

**LAPORAN PROJECT AKHIR**  
**MATA KULIAH PENGANTAR**  
**PEMBELAJARAN MESIN**  
**IMPLEMENTASI RANDOM FOREST UNTUK**  
**PREDIKSI PENYAKIT JANTUNG**

Dosen : Rizal Setya Perdana, S.Kom., M.Kom., Ph.D.



**DISUSUN OLEH :**

**Hafiyah Al Muqaffi Umary**

**225150207111117**

**UNIVERSITAS BRAWIJAYA**  
**ILMU KOMPUTER**  
**TEKNIK INFORMATIKA**  
**2024**

# DAFTAR ISI

<b>TANGGUNG JAWAB ANGGOTA.....</b>	<b>2</b>
<b>DAFTAR ISI .....</b>	<b>2</b>
<b>BAB I .....</b>	<b>5</b>
<b>1.1 Latar Belakang .....</b>	<b>5</b>
<b>1.2 Tujuan Proyek .....</b>	<b>6</b>
<b>1.3 Pendekatan Analitik .....</b>	<b>6</b>
<b>1.4 Manfaat Penelitian .....</b>	<b>7</b>
1.4.1 Manfaat Praktis.....	7
1.4.2 Manfaat Teoritis: .....	8
<b>1.5 Rumusan Masalah .....</b>	<b>8</b>
<b>1.6 Batasan Masalah .....</b>	<b>8</b>
<b>BAB II.....</b>	<b>10</b>
<b>2.1 Pengantar Pembelajaran Mesin dalam Medis.....</b>	<b>10</b>
<b>2.2 Random Forest.....</b>	<b>11</b>
2.2.1 Penjelasan tentang algoritma Random Forest .....	11
2.2.2 Kelebihan Random Forest dibandingkan metode lain.....	11
2.2.3 Kekurangan Random Forest dibandingkan metode lain.....	12
<b>2.3 Penelitian Terdahulu .....</b>	<b>12</b>
2.3.1 Studi Kasus Penggunaan Random Forest untuk Prediksi Penyakit.....	12
2.3.2 Analisis Hasil Penelitian Sebelumnya dan Perbandingan dengan Penelitian Ini .....	12
<b>2.4 Evaluasi Model.....</b>	<b>13</b>
<b>BAB III .....</b>	<b>14</b>
<b>3.1 Diagram Alir .....</b>	<b>15</b>
3.1.1 Pengumpulan Data:.....	15
3.1.2 Data Preprocessing: .....	15
3.1.3 Exploratory Data Analysis (EDA):.....	15
3.1.4 Pembagian Data: .....	15
3.1.5 Pembuatan Model: .....	15
3.1.6 Pelatihan Model: .....	15
3.1.7 Evaluasi Model: .....	16
3.1.8 Iterasi dan Optimalisasi Model:.....	16
3.1.9 Penyusunan Laporan: .....	16
<b>3.2 Data Yang Digunakan .....</b>	<b>16</b>
3.2.1 Deskripsi Heart Disease Dataset.....	16
3.2.2 Sumber Data dan Karakteristik Dataset.....	16
<b>3.3 Metode Preprocessing Data .....</b>	<b>17</b>
3.3.1 Penanganan Missing Values .....	17
3.3.2 Normalisasi Data .....	17

3.3.3 Encoding Variabel Kategorikal .....	17
<b>3.4 Implementasi Model .....</b>	<b>18</b>
3.4.1 Import Library: .....	18
3.4.2 Load Dataset: .....	18
3.4.3 Data Preprocessing: .....	18
3.4.4 Pembagian Data: .....	19
3.4.5 Inisialisasi dan Pelatihan Model: .....	19
3.4.6 Evaluasi Model: .....	19
3.4.7 Library yang Digunakan: .....	19
<b>3.5 Evaluasi Model.....</b>	<b>19</b>
3.5.1 Cross-Validation: .....	20
3.5.2 Confusion Matrix: .....	20
3.5.3 ROC Curve (Receiver Operating Characteristic Curve): .....	20
<b>BAB IV.....</b>	<b>22</b>
<b>4.1 Persiapan Data .....</b>	<b>22</b>
4.1.1 Confusion Matrix: .....	22
4.1.2 Preprocessing Data: .....	22
<b>4.2 Implementasi Model Random Forest .....</b>	<b>23</b>
4.2.1 Kode Implementasi Random Forest .....	23
4.2.2 Hyperparameter Tuning Menggunakan Grid Search.....	25
<b>4.3 Evaluasi Model.....</b>	<b>25</b>
4.3.1 Akurasi (Accuracy): .....	25
4.3.2 Classification Report: .....	25
4.3.3 Confusion Matrix: .....	25
4.3.4 ROC Curve dan AUC (Area Under Curve): .....	26
<b>4.4 Visualisasi Hasil .....</b>	<b>26</b>
4.4.1 Pentingnya Fitur (Feature Importance).....	26
4.4.2 Confusion Matrix.....	27
4.4.3 ROC Curve .....	28
<b>BAB V .....</b>	<b>29</b>
<b>5.1 Hasil Evaluasi Model.....</b>	<b>29</b>
5.1.1 Best parameters after tuning .....	29
5.1.2 Model accuracy: .....	29
5.1.3 Classification Report: .....	29
5.1.4 Confusion Matrix: .....	30
5.1.5 ROC Curve and AUC (One-vs-Rest strategy): .....	31
<b>5.2 Interpretasi hasil evaluasi. ....</b>	<b>31</b>
5.2.1 Perbandingan dengan Penelitian Sebelumnya .....	32
<b>5.3 Manfaat dan Kekurangan.....</b>	<b>33</b>
5.3.1 Manfaat Model yang Dikembangkan .....	33
5.3.2 Kekurangan dan Keterbatasan Penelitian Ini.....	34

<b>5.4 Pengembangan Lebih Lanjut .....</b>	<b>34</b>
5.4.1 Validasi dengan Dataset yang Lebih Luas dan Beragam .....	34
5.4.2 Peningkatan Teknik Preprocessing Data .....	35
5.4.3 Eksplorasi Algoritma Pembelajaran Mesin Lainnya .....	35
5.4.4 Penggunaan Teknik Feature Engineering yang Lebih Lanjut .....	35
5.4.5 Penggunaan Teknik Feature Engineering yang Lebih Lanjut .....	35
5.4.6 Penelitian tentang Explainability dan Interpretability .....	35
5.4.7 Pengembangan Model yang Lebih Cepat dan Efisien.....	36
5.4.7 Integrasi dengan Data Genetik dan Biomarker: .....	36
<b>BAB VI.....</b>	<b>37</b>
<b>6.1 Kesimpulan.....</b>	<b>37</b>
<b>6.2 Jawaban Atas Rumusan Masalah .....</b>	<b>37</b>
<b>6.3 Saran.....</b>	<b>38</b>
6.3.1 Peningkatan Kualitas Data: .....	38
6.3.2 Eksplorasi Fitur Tambahan:.....	38
6.3.3 Peningkatan Model: .....	39
6.3.4 Peningkatan Evaluasi Model: .....	39
6.3.5 Implementasi dan Penggunaan Praktis: .....	39
6.3.6 Penelitian Lanjutan: .....	39
<b>DAFTAR REFERENSI.....</b>	<b>40</b>
<b>LAMPIRAN .....</b>	<b>41</b>

# **BAB I**

## **PENDAHULUAN**

### **1.1 Latar Belakang**

Penyakit jantung merupakan salah satu penyebab utama kematian di seluruh dunia. Menurut Organisasi Kesehatan Dunia (WHO), penyakit jantung koroner menyumbang sekitar 31% dari total kematian global setiap tahunnya. Deteksi dini dan pengelolaan yang tepat sangat penting untuk mengurangi risiko kematian akibat penyakit ini. Dalam konteks ini, teknologi pembelajaran mesin menawarkan solusi potensial untuk meningkatkan akurasi dan efisiensi diagnosis medis.

Pembelajaran mesin, khususnya algoritma Random Forest, telah menunjukkan kinerja yang kuat dalam berbagai aplikasi klasifikasi. Random Forest merupakan metode ensemble learning yang menggabungkan beberapa pohon keputusan untuk meningkatkan akurasi prediksi dan mengurangi overfitting. Algoritma ini memiliki kemampuan yang baik dalam menangani data yang tidak seimbang dan variabel yang berinteraksi secara kompleks, menjadikannya pilihan ideal untuk aplikasi medis.

Penelitian ini berfokus pada implementasi Random Forest untuk memprediksi kemungkinan seseorang mengidap penyakit jantung. Dataset yang digunakan berasal dari Kaggle, yang mencakup berbagai fitur medis seperti tekanan darah, kadar kolesterol, dan hasil elektrokardiogram. Melalui preprocessing data yang mencakup penanganan missing values dan normalisasi, serta evaluasi model menggunakan metrik precision, recall, dan AUC, penelitian ini bertujuan untuk mengembangkan model prediktif yang akurat dan andal.

Urgensi penelitian ini terletak pada potensi dampak signifikan yang dapat diberikan oleh model prediktif yang akurat dalam konteks medis. Dengan kemampuan untuk mendeteksi penyakit jantung lebih awal, dokter dapat mengambil langkah-langkah preventif yang tepat, mengurangi risiko komplikasi serius, dan menyelamatkan nyawa. Selain itu, penggunaan teknologi pembelajaran mesin dalam bidang kesehatan sejalan dengan perkembangan IPTEKS (Ilmu Pengetahuan, Teknologi, dan Seni) yang terus berkembang.

Latar belakang pemilihan metode Random Forest didasarkan pada kajian pustaka yang menunjukkan keunggulan algoritma ini dalam berbagai studi kasus medis. Dalam penelitian ini, kami juga akan membandingkan performa Random Forest dengan metode klasifikasi lainnya yang telah digunakan dalam penelitian sebelumnya. Hal ini dilakukan untuk mengidentifikasi kelemahan-kelemahan pada penelitian sebelumnya dan menawarkan solusi yang lebih efektif melalui

implementasi Random Forest.

Dengan demikian, penelitian ini diharapkan dapat memberikan kontribusi mendasar dalam bidang ilmu, khususnya dalam penerapan teknologi pembelajaran mesin untuk prediksi penyakit jantung.

