

# Algoritma Pohon Keputusan Termodifikasi sebagai Alat Imputasi Multivariat Data Nonlinear

Muhammad Fathur Rizky - 13523105<sup>1</sup>

Program Studi Teknik Informatika

Sekolah Teknik Elektro dan Informatika

Institut Teknologi Bandung, Jl. Ganesha 10 Bandung 40132, Indonesia

[rizkyfathur326@gmail.com](mailto:rizkyfathur326@gmail.com), [13523105@std.stei.itb.ac.id](mailto:13523105@std.stei.itb.ac.id)

**Abstract**—Data imputation plays a crucial role in data processing to address missing or incomplete data effectively. Missing data can compromise the accuracy of analyses and predictive models, highlighting the need for reliable imputation techniques to ensure robust results. Decision trees, widely used in predictive modeling, excel in capturing nonlinear relationships. This study demonstrates that decision tree algorithms are effective in handling missing data, achieving an average RMSE of 4.98, while offering a competitive time complexity of  $O(d^2n^2 \log n)$  compared to simpler imputation methods such as mean, median, and mode.

**Keywords**—Data Imputation, Decision Tree, Predictive Modeling.

## I. PENDAHULUAN

Imputasi data adalah salah satu langkah penting dalam pengolahan data yang digunakan untuk mengatasi masalah data yang hilang atau tidak lengkap. Kehilangan data dapat mempengaruhi akurasi analisis dan model prediksi sehingga imputasi yang efektif diperlukan untuk memastikan validitas hasil analisis. Salah satu metode yang sering digunakan dalam analisis prediktif adalah pohon keputusan yang dapat menangani data dengan hubungan nonlinear dan kompleks. Meskipun demikian, pada implementasi standar pohon keputusan, nilai *null* atau data yang hilang tidak diimplementasikan sehingga diperlukan teknik imputasi di luar proses pemodelan.

Pada makalah ini, dilakukan modifikasi pada algoritma pohon keputusan untuk memperhitungkan nilai *null* dalam proses pemisahan data. Modifikasi ini bertujuan untuk mengatasi kelemahan dalam pemisahan yang tidak memperhitungkan data hilang sehingga proses imputasi dapat dilakukan dengan lebih efektif. Dalam modifikasi ini, nilai *null* diperlakukan secara eksplisit dalam pemisahan data. Nilai hilang dipertimbangkan bersama dengan data yang memenuhi kriteria pemisahan tertentu. Dengan pendekatan ini, pohon keputusan dapat menghasilkan pemisahan yang lebih baik, termasuk untuk data yang memiliki nilai hilang, dan memungkinkan imputasi yang lebih akurat.

Pendekatan ini diharapkan dapat meningkatkan kinerja pohon keputusan dalam imputasi data sehingga model dapat mengisi nilai yang hilang dengan lebih tepat. Dengan demikian, modifikasi ini tidak hanya meningkatkan akurasi imputasi, tetapi juga memperluas kemampuan pohon keputusan dalam menangani dataset yang memiliki sejumlah nilai hilang dan

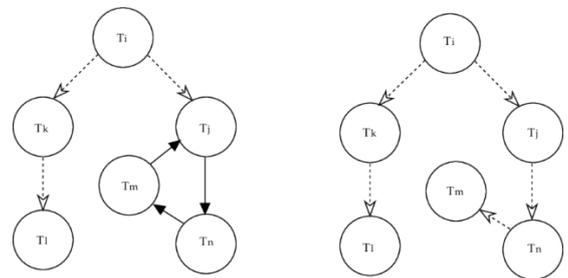
meningkatkan kualitas data yang digunakan dalam analisis lebih lanjut.

## II. DASAR TEORI

### A. Graf

Graf adalah suatu struktur data diskret yang terdiri atas simpul-simpul dan sisi-sisi. Graf dinotasikan sebagai  $G = (V, E)$  dengan  $V$  merupakan simpul (*vertex*) dan  $E$  merupakan tepi (*edge*) [1]. Graf digunakan untuk menggambarkan hubungan antara objek-objek yang memiliki karakteristik sebagai berikut.

- **Simpul:** Simpul adalah elemen dasar dalam graf yang merepresentasikan titik atau objek dalam suatu sistem. Setiap simpul biasanya diberi label atau nomor untuk membedakan satu simpul dengan yang lainnya.
- **Tepi:** Tepi adalah hubungan atau penghubung antara dua simpul. Setiap sisi menghubungkan sepasang simpul dan dapat berupa sisi tak terarah (*undirected*) atau sisi berarah (*directed*) tergantung pada jenis graf.
- **Siklus:** Siklus adalah jalur yang dimulai dan berakhir pada simpul yang sama, tanpa melewati sisi atau simpul lebih dari sekali. Graf yang mengandung siklus disebut graf berbentuk siklus, sedangkan graf yang tidak mengandung siklus disebut graf asiklik.



