

**IMPLEMENTASI ALGORITMA RANDOM FOREST UNTUK
KLASIFIKASI TINGKAT RISIKO MENGIDAP KANKER PARU-PARU**



SKRIPSI

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

Nehemia Gueldi
00000043202

**PROGRAM STUDI INFORMATIKA
FAKULTAS TEKNIK DAN INFORMATIKA
UNIVERSITAS MULTIMEDIA NUSANTARA
TANGERANG
2024**

**IMPLEMENTASI ALGORITMA RANDOM FOREST UNTUK
KLASIFIKASI TINGKAT RISIKO MENGIDAP KANKER PARU-PARU**



Nehemia Gueldi

00000043202

UMMN

UNIVERSITAS

MULTIMEDIA

NUSANTARA

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

TANGERANG

2024

HALAMAN PERNYATAAN TIDAK PLAGIAT

Dengan ini saya,

Nama : Nehemia Gueldi
Nomor Induk Mahasiswa : 00000043202
Program Studi : Informatika

Skripsi dengan judul:

Implementasi Algoritma Random Forest Untuk Klasifikasi Tingkat Risiko Mengidap Kanker Paru-Paru

merupakan hasil karya saya sendiri bukan plagiat