

**PERBANDINGAN PERFORMA ALGORITMA RANDOM FOREST DAN
XGBOOST DALAM MENDETEKSI DIABETES BERDASARKAN DATA
REKAM MEDIS**



UMN
UNIVERSITAS
MULTIMEDIA
NUSANTARA

SKRIPSI

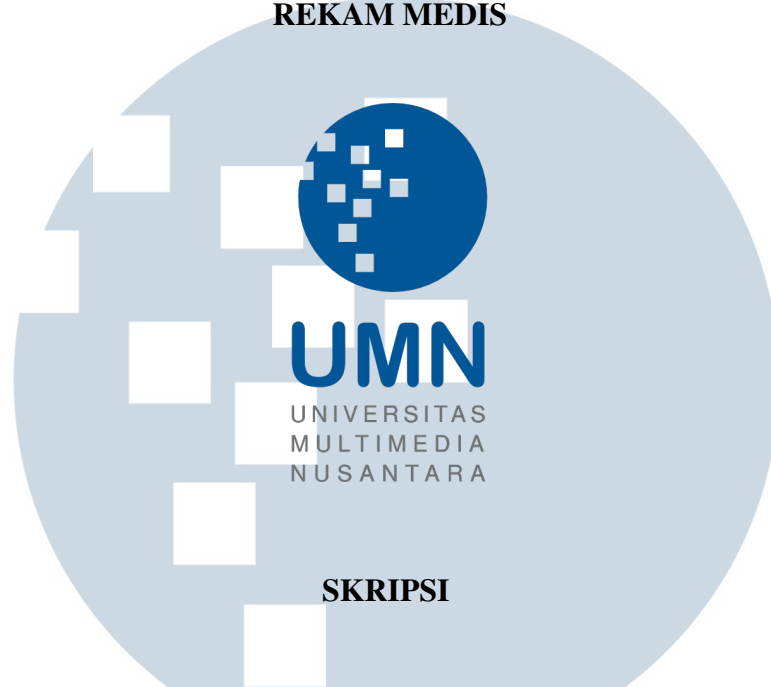
Shyehan Rafael Adlinugroho

00000052738

**PROGRAM STUDI INFORMATIKA
FAKULTAS TEKNIK DAN INFORMATIKA
UNIVERSITAS MULTIMEDIA NUSANTARA
TANGERANG**

2025

**PERBANDINGAN PERFORMA ALGORITMA RANDOM FOREST DAN
XGBOOST DALAM MENDETEKSI DIABETES BERDASARKAN DATA
REKAM MEDIS**



Diajukan sebagai salah satu syarat untuk memperoleh
Gelar Sarjana Komputer (S.Kom.)

Shyehan Rafael Adlinugroho

00000052738

UMN

UNIVERSITAS

MULTIMEDIA

NUSANTARA

**PROGRAM STUDI INFORMATIKA
FAKULTAS TEKNIK DAN INFORMATIKA
UNIVERSITAS MULTIMEDIA NUSANTARA**

TANGERANG

2025

HALAMAN PERNYATAAN TIDAK PLAGIAT

Dengan ini saya,

Nama : Shyehan Rafael Adlinugroho
Nomor Induk Mahasiswa : 00000052738
Program Studi : Informatika

Skripsi dengan judul:

Perbandingan Performa Algoritma Random Forest dan XGBoost dalam Mendeteksi Diabetes Berdasarkan Data Rekam Medis

merupakan hasil karya saya sendiri bukan plagiat dari laporan karya tulis ilmiah yang ditulis oleh orang lain, dan semua sumber, baik yang dikutip maupun dirujuk, telah saya nyatakan dengan benar serta dicantumkan di Daftar Pustaka.

Jika di kemudian hari terbukti ditemukan kecurangan/penyimpangan, baik dalam pelaksanaan maupun dalam penulisan laporan karya tulis ilmiah, saya bersedia menerima konsekuensi dinyatakan TIDAK LULUS untuk mata kuliah yang telah saya tempuh.