Gambar 1. Graf Siklis dan Asiklis

Sumber:

- **Arah:** Pada graf berarah (*digraph*), sisi memiliki arah, yang berarti hubungan antar simpul bersifat satu arah. Sedangkan pada graf tak terarah, sisi tidak memiliki arah, yang berarti hubungan antar simpul bersifat dua arah.

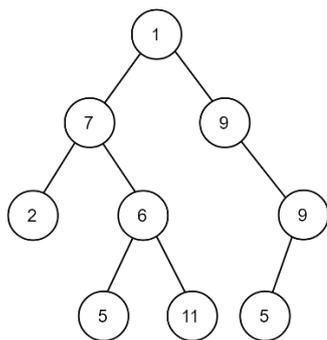
### B. Pohon

Pohon adalah salah satu struktur data yang berbasis graf dengan karakteristik khusus berupa hubungan hierarkis antara simpul-simpulnya. Secara formal, pohon dapat didefinisikan sebagai graf berarah tanpa siklus (*Directed Acyclic*

Graph/DAG) yang memiliki satu simpul utama sebagai akar (root) dan cabang-cabang yang menghubungkan simpul lainnya. Simpul-simpul dalam pohon dibedakan menjadi beberapa jenis, yaitu akar (simpul paling atas yang menjadi titik awal struktur pohon), simpul daun (*leaf*) yang tidak memiliki anak, dan simpul internal yang memiliki anak namun bukan akar [2].

Komponen utama dari pohon meliputi simpul (*node*), tepi (*edge*), tingkat (*level*), kedalaman (*depth*), tinggi (*height*), dan subpohon (*subtree*). Tepi adalah hubungan antara dua simpul yang menggambarkan keterkaitan hierarkis antara orang tua (*parent*) dan anak (*child*). Tingkat (*level*) dihitung berdasarkan jarak suatu simpul dari akar, sedangkan kedalaman (*depth*) menunjukkan jarak dari akar ke simpul tertentu.

Salah satu varian dari pohon yang banyak digunakan adalah pohon biner. Pohon biner adalah jenis pohon dengan setiap simpul hanya dapat memiliki paling banyak dua anak, yang disebut sebagai anak kiri (*left child*) dan anak kanan (*right child*).



Gambar 2. Pohon Biner

### C. Pohon Keputusan (Decision Tree)

Pohon keputusan adalah salah satu algoritma yang digunakan dalam pembelajaran mesin (*machine learning*) untuk tugas klasifikasi dan regresi. Algoritma ini membangun model berbentuk pohon yang digunakan untuk membuat keputusan berdasarkan fitur-fitur input. Setiap simpul dalam pohon mewakili keputusan atau pertanyaan mengenai atribut atau fitur, dan cabang-cabangnya mewakili hasil dari keputusan tersebut. Pada akhirnya, setiap daun pohon mewakili kelas atau nilai yang diprediksi [3].

Proses pembangkitan pohon keputusan dimulai dengan pemilihan atribut yang paling baik untuk membagi data pada setiap simpul. Pembagian ini dilakukan secara rekursif hingga seluruh data terklasifikasi atau mencapai kriteria berhenti tertentu. Pohon keputusan memiliki komponen utama sebagai berikut.

- Variansi: Variansi digunakan untuk mengukur variabilitas dalam data. Semakin tinggi variansi, semakin luas persebaran data tersebut. Variansi dapat dirumuskan sebagai berikut.

$$\text{Var}(D) = \frac{1}{|D|} \sum_{i=1}^{|D|} (y_i - \bar{y})^2 \quad (1)$$

dengan  $y_i$  adalah nilai target pada data ke- $i$ ,  $|D|$  adalah jumlah sampel dalam dataset  $D$ , dan  $\bar{y}$  adalah nilai rata-rata dalam dataset.

- Information gain: Information gain adalah ukuran yang

digunakan untuk memilih atribut yang paling baik dalam membagi dataset pada simpul tertentu. Information Gain dihitung berdasarkan pengurangan *parent variance* dengan *child variance*. Atribut dengan Information Gain tertinggi akan dipilih sebagai pembagi pada simpul tersebut. Secara matematis, information gain dirumuskan sebagai berikut.

$$\text{IG}_{\text{var}} = \text{Var}(D) - \left( \frac{|D_L|}{|D|} \text{Var}(D_L) + \frac{|D_R|}{|D|} \text{Var}(D_R) \right) \quad (2)$$

Dengan  $\text{Var}(D)$  adalah varians sebelum pembagian,  $\frac{|D_L|}{|D|} \text{Var}(D_L)$  dan  $\frac{|D_R|}{|D|} \text{Var}(D_R)$  adalah varians rata-rata tertimbang dari subset data setelah pemisahan, dan  $\text{IG}_{\text{var}}$  adalah selisih antara varians sebelum pemisahan dan setelah pemisahan. Semakin besar nilai IG, semakin baik pemisahan tersebut dalam mengurangi ketidakpastian dalam dataset.