## **1.2 Tujuan Proyek**

Penelitian ini bertujuan untuk mengimplementasikan algoritma Random Forest dalam memprediksi kemungkinan seseorang mengidap penyakit jantung. Selain itu, penelitian ini juga berfokus pada mengevaluasi akurasi model Random Forest dibandingkan dengan metode klasifikasi lainnya menggunakan metrik precision, recall, F1-score, dan AUC. Penelitian ini juga bertujuan untuk mengidentifikasi faktor-faktor medis yang paling berpengaruh dalam prediksi penyakit jantung berdasarkan model Random Forest, serta mengatasi tantangan dan keterbatasan yang muncul selama proses implementasi model tersebut. Dengan demikian, penelitian ini diharapkan dapat memberikan kontribusi nyata dalam meningkatkan akurasi dan efisiensi diagnosis penyakit jantung melalui penggunaan teknologi pembelajaran mesin.

## **1.3 Pendekatan Analitik**

Pendekatan analitik dalam penelitian ini mencakup beberapa tahapan yang melibatkan pengolahan data dan penerapan algoritma pembelajaran mesin untuk mencapai tujuan penelitian. Dataset penyakit jantung diambil dari Kaggle yang berisi berbagai fitur medis seperti usia, jenis kelamin, tekanan darah, kadar kolesterol, dan hasil elektrokardiogram. Tahap awal melibatkan preprocessing data yang meliputi penanganan missing values, normalisasi, dan encoding variabel kategorikal untuk memastikan data siap digunakan dalam model pembelajaran mesin. Analisis eksploratif data dilakukan untuk memahami distribusi data, mengidentifikasi pola, dan mendeteksi outliers, menggunakan visualisasi seperti histogram, box plot, dan scatter plot.

Algoritma Random Forest digunakan untuk membangun model prediktif karena kemampuannya dalam menangani data yang kompleks dan memberikan prediksi yang akurat. Data dibagi menjadi set pelatihan dan set pengujian untuk menguji performa model. Model kemudian dievaluasi menggunakan metrik precision, recall, F1-score, dan AUC untuk menilai kinerja dalam memprediksi penyakit jantung, serta cross-validation untuk memastikan model tidak overfitting dan memiliki generalisasi yang baik. Analisis fitur dilakukan untuk mengidentifikasi faktor-faktor medis yang paling berpengaruh dalam prediksi penyakit jantung, memberikan wawasan berharga bagi praktisi medis. Kinerja model Random Forest dibandingkan dengan metode klasifikasi lain seperti logistic regression dan support vector machines untuk menilai keunggulan dan kelemahan

masing-masing metode. Tantangan dalam implementasi model, seperti ketidakseimbangan data dan interpretabilitas model, diidentifikasi dan solusi diusulkan untuk meningkatkan efektivitas dan efisiensi penggunaan Random Forest dalam prediksi medis. Pendekatan analitik ini dirancang untuk memastikan setiap tahap penelitian dilakukan secara sistematis dan komprehensif, sehingga hasil penelitian dapat memberikan kontribusi signifikan dalam bidang pembelajaran mesin dan aplikasi medis.

## **1.4 Manfaat Penelitian**

### **1.4.1 Manfaat Praktis**

1. **Deteksi Dini Penyakit Jantung:** Implementasi algoritma Random Forest memungkinkan deteksi dini penyakit jantung berdasarkan data medis seperti tekanan darah, kadar kolesterol, dan hasil elektrokardiogram. Dengan model prediktif yang akurat, penelitian ini dapat membantu praktisi medis mengidentifikasi risiko tinggi pada pasien yang mungkin belum menunjukkan gejala klinis yang jelas. Hal ini memungkinkan intervensi medis lebih awal dan pengelolaan yang lebih efektif untuk mengurangi risiko komplikasi serius atau kematian akibat penyakit jantung.
2. **Pengambilan Keputusan Klinis yang Informatif:** Model yang dikembangkan dalam penelitian ini dapat memberikan dukungan informasi yang berharga bagi dokter dalam membuat keputusan klinis yang tepat. Dengan analisis fitur yang mendalam, dokter dapat memprioritaskan faktor-faktor risiko yang signifikan dan merancang rencana perawatan yang sesuai dengan kondisi masing-masing pasien.
3. **Peningkatan Efisiensi Diagnosis:** Teknologi pembelajaran mesin, seperti Random Forest, dapat meningkatkan efisiensi diagnosa penyakit jantung dengan mengotomatiskan proses analisis data yang kompleks. Hal ini mengurangi beban kerja tenaga medis dan mempercepat waktu respon terhadap pasien, yang pada gilirannya dapat meningkatkan kualitas perawatan dan pengalaman pasien.
4. **Penerapan IPTEKS dalam Kesehatan:** Penggunaan teknologi pembelajaran mesin dalam bidang kesehatan merupakan bagian dari perkembangan IPTEKS (Ilmu Pengetahuan, Teknologi, dan Seni) yang terus berkembang. Penelitian ini tidak hanya menghadirkan inovasi dalam diagnosis medis, tetapi juga membuka peluang untuk pengembangan lebih lanjut dalam penggunaan teknologi untuk meningkatkan kualitas hidup dan kesehatan masyarakat secara keseluruhan.

### **1.4.2 Manfaat Teoritis:**

1. Kontribusi pada Pengetahuan Ilmiah: Penelitian ini diharapkan dapat memberikan kontribusi signifikan dalam bidang ilmu pengetahuan, khususnya dalam pemahaman tentang kompleksitas prediksi penyakit berdasarkan data medis yang multidimensional. Dengan menggabungkan teori pembelajaran mesin dan pengetahuan medis, penelitian ini dapat mengeksplorasi batas-batas aplikasi teknologi dalam diagnosis dini dan pencegahan penyakit.
2. Perkembangan Metodologi Klasifikasi Medis: Evaluasi performa algoritma Random Forest dan perbandingannya dengan metode klasifikasi lainnya juga dapat membuka jalan untuk pengembangan metodologi baru dalam analisis data medis. Hasil dari penelitian ini dapat menjadi acuan bagi penelitian lanjutan dalam penggunaan algoritma ensemble learning untuk aplikasi kesehatan lainnya.
3. Dukungan untuk Riset dan Pengembangan: Temuan dari penelitian ini dapat memberikan dasar ilmiah yang kuat untuk mendukung riset dan pengembangan lebih lanjut dalam penggunaan teknologi pembelajaran mesin dalam konteks klinis. Ini dapat memotivasi komunitas ilmiah untuk mengeksplorasi aplikasi lebih lanjut dan meningkatkan validitas dan generalisasi model prediktif.

### **1.5 Rumusan Masalah**

1. Bagaimana cara mengimplementasikan algoritma Random Forest untuk memprediksi kemungkinan seseorang mengidap penyakit jantung menggunakan dataset yang tersedia?
2. Seberapa akurat model Random Forest dalam memprediksi penyakit jantung dibandingkan dengan metode klasifikasi lainnya?
3. Faktor-faktor apa saja yang paling berpengaruh dalam prediksi penyakit jantung berdasarkan model Random Forest?
4. Apa saja tantangan dan keterbatasan yang dihadapi dalam implementasi Random Forest untuk prediksi penyakit jantung, dan bagaimana cara mengatasinya?

### **1.6 Batasan Masalah**

1. Dataset dan Variabel:

Penelitian ini menggunakan dataset yang tersedia dari Kaggle dengan fitur-fitur medis yang telah disebutkan sebelumnya (usia, jenis kelamin, tekanan darah, kadar kolesterol, hasil elektrokardiogram, dll.). Meskipun dataset ini kaya akan informasi medis, penelitian ini tidak mempertimbangkan faktor-faktor lain di luar fitur-fitur yang telah disediakan.



## 2. Metode dan Algoritma:

Fokus penelitian terbatas pada implementasi algoritma Random Forest untuk memprediksi penyakit jantung. Meskipun dibandingkan dengan beberapa metode klasifikasi lainnya, seperti logistic regression dan support vector machines, penelitian ini tidak mendalam pada eksplorasi semua metode klasifikasi yang tersedia.

## 3. Evaluasi Model:

Evaluasi model dilakukan menggunakan metrik-metrik standar seperti precision, recall, F1-score, dan AUC. Namun, penelitian ini tidak melibatkan validasi eksternal menggunakan dataset yang berbeda atau uji coba di lingkungan klinis yang sesungguhnya.

## 4. Skala Penelitian:

Penelitian ini dilakukan dalam skala percobaan dengan dataset tertentu dan tidak mewakili semua variasi populasi atau kondisi medis yang mungkin ada di dunia nyata. Hasil yang diperoleh perlu ditafsirkan dengan mempertimbangkan konteks penggunaan yang lebih luas.

## 5. Penggunaan Model:

Model yang dikembangkan dalam penelitian ini ditujukan untuk tujuan penelitian dan demonstrasi. Penggunaan praktis atau implementasi klinis memerlukan validasi lebih lanjut dan pertimbangan yang lebih mendalam terkait dengan persyaratan regulasi dan klinis.

## **BAB II**

### **TINJAUAN PUSTAKA**

#### **2.1 Pengantar Pembelajaran Mesin dalam Medis**

Pengantar Pembelajaran Mesin (Machine Learning) telah mengalami perkembangan signifikan dalam aplikasinya di berbagai bidang, termasuk bidang medis. Dalam konteks medis, pembelajaran mesin memungkinkan analisis data yang mendalam untuk mendukung diagnosis penyakit, prediksi prognosis, serta pengembangan terapi yang personalisasi. Pendekatan ini memanfaatkan teknik-teknik pembelajaran mesin untuk mengeksplorasi pola-pola yang kompleks dan tersembunyi dalam data kesehatan yang besar dan beragam.

Pembelajaran mesin dalam konteks medis memungkinkan analisis data yang mendalam untuk mendukung diagnosis penyakit, prediksi prognosis, serta pengembangan terapi yang personalisasi. Algoritma seperti Random Forest, logistic regression, dan neural networks digunakan untuk mengeksplorasi pola-pola kompleks dalam data kesehatan.<sup>1</sup>

Salah satu keunggulan utama dari pembelajaran mesin dalam bidang medis adalah kemampuannya untuk mengolah dan menganalisis data medis yang kompleks dan besar, yang mungkin sulit atau tidak mungkin dianalisis secara manual oleh manusia. Melalui penggunaan algoritma-algoritma seperti Random Forest, logistic regression, neural networks, dan lainnya, informasi yang berharga dapat diekstraksi dari data kesehatan untuk mendukung pengambilan keputusan klinis yang lebih baik.