Tangerang, 3 Januari 2025



(Shyehan Rafael Adlinugroho)

UNIVERSITAS
MULTIMEDIA
NUSANTARA

HALAMAN PENGESAHAN

Skripsi dengan judul

PERBANDINGAN PERFORMA ALGORITMA RANDOM FOREST DAN XGBOOST DALAM MENDETEKSI DIABETES BERDASARKAN DATA REKAM MEDIS

oleh

Nama : Shyehan Rafael Adlinugroho
NIM : 00000052738
Program Studi : Informatika
Fakultas : Fakultas Teknik dan Informatika

Telah diujikan pada hari Rabu, 8 Januari 2025
Pukul 13.00 s/s 15.00 dan dinyatakan


LULUS

Dengan susunan penguji sebagai berikut

Ketua Sidang


(Arya Wicaksana, S.Kom., M.Eng.Sc.
(OCA, CEH, CEI))
NIDN: 0315109103


Penguji


(Adhi Kusnadi, S.T, M.Si.)
NIDN: 0303037304

Pembimbing


(David Agustriawan, S.Kom., M.Sc., Ph.D.)
NIDN: 0525088601

Ketua Program Studi Informatika,


(Arya Wicaksana, S.Kom., M.Eng.Sc. (OCA, CEH, CEI))
NIDN: 0315109103

iii

**HALAMAN PERSETUJUAN PUBLIKASI KARYA ILMIAH UNTUK
KEPENTINGAN AKADEMIS**

Yang bertanda tangan di bawah ini:

Nama : Shyehan Rafael Adlinugroho
NIM : 00000052738
Program Studi : Informatika
Jenjang : S1
Judul Karya Ilmiah : Perbandingan Performa Algoritma
Random Forest dan XGBoost dalam
Mendeteksi Diabetes Berdasarkan
Data Rekam Medis

Menyatakan dengan sesungguhnya bahwa saya bersedia (**pilih salah satu**):

- Saya bersedia memberikan izin sepenuhnya kepada Universitas Multimedia Nusantara untuk mempublikasikan hasil karya ilmiah saya ke dalam repositori Knowledge Center sehingga dapat diakses oleh Sivitas Akademika UMN/Publik. Saya menyatakan bahwa karya ilmiah yang saya buat tidak mengandung data yang bersifat konfidensial.
- Saya tidak bersedia mempublikasikan hasil karya ilmiah ini ke dalam repositori Knowledge Center, dikarenakan: dalam proses pengajuan publikasi ke jurnal/konferensi nasional/internasional (dibuktikan dengan *letter of acceptance*) **.
- Lainnya, pilih salah satu:
- Hanya dapat diakses secara internal Universitas Multimedia Nusantara
 - Embargo publikasi karya ilmiah dalam kurun waktu tiga tahun.

Tangerang, 3 Januari 2025

Yang menyatakan



Shyehan Rafael Adlinugroho

**Jika tidak bisa membuktikan LoA jurnal/HKI, saya bersedia mengizinkan penuh karya ilmiah saya untuk dipublikasikan ke KC UMN dan menjadi hak institusi UMN.



UMMN

UNIVERSITAS
MULTIMEDIA
NUSANTARA

Halaman Persembahan / Motto

”Discipline is the bridge between goals and accomplishment.”

Jim Rohn



KATA PENGANTAR

Segala puji dan syukur penulis panjatkan ke hadirat Tuhan Yang Maha Esa atas rahmat dan karunia-Nya, sehingga penulis dapat menyelesaikan tugas akhir ini yang berjudul “PERBANDINGAN PERFORMA ALGORITMA RANDOM FOREST DAN XGBOOST DALAM MENDETEKSI DIABETES BERDASARKAN DATA REKAM MEDIS”. Penulisan tugas akhir ini bertujuan untuk memberikan kontribusi dalam pengembangan teknologi prediktif di bidang kesehatan, khususnya dalam deteksi dini diabetes.

Penulis menyadari bahwa penyelesaian tugas akhir ini tidak terlepas dari dukungan dan bantuan berbagai pihak. Oleh karena itu, penulis mengucapkan terima kasih yang sebesar-besarnya kepada semua pihak yang telah memberikan dukungan selama proses penelitian ini.

