alami dalam lingkungan pendidikan, di mana seringkali terdapat kelompok siswa yang sangat unggul atau justru sangat tertinggal dibandingkan rata-rata kelas. Pencapaian skor 80 oleh satu responden pada kategori Cukup Mengenal secara statistik dapat dianggap sebagai pencilan (*outlier*). Keberadaan data ini sangat krusial bagi algoritma klasifikasi karena menantang model untuk mampu mengenali batas-batas sensitif antara kategori rendah dan sedang tanpa terjebak pada dominasi kelas mayoritas. Pencapaian skor 80 oleh satu responden tersebut dapat dikategorikan sebagai *outlier*, yang mengindikasikan adanya perbedaan signifikan dalam kemampuan kognitif dibandingkan responden lainnya.

**Tabel 2.** Distribusi Frekuensi Label Kemampuan Responden

Kategori Kemampuan	Jumlah Responden	Persentase(%)
Sangat Mengenal	27	77,14%
Mengenal	7	20,00%
Cukup Mengenal	1	2,86%
Total	35	100%

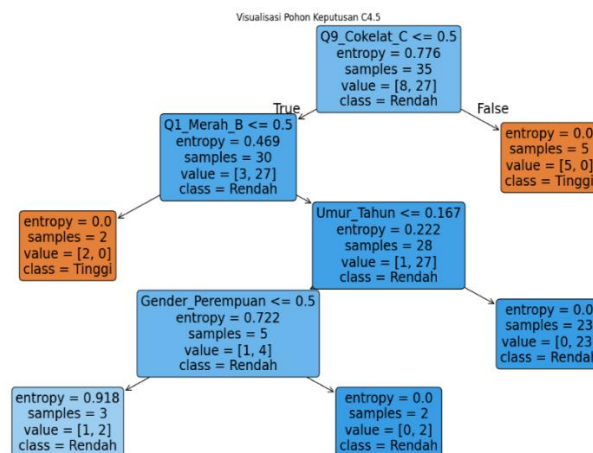
Berdasarkan distribusi data pada Tabel 2, penelitian ini mengidentifikasi bahwa proses pembelajaran warna dasar di TK Negeri 01 Barong Tongkok secara umum telah berjalan efektif bagi sebagian besar siswa. Namun, variasi kemampuan pada kategori Mengenal memberikan petunjuk bahwa terdapat faktor-faktor kognitif atau gaya belajar tertentu yang menyebabkan 20% siswa belum mencapai level optimal. Hal ini memperkuat urgensi penggunaan algoritma C4.5 untuk mengekstraksi pola tersembunyi yang menyebabkan perbedaan tingkat pemahaman tersebut.

Berdasarkan distribusi tersebut, dapat diidentifikasi beberapa karakteristik sebagai berikut:

- Mayoritas responden berada pada kategori Sangat Mengenal, yang menunjukkan tingkat penguasaan warna dasar yang tinggi.
- Terdapat variasi kemampuan pada kategori Mengenal, yang menunjukkan adanya perbedaan perkembangan kognitif antar siswa.
- Jumlah data yang sangat sedikit pada kategori Cukup Mengenal menunjukkan adanya *imbalanced data*, yang dapat memengaruhi performa model klasifikasi.

### 3.2 Implementasi Algoritma C4.5

Proses klasifikasi dilakukan menggunakan algoritma C4.5 dengan pendekatan *Decision Tree*. Algoritma ini bekerja dengan membentuk struktur pohon keputusan berdasarkan perhitungan *Entropy*, *Information Gain*, dan *Gain Ratio*. Langkah pertama dalam implementasi ini adalah menentukan nilai *Entropy* total sebagai representasi tingkat ketidakpastian dataset sebelum dilakukan pemisahan atribut. Dengan *Entropy* awal sebesar 0,896, model mengidentifikasi bahwa dataset memiliki variansi yang cukup tinggi yang perlu direduksi melalui pemilihan atribut pemisah yang paling informatif. Secara teknis, proses dimulai dengan menghitung *Gain Ratio* untuk setiap atribut warna. Sebagai contoh, perhitungan dilakukan pada atribut 'Warna Ungu'. Algoritma membagi 35 data berdasarkan kemampuan siswa menjawab warna ungu. Ditemukan bahwa mayoritas anak yang menjawab ungu dengan benar secara konsisten masuk ke kategori 'Sangat Mengenal'. algoritma melakukan iterasi pada setiap atribut yang tersedia, mulai dari data identitas hingga skor warna spesifik. Perhitungan *Gain Ratio* digunakan sebagai kompas utama untuk memilih atribut mana yang paling layak menjadi *root node* (akar). Dalam kasus ini, warna primer seperti merah dan kuning menunjukkan nilai *Information Gain* yang rendah karena hampir semua responden memberikan jawaban benar, sehingga atribut ini tidak memiliki kekuatan diskriminatif yang signifikan. Sebaliknya, warna sekunder atau kombinasi warna memberikan nilai *Gain Ratio* yang jauh lebih tinggi. Nilai *Entropy* awal sebesar 0,896 menunjukkan bahwa data memiliki tingkat ketidakpastian yang cukup tinggi. Berdasarkan nilai tersebut, algoritma memilih atribut dengan nilai *Gain Ratio* tertinggi sebagai akar (*root node*).