Kompleksitas algoritma decision tree dapat dianalisis berdasarkan proses pembentukannya. Dalam mencari pemisah terbaik pada setiap simpul, algoritma melakukan iterasi sebanyak  $d$  kali untuk setiap fitur. Pada setiap iterasi tersebut, algoritma akan mengiterasi nilai unik pada sampel sebanyak  $n$  kali, dimana untuk setiap nilai unik dibutuhkan perhitungan entropi dengan kompleksitas  $O(n)$ . Dengan demikian, total kompleksitas untuk mencari pemisah terbaik pada satu simpul adalah  $O(dn^2)$ . Proses ini kemudian dilakukan secara rekursif dengan membagi sampel menjadi dua bagian pada setiap level pohon. Dalam kasus tree yang seimbang, jumlah level maksimum yang terbentuk adalah  $\log_2 n$ . Oleh karena itu, total kompleksitas waktu untuk pembangkitan pohon keputusan adalah  $O(dn^2 \log_2 n)$  dengan  $d$  merepresentasikan jumlah fitur dan  $n$  merepresentasikan jumlah sampel dalam dataset.

### D. Root Mean Squared Error

Root Mean Squared Error (RMSE) adalah salah satu metrik evaluasi yang sering digunakan untuk mengukur kinerja model regresi. RMSE menghitung rata-rata dari kuadrat selisih antara nilai yang sebenarnya ( $y_i$ ) dan nilai yang diprediksi oleh model ( $\hat{y}_i$ ). Rumus RMSE secara matematis dengan  $n$  adalah jumlah sampel dinyatakan sebagai berikut.

$$\text{RMSE} = \sqrt{\frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n}} \quad (3)$$

## III. IMPLEMENTASI ALAT IMPUTASI DATA NONLINEAR

### A. Penjelasan dan Penyiapan Dataset

Digunakan sebuah dataset performa murid dalam menghadapi ujian. Dataset ini merupakan dataset generatif dan ditujukan untuk kepentingan edukasi. Dataset ini terdiri atas 6607 baris dan 7 kolom tanpa nilai hilang (*null*). Berikut merupakan penjelasan mengenai tiap-tiap kolom pada dataset.

TABEL I  
PENJELASAN DATASET

Nama kolom	Penjelasan
Hours_Studied	Jumlah jam yang dihabiskan untuk belajar per minggu.
Attendance	Persentase jumlah kehadiran

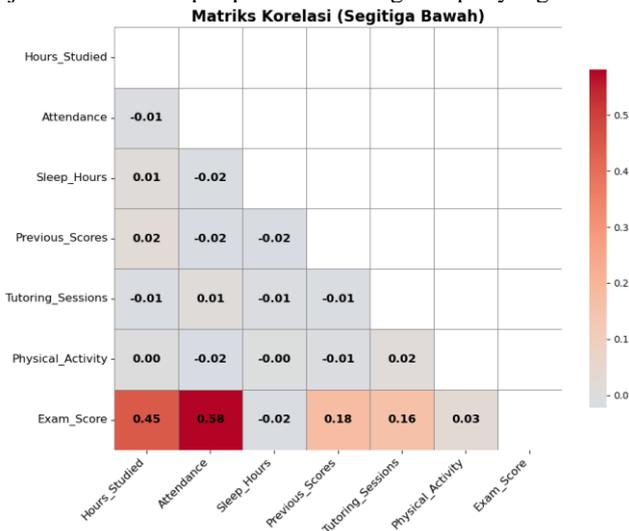
Sleep_Hours	Lama murid tidur (jam)
Previous_Scores	Skor ujian sebelumnya
Tutoring_Sessions	Jumlah sesi bimbingan belajar yang dihadiri per bulan.
Physical_Activity	Rata-rata jumlah aktivitas fisik per minggu
Exam_Score	Nilai ujian akhir

Selain penjelasan, diberikan juga statistik dasar dari setiap kolom yang meliputi nilai minimum, mean, maksimum, dan standar deviasi. Dari tabel di bawah, terlihat bahwa kolom Previous\_Scores mempunyai standar deviasi paling tinggi sehingga ini akan menjadi tantangan yang besar bagi model dalam melakukan imputasi.