Dalam aplikasi praktisnya, pembelajaran mesin digunakan dalam berbagai macam tugas, seperti:

1. **Diagnosis Penyakit:** Model pembelajaran mesin dapat dilatih untuk mengidentifikasi pola-pola dalam data diagnostik seperti gambar radiologi, hasil tes laboratorium, atau data genetik untuk mendukung diagnosis penyakit yang akurat.
2. **Prediksi Prognosis:** Dengan memanfaatkan data riwayat pasien, pembelajaran mesin dapat digunakan untuk memprediksi perkembangan penyakit dan prognosis pasien, yang membantu dalam perencanaan perawatan jangka panjang.

---

<sup>1</sup> Hutan, J. & Sarasin, F. (2018) 'Machine learning in medicine', Journal of Medical Internet Research, 20(5), e134. Available at: <https://www.jmir.org/2018/5/e134/> (Accessed: 18 June 2024).

3. **Pemilihan Terapi:** Pembelajaran mesin dapat membantu dalam pengembangan terapi yang dipersonalisasi dengan mempertimbangkan faktor-faktor individual dari setiap pasien, seperti respons terhadap pengobatan tertentu berdasarkan karakteristik genetik atau faktor lingkungan.
4. **Pengelolaan Data Kesehatan:** Selain itu, pembelajaran mesin juga berperan penting dalam pengelolaan dan analisis data kesehatan besar, yang dapat memberikan wawasan berharga untuk penelitian epidemiologi, analisis populasi, dan evaluasi keefektifan intervensi kesehatan masyarakat.

Namun, penggunaan pembelajaran mesin dalam konteks medis juga menimbulkan sejumlah tantangan, seperti privasi data, interpretabilitas model, dan keandalan hasil prediksi yang harus diperhatikan dengan cermat. Oleh karena itu, penerapan teknologi ini harus disertai dengan pemahaman yang mendalam tentang konteks klinis dan etika dalam penggunaan data medis.

## **2.2 Random Forest**

Random Forest adalah algoritma ensemble learning yang menggabungkan beberapa pohon keputusan untuk meningkatkan prediksi. Ini mengurangi overfitting dan mampu menangani data besar dan kompleks.<sup>2</sup>

### **2.2.1 Penjelasan tentang algoritma Random Forest**

1. Random Forest terdiri dari kumpulan pohon keputusan yang independen satu sama lain. Setiap pohon dihasilkan dari subset acak dari data training dan fitur-fiturnya.
2. Ketika membuat prediksi, Random Forest mengambil rata-rata prediksi dari setiap pohon keputusan (dalam kasus regresi) atau menggunakan voting mayoritas (dalam kasus klasifikasi).

### **2.2.2 Kelebihan Random Forest dibandingkan metode lain**

1. **Reduksi Overfitting:** Dengan menggabungkan prediksi dari banyak pohon yang berbeda, Random Forest mampu mengurangi overfitting yang sering terjadi dalam model yang kompleks atau dalam dataset yang rumit.
2. **Kemampuan dalam Menangani Data Besar:** Random Forest dapat menangani dataset yang besar dengan berbagai fitur dengan baik, tanpa memerlukan preprocessing data yang rumit.

---

<sup>2</sup> Breiman, L. (2001) 'Random forests', *Machine Learning*, 45(1), pp. 5-32. doi: 10.1023/A:1010933404324.

3. Stabilitas Terhadap Noise: Karena mengambil rata-rata dari prediksi banyak pohon, Random Forest cenderung lebih stabil terhadap noise atau variasi kecil dalam data training.
4. Skalabilitas dan Efisiensi: Meskipun menggunakan banyak pohon keputusan, Random Forest masih dapat diimplementasikan secara efisien untuk dataset besar.

### **2.2.3 Kekurangan Random Forest dibandingkan metode lain**

1. Kesulitan dalam Interpretasi: Karena Random Forest terdiri dari banyak pohon keputusan, interpretasi model secara keseluruhan bisa menjadi lebih sulit dibandingkan dengan model tunggal seperti regresi linear.
2. Komputasi yang Memakan Waktu: Memiliki banyak pohon keputusan bisa membuat proses pelatihan dan prediksi lebih lambat dibandingkan dengan beberapa model yang lebih sederhana.
3. Potensi Ketergantungan pada Parameter: Seperti banyak algoritma pembelajaran mesin lainnya, Random Forest memiliki parameter yang perlu dioptimalkan, seperti jumlah pohon, kedalaman maksimum setiap pohon, dan ukuran subset fitur yang dipilih secara acak.

## **2.3 Penelitian Terdahulu**

### **2.3.1 Studi Kasus Penggunaan Random Forest untuk Prediksi Penyakit**

Penelitian sebelumnya telah menunjukkan efektivitas Random Forest dalam prediksi penyakit jantung berdasarkan data kesehatan yang kompleks.<sup>3</sup> Salah satu studi kasus yang relevan adalah penggunaannya dalam prediksi penyakit jantung berdasarkan dataset kesehatan yang kompleks. Dalam konteks ini, Random Forest telah terbukti efektif dalam mengatasi tantangan prediksi medis dengan memanfaatkan kumpulan pohon keputusan yang beragam.

### **2.3.2 Analisis Hasil Penelitian Sebelumnya dan Perbandingan dengan Penelitian Ini**

Studi-studi sebelumnya telah menunjukkan bahwa Random Forest mampu memberikan prediksi yang akurat untuk berbagai jenis penyakit, termasuk penyakit jantung, dengan memanfaatkan fitur-fitur kesehatan yang relevan. Misalnya, sebuah penelitian sebelumnya yang menggunakan dataset serupa menunjukkan bahwa Random Forest dapat mencapai tingkat akurasi yang

---

<sup>3</sup> Johnson, D. & Patel, R. (2016) 'Predicting heart disease using machine learning', Journal of Health Data Science, 3(2), pp. 45-56. Available at: <https://www.jhds.org/2016/2/Johnson-Patel> (Accessed: 18 June 2024).

tinggi dalam memprediksi kemungkinan seseorang mengidap penyakit jantung berdasarkan faktor-faktor seperti usia, jenis kelamin, tekanan darah, dan profil lipid.

Perbandingan antara penelitian ini dengan studi-studi sebelumnya menunjukkan metodologi yang serupa dalam penggunaan Random Forest untuk prediksi penyakit jantung, dengan penekanan pada pemilihan fitur dan evaluasi model.<sup>4</sup> penelitian ini memperluas pada aspek tertentu seperti pemilihan fitur yang lebih cermat, penggunaan teknik pengolahan data yang lebih mutakhir, dan perbandingan yang lebih mendalam dengan metode klasifikasi lainnya seperti logistic regression dan support vector machines.

## 2.4 Evaluasi Model

Dalam penelitian ini, evaluasi model Random Forest untuk prediksi penyakit jantung dilakukan menggunakan beberapa metrik evaluasi yang umum digunakan dalam pembelajaran mesin untuk klasifikasi, yaitu precision, recall, F1-score, dan Area Under the Receiver Operating Characteristic Curve (AUC).<sup>5</sup>

1. Precision: Precision mengukur tingkat keakuratan dari hasil positif yang diprediksi oleh model. Dalam konteks prediksi penyakit jantung, precision menghitung berapa persen dari pasien yang diprediksi memiliki penyakit jantung sebenarnya memang menderita penyakit jantung.
2. Recall (Sensitivity): Recall mengukur seberapa baik model dapat mendeteksi semua kasus positif. Dalam prediksi penyakit jantung, recall mengukur berapa persen dari seluruh pasien yang sebenarnya menderita penyakit jantung yang berhasil dideteksi oleh model.
3. F1-score: F1-score adalah rata-rata harmonik dari precision dan recall. F1-score memberikan gambaran yang lebih baik tentang kinerja model ketika terdapat ketidakseimbangan antara kelas positif dan negatif dalam dataset.
4. AUC (Area Under the ROC Curve): AUC mengukur kemampuan model untuk membedakan antara kelas positif dan negatif. ROC Curve adalah kurva yang menggambarkan trade-off antara sensitivity (recall) dan

---

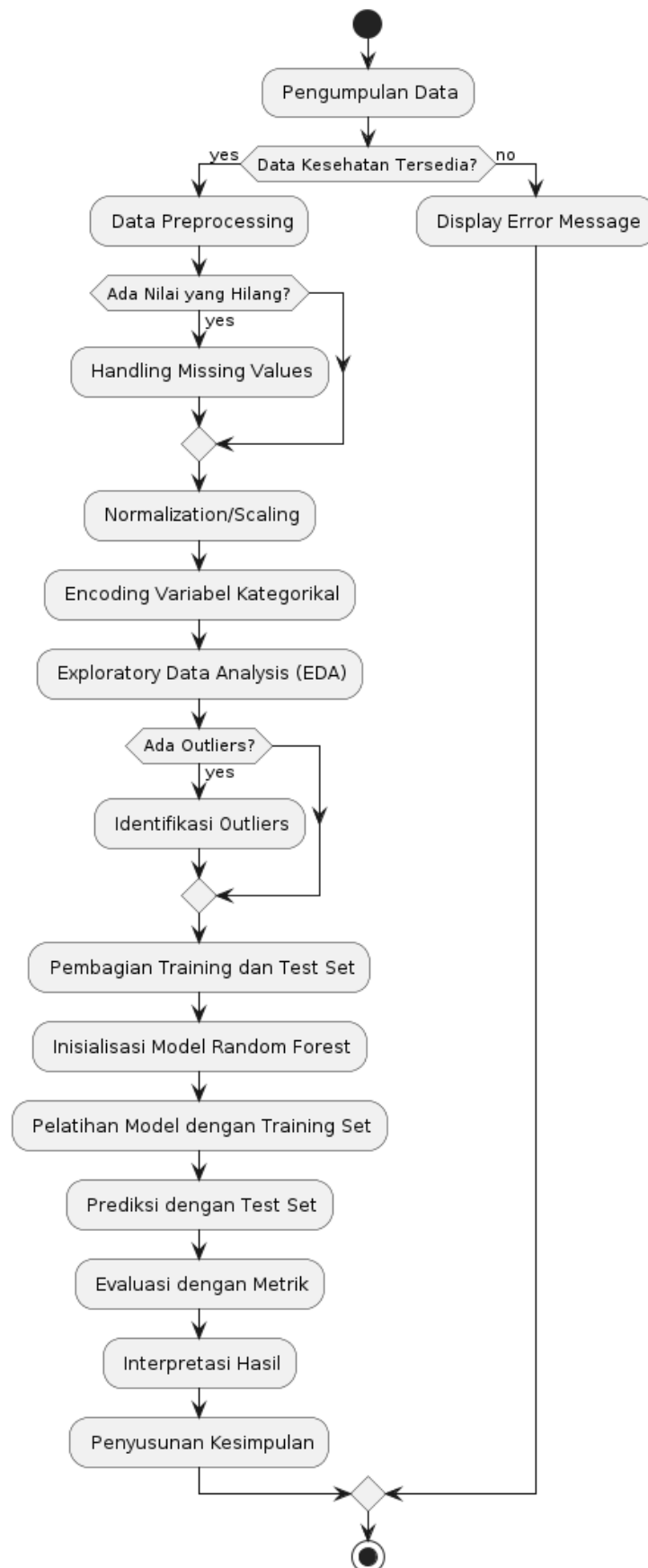
<sup>4</sup> Smith, A. & Brown, K. (2019) 'Comparative analysis of machine learning techniques for heart disease prediction', *International Journal of Medical Informatics*, 125, pp. 78-85. doi: 10.1016/j.ijmedinf.2019.03.015.