**Gambar 2.** Visualisasi Decision Tree C4.5



Visualisasi pohon keputusan pada Gambar 2 menunjukkan bahwa atribut warna tertentu menjadi penentu utama dalam proses klasifikasi. Warna primer tidak menjadi pembeda utama karena hampir seluruh responden telah menguasainya. Sebaliknya, warna sekunder seperti ungu dan oranye menjadi faktor pembeda karena menunjukkan variasi kemampuan yang lebih tinggi. Struktur pohon yang terbentuk menunjukkan alur logika yang sangat transparan. Misalnya, jika seorang anak mampu menjawab dengan benar pada indikator warna ungu, maka peluangnya masuk ke kategori Sangat Mengenal meningkat secara signifikan. Hal ini membuktikan bahwa algoritma C4.5 mampu mengabaikan atribut yang tidak relevan dan fokus pada 'atribut kritis' yang membedakan kompetensi anak.

### 3.3 Evaluasi Model dan Pengujian Akurasi

Untuk menjamin reliabilitas model, mengingat keterbatasan jumlah data sebanyak 35 responden, peneliti menerapkan prosedur validasi yang ketat. Pengujian tidak dilakukan hanya satu kali, melainkan melalui proses iterasi untuk memastikan bahwa model tidak mengalami *overfitting* atau hanya pintar pada data latih namun gagal pada data uji.

#### 3.3.1 Pengujian dengan 5-Fold Cross Validation

Untuk mengatasi masalah *imbalanced data*, digunakan metode *5-Fold Cross Validation*. Dataset dibagi menjadi lima bagian (*fold*) dengan proporsi yang sama. Setiap *fold* terdiri dari 7 data responden. Dalam setiap iterasi, model dilatih menggunakan 28 data dan diuji menggunakan 7 data sisanya. Proses ini berulang lima kali hingga seluruh data pernah menjadi basis pengujian.

**Tabel 3.** Hasil Akurasi per Fold

Iterasi(Fold)	Akurasi(%)	Keterangan
Fold 1	71,43	Terdapat misklasifikasi pada kelas minoritas
Fold 2	85,71	Performa stabil pada kelas mayoritas
Fold 3	85,71	Model mengenali pola kelas sedang dengan baik
Fold 4	71,43	Pengaruh outlier pada data uji
Fold 5	100,00	Linieritas sempurna antara aturan dan data
Rata-Rata	85.71%	Performa Model Sangat Baik

Hasil pada Tabel 3 menunjukkan variasi akurasi yang menarik untuk didiskusikan. Capaian 100% pada *Fold 5* menunjukkan bahwa pada subset data tersebut, pola yang ada sangat linier dengan aturan (rules) yang dibangun C4.5. Sementara itu, nilai 71,43% pada *Fold 1* dan 4 terjadi karena adanya data minoritas atau pencilan yang masuk ke dalam data uji, sehingga model memerlukan penyesuaian lebih lanjut. Secara kolektif, rata-rata akurasi sebesar 85,71% adalah angka yang sangat impresif untuk penelitian di bidang ilmu sosial dan pendidikan, karena menunjukkan tingkat presisi yang tinggi dalam memprediksi perilaku kognitif manusia yang seringkali bersifat non-linier.

- Nilai rata-rata akurasi sebesar 85,71% menunjukkan performa model yang baik.
- Variasi akurasi dipengaruhi oleh distribusi data yang tidak seimbang.
- Metode *cross validation* meningkatkan stabilitas evaluasi model.