TABEL II  
STATISTIK DATASET

Nama kolom	Statistik			
	Min	Mean	Max	Std
Hours_Studied	1	19.97	44	5.99
Attendance	60	79.98	100	11.55
Sleep_Hours	4	7.03	10	1.47
Previous_Scores	50	75.07	100	14.4
Tutoring_Sessions	0	1.49	8	1.23
Physical_Activity	0	2.97	6	1.03
Exam_Score	55	67.24	101	3.89

Untuk membuktikan bahwa data yang digunakan adalah data dengan hubungan nonlinear, maka ditampilkan gambar matriks korelasi menggunakan koefisien pearson sebagai penghitung korelasi. Terlihat bahwa terdapat hubungan cukup linear antara Exam\_Score dengan Hours\_Studied dan Attendance, tetapi tidak linear dengan yang lainnya. Oleh karena itu, data yang diujikan sudah cukup representatif dengan topik yang ditulis.



Gambar 3. Matriks Korelasi (Segitiga Bawah)

Karena topik makalah ini bertujuan menangani kasus data hilang, maka akan diimplementasikan fungsi untuk merandomisasi indeks ( $i, j$ ) agar dipilih sebagai nilai hilang. Penulis mengubah nilai menjadi *null* pada dua fitur untuk setiap barisnya dengan melakukan *setting* parameter `random_state = 42`. Berikut kode untuk penghapusan dalam bahasa Python.

```
def random_nan(df: pd.DataFrame) -> pd.DataFrame:
    df = df.copy()
```

```
for i in range(df.shape[0]):
    j = np.random.choice(df.shape[1], 2,
        replace=False)
    df.iloc[i, j] = np.nan
return df
```

Dari proses tersebut, dihasilkan data terbaru dengan 2/7 merupakan nilai hilang.

### B. Implementasi Pohon Keputusan untuk Regresi

Setelah melakukan penyiapan dataset, selanjutnya adalah membuat algoritma pohon keputusan. Dalam membuat pohon keputusan, penulis menggunakan *library* tambahan, yakni *pandas*, *numpy*, dan *matplotlib* untuk visualisasi. Pertama-tama, dibuat kelas *Node* terlebih dahulu.

```
class Node:
    def __init__(self, f=None, t=None, l=None,
        r=None, v=None):
        self.f = f # fitur
        self.t = t # batas pemisah
        self.l = l # anak kiri
        self.r = r # anak kanan
        self.v = v # nilai pada daun
```

Setelah membuat kelas *Node*, selanjutnya adalah membuat kelas *DecisionTreeRegressor*. Kelas ini memiliki metode inialisasi (`__init__`), `fit` untuk melakukan pembangkitan pohon, dan `predict` yang digunakan dalam melakukan prediksi/inferensi.

```
class DecisionTreeRegressor:
    def __init__(self, min_split=2, min_leaf=1,
        max_depth=None):
        self.min_split = min_split
        self.min_leaf = min_leaf
        self.max_depth = max_depth
        self.root = None

    def fit(self, X, y):
        self.root = self._grow(X, y, d=0)
        return self

    def predict(self, X):
        return np.array([self._traverse(row,
            self.root) for _, row in X.iterrows()])
```

Fungsi berikutnya adalah `_grow` yang bertanggung jawab untuk membangun pohon keputusan secara rekursif. Proses dimulai dengan memeriksa apakah jumlah data pada simpul saat ini kurang dari batas minimum (`min_split`), yang jika benar, fungsi akan berhenti dan mengembalikan simpul daun dengan nilai yang dihitung menggunakan `_leaf_val(y)`. Selanjutnya, fungsi mencari pemisahan terbaik pada dataset dengan memanggil `_best_split(X, y)` yang akan mengembalikan fitur ( $f$ ) dan nilai pemisahan ( $t$ ). Jika tidak ada pemisahan yang dapat dilakukan, fungsi akan mengembalikan simpul daun. Data kemudian dibagi menjadi dua subset: satu untuk nilai fitur yang lebih kecil dari  $t$  atau nilai fitur hilang dan satu lagi untuk nilai yang lebih besar atau sama dengan  $t$ . Setelah pemisahan, fungsi memeriksa apakah jumlah data pada salah satu sisi (kiri atau kanan) kurang dari batas minimal (`min_leaf`). Jika salah satu subset memiliki data yang terlalu sedikit, maka simpul tersebut menjadi daun. Jika kedua subset memiliki cukup data, fungsi `_grow` dipanggil kembali secara rekursif untuk membangun pohon keputusan pada kedua sisi. Pada akhirnya, fungsi ini mengembalikan simpul yang menyimpan informasi tentang fitur pemisah, nilai pemisah, serta anak-anak simpul kiri dan kanan

yang telah dibangun melalui rekursi.

```
def _grow(self, X: pd.DataFrame, y: pd.Series,
d: int):
    n = len(y)

    if n < self.min_split or (self.max_depth is
not None and d >= self.max_depth):
        return Node(v=self._leaf_val(y))

    f, t = self._best_split(X, y)

    if f is None:
        return Node(v=self._leaf_val(y))

    l_idx = X[f].isna() | (X[f] < t)
    r_idx = ~l_idx

    if l_idx.sum() < self.min_leaf or
r_idx.sum() < self.min_leaf:
        return Node(v=self._leaf_val(y))

    l = self._grow(X[l_idx], y[l_idx], d + 1)
    r = self._grow(X[r_idx], y[r_idx], d + 1)
    return Node(f, t, l, r)
```