Mengucapkan terima kasih

1. Bapak Dr. Andrey Andoko, selaku Rektor Universitas Multimedia Nusantara.
2. Bapak Dr. Eng. Niki Prastomo, S.T., M.Sc., selaku Dekan Fakultas Teknik dan Informatika Universitas Multimedia Nusantara.
3. Bapak Arya Wicaksana, S.Kom., M.Eng.Sc. (OCA, CEH, CEI), selaku Ketua Program Studi Informatika Universitas Multimedia Nusantara.
4. Bapak David Agustriawan, S.Kom., M.Sc., Ph.D., sebagai Pembimbing pertama yang telah memberikan bimbingan, arahan, dan motivasi atas terselesainya tugas akhir ini.
5. dr. Ria Novitasari dan dr. Gita Permatasari, MKK, selaku narasumber yang telah memberikan wawasan dan informasi yang sangat berguna untuk mendukung penelitian ini.
6. Keluarga saya yang telah memberikan bantuan dukungan material dan moral, sehingga penulis dapat menyelesaikan tugas akhir ini.
7. Teman-teman yang telah membantu baik secara langsung maupun tidak langsung dalam proses penyelesaian tugas akhir ini, atas segala dukungan, diskusi, serta motivasinya.

Semoga karya ilmiah ini dapat bermanfaat bagi pengembangan ilmu pengetahuan, khususnya dalam bidang deteksi dini diabetes, serta menjadi referensi bagi penelitian selanjutnya.

Tangerang, 3 Januari 2025



Shyehan Rafael Adlinugroho



PERBANDINGAN PERFORMA ALGORITMA RANDOM FOREST DAN XGBOOST DALAM MENDETEKSI DIABETES BERDASARKAN DATA REKAM MEDIS

Shyehan Rafael Adlinugroho

ABSTRAK

Diabetes merupakan kondisi kronis dengan kadar gula darah tinggi yang dapat menyebabkan komplikasi serius. Faktor risiko meliputi pola makan tidak sehat, kurang aktivitas fisik, faktor genetik, dan gaya hidup yang buruk. Berdasarkan data WHO tahun 2022, terdapat 537 juta orang dewasa dengan diabetes di dunia, dengan prevalensi yang terus meningkat, terutama di negara berkembang seperti Indonesia. Penelitian ini membandingkan algoritma Random Forest dan XGBoost dalam mendeteksi diabetes berdasarkan data medical check-up dari Rumah Sakit Pusat Pertamina dengan rentang usia 46-65 tahun. Pemilihan algoritma Random Forest dan XGBoost didasarkan pada kemampuan keduanya dalam menganalisis data dengan banyak fitur dan menangani klasifikasi yang kompleks, termasuk dalam kasus data kesehatan yang sering kali memiliki distribusi tidak seimbang. Random Forest menggunakan pendekatan bagging dengan membangun banyak pohon keputusan secara paralel, sedangkan XGBoost menggunakan pendekatan boosting yang membangun pohon secara bertahap dengan memperbaiki kesalahan prediksi sebelumnya. Perbandingan ini penting dilakukan untuk mengevaluasi akurasi dan stabilitas prediksi dalam kasus deteksi diabetes yang kompleks. Hasil percobaan menunjukkan bahwa Random Forest menghasilkan akurasi pengujian sebesar 90% dengan F1-Score 0,87, sedangkan XGBoost mencapai akurasi 88% dengan F1-Score 0,77. Random Forest menunjukkan performa yang lebih baik dalam mendeteksi diabetes dengan distribusi data yang lebih seimbang, sementara XGBoost cenderung lebih selektif pada fitur dengan kontribusi signifikan. Distribusi data yang lebih terfokus dan seimbang terbukti meningkatkan performa model dalam mendeteksi diabetes, yang diharapkan dapat membantu upaya deteksi dini dan pengambilan keputusan klinis yang lebih efektif di masa depan.