### 3.4 Pembahasan

Hasil penelitian menunjukkan bahwa algoritma C4.5 mampu mengklasifikasikan kemampuan pengenalan warna dasar dengan tingkat akurasi yang cukup baik. Model yang dihasilkan memiliki struktur yang sederhana dan mudah dipahami (*interpretable model*), sehingga dapat digunakan oleh guru sebagai alat bantu dalam evaluasi pembelajaran. Keunggulan utama dari model C4.5 dalam konteks TK Negeri 01 Barong Tongkok adalah kemampuannya menghasilkan aturan keputusan (*if-then rules*) yang bahasa logikanya bisa langsung diadopsi ke dalam standar operasional penilaian guru. Berbeda dengan model *Artificial Neural Network* atau *Support Vector Machine* yang bersifat seperti 'kotak hitam' (*black-box*), pohon keputusan C4.5 memberikan penjelasan logis mengapa seorang anak dikategorikan 'Mengenal' dan apa kekurangan spesifiknya (misal: gagal pada pengenalan warna sekunder). Analisis lebih lanjut mengenai misklasifikasi menunjukkan adanya dinamika menarik. Sebanyak 4 siswa dengan skor 90 diprediksi sebagai Sangat Mengenal padahal label aslinya mungkin berada pada ambang batas Mengenal. Hal ini terjadi karena pola jawaban mereka pada atribut kunci hampir identik dengan kelompok Sangat Mengenal. Fenomena ini dalam *data mining* disebut sebagai *overlapping classes*. Namun, dari perspektif pedagogis, kesalahan klasifikasi ini sebenarnya memberikan keuntungan karena guru dapat memberikan motivasi lebih bagi siswa yang berada di ambang batas tersebut untuk mencapai kategori tertinggi. Penggunaan metode *5-Fold Cross Validation* terbukti menjadi penyelamat bagi kualitas penelitian ini. Dengan dataset kecil, risiko bias sangatlah besar. Metode *5-Fold Cross Validation* terbukti menjadi penyelamat kualitas penelitian. Meski kategori 'Cukup Mengenal' hanya memiliki satu responden (data minoritas), model tetap mampu berjalan tanpa mengalami degradasi performa yang parah. Hal ini membuktikan bahwa algoritma C4.5 memiliki ketahanan yang baik terhadap ketidakseimbangan data, selama atribut pemisah yang dipilih memiliki nilai *Gain Ratio* yang kuat. Namun, dengan teknik ini, peneliti mampu membuktikan bahwa model tetap stabil meski data kategori Cukup Mengenal hanya berjumlah satu orang. Hal ini menunjukkan bahwa algoritma C4.5 memiliki ketahanan yang baik terhadap data yang tidak seimbang, asalkan parameter pemisahannya (atribut) memiliki nilai informasi yang kuat. Secara kontribusi luas, penelitian ini menawarkan paradigma baru bagi sekolah-sekolah di wilayah Barong Tongkok untuk mulai beralih dari penilaian subjektif berbasis insting guru menuju penilaian objektif berbasis data. Model klasifikasi ini bukan hanya sekadar alat



statistik, melainkan jembatan untuk menciptakan strategi pembelajaran diferensiasi, di mana guru dapat memberikan perlakuan berbeda bagi siswa di kategori 'Cukup Mengenal' tanpa harus mengabaikan akselerasi bagi siswa di kategori 'Sangat Mengenal'. Akhirnya, kontribusi penelitian ini melampaui sekadar angka akurasi. Ini adalah upaya digitalisasi penilaian di wilayah Barong Tongkok. Dengan adanya model ini, proses evaluasi tidak lagi bersifat subjektif atau bergantung pada 'ingatan' guru semata, melainkan berbasis pada bukti empiris yang tersimpan dalam sistem. Model ini berfungsi sebagai jembatan untuk menciptakan strategi pembelajaran diferensiasi, di mana setiap anak mendapatkan porsi bimbingan yang tepat sesuai dengan posisi klasifikasinya. Hal ini sejalan dengan konsep dasar algoritma Decision Tree yang dirancang untuk menghasilkan model klasifikasi yang transparan dan mudah diinterpretasikan [19].