Fungsi selanjutnya adalah `_traverse` yang digunakan untuk memprediksi nilai `X` berdasarkan pohon keputusan yang sudah dibangkitkan. Proses ini dilakukan dengan membandingkan apakah nilai yang dibandingkan lebih besar, kecil atau tidak mempunyai nilai.

```
def _traverse(self, row: pd.Series, n: Node):
    if n.v is not None: return n.v
    if pd.isna(row[n.f]): # jika nilai hilang
        return self._traverse(row, n.l)
    if row[n.f] < n.t:
        return self._traverse(row, n.l)
    return self._traverse(row, n.r)
```

Fungsi selanjutnya adalah `_best_split` yang digunakan untuk mencari pemisah mana yang terbaik. Fungsi ini mempertimbangkan fungsi `_information_gain` sebagai parameter bagus atau tidaknya. Semakin besar nilai `information gain` untuk fitur yang dipilih, maka fitur tersebut memiliki bobot terbesar dalam pemisahan.

```
def _best_split(self, X: pd.DataFrame, y:
pd.Series):
    best_gain = -np.inf
    best_f = None
    best_t = None
    for f in X.columns:
        thresholds = X[f].dropna().unique()
        for t in thresholds:
            gain = self._gain(y, X[f], t)
            if gain > best_gain:
                best_gain = gain
                best_f = f
                best_t = t
    return best_f, best_t
```

Fungsi berikutnya adalah `_leaf_val`. Fungsi ini mengeluarkan rata-rata nilai target di daun pada saat itu.

```
def _leaf_val(self, y: pd.Series):
    return y.mean()
```

Fungsi yang diberikan digunakan untuk menghitung `information gain` atau peningkatan informasi dalam pembentukan pohon keputusan. Fungsi ini terdiri dari dua bagian utama, yakni `_gain` dan `_variance`.

Fungsi `_variance` menghitung variansi dari data target, yang menggambarkan sejauh mana nilai-nilai target tersebar. Semakin besar variansi, semakin besar ketidakpastian dalam data tersebut.

Fungsi `_gain` menghitung peningkatan informasi yang diperoleh dengan membagi data berdasarkan suatu `threshold (t)` pada salah satu fitur (`X_col`). Langkah pertama adalah menghitung variansi seluruh data target (`p_var`). Data kemudian dibagi menjadi dua subset, berdasarkan apakah nilai fitur lebih kecil atau lebih besar dari `threshold`. Setelah itu, dihitung variansi masing-masing subset. Jika salah satu subset kosong, fungsi mengembalikan nilai `gain 0` karena pembagian tersebut tidak memberikan informasi yang berguna. Jika kedua subset berisi data, dihitung rata-rata berimbang dari variansi kedua subset. `Information gain` dihitung dengan mengurangi variansi gabungan subset dari variansi awal. Nilai `gain` menunjukkan seberapa baik pembagian data tersebut dapat mengurangi ketidakpastian dalam target sehingga membantu dalam memilih pembagian terbaik dalam pohon keputusan.

```
def _gain(self, y: pd.Series, X_col:
pd.Series, t):
    p_var = self._variance(y)

    l_idx = X_col.isna() | (X_col < t)
    r_idx = ~l_idx

    if l_idx.sum() == 0 or r_idx.sum() == 0:
        return 0

    n = len(y)
    n_l, n_r = l_idx.sum(), r_idx.sum()
    v_l, v_r = self._variance(y[l_idx]),
self._variance(y[r_idx])
    c_var = (n_l / n) * v_l + (n_r / n) * v_r

    return p_var - c_var

def _variance(self, y: pd.Series):
    return y.var()
```

### C. Pipelining dan Pengujian

Pada gambar di bawah ini, terlihat proses pembagian dataset menjadi dua subset: data pelatihan dan data pengujian. Proses dimulai dengan pengacakan data secara acak untuk memastikan bahwa pembagian data tidak mengandung bias. Selanjutnya, 80% dari total dataset dialokasikan untuk data pelatihan, yang digunakan untuk melatih model, sementara 20% sisanya digunakan sebagai data pengujian untuk mengevaluasi kinerja model yang telah dilatih. Pembagian ini penting untuk memastikan bahwa model dapat dievaluasi menggunakan data yang tidak digunakan dalam proses pelatihan, sehingga dapat memberikan gambaran yang lebih akurat mengenai kemampuan model dalam melakukan generalisasi pada data baru.

```
def _train_test_split(self, train:
pd.DataFrame, test_size=0.2) -> tuple:
    n = len(train)
    idx = np.arange(n)
    np.random.shuffle(idx)
    train_idx = idx[int(n * test_size):]
    test_idx = idx[:int(n * test_size)]
    return train.iloc[train_idx],
train.iloc[test_idx]
```