Kata Kunci: Deteksi Dini, Diabetes, Medical Check-Up, Random Forest, XGBoost

**PERFORMANCE COMPARISON OF RANDOM FOREST AND XGBOOST
ALGORITHMS IN DETECTING DIABETES BASED ON MEDICAL
RECORD DATA**

Shyehan Rafael Adlinugroho

ABSTRACT

Diabetes is a chronic condition characterized by high blood sugar levels, potentially leading to severe complications. Risk factors include unhealthy eating habits, physical inactivity, genetic predisposition, and poor lifestyle choices. According to WHO data in 2022, there are 537 million adults living with diabetes worldwide, with a rising prevalence, particularly in developing countries like Indonesia. This study compares the Random Forest and XGBoost algorithms in detecting diabetes based on medical check-up data from Rumah Sakit Pusat Pertamina for individuals aged 46-65 years. The selection of Random Forest and XGBoost is based on their ability to analyze datasets with numerous features and handle complex classification tasks, including imbalanced health data. Random Forest employs a bagging approach, building multiple decision trees in parallel, whereas XGBoost uses a boosting approach, constructing trees iteratively to correct previous prediction errors. This comparison is crucial for evaluating prediction accuracy and stability in complex diabetes detection cases. Experimental results indicate that Random Forest achieves a testing accuracy of 90% with an F1-score of 0.87, while XGBoost achieves an accuracy of 88% with an F1-score of 0.77. Random Forest demonstrates better performance in detecting diabetes with more balanced data distribution, whereas XGBoost tends to focus on features with significant contributions. A more focused and balanced data distribution has been proven to enhance model performance in detecting diabetes, offering potential support for early detection efforts and more effective clinical decision-making in the future.

Keywords: Diabetes, Early Detection, Medical Check-Up, Random Forest, XGBoost

U N I V E R S I T A S
M U L T I M E D I A
N U S A N T A R A

DAFTAR ISI

HALAMAN JUDUL	i
PERNYATAAN TIDAK MELAKUKAN PLAGIAT	ii
HALAMAN PENGESAHAN	iii
HALAMAN PERSETUJUAN PUBLIKASI ILMIAH	iv
HALAMAN PERSEMBAHAN/MOTO	vi
KATA PENGANTAR	vii
ABSTRAK	ix
ABSTRACT	x
DAFTAR ISI	xi
DAFTAR TABEL	xiii
DAFTAR GAMBAR	xiv
DAFTAR KODE	xv
DAFTAR RUMUS	xvi
DAFTAR LAMPIRAN	xvii
BAB 1 PENDAHULUAN	1
1.1 Latar Belakang Masalah	1
1.2 Rumusan Masalah	4
1.3 Batasan Permasalahan	4
1.4 Tujuan Penelitian	4
1.5 Manfaat Penelitian	5
1.6 Sistematika Penulisan	5
BAB 2 LANDASAN TEORI	8
BAB 3 METODOLOGI PENELITIAN	15
3.0.1 Alur penelitian	15
3.0.2 Alur Preprocessing Data	17
3.1 Alur Preprocessing Data	18
3.2 Alur Penggabungan Dataset	19
3.3 Alur Pembuatan Model	20
3.4 Spesifikasi Sistem	21
3.4.1 Spesifikasi Perangkat Keras	21
3.4.2 Spesifikasi Perangkat Lunak	22
BAB 4 HASIL DAN DISKUSI	23
4.1 Dataset	23
4.2 Preprocessing Data	24
4.2.1 Penghapusan Kolom Tidak Diperlukan	24
4.2.2 Transformasi Data	25
4.2.3 Pembersihan Data	27
4.2.4 Pemberian Label Diabetes	27
4.2.5 Penanganan Data Duplikat	28
4.2.6 Pemberian Label Penyakit Lain	29
4.2.7 Penyimpanan Data	32
4.3 Pembuatan Model	33
4.3.1 Persiapan Dataset	33
4.3.2 Pengecekan Total Nilai Label Penyakit	33
4.3.3 Penentuan Fitur dan Target	35
4.3.4 Encoding Data Menggunakan Dummy	35
4.3.5 Normalisasi Data	36
4.3.6 Pembagian Data: Training dan Testing	36