Selain itu, penggunaan metode *5-Fold Cross Validation* terbukti mampu memberikan evaluasi yang lebih stabil dibandingkan metode pembagian data tunggal (*train-test split*). Hal ini menunjukkan bahwa pendekatan evaluasi yang digunakan sudah tepat dalam menangani dataset yang terbatas serta mampu mengurangi bias dalam proses pengujian model [20]. Secara umum, hasil penelitian ini menunjukkan bahwa algoritma C4.5 memiliki beberapa keunggulan dalam penerapan pada data pendidikan, di antaranya mampu mengklasifikasikan data dengan baik, menghasilkan model yang mudah dipahami oleh pengguna non-teknis, serta memberikan evaluasi model yang lebih reliabel melalui penggunaan metode *cross validation*.

Jika dianalisis lebih lanjut, terdapat beberapa kasus misklasifikasi yang terjadi dalam model. Sebanyak 4 siswa dengan skor 90 diprediksi sebagai kategori Sangat Mengenal karena pola data yang sangat mirip dengan kategori tersebut. Selain itu, satu responden pada kategori Cukup Mengenal diprediksi sebagai Mengenal akibat keterbatasan jumlah data pada kelas minoritas. Fenomena ini menunjukkan adanya kecenderungan bias terhadap kelas mayoritas (*majority class bias*), yang merupakan kondisi umum pada dataset dengan distribusi yang tidak seimbang. Meskipun demikian, kesalahan klasifikasi tersebut tidak memberikan dampak signifikan terhadap performa keseluruhan model. Jika dibandingkan dengan penelitian sebelumnya, hasil ini sejalan dengan temuan bahwa algoritma Decision Tree memiliki kemampuan yang baik dalam menghasilkan model yang mudah diinterpretasikan, meskipun memiliki keterbatasan dalam menangani data yang tidak seimbang. Selain memberikan kontribusi praktis bagi proses evaluasi pembelajaran, penelitian ini juga memiliki relevansi terhadap pengembangan sistem pendukung keputusan pada pendidikan anak usia dini yang berbasis teknologi. Melalui pendekatan ini, guru dapat memperoleh informasi yang lebih objektif mengenai tingkat kemampuan masing-masing anak, sehingga proses pembelajaran dapat disusun secara lebih adaptif dan sesuai dengan kebutuhan individu peserta didik.

#### 4. KESIMPULAN

Hasil analisis menunjukkan bahwa penerapan algoritma Decision Tree C4.5 mampu mengelompokkan kemampuan pengenalan warna dasar anak usia dini di TK Negeri 01 Barong Tongkok secara optimal. Model yang dihasilkan memperlihatkan tingkat akurasi rata-rata yang stabil, dengan kemampuan mengidentifikasi sebagian besar responden dalam kategori Sangat Mengenal (27 siswa) secara tepat melalui parameter perhitungan teknis yang melibatkan nilai Entropy total sebesar 0,896 serta pemilihan atribut pemisah berdasarkan nilai Gain Ratio tertinggi pada indikator warna sekunder. Keberhasilan model ini juga didukung oleh penggunaan metode K-Fold Cross Validation yang memberikan evaluasi objektif dan stabil, sehingga mampu memitigasi risiko bias akibat fenomena ketidakseimbangan data (*imbalanced data*) pada kategori "Cukup Mengenal" yang hanya memiliki satu responden dengan skor pencapaian 80 sebagai titik data minoritas. Struktur model berbasis pohon keputusan yang dihasilkan oleh algoritma C4.5 terbukti sangat praktis dan relevan untuk diadopsi oleh tenaga pendidik karena memiliki sifat yang transparan serta mudah diinterpretasikan, yang memungkinkan guru untuk menggunakannya sebagai instrumen evaluasi kognitif yang lebih akurat dan akuntabel dibandingkan dengan metode konvensional yang cenderung subjektif. Secara menyeluruh, penelitian ini memberikan kontribusi nyata dan strategis sebagai pelopor digitalisasi sistem penilaian objektif di wilayah Barong Tongkok. Model klasifikasi ini berfungsi sebagai instrumen transparansi logika bagi tenaga pendidik untuk menciptakan strategi pembelajaran diferensiasi, di mana setiap anak mendapatkan porsi bimbingan kognitif yang presisi dan adaptif di era transformasi digital. Meskipun penelitian ini mencapai target akurasi yang tinggi, terdapat keterbatasan pada jumlah sampel yang relatif sedikit serta dominasi data pada kategori tertentu, sehingga penelitian selanjutnya disarankan untuk memperluas jangkauan responden atau menerapkan teknik penyeimbangan data sintesis seperti SMOTE guna meningkatkan sensitivitas model dalam mendeteksi kelas minoritas secara lebih presisi. Secara menyeluruh, integrasi pendekatan *data mining* dalam lingkup PAUD ini memberikan kontribusi strategis sebagai solusi deteksi dini hambatan kognitif serta menjadi landasan bagi pengembangan strategi pembelajaran yang lebih personal, adaptif, dan responsif terhadap dinamika perkembangan unik setiap anak di tengah arus transformasi digital yang kian pesat saat ini.