Pada bagian pengujian, penulis menggunakan metrik *root mean squared error (RMSE)* untuk setiap fitur sebagai indikator utama dalam evaluasi model. Metrik ini dipilih karena memiliki

kecocokan yang tinggi dan merupakan salah satu metode yang paling banyak digunakan dalam analisis galat dan mengingat kemampuannya untuk memberikan penalti yang lebih besar pada galat yang lebih besar, sehingga memberikan gambaran yang jelas tentang seberapa jauh prediksi menyimpang dari nilai sebenarnya.

```

error_list = [[] for _ in range(7)]
for col_idx, col in enumerate(test.columns):
    X_train, y_train =
train.drop(columns=col), train[col]
    y_train = y_train.dropna()
    X_train = X_train.loc[y_train.index]

    X_test, y_test = test.drop(columns=col),
test[col]
    nan_indices = y_test.index[y_test.isna()]
    X_test = X_test.loc[nan_indices]

    model = DecisionTreeRegressor()
    model.fit(X_train, y_train)

    y_pred = model.predict(X_test)
    y_true = df.loc[nan_indices, col]

    squared_error = (y_pred - y_true.values)
    ** 2
    error_list[col_idx] = squared_error

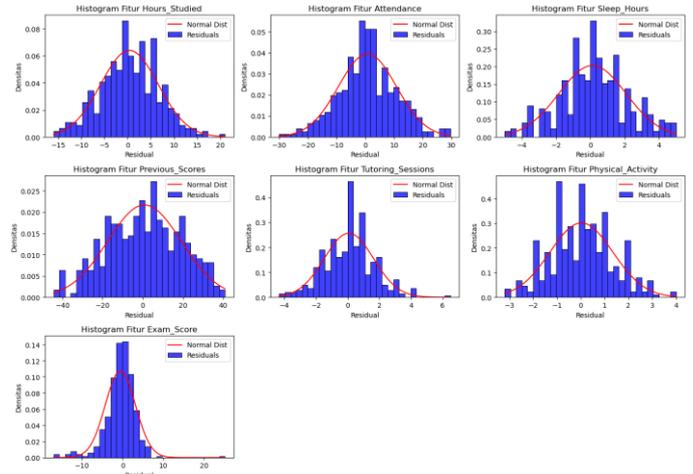
rmse_per_feature = {
test.columns[j]:
np.sqrt(np.mean(error_list[j])).astype(float)
for j in range(len(test.columns))
}

```

#### D. Evaluasi dan Analisis Galat

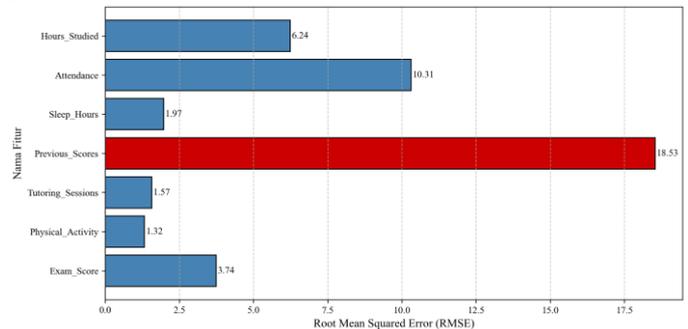
Pada Gambar 4, terlihat tujuh histogram yang menggambarkan distribusi residual dari Persamaan (3) pada model Decision Tree untuk berbagai fitur yang berbeda, seperti Hours\_Studied, Attendance, Sleep\_Hours, Previous\_Scores, Tutoring\_Sessions, Physical\_Activity, dan Exam\_Score. Pada setiap histogram, dapat terlihat distribusi residual yang diwakili oleh batang berwarna biru dengan kurva distribusi normal yang di-overlay menggunakan garis merah. Distribusi residual ini mengindikasikan bagaimana perbedaan antara nilai prediksi dan nilai aktual untuk setiap fitur.

Sebagian besar fitur menunjukkan distribusi residual yang cukup mirip dengan distribusi normal. Hal ini mengindikasikan bahwa model pohon keputusan dapat memberikan prediksi yang cukup baik untuk fitur-fitur tersebut. Namun, ada beberapa histogram yang menunjukkan ketidaksimetrisan atau puncak tajam, seperti pada fitur Tutoring\_Sessions, yang mengindikasikan adanya masalah dalam kemampuan model untuk menangkap variabilitas data dengan baik, atau menunjukkan adanya *overfitting*. Pada fitur-fitur seperti Attendance dan Sleep\_Hours, distribusi residual menunjukkan variasi yang lebih besar, yang dapat menunjukkan bahwa model kesulitan menangkap pola yang lebih kompleks dalam data. Secara keseluruhan, distribusi residual ini memberikan wawasan penting mengenai seberapa baik model Decision Tree dalam memprediksi nilai pada masing-masing fitur, serta menunjukkan di mana model mungkin perlu disempurnakan.