4.3.7	Pemeriksaan Jumlah Data pada Setiap Bagian	37
4.3.8	Pembangunan Model	37
4.3.9	Evaluasi Model	40
4.4	Diskusi	42
4.5	Scenario Percobaan	42
4.6	Evaluasi Algoritma Optimal	44
BAB 5	SIMPULAN DAN SARAN	49
5.1	Simpulan	49
5.2	Saran	49
DAFTAR PUSTAKA	51



DAFTAR TABEL

Tabel 2.1	Confusion Matrix untuk Prediksi Diabetes	13
Tabel 4.1	10 Fitur dengan Korelasi Tertinggi terhadap Diabetes	23
Tabel 4.2	Performa Algoritma Random Forest dan XGBoost pada Percobaan 1–11 (Hasil Testing)	43
Tabel 4.3	Performa Algoritma Random Forest dan XGBoost pada Percobaan 11 (Hasil Testing)	45
Tabel 4.4	Confusion Matrix untuk Random Forest dan XGBoost	45
Tabel 4.5	Perbedaan Data pada Percobaan 8 dan Percobaan 11 (46-65)	46



DAFTAR GAMBAR

Gambar 2.1	Ilustrasi Pohon Keputusan dalam Random Forest	10
Gambar 3.1	Flowchart penelitian	15
Gambar 3.2	Flowchart Data Preparation	17
Gambar 3.3	Flowchart Preprocessing Data	18
Gambar 3.4	Flowchart Penggabungan Data	19
Gambar 3.5	Flowchart Pembuatan Model	20
Gambar 4.1	Data Mentah	23
Gambar 4.2	Contoh Data Sesudah Pembersihan	27
Gambar 4.3	Distribusi Label Diabetes dalam Dataset	28
Gambar 4.4	Distribusi Label Penyakit Hipertensi	30
Gambar 4.5	Distribusi Label Penyakit Jantung	31
Gambar 4.6	Distribusi Label Penyakit Ginjal	32
Gambar 4.7	Distribusi Label Penyakit Lain dalam Dataset	32
Gambar 4.8	Data siap digunakan	33
Gambar 4.9	Data siap digunakan	34
Gambar 4.10	Hasil Evaluasi Random Forest	41
Gambar 4.11	Hasil Evaluasi XGBoost	41
Gambar 4.12	Fitur Importance Random Forest	46
Gambar 4.13	Fitur Importance XGBoost	47