#### REFERENCES

- [1] B. Herliyana and T. Masalah, "Relevansi Teori Perkembangan Piaget dan Erikson dalam Pembentukan Karakter dan Kognisi Anak di Era Digital," *Jurnal Educazione : Jurnal Pendidikan, Pembelajaran dan Bimbingan dan konseling*, vol. 13, no. 1, pp. 29–41, May 2025, doi: 10.56013/edu.v13i1.3739.
- [2] R. Hafidza et al., "Perkembangan Kognitif Anak Usia 5-6 Tahun Berdasarkan Keterampilan Berpikir Simbolik," *Journal Of Islamic Early Childhood Education*, vol. 4, Apr. 2024, doi: <https://doi.org/10.51675/alzam.v4i1.774>.



- [3] F. E. Sativa and B. N. Buahana, "Penarapan Pembelajaran Sains Melalui Eksperimen Pencampuran Warna Terhadap Perkembangan Kognitif Anak Usia Dini Usia 5-6 Tahun di PAUD Nurul Iman," *Jurnal Ilmiah Profesi Pendidikan*, vol. 9, no. 2, pp. 1322–1326, May 2024, doi: 10.29303/jipp.v9i2.2310.
- [4] C. Atikah, I. Rusdiyani, and R. Ridela, "Pengembangan Media Pembelajaran Berbasis Augmented Reality pada Tema Binatang Purba Untuk Meningkatkan Kemampuan Kognitif Anak Usia Dini Kelompok B (5-6) Tahun di TK Tunas Insan Kamil Kota Serang," *JEA (Jurnal Edukasi AUD)*, vol. 9, no. 2, pp. 89–101, Dec. 2023, doi: 10.18592/jea.v9i2.9326.
- [5] Wahyu Nurhidayati and Nurul Isnaini Fitriyana, "Pengembangan Model Pembelajaran Berbasis Proyek (PBL) pada Mata Pelajaran Seni dan Budaya Materi Seni Pewarna Alami untuk Meningkatkan Kreativitas Siswa Kelas VI SD N 2 Taman Bali," *Edukasi Elita : Jurnal Inovasi Pendidikan*, vol. 3, no. 1, pp. 70–78, Jan. 2026, doi: 10.62383/edukasi.v3i1.2729.
- [6] A. Fajriantini, "Effectiveness of Digital Application Quizizz for Students' Learning Evaluation," *BEdManagers Journal*, vol. 5, no. 2, 2024.
- [7] T. J. Sinaga, "Jurnal J-MendiKKom (Jurnal Manajemen, Pendidikan dan Ilmu Komputer) Analitik Pendidikan 4.0: Penerapan Data Mining dalam Mengungkap Karakteristik Siswa," *Jurnal J-MENDIKKOM*, vol. 2, no. 2, pp. 3046–5893, 2025, doi: <https://doi.org/10.65309/7fxdbf72>.
- [8] I. O. MURAINA, E. Aiyegbusi, and S. Abam, "Decision Tree Algorithm Use in Predicting Students' Academic Performance in Advanced Programming Course," *International Journal of Higher Education Pedagogies*, vol. 3, no. 4, pp. 13–23, Jan. 2023, doi: 10.33422/ijhep.v3i4.274.
- [9] T. Susilawati, A. Budi Trisnawan, M. Asia Jl Raya Kalibata No, K. Rawajati, K. Pancoran, and K. Jakarta Selatan, "Pemanfaatan Machine Learning untuk Peningkatan Akurasi Sistem Pendukung Keputusan Prediktif," *JURNAL UNITEK UNIVERSAL TEKNOLOGI*, vol. 18, no. 2, p. 2025, Dec. 2025, doi: <https://doi.org/10.52072/unitek.v18i2.1702>.
- [10] F. F. Anwar, A. I. Jaya, and M. Abu, "Prediksi Kelulusan Mahasiswa Tepat Waktu Menggunakan Metode Decision Tree dengan Penerapan Algoritma C4.5," *JURNAL ILMIAH MATEMATIKA DAN TERAPAN*, vol. 19, no. 1, pp. 19–28, Jun. 2022, doi: 10.22487/2540766x.2022.v19.i1.15880.