Gambar 4. Histogram residual setiap fitur

Selanjutnya pada gambar 5, ditampilkan RMSE setiap fitur agar dapat mengevaluasi lebih jelas. Terlihat bahwa fitur Previous\_Scores memiliki RMSE tertinggi. Hal ini bisa saja terjadi mengingat nilai yang dihilangkan adalah 28% dari keseluruhan data sehingga akan mengurangi variabilitas prediksi model.



Gambar 5. RMSE baseline setiap fitur

Walaupun terdapat kekurangan pada satu fitur yang sudah dijelaskan di atas, model ini sudah terbilang cukup baik berdasarkan skor RMSE dari fitur-fitur yang lain.

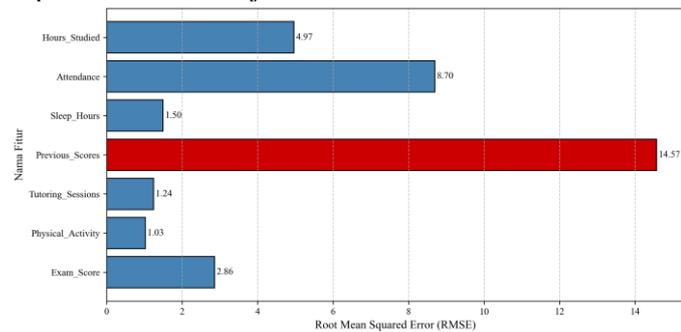
#### E. Hyperparameter Tuning

Pada implementasi ini, dilakukan proses *hyperparameter tuning* untuk menyempurnakan kinerja model dengan menggunakan pendekatan Grid Search. *Hyperparameter tuning* bertujuan untuk menemukan kombinasi parameter terbaik yang menghasilkan galat terkecil dalam proses prediksi. Parameter yang disesuaikan meliputi *max\_depth*, *min\_split*, dan *min\_leaf*. Pada eksperimen ini, nilai *max\_depth* divariasikan dari 1 hingga 8, *min\_split* dari 2 hingga 4, dan *min\_leaf* dari 1 hingga 4.

Grid search adalah algoritma *brute force* yang melakukan pencarian secara menyeluruh pada subset ruang hyperparameter yang diberikan [4]. Proses ini dilakukan sebanyak  $8 \times 3 \times 4 = 96$  kali dengan melakukan *fitting* dan prediksi berulang.

Berdasarkan Gambar 6, dapat terlihat bahwa kombinasi *hyperparameter* terbaik untuk model yang diuji adalah dengan *max\_depth* sebesar 5, *min\_split* sebesar 2, dan *min\_leaf* sebesar 4. Kombinasi ini menghasilkan nilai rata-rata Root Mean Squared Error (RMSE) keseluruhan terkecil, yaitu 4.98—semula 6.24. Dengan menggunakan kombinasi parameter ini, model dapat dioptimalkan untuk memberikan hasil prediksi

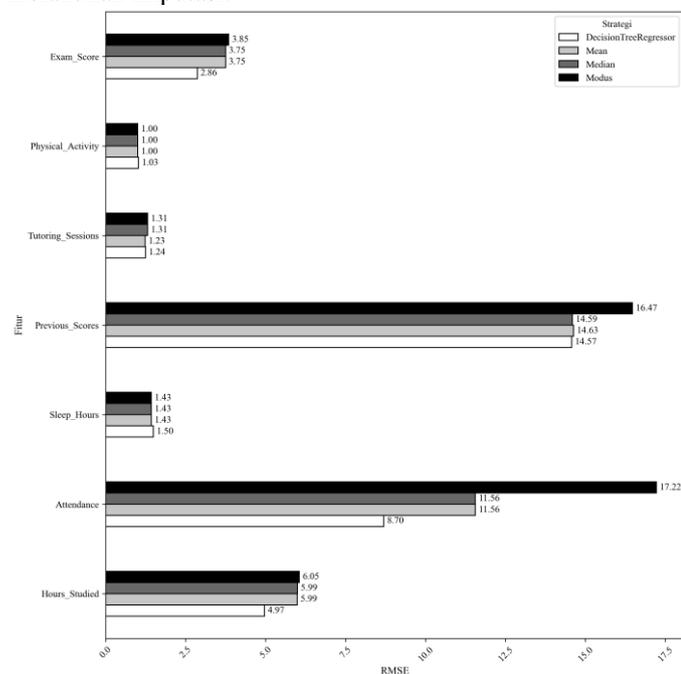
yang lebih akurat, menjadikannya konfigurasi yang direkomendasikan untuk digunakan dalam pengujian atau implementasi lebih lanjut.