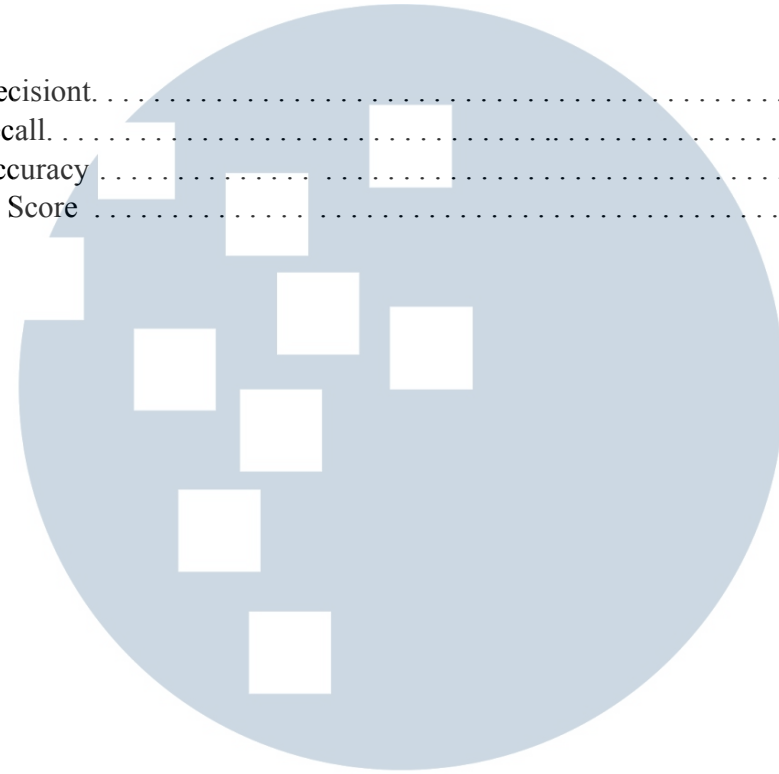
DAFTAR KODE

Kode 4.1	Kolom yang Dihapus	24
Kode 4.2	Kolom yang Ditransformasi	25
Kode 4.3	Pembersihan kolom Riwayat Penyakit Terdahulu	27
Kode 4.4	Pembersihan kolom Riwayat Penyakit Terdahulu	28
Kode 4.5	Hapus Duplikat dengan label ganda	28
Kode 4.6	Hapus sisa Duplikat pada dataset	29
Kode 4.7	Label Hipertensi	29
Kode 4.8	Label Jantung	30
Kode 4.9	Label Ginjal	31
Kode 4.10	Total Nilai Label Penyakit	33
Kode 4.11	Fitur Selection	35
Kode 4.12	Encoding data menggunakan Dummy	36
Kode 4.13	Normalisasi Data	36
Kode 4.14	Fitur dan Target	37
Kode 4.15	Training & testing	37
Kode 4.16	Jumlah data Train dan Test	37
Kode 4.17	Pembuatan Model Random Forest	37
Kode 4.18	Pembuatan Model Random Forest menggunakan hyper parameter Tunning	38
Kode 4.19	Pembuatan Model XGBoost	39
Kode 4.20	Pembuatan Model XGBoost menggunakan hyper parameter Tunning	39



DAFTAR RUMUS

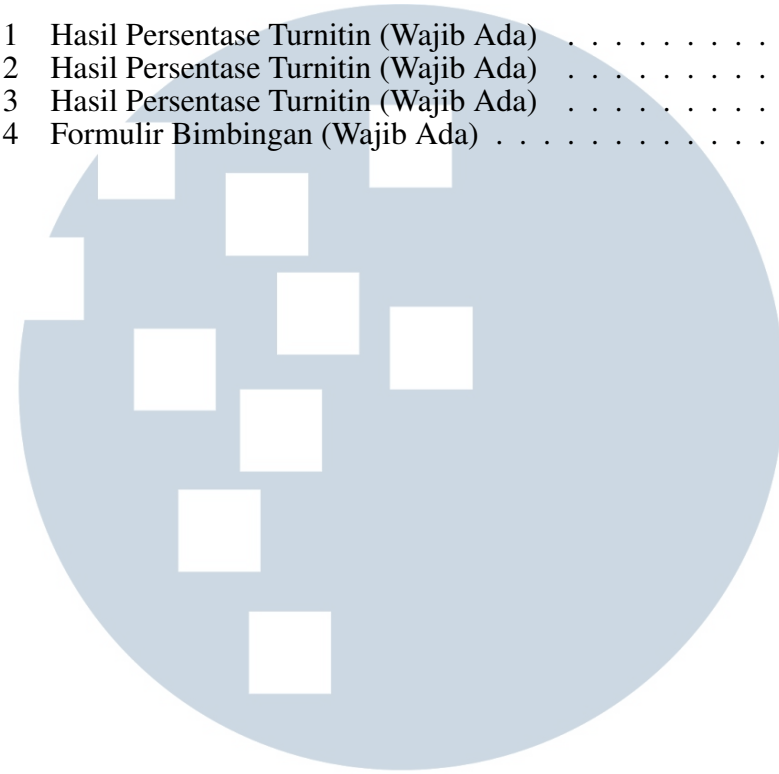
Rumus 2.1 Precisiont.....	14
Rumus 2.2 Recall.....	14
Rumus 2.3 Accuracy.....	14
Rumus 2.4 F1 Score.....	14



UMMN
UNIVERSITAS
MULTIMEDIA
NUSANTARA

DAFTAR LAMPIRAN

Lampiran 1	Hasil Persentase Turnitin (Wajib Ada)	53
Lampiran 2	Hasil Persentase Turnitin (Wajib Ada)	54
Lampiran 3	Hasil Persentase Turnitin (Wajib Ada)	55
Lampiran 4	Formulir Bimbingan (Wajib Ada)	56



UMN

UNIVERSITAS
MULTIMEDIA
NUSANTARA