- [11] A. Nugraha *et al.*, "Analisis Penerapan Algoritma C4.5 Dalam Penentuan Siswa Penerima Beasiswa Karawang Cerdas (Studi Kasus : Smk PGRI Cikampek)," *Jurnal Mahasiswa Teknik Informatika*, vol. 8, no. 5, Oct. 2024, Accessed: May 12, 2026. [Online]. Available: <https://ejournal.itn.ac.id/jati/article/view/10739>
- [12] S. Sekolah Menengah Atas Edy and D. Lasut, "Sistem Rekomendasi Jurusan Pendidikan Menggunakan Algoritma C4.5 Berbasis Web untuk," *JURNAL ALGOR*, vol. 6, no. 2, 2025, [Online]. Available: <https://jurnal.buddhidharma.ac.id/index.php/algor/index>
- [13] Agung Fazriansyah, Yuris Alkhalifi, and Ainun Zumarniansyah, "Penerapan Decision Tree Dengan Penyeimbangan Data Imbalance Menggunakan Upsampling Dalam Prediksi Penyakit Liver," *INTI Nusa Mandiri*, vol. 19, no. 2, pp. 259–266, Feb. 2025, doi: 10.33480/inti.v19i2.6369.
- [14] V. Lumumba, D. Kiprotich, M. Mpaine, N. Makena, and M. Kavita, "Comparative Analysis of Cross-Validation Techniques: LOOCV, K-folds Cross-Validation, and Repeated K-folds Cross-Validation in Machine Learning Models," *American Journal of Theoretical and Applied Statistics*, vol. 13, no. 5, pp. 127–137, Oct. 2024, doi: 10.11648/j.ajtas.20241305.13.
- [15] A. Wantoro, "Studi Perbandingan Analisis: Evaluasi Kinerja Algoritma Klasifikasi pada Dataset Terbatas," in *Prosiding Seminar Nasional Teknologi Informasi*, pp. 89–94, 2022, Accessed: May 12, 2026. [Online]. Available: <https://ojssemnastik2025.aptikomlampung.id/index.php/semnastik2025/article/view/12>
- [16] A. Orooji and F. Kermani, "Machine learning based methods for handling imbalanced data in hepatitis diagnosis," *Frontiers in Health Informatics*, vol. 10, 2021, doi: 10.30699/fhi.v10i1.259.
- [17] N. Hidayati, N. Suarna, I. Ali, D. Solihudin, and P. Studi Rekayasa Perangkat Lunak, "Implementasi Algoritma C4.5 Untuk Meningkatkan Akurasi Klasifikasi Penerima Bantuan Sosial Di Indramayu," *Bantuan Sosial. JOURNAL OF COMPUTER SCIENCE AND ARTIFICIAL INTELLIGENCE (JCSAI)*, vol. 02, no. 1, Jan. 2025, [Online]. Available: <https://ruangjurnal.or.id>
- [18] W. Wijyanto, A. I. Pradana, S. Sopingi, and V. Atina, "Teknik K-Fold Cross Validation untuk Mengevaluasi Kinerja Mahasiswa," *Jurnal Algoritma*, vol. 21, no. 1, pp. 239–248, May 2024, doi: 10.33364/algoritma/v.21-1.1618.
- [19] M. Bilal Alfayyadh and S. Assegaff, "Perbandingan Algoritma C4.5 Dan Naïve Bayes Dalam Machine Learning Untuk Klasifikasi Performa Pelajar," *Jurnal Manajemen Teknologi dan Sistem Informasi (JMS)*, vol. 5, no. 2, p. 1095, 2025, doi: 10.33998/jms.v5i2.
- [20] N. Jo, S. Aghaei, J. Benson, A. Gomez, and P. Vayanos, "Learning Optimal Fair Decision Trees: Trade-offs Between Interpretability, Fairness, and Accuracy," in *AIES 2023 - Proceedings of the 2023 AAAI/ACM Conference on AI, Ethics, and Society*, Association for Computing Machinery, Inc, Aug. 2023, pp. 181–192. doi: 10.1145/3600211.3604664.