Gambar 6. RMSE setelah Hyperparameter Tuning

### F. Benchmarking

Benchmarking digunakan untuk mengukur seberapa baik suatu model dalam melakukan prediksi dibandingkan dengan model lainnya. Dilakukan sebuah *benchmarking* dengan algoritma imputasi tradisional, seperti mean, median, dan modus. Berdasarkan Tabel III, terlihat bahwa—secara keseluruhan—algoritma Tuned DecisionTreeRegressor (TDTR) memiliki nilai RMSE terendah dibandingkan model-model tradisional. Hal ini dapat disimpulkan bahwa pada kasus ini, DecisionTreeRegressor bisa dijadikan salah satu opsi dalam melakukan imputasi.



Gambar 7. RMSE pada empat strategi

TABEL III  
SKOR BENCHMARKING

Fitur	RMSE			
	TDTR	Mean	Median	Modus
Hours_Studied	4.97	5.99	5.99	6.05
Attendance	8.7	11.56	11.56	17.22
Sleep_Hours	1.5	1.43	1.43	1.43
Previous_Scores	14.57	14.63	14.59	16.47
Tutoring_Sessions	1.24	1.23	1.31	1.31

Physical_Activity	1.03	1	1	1
Exam_Score	2.86	3.75	3.75	3.85
Overall	<b>4.98</b>	5.66	5.66	6.76

## V. KESIMPULAN

Imputasi data nonlinear menggunakan pohon keputusan yang dimodifikasi dapat menjadi salah satu opsi dari sekian banyaknya teknik imputasi yang sudah ada. Dengan mempertimbangkan kompleksitas dan secara kualitatif, pohon keputusan merupakan algoritma yang cukup efisien dengan kompleksitas total untuk mengimputasi semua fitur adalah  $O(d^2 n^2 \log n)$  dan cukup andal dalam menghadapi kasus hubungan data nonlinear dengan rata-rata nilai galat secara keseluruhan 4.98.

## VI. LAMPIRAN

Sebagai bagian dari transparansi dan untuk memudahkan replikasi penelitian ini, seluruh kode sumber, dataset, dan dokumentasi yang digunakan dalam penelitian ini dapat diakses melalui repositori <https://github.com/fathurwithyou/Matdis-MultivariateImputation>.

Selain itu, untuk mempermudah pemahaman terhadap proses dan hasil penelitian ini, sebuah video penjelasan telah disediakan. Video tersebut dapat diakses melalui tautan <https://youtu.be/zNDRNfv6id4>. Video ini memberikan penjelasan rinci tentang langkah-langkah penelitian, implementasi algoritma, serta analisis hasil yang diperoleh.

## VII. PENUTUP

Ucapan terima kasih penulis berikan kepada Tuhan Yang Maha Esa, Institut Teknologi Bandung, Prodi Teknik Informatika, dan dosen mata kuliah IF2120 Matematika Diskrit, Bapak Rinaldi Munir. Diharapkan makalah ini dapat bermanfaat bagi pembaca dan dapat dikembangkan lebih lanjut. Seperti kata Musashi dalam blog Rinaldi Munir, “*Hidup lebih dari harus sekedar tetap hidup. Masalahnya, bagaimana membuat hidup itu bermakna, memancarkan cahaya gemilang ke masa depan, sekalipun harus mengorbankan hidup itu sendiri. Bila ini sudah tercapai, maka tidak ada bedanya dengan berapa lama hidup itu.*”

## REFERENSI

- [1] Rinaldi Munir, “Graf (Bag.1).” Accessed: Dec. 25, 2024. [Online]. Available: <https://informatika.stei.itb.ac.id/~rinaldi.munir/Matdis/2024-2025/20-Graf-Bagian1-2024.pdf>
- [2] Rinaldi Munir, “Pohon (Bag. 1).” Accessed: Dec. 25, 2024. [Online]. Available: <https://informatika.stei.itb.ac.id/~rinaldi.munir/Matdis/2024-2025/23-Pohon-Bag1-2024.pdf>
- [3] J. R. Quinlan, “Induction of decision trees,” *Mach Learn*, vol. 1, no. 1, pp. 81–106, 1986, doi: 10.1007/BF00116251.
- [4] P. Liashchynskyi and P. Liashchynskyi, “Grid Search, Random Search, Genetic Algorithm: A big comparison for NAS,” 2019.

## PERNYATAAN

Dengan ini saya menyatakan bahwa makalah yang saya tulis ini adalah tulisan saya sendiri, bukan saduran, atau terjemahan dari makalah orang lain, dan bukan plagiasi.

Bandung, 5 Januari 2024



Muhammad Fathur Rizky (13523105)