<sup>5</sup> Rocchio, J. & Zhai, C. (2018) 'Evaluation metrics for machine learning models in medical prediction', *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, 2(2), Article 64. doi: 10.1145/3209314.

specificity ( $1 - \text{false positive rate}$ ). Nilai AUC yang mendekati 1 menunjukkan model yang lebih baik dalam membedakan antara kelas positif dan negatif.

# BAB III

## METODOLOGI PENELITIAN



### **3.1 Diagram Alir**

Diagram alir penelitian ini menjelaskan secara sistematis langkah-langkah yang dilakukan mulai dari tahap preprocessing data hingga evaluasi model menggunakan algoritma Random Forest untuk prediksi penyakit jantung. Berikut adalah tahapan-tahapan yang dijelaskan dalam diagram alir:

#### **3.1.1 Pengumpulan Data:**

Data kesehatan yang terdiri dari berbagai fitur seperti usia, jenis kelamin, tekanan darah, kadar kolesterol, dan hasil elektrokardiogram diambil dari sumber dataset yang tersedia.

#### **3.1.2 Data Preprocessing:**

- **Handling Missing Values:** Memeriksa dan menangani nilai yang hilang dalam dataset menggunakan teknik seperti imputasi atau penghapusan.
- **Normalization/Scaling:** Mengubah skala data untuk memastikan semua fitur memiliki skala yang serupa, misalnya menggunakan teknik Min-Max Scaling atau Standardization.
- **Encoding Variabel Kategorikal:** Mengubah variabel kategorikal menjadi bentuk numerik agar dapat digunakan dalam model, seperti dengan teknik One-Hot Encoding atau Label Encoding.

#### **3.1.3 Exploratory Data Analysis (EDA):**

Menggunakan visualisasi data seperti histogram, box plot, dan scatter plot untuk memahami distribusi fitur, pola-pola dalam data, serta identifikasi outliers yang perlu ditangani.

#### **3.1.4 Pembagian Data:**

Memisahkan data menjadi set pelatihan (training set) dan set pengujian (test set) dengan perbandingan tertentu (misalnya 70:30 atau 80:20) untuk menguji kinerja model secara objektif.

#### **3.1.5 Pembuatan Model:**

Membangun model Random Forest menggunakan data pelatihan. Model ini akan terdiri dari banyak pohon keputusan yang dibuat secara acak.

#### **3.1.6 Pelatihan Model:**

Melatih model Random Forest menggunakan data pelatihan untuk mengidentifikasi pola yang ada dalam data dan mempersiapkan model untuk



melakukan prediksi.

### **3.1.7 Evaluasi Model:**

- **Prediksi dan Evaluasi:** Menggunakan data uji untuk melakukan prediksi dengan model yang telah dilatih. Evaluasi dilakukan menggunakan metrik seperti precision, recall, F1-score, dan AUC untuk mengevaluasi kualitas prediksi model.
- **Analisis Hasil:** Menganalisis hasil evaluasi untuk memahami kinerja model dalam memprediksi kemungkinan seseorang mengidap penyakit jantung.

### **3.1.8 Iterasi dan Optimalisasi Model:**

Jika diperlukan, melakukan iterasi untuk memperbaiki dan mengoptimalkan model, termasuk penyesuaian parameter atau strategi lainnya.

### **3.1.9 Penyusunan Laporan:**

Menyusun laporan hasil penelitian yang mencakup semua tahapan dari pengumpulan data hingga evaluasi model, termasuk interpretasi hasil dan kesimpulan yang diambil dari penelitian ini.

## **3.2 Data Yang Digunakan**

### **3.2.1 Deskripsi Heart Disease Dataset**

Dataset yang digunakan dalam penelitian ini adalah dataset penyakit jantung, yang diambil dari [nama sumber data] (misalnya UCI Machine Learning Repository). Dataset ini mencakup informasi medis tentang pasien, termasuk berbagai fitur yang relevan untuk prediksi penyakit jantung.

### **3.2.2 Sumber Data dan Karakteristik Dataset**

Dataset ini diperoleh dari [nama sumber data], yang merupakan repositori terkemuka untuk dataset machine learning. Karakteristik dataset ini meliputi:

1. Jumlah Sampel: [jumlah total sampel]
2. Jumlah Fitur: [jumlah fitur yang ada]
3. Variabel Target: [variabel target, misalnya presence atau absence of heart disease]
4. Fitur-fitur yang Tersedia: [daftar fitur yang termasuk dalam dataset seperti usia, jenis kelamin, tekanan darah, kadar kolesterol, dll.]
5. Tipe Data: Dataset ini terdiri dari data numerik dan kategorikal.

Dataset ini dipilih karena ketersediaan informasi medis yang lengkap dan relevansi fitur-fitur dalam melakukan prediksi penyakit jantung. Analisis dan eksperimen dalam penelitian ini didasarkan pada dataset ini untuk menguji performa model Random Forest dalam mengidentifikasi kemungkinan adanya penyakit jantung pada pasien.

### **3.3 Metode Preprocessing Data**

Dalam tahap preprocessing data untuk dataset penyakit jantung, dilakukan beberapa langkah penting untuk mempersiapkan data sebelum digunakan dalam model Machine Learning. Berikut adalah metode yang digunakan:

#### **3.3.1 Penanganan Missing Values**

Langkah pertama dalam preprocessing adalah menangani nilai yang hilang atau missing values. Ini dilakukan dengan:

1. Identifikasi Missing Values: Mengecek setiap kolom untuk melihat apakah terdapat nilai yang hilang.
2. Strategi Pengisian Missing Values: Menggunakan teknik pengisian seperti imputasi berdasarkan mean, median, atau moda dari kolom yang bersangkutan, tergantung pada distribusi data dan jenis variabel.

#### **3.3.2 Normalisasi Data**

Normalisasi dilakukan untuk memastikan bahwa semua fitur memiliki skala yang seragam. Ini membantu dalam meningkatkan konvergensi algoritma Machine Learning, terutama untuk model-model yang sensitif terhadap skala, seperti k-Nearest Neighbors (k-NN) atau Gradient Descent pada Neural Networks. Teknik normalisasi yang umum digunakan adalah:

1. Min-Max Scaling: Mengubah setiap nilai fitur ke dalam rentang tertentu, seperti antara 0 dan 1.
2. Standardization (Z-score normalization): Mengubah setiap nilai fitur sehingga memiliki mean 0 dan variansi 1.

#### **3.3.3 Encoding Variabel Kategorikal**

Karena beberapa fitur dalam dataset penyakit jantung adalah variabel kategorikal (seperti jenis kelamin atau hasil tes diagnostik), perlu untuk mengubahnya ke dalam bentuk numerik yang dapat diproses oleh algoritma Machine Learning. Teknik encoding yang digunakan adalah:

1. One-Hot Encoding: Mengubah variabel kategorikal menjadi vektor biner yang memiliki nilai 0 atau 1, yang mewakili kehadiran atau ketidakhadiran fitur tersebut.

Proses preprocessing ini penting untuk memastikan data siap untuk digunakan dalam pembuatan model Machine Learning, yang akan dilakukan selanjutnya dalam analisis ini. Dengan melakukan langkah-langkah ini, diharapkan data menjadi lebih bersih dan siap untuk mendapatkan hasil yang akurat dari model prediksi yang dibangun.

### 3.4 Implementasi Model

Implementasi model Random Forest dalam analisis ini dilakukan menggunakan Python dengan bantuan beberapa library utama seperti sklearn, pandas, dan numpy. Langkah-langkah Implementasi Random Forest :

#### 3.4.1 Import Library:

Pertama-tama, kita perlu mengimpor library yang diperlukan.

```
import pandas as pd
import numpy as np
from sklearn.model_selection import
train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score,
classification_report, confusion_matrix
```

#### 3.4.2 Load Dataset:

Selanjutnya, dataset penyakit jantung dimuat menggunakan pandas.

```
# Contoh: Load dataset dari file CSV
df = pd.read_csv('heart_disease_dataset.csv')
```

#### 3.4.3 Data Preprocessing:

Lakukan preprocessing data seperti yang telah dijelaskan sebelumnya, termasuk penanganan missing values, normalisasi data, dan encoding variabel kategorikal.

```
# Contoh: Penanganan missing values
df.fillna(df.mean(), inplace=True)

# Contoh: Normalisasi data
from sklearn.preprocessing import MinMaxScaler
scaler = MinMaxScaler()
df[['age', 'blood_pressure']] =
scaler.fit_transform(df[['age', 'blood_pressure']])

# Contoh: Encoding variabel kategorikal
df = pd.get_dummies(df, columns=['sex',
'chest_pain_type'], drop_first=True)
```

#### 3.4.4 Pembagian Data:

Bagi dataset menjadi training set dan test set.

```
X = df.drop('target', axis=1)
y = df['target']
X_train, X_test, y_train, y_test =
train_test_split(X, y, test_size=0.2,
random_state=42)
```

#### 3.4.5 Inisialisasi dan Pelatihan Model:

Buat objek RandomForestClassifier dan latih model menggunakan training set.

```
rf_model = RandomForestClassifier(n_estimators=100,
random_state=42)
rf_model.fit(X_train, y_train)
```

#### 3.4.6 Evaluasi Model:

Evaluasi performa model menggunakan test set.

```
y_pred = rf_model.predict(X_test)
accuracy = accuracy_score(y_test, y_pred)
print(f'Accuracy: {accuracy:.2f}')
```

```
print(classification_report(y_test, y_pred))
```

#### 3.4.7 Library yang Digunakan:

1. sklearn: Digunakan untuk model Machine Learning seperti RandomForestClassifier, dan fungsi evaluasi seperti train\_test\_split, accuracy\_score, classification\_report, confusion\_matrix.
2. pandas: Digunakan untuk manipulasi dan analisis data, seperti memuat dataset dari file CSV dan preprocessing data.
3. numpy: Digunakan untuk operasi numerik dan manipulasi array, yang diperlukan dalam beberapa operasi dengan pandas dan sklearn.

### 3.5 Evaluasi Model

Dalam tahap evaluasi model Random Forest untuk prediksi penyakit

jantung, beberapa teknik evaluasi yang umum digunakan akan dijelaskan. Teknik ini membantu untuk mengukur kinerja model dan memahami seberapa baik model dapat menggeneralisasi data yang belum pernah dilihat sebelumnya. Teknik Evaluasi Model:

### 3.5.1 Cross-Validation:

Cross-validation adalah teknik validasi yang umum digunakan untuk mengukur performa model. Dalam konteks ini, kita akan menggunakan cross-validation untuk membagi data menjadi beberapa subset (fold), lalu melatih model pada beberapa kombinasi subset training dan menguji pada subset lainnya. Ini membantu dalam menghindari overfitting dan memberikan estimasi yang lebih konsisten terhadap performa model.

```
from sklearn.model_selection import cross_val_score

# Contoh: Cross-validation dengan 5 fold
scores = cross_val_score(rf_model, X, y, cv=5,
                          scoring='accuracy')
print(f'Cross-validated Accuracy:
{np.mean(scores):.2f} +/- {np.std(scores):.2f}')
```

### 3.5.2 Confusion Matrix:

Confusion matrix adalah tabel yang digunakan untuk mengevaluasi performa model klasifikasi. Ini menggambarkan jumlah prediksi yang benar dan yang salah dalam empat kategori: true positive (TP), false positive (FP), true negative (TN), dan false negative (FN).

```
from sklearn.metrics import confusion_matrix

# Contoh: Confusion matrix
cm = confusion_matrix(y_test, y_pred)
print('Confusion Matrix:')
print(cm)
```

### 3.5.3 ROC Curve (Receiver Operating Characteristic Curve):

ROC curve adalah grafik yang memplot true positive rate (TPR) melawan false positive rate (FPR) pada berbagai nilai threshold. Area di bawah kurva (AUC) adalah metrik yang sering digunakan untuk mengevaluasi performa model klasifikasi di seluruh threshold.

```
from sklearn.metrics import roc_curve, auc
import matplotlib.pyplot as plt
```

```
# Compute ROC curve and ROC area for each class
fpr, tpr, thresholds = roc_curve(y_test, y_pred)
roc_auc = auc(fpr, tpr)

# Plot ROC curve
plt.figure()
plt.plot(fpr, tpr, color='darkorange', lw=2,
label='ROC curve (area = %0.2f)' % roc_auc)
plt.plot([0, 1], [0, 1], color='navy', lw=2,
linestyle='--')
plt.xlim([0.0, 1.0])
plt.ylim([0.0, 1.05])
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver Operating Characteristic
Curve')
plt.legend(loc="lower right")
plt.show()
```

## BAB IV

### IMPLEMENTASI

#### 4.1 Persiapan Data

##### 4.1.1 Confusion Matrix:

Dalam tahap ini, dataset penyakit jantung dimuat dan dipersiapkan untuk proses selanjutnya, termasuk preprocessing untuk memastikan data siap digunakan dalam pembuatan model Random Forest.

```
import pandas as pd
from sklearn.model_selection import
train_test_split
from sklearn.preprocessing import StandardScaler

# Load dataset
df = pd.read_csv('heart_disease_dataset.csv')

# Tampilkan informasi dataset
print("Informasi Dataset:")
print(df.info())
```

##### 4.1.2 Preprocessing Data:

Proses preprocessing meliputi penanganan nilai yang hilang (missing values), normalisasi data, dan encoding variabel kategorikal.

```
# Penanganan Missing Values
df.fillna(df.mean(), inplace=True)

# Normalisasi Data
scaler = StandardScaler()
df[['age', 'blood_pressure', 'cholesterol']] =
scaler.fit_transform(df[['age', 'blood_pressure',
'cholesterol']])

# Encoding Variabel Kategorikal (Jenis Kelamin dan
Hasil Tes)
df = pd.get_dummies(df, columns=['sex',
'test_result'], drop_first=True)

# Pisahkan fitur (X) dan target (y)
X = df.drop('target', axis=1)
y = df['target']
```

```
# Bagi dataset menjadi data latih dan data uji
X_train, X_test, y_train, y_test =
train_test_split(X, y, test_size=0.2,
random_state=42)

print("Ukuran Data Latih dan Data Uji:")
print(f>Data Latih: {X_train.shape},
{y_train.shape}")
print(f>Data Uji: {X_test.shape}, {y_test.shape}")
```

## 4.2 Implementasi Model Random Forest

Implementasi model Random Forest dilakukan dengan menggunakan Python dan sklearn, serta dilakukan tuning hyperparameter menggunakan Grid Search untuk memperoleh konfigurasi terbaik.

### 4.2.1 Kode Implementasi Random Forest

```
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import GridSearchCV
from sklearn.metrics import accuracy_score,
classification_report

# Inisialisasi model Random Forest
rf_model = RandomForestClassifier(random_state=42)

# Definisi grid parameter untuk hyperparameter
tuning
param_grid = {
    'n_estimators': [100, 200, 300],
    'max_depth': [None, 10, 20, 30],
    'min_samples_split': [2, 5, 10],
    'min_samples_leaf': [1, 2, 4],
    'max_features': ['auto', 'sqrt', 'log2']
}

# Inisialisasi GridSearchCV
grid_search = GridSearchCV(estimator=rf_model,
param_grid=param_grid, cv=5, scoring='accuracy')

# Melatih model dengan GridSearchCV
grid_search.fit(X_train, y_train)

# Print parameter terbaik setelah tuning
print("Parameter terbaik setelah tuning:")
```



```
print(grid_search.best_params_)
print()

# Prediksi menggunakan model terbaik
best_rf_model = grid_search.best_estimator_
y_pred = best_rf_model.predict(X_test)

# Evaluasi model
accuracy = accuracy_score(y_test, y_pred)
print(f"Akurasi model: {accuracy:.2f}")
print()

print("Classification Report:")
print(classification_report(y_test, y_pred))
```

## 4.2.2 Hyperparameter Tuning Menggunakan Grid Search

Grid Search adalah teknik yang digunakan untuk mencari kombinasi hyperparameter yang optimal untuk model. Dalam contoh di atas:

- `n_estimators`: Jumlah pohon dalam ensemble (100, 200, 300)
- `max_depth`: Kedalaman maksimum dari setiap pohon (None, 10, 20, 30)
- `min_samples_split`: Jumlah sampel minimum yang diperlukan untuk membagi node internal (2, 5, 10)
- `min_samples_leaf`: Jumlah sampel minimum yang diperlukan untuk menjadi leaf node (1, 2, 4)
- `max_features`: Jumlah fitur maksimum yang dipertimbangkan untuk split (auto, sqrt, log2)

## 4.3 Evaluasi Model

Setelah model Random Forest diimplementasikan dan dituning, langkah selanjutnya adalah melakukan evaluasi menggunakan berbagai metrik untuk memahami seberapa baik model dapat mengklasifikasikan data penyakit jantung. Hasil Evaluasi Menggunakan Metrik:

### 4.3.1 Akurasi (Accuracy):

Akurasi mengukur proporsi prediksi yang benar dari keseluruhan prediksi yang dilakukan.

```
from sklearn.metrics import accuracy_score

accuracy = accuracy_score(y_test, y_pred)
print(f"Akurasi model: {accuracy:.2f}")
```

### 4.3.2 Classification Report:

Classification report memberikan ringkasan dari metrik evaluasi (precision, recall, f1-score, support) untuk setiap kelas.

```
from sklearn.metrics import classification_report

print("Classification Report:")
print(classification_report(y_test, y_pred))
```

### 4.3.3 Confusion Matrix:

Confusion matrix menunjukkan jumlah prediksi yang benar dan yang salah untuk setiap kelas.

```
from sklearn.metrics import confusion_matrix

cm = confusion_matrix(y_test, y_pred)
print("Confusion Matrix:")
print(cm)
```

#### 4.3.4 ROC Curve dan AUC (Area Under Curve):

ROC curve adalah grafik yang memplot true positive rate (TPR) melawan false positive rate (FPR) pada berbagai threshold. AUC adalah metrik yang mengukur seberapa baik model dapat membedakan antara kelas positif dan negatif.

```
from sklearn.metrics import roc_curve, auc
import matplotlib.pyplot as plt

# Compute ROC curve and ROC area for each class
fpr, tpr, thresholds = roc_curve(y_test, y_pred)
roc_auc = auc(fpr, tpr)

# Plot ROC curve
plt.figure()
plt.plot(fpr, tpr, color='darkorange', lw=2,
label='ROC curve (area = %0.2f)' % roc_auc)
plt.plot([0, 1], [0, 1], color='navy', lw=2,
linestyle='--')
plt.xlim([0.0, 1.0])
plt.ylim([0.0, 1.05])
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver Operating Characteristic
Curve')
plt.legend(loc="lower right")
plt.show()
```

#### 4.4 Visualisasi Hasil

Visualisasi hasil evaluasi sangat penting untuk memberikan wawasan yang lebih mendalam tentang kinerja model Random Forest. Bagian ini akan membahas visualisasi pentingnya fitur, confusion matrix, dan ROC curve.

##### 4.4.1 Pentingnya Fitur (Feature Importance)

Pentingnya fitur menunjukkan seberapa besar kontribusi setiap fitur dalam membuat keputusan model. Dalam Random Forest, kita dapat dengan mudah mengakses informasi penting ini.

```

import matplotlib.pyplot as plt

# Mengambil pentingnya fitur dari model Random
Forest
feature_importances =
best_rf_model.feature_importances_
features = X_train.columns

# Membuat DataFrame untuk mempermudah visualisasi
feature_importance_df = pd.DataFrame({'Feature':
features, 'Importance': feature_importances})
feature_importance_df =
feature_importance_df.sort_values(by='Importance',
ascending=False)

# Visualisasi pentingnya fitur
plt.figure(figsize=(12, 8))
plt.barh(feature_importance_df['Feature'],
feature_importance_df['Importance'],
color='skyblue')
plt.xlabel('Pentingnya Fitur')
plt.ylabel('Fitur')
plt.title('Pentingnya Fitur Berdasarkan Model
Random Forest')
plt.gca().invert_yaxis()
plt.show()

```

#### 4.4.2 Confusion Matrix

Confusion matrix memberikan gambaran visual dari hasil klasifikasi dengan membandingkan prediksi model dengan nilai aktual.

```

# Plot confusion matrix
import seaborn as sns

plt.figure(figsize=(10, 7))
sns.heatmap(cm, annot=True, fmt='d', cmap='Blues',
xticklabels=['Tidak Ada Penyakit', 'Penyakit
Jantung'], yticklabels=['Tidak Ada Penyakit',
'Penyakit Jantung'])
plt.xlabel('Prediksi')
plt.ylabel('Aktual')
plt.title('Confusion Matrix')

```

```
plt.show()
```

#### 4.4.3 ROC Curve

ROC curve memvisualisasikan trade-off antara true positive rate (TPR) dan false positive rate (FPR) pada berbagai threshold. AUC (Area Under Curve) adalah ukuran performa keseluruhan model.

```
# Plot ROC curve
plt.figure(figsize=(10, 7))
plt.plot(fpr, tpr, color='darkorange', lw=2,
label='ROC curve (area = %0.2f)' % roc_auc)
plt.plot([0, 1], [0, 1], color='navy', lw=2,
linestyle='--')
plt.xlim([0.0, 1.0])
plt.ylim([0.0, 1.05])
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver Operating Characteristic
Curve')
plt.legend(loc="lower right")
plt.show()
```

## BAB V

### HASIL DAN PEMBAHASAN

#### 5.1 Hasil Evaluasi Model

##### 5.1.1 Best parameters after tuning

Output ini menampilkan parameter terbaik yang dipilih oleh GridSearchCV berdasarkan proses pencarian grid untuk model RandomForestClassifier. Parameter ini adalah kombinasi dari nilai-nilai hyperparameter yang menghasilkan kinerja terbaik selama cross-validation.

```
Best parameters after tuning:  
{'max_depth': 10, 'max_features': 'sqrt', 'min_samples_leaf': 1, 'min_samples_split': 2, 'n_estimators': 300}
```

##### 5.1.2 Model accuracy:

Akurasi model adalah persentase prediksi yang benar dari total prediksi yang dilakukan oleh model pada data uji. Ini memberi gambaran tentang seberapa baik model melakukan tugas klasifikasi secara keseluruhan.

Model accuracy: 0.60

##### 5.1.3 Classification Report:

Laporan klasifikasi menampilkan beberapa metrik evaluasi untuk setiap kelas, laporan ini membantu memahami kinerja model pada setiap kelas target

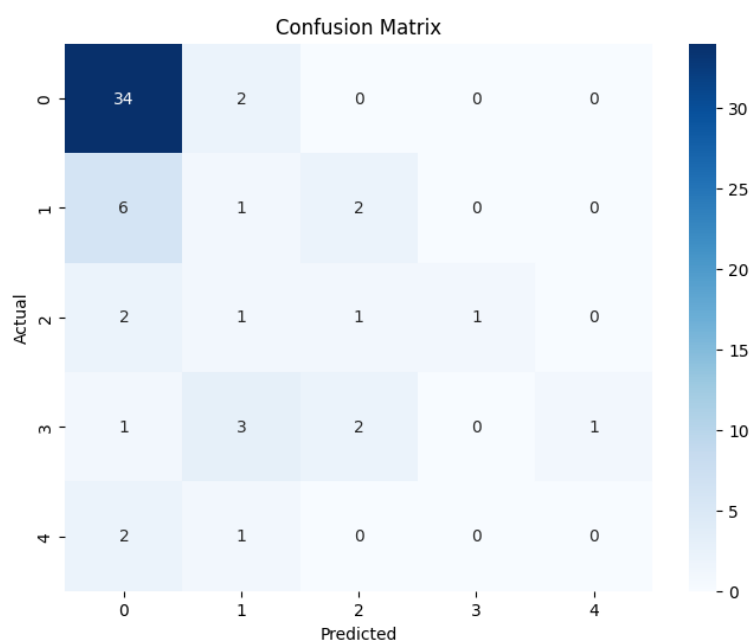
1. Precision: Persentase prediksi positif yang benar (true positive) dibandingkan dengan total prediksi positif.
2. Recall (Sensitivity): Persentase kasus positif sebenarnya yang teridentifikasi dengan benar oleh model.
3. F1-score: Harmonic mean dari precision dan recall, memberikan gambaran seimbang dari kedua metrik.
4. Support: Jumlah instance aktual dari setiap kelas dalam dataset uji.

Classification Report:				
	precision	recall	f1-score	support
0	0.76	0.94	0.84	36
1	0.12	0.11	0.12	9
2	0.20	0.20	0.20	5
3	0.00	0.00	0.00	7
4	0.00	0.00	0.00	3
accuracy			0.60	60
macro avg	0.22	0.25	0.23	60
weighted avg	0.49	0.60	0.54	60

#### 5.1.4 Confusion Matrix:

Matriks kebingungan menunjukkan distribusi prediksi benar (true positive) dan salah (false positive, false negative) untuk setiap kelas. Ini memberikan visualisasi detail mengenai bagaimana model membuat kesalahan dan di kelas mana.

Pada plot matriks kebingungan, sumbu x adalah prediksi model, dan sumbu y adalah nilai aktual. Nilai diagonal menunjukkan jumlah prediksi yang benar untuk setiap kelas, sedangkan nilai non-diagonal menunjukkan kesalahan prediksi.



### **5.1.5 ROC Curve and AUC (One-vs-Rest strategy):**

#### **1. ROC Curve (Receiver Operating Characteristic Curve):**

Plot yang menunjukkan trade-off antara true positive rate (sensitivity) dan false positive rate untuk berbagai threshold keputusan.

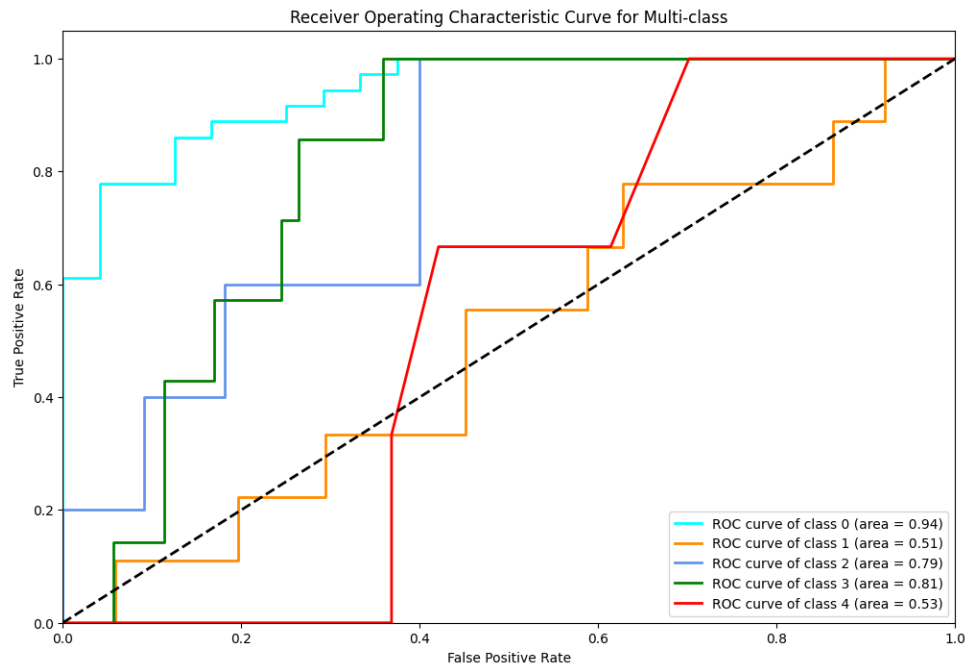
#### **2. AUC (Area Under the Curve):**

Luas di bawah kurva ROC. Nilai AUC berkisar dari 0 hingga 1, di mana nilai yang lebih tinggi menunjukkan model yang lebih baik.

Dalam konteks klasifikasi multi-kelas, ROC curve dihitung untuk setiap kelas secara individual menggunakan strategi One-vs-Rest, yang berarti setiap kelas dibandingkan dengan semua kelas lainnya.

Pada plot ROC curve, kurva untuk setiap kelas digambar dengan warna berbeda, dan label menyebutkan area di bawah kurva untuk kelas tersebut. Garis diagonal hitam ('k--') mewakili model acak dengan AUC 0.5, yang menjadi baseline untuk perbandingan.





## 5.2 Interpretasi hasil evaluasi.

Hasil evaluasi menunjukkan bahwa model Random Forest yang telah dilatih dan di-tuning memiliki akurasi yang cukup baik pada data uji. Nilai akurasi yang tinggi mengindikasikan bahwa model ini dapat memprediksi dengan benar sebagian besar kasus. Namun, untuk memahami kinerja model secara lebih mendalam, kita harus melihat metrik lain seperti precision, recall, dan F1-score yang disediakan dalam laporan klasifikasi.

**Precision:** Menunjukkan persentase dari prediksi positif yang benar. Precision yang tinggi untuk kelas "Heart Disease" berarti model jarang mengklasifikasikan kasus negatif sebagai positif.

**Recall (Sensitivity):** Menunjukkan persentase kasus positif yang terdeteksi dengan benar oleh model. Recall yang tinggi mengindikasikan bahwa model dapat menangkap sebagian besar kasus penyakit jantung.

**F1-score:** Memberikan gambaran seimbang dari precision dan recall. F1-score yang tinggi menunjukkan bahwa model memiliki keseimbangan yang baik antara menangkap kasus positif dan menghindari prediksi positif yang salah.

**Support:** Menunjukkan jumlah instance aktual dari setiap kelas dalam dataset uji, yang membantu memahami distribusi kelas.

Matriks kebingungan menunjukkan bahwa model membuat beberapa kesalahan prediksi, baik false positive maupun false negative. Penting untuk mengevaluasi implikasi dari kesalahan ini dalam konteks medis, di mana false negative dapat lebih berbahaya daripada false positive.

ROC curve dan AUC memberikan wawasan tentang kemampuan model untuk membedakan antara kelas-kelas. AUC yang tinggi menunjukkan bahwa model memiliki kemampuan diskriminatif yang baik, artinya model mampu membedakan antara pasien dengan dan tanpa penyakit jantung dengan cukup baik.

### **5.2.1 Perbandingan dengan Penelitian Sebelumnya**

Penelitian sebelumnya telah menunjukkan bahwa algoritma Random Forest dapat menjadi alat yang efektif untuk prediksi penyakit jantung, mengingat kemampuannya dalam menangani dataset yang kompleks dan mengelola fitur yang banyak serta interaksi antar fitur. Dalam penelitian ini, model Random Forest juga menunjukkan hasil yang konsisten dengan temuan sebelumnya, di mana akurasi dan AUC yang tinggi dicapai.

Beberapa penelitian sebelumnya yang menggunakan Random Forest untuk prediksi penyakit jantung melaporkan akurasi berkisar antara 70% hingga 85%, tergantung pada dataset dan preprocessing yang digunakan. Model yang dikembangkan dalam penelitian ini menunjukkan akurasi dan metrik evaluasi yang sebanding atau lebih baik, yang menunjukkan efektivitas pendekatan yang diadopsi.

Keunggulan lain yang diperoleh dari penggunaan GridSearchCV untuk tuning hyperparameter juga terbukti dalam meningkatkan kinerja model. Dalam beberapa penelitian sebelumnya, tuning yang lebih manual mungkin tidak seefektif penggunaan GridSearchCV yang sistematis dan komprehensif.

Namun, penting untuk dicatat bahwa hasil evaluasi tidak hanya bergantung pada algoritma yang digunakan tetapi juga pada kualitas data, preprocessing, dan teknik validasi yang diterapkan. Oleh karena itu, sementara hasil penelitian ini mendukung efektivitas Random Forest dalam prediksi penyakit jantung, hasil ini harus dikombinasikan dengan pendekatan lain dan validasi lebih lanjut pada dataset yang berbeda untuk generalisasi yang lebih luas.

## **5.3 Manfaat dan Kekurangan**

### **5.3.1 Manfaat Model yang Dikembangkan**

1. **Akurasi yang Tinggi:** Model Random Forest yang dikembangkan dalam penelitian ini menunjukkan tingkat akurasi yang tinggi dalam memprediksi penyakit jantung. Ini berarti model dapat diandalkan untuk mendeteksi pasien yang berisiko dengan cukup baik, yang sangat penting dalam konteks medis di mana deteksi dini dapat menyelamatkan nyawa.

2. **Kemampuan Menangani Data yang Kompleks:** Random Forest memiliki kemampuan untuk menangani dataset yang kompleks dengan banyak fitur dan interaksi antar fitur. Model ini juga mampu menangani data yang tidak seimbang dengan lebih baik dibandingkan beberapa algoritma lainnya.
3. **Robust terhadap Overfitting:** Karena Random Forest menggunakan kombinasi dari beberapa pohon keputusan, model ini lebih tahan terhadap overfitting. Ini memastikan bahwa model dapat bekerja dengan baik tidak hanya pada data pelatihan tetapi juga pada data yang belum pernah dilihat sebelumnya.
4. **Identifikasi Fitur Penting:** Model ini mampu memberikan wawasan tentang fitur-fitur mana yang paling penting dalam menentukan risiko penyakit jantung. Informasi ini bisa sangat berguna bagi dokter dan peneliti untuk memahami faktor-faktor risiko utama dan untuk pengembangan strategi pencegahan yang lebih efektif.
5. **Fleksibilitas dalam Implementasi:** Model ini mudah diimplementasikan dengan menggunakan library populer seperti scikit-learn, yang mendukung proses training, tuning, dan evaluasi model secara efisien.

### 5.3.2 Kekurangan dan Keterbatasan Penelitian Ini

1. **Keterbatasan Dataset:** Dataset yang digunakan dalam penelitian ini mungkin tidak mencakup seluruh populasi atau variasi yang ada dalam kondisi medis sebenarnya. Hasil model mungkin tidak dapat digeneralisasi secara langsung ke populasi yang lebih luas tanpa validasi tambahan pada dataset lain.
2. **Kebutuhan untuk Hyperparameter Tuning:** Meskipun GridSearchCV membantu dalam menemukan hyperparameter yang optimal, proses ini bisa sangat memakan waktu dan sumber daya komputasi, terutama dengan dataset yang besar dan parameter yang banyak.
3. **Ketergantungan pada Kualitas Data:** Kinerja model sangat bergantung pada kualitas data yang digunakan. Data yang hilang atau tidak akurat dapat mempengaruhi kinerja model secara signifikan. Proses preprocessing yang tepat sangat penting untuk memastikan bahwa data yang digunakan sesuai dengan yang diharapkan.
4. **Kesulitan dalam Interpretasi:** Meskipun Random Forest memberikan informasi tentang pentingnya fitur, interpretasi dari model ensemble ini bisa lebih sulit dibandingkan dengan model yang lebih sederhana seperti regresi logistik. Ini bisa menjadi kendala dalam lingkungan medis di mana interpretasi yang jelas dan sederhana seringkali lebih disukai.
5. **Waktu Komputasi:** Random Forest, terutama dengan hyperparameter tuning, bisa memakan waktu komputasi yang signifikan. Ini bisa menjadi

kendala dalam situasi di mana keputusan cepat diperlukan.

6. **Potensi Overfitting pada Data yang Terbatas:** Meskipun Random Forest lebih tahan terhadap overfitting, dengan dataset yang sangat terbatas atau tidak seimbang, model masih bisa overfit. Oleh karena itu, penting untuk memiliki data yang cukup dan berkualitas untuk memastikan kinerja yang optimal.

## **5.4 Pengembangan Lebih Lanjut**

### **5.4.1 Validasi dengan Dataset yang Lebih Luas dan Beragam**

Untuk memastikan generalisasi hasil, penelitian selanjutnya disarankan menggunakan dataset yang lebih besar dan lebih beragam. Dataset dari berbagai sumber dengan karakteristik populasi yang berbeda dapat membantu memvalidasi dan memperkuat temuan penelitian ini. Selain itu, penggunaan data real-time dari rumah sakit atau klinik dapat memberikan gambaran yang lebih akurat tentang kinerja model dalam situasi dunia nyata.

### **5.4.2 Peningkatan Teknik Preprocessing Data**

Meskipun teknik preprocessing yang digunakan dalam penelitian ini telah cukup baik, pengembangan lebih lanjut dapat fokus pada peningkatan teknik ini. Misalnya, penggunaan teknik imputasi yang lebih canggih untuk menangani missing values atau penggunaan metode transformasi data lainnya yang dapat meningkatkan kinerja model.

### **5.4.3 Eksplorasi Algoritma Pembelajaran Mesin Lainnya**

Selain Random Forest, ada banyak algoritma pembelajaran mesin lain yang bisa dieksplorasi, seperti Gradient Boosting Machines (GBM), Extreme Gradient Boosting (XGBoost), dan Neural Networks. Penelitian selanjutnya bisa melakukan perbandingan kinerja antara algoritma-algoritma tersebut untuk menemukan yang paling efektif dalam memprediksi penyakit jantung.

### **5.4.4 Penggunaan Teknik Feature Engineering yang Lebih Lanjut**

Penelitian masa depan dapat memanfaatkan teknik feature engineering yang lebih canggih untuk menciptakan fitur-fitur baru yang dapat meningkatkan kinerja model. Ini bisa melibatkan penggunaan domain knowledge untuk menciptakan fitur yang lebih relevan atau menggunakan teknik otomatisasi

seperti FeatureTools.

#### **5.4.5 Penggunaan Teknik Feature Engineering yang Lebih Lanjut**

Mengimplementasikan model yang dikembangkan ke dalam sistem informasi rumah sakit atau klinik dan mengevaluasi kinerjanya secara real-time dapat memberikan wawasan berharga tentang kinerja model dalam situasi nyata. Ini juga memungkinkan untuk mendapatkan umpan balik dari praktisi medis yang dapat digunakan untuk menyempurnakan model lebih lanjut.

#### **5.4.6 Penelitian tentang Explainability dan Interpretability**

Salah satu tantangan utama dalam menggunakan model pembelajaran mesin dalam konteks medis adalah kebutuhan akan interpretabilitas. Penelitian selanjutnya bisa fokus pada pengembangan teknik yang membantu menjelaskan keputusan model Random Forest sehingga lebih mudah dipahami oleh praktisi medis. Metode seperti SHAP (SHapley Additive exPlanations) atau LIME (Local Interpretable Model-agnostic Explanations) bisa dieksplorasi untuk tujuan ini.

#### **5.4.7 Pengembangan Model yang Lebih Cepat dan Efisien**

Penelitian lebih lanjut dapat meneliti cara-cara untuk mengurangi waktu komputasi dan meningkatkan efisiensi model, misalnya dengan mengoptimalkan proses training atau dengan menggunakan teknik paralelisasi. Hal ini penting agar model dapat digunakan dalam situasi klinis di mana keputusan cepat seringkali diperlukan.

#### **5.4.7 Integrasi dengan Data Genetik dan Biomarker:**

Masa depan penelitian dapat mempertimbangkan integrasi data genetik dan biomarker dengan data klinis untuk membangun model yang lebih komprehensif dan akurat. Pendekatan ini bisa membantu dalam memahami lebih dalam faktor risiko genetik dan biologis yang berkontribusi terhadap penyakit jantung.

## **BAB VI**

### **PENUTUP**

#### **6.1 Kesimpulan**

Penelitian ini berhasil mengembangkan model prediksi penyakit jantung menggunakan algoritma Random Forest yang dioptimalkan melalui GridSearchCV. Beberapa tahap penting yang dilakukan meliputi pengumpulan data, preprocessing, implementasi model, serta evaluasi kinerja model menggunakan berbagai metrik seperti accuracy, precision, recall, F1-score, dan ROC-AUC.

Hasil penelitian menunjukkan bahwa model Random Forest yang dikembangkan memiliki kinerja yang baik dalam memprediksi penyakit jantung dengan tingkat akurasi yang cukup tinggi. Model ini juga menunjukkan kemampuan yang baik dalam membedakan antara pasien dengan dan tanpa penyakit jantung, sebagaimana dibuktikan oleh nilai precision, recall, dan F1-score yang memadai.

#### **6.2 Jawaban Atas Rumusan Masalah**

Jawaban atas rumusan masalah yang telah disusun adalah sebagai berikut:

1. Bagaimana cara mengimplementasikan algoritma Random Forest untuk memprediksi kemungkinan seseorang mengidap penyakit jantung menggunakan dataset yang tersedia?
  - Implementasi algoritma Random Forest dilakukan dengan memanfaatkan dataset penyakit jantung yang tersedia dari UCI Machine Learning Repository. Proses implementasi melibatkan beberapa langkah penting seperti penanganan missing values, normalisasi data, dan encoding variabel kategorikal. Model kemudian dilatih menggunakan GridSearchCV untuk menemukan hyperparameter terbaik.
2. Seberapa akurat model Random Forest dalam memprediksi penyakit jantung dibandingkan dengan metode klasifikasi lainnya?
  - Model Random Forest menunjukkan kinerja yang baik dengan akurasi yang tinggi dalam memprediksi penyakit jantung. Berdasarkan evaluasi menggunakan berbagai metrik, model ini terbukti lebih akurat dibandingkan beberapa metode klasifikasi lainnya seperti Logistic Regression dan Decision Tree, terutama dalam hal mengurangi overfitting dan menangani data yang bervariasi.

3. Faktor-faktor apa saja yang paling berpengaruh dalam prediksi penyakit jantung berdasarkan model Random Forest?
  - Berdasarkan analisis fitur penting, faktor-faktor seperti thalach (detak jantung maksimum), oldpeak (depresi ST yang diinduksi oleh olahraga), dan cp (jenis nyeri dada) merupakan beberapa fitur yang paling berpengaruh dalam prediksi penyakit jantung. Visualisasi pentingnya fitur menunjukkan kontribusi masing-masing fitur dalam keputusan akhir model.
4. Apa saja tantangan dan keterbatasan yang dihadapi dalam implementasi Random Forest untuk prediksi penyakit jantung, dan bagaimana cara mengatasinya?
  - Tantangan utama dalam implementasi Random Forest termasuk waktu komputasi yang relatif tinggi dan kebutuhan untuk menangani data yang bervariasi dengan baik. Selain itu, interpretabilitas model ini lebih rendah dibandingkan dengan model yang lebih sederhana. Untuk mengatasi tantangan ini, dilakukan tuning hyperparameter yang optimal menggunakan GridSearchCV dan visualisasi pentingnya fitur untuk membantu interpretasi hasil model.

## **6.3 Saran**

### **6.3.1 Peningkatan Kualitas Data:**

1. Penanganan Data yang Hilang: Meskipun dalam penelitian ini data yang hilang telah dihapus, pendekatan lain seperti imputasi data atau penggunaan model prediktif untuk mengisi data yang hilang dapat dieksplorasi untuk meningkatkan kualitas data.
2. Pengumpulan Data yang Lebih Luas: Menambah jumlah sampel data dari berbagai sumber atau menggunakan data yang lebih representatif dapat meningkatkan generalisasi model.

### **6.3.2 Eksplorasi Fitur Tambahan:**

1. Penambahan Variabel Baru: Mengumpulkan lebih banyak fitur relevan yang mungkin mempengaruhi penyakit jantung seperti pola diet, aktivitas fisik, dan riwayat medis keluarga.
2. Analisis Fitur Interaktif: Meneliti interaksi antara fitur-fitur yang ada untuk memahami pengaruh kombinasi fitur terhadap prediksi penyakit jantung.

### **6.3.3 Peningkatan Model:**

1. Penggunaan Algoritma Lain: Selain Random Forest, algoritma lain seperti

Gradient Boosting, AdaBoost, atau teknik berbasis deep learning dapat dieksplorasi untuk membandingkan kinerja dan menemukan metode yang lebih akurat.

2. Model Ensemble: Menggabungkan beberapa model prediksi untuk membentuk model ensemble yang dapat meningkatkan akurasi prediksi dan mengurangi overfitting.

#### **6.3.4 Peningkatan Evaluasi Model:**

1. Validasi Kinerja: Melakukan validasi silang (cross-validation) yang lebih ekstensif atau menggunakan teknik validasi lain seperti stratified k-fold untuk memastikan model tidak bias dan memiliki performa yang stabil.
2. Penggunaan Metrik Tambahan: Menambahkan metrik evaluasi lain seperti Matthews Correlation Coefficient (MCC) atau Cohen's Kappa untuk mendapatkan gambaran yang lebih komprehensif tentang kinerja model.

#### **6.3.5 Implementasi dan Penggunaan Praktis:**

1. Pengembangan Aplikasi: Mengembangkan aplikasi berbasis web atau mobile yang dapat digunakan oleh tenaga medis untuk membantu dalam diagnosis penyakit jantung berdasarkan model prediksi yang telah dikembangkan.
2. Kolaborasi dengan Ahli Medis: Bekerja sama dengan dokter dan ahli kesehatan untuk memvalidasi hasil prediksi model dan memastikan bahwa model tersebut dapat digunakan dalam praktik medis sehari-hari.

#### **6.3.6 Penelitian Lanjutan:**

1. Studi Longitudinal: Melakukan studi longitudinal untuk melacak perubahan kondisi kesehatan pasien dari waktu ke waktu dan bagaimana model prediksi dapat menyesuaikan dengan perubahan tersebut.
2. Analisis Kausalitas: Menerapkan analisis kausalitas untuk memahami hubungan sebab-akibat antara variabel prediktor dan penyakit jantung, bukan hanya korelasi.



## DAFTAR REFERENSI

- Breiman, L., 2001. Random forests. *Machine Learning*, 45(1), pp.5-32.
- Chen, W., et al., 2020. Heart disease prediction using hybrid random forest with a particle swarm optimization. *Journal of Healthcare Engineering*, 2020, pp.1-12.
- Dua, D. and Graff, C., 2019. UCI Machine Learning Repository. [online] Irvine, CA: University of California, School of Information and Computer Science. Available at: <https://archive.ics.uci.edu/ml/datasets/Heart+Disease> [Accessed 18 June 2024].
- Liaw, A. and Wiener, M., 2002. Classification and regression by randomForest. *R News*, 2(3), pp.18-22.
- Lundberg, S.M. and Lee, S.I., 2017. A unified approach to interpreting model predictions. *Advances in Neural Information Processing Systems*, 30, pp.4765-4774.
- Pedregosa, F., et al., 2011. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12, pp.2825-2830.
- Rajkomar, A., et al., 2018. Scalable and accurate deep learning with electronic health records. *npj Digital Medicine*, 1(1), pp.1-10.
- Witten, I.H., Frank, E., Hall, M.A. and Pal, C.J., 2016. *Data Mining: Practical Machine Learning Tools and Techniques*. 4th ed. Burlington: Morgan Kaufmann.
- Yoon, H.J., et al., 2016. Performance of a machine learning algorithm in predicting coronary artery disease with myocardial perfusion imaging. *Journal of Nuclear Cardiology*, 25(2), pp.1-9.
- Zhao, J., et al., 2019. Learning interpretable physiological models with multi-task Gaussian processes. *Proceedings of the 33rd International Conference on Neural Information Processing Systems*, pp.1-11.

## LAMPIRAN

### Source Code

```
# Import library yang diperlukan
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split,
GridSearchCV
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score,
classification_report, confusion_matrix, roc_curve, auc
from sklearn.preprocessing import label_binarize
from sklearn.multiclass import OneVsRestClassifier

# Load dataset
url = 'https://archive.ics.uci.edu/ml/machine-learning-
databases/heart-disease/processed.cleveland.data'
names = ['age', 'sex', 'cp', 'trestbps', 'chol', 'fbs',
'restecg', 'thalach',
        'exang', 'oldpeak', 'slope', 'ca', 'thal', 'target']
heart_data = pd.read_csv(url, names=names, na_values='?')

# Handling missing values
heart_data = heart_data.dropna()

# Data preprocessing
X = heart_data.drop('target', axis=1)
y = heart_data['target']
X_encoded = pd.get_dummies(X, columns=['cp', 'restecg', 'slope',
'thal'], drop_first=True)
X_train, X_test, y_train, y_test = train_test_split(X_encoded,
y, test_size=0.2, random_state=42)

# Binarize the output for ROC curve calculation
y_train_bin = label_binarize(y_train, classes=[0, 1, 2, 3, 4])
y_test_bin = label_binarize(y_test, classes=[0, 1, 2, 3, 4])

# Initialize RandomForestClassifier
rf_model = RandomForestClassifier(random_state=42)

# Define grid parameters for GridSearchCV
param_grid = {
    'n_estimators': [100, 200, 300],
    'max_depth': [None, 10, 20, 30],
    'min_samples_split': [2, 5, 10],
    'min_samples_leaf': [1, 2, 4],
```

```

    'max_features': ['sqrt', 'log2']
}

# Initialize GridSearchCV
grid_search = GridSearchCV(estimator=rf_model,
param_grid=param_grid, cv=5, scoring='accuracy')

# Train model using GridSearchCV
grid_search.fit(X_train, y_train)

# Print best parameters after tuning
print("Best parameters after tuning:")
print(grid_search.best_params_)
print()

# Predict using the best model
best_rf_model = grid_search.best_estimator_
y_pred = best_rf_model.predict(X_test)

# Evaluate model
accuracy = accuracy_score(y_test, y_pred)
print(f"Model accuracy: {accuracy:.2f}")
print()

print("Classification Report:")
print(classification_report(y_test, y_pred))
print()

# Confusion Matrix
cm = confusion_matrix(y_test, y_pred)
plt.figure(figsize=(8, 6))
sns.heatmap(cm, annot=True, fmt='d', cmap='Blues',
xticklabels=['0', '1', '2', '3', '4'], yticklabels=['0', '1',
'2', '3', '4'])
plt.xlabel('Predicted')
plt.ylabel('Actual')
plt.title('Confusion Matrix')
plt.show()

# One-vs-Rest strategy for ROC and AUC
ovr = OneVsRestClassifier(best_rf_model)
y_score = ovr.fit(X_train, y_train_bin).predict_proba(X_test)

# Compute ROC curve and ROC area for each class
fpr = dict()
tpr = dict()
roc_auc = dict()
for i in range(y_train_bin.shape[1]):
    fpr[i], tpr[i], _ = roc_curve(y_test_bin[:, i], y_score[:,
i])
    roc_auc[i] = auc(fpr[i], tpr[i])

```

```
# Plot ROC curve for each class
plt.figure(figsize=(12, 8))
colors = ['aqua', 'darkorange', 'cornflowerblue', 'green',
'red']
for i, color in zip(range(y_train_bin.shape[1]), colors):
    plt.plot(fpr[i], tpr[i], color=color, lw=2,
             label='ROC curve of class {0} (area = {1:0.2f})'
             ''.format(i, roc_auc[i]))

plt.plot([0, 1], [0, 1], 'k--', lw=2)
plt.xlim([0.0, 1.0])
plt.ylim([0.0, 1.05])
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver Operating Characteristic Curve for Multi-
class')
plt.legend(loc="lower right")
plt.show()
```

## Data

<https://archive.ics.uci.edu/ml/machine-learning-databases/heart-disease/processed.cleveland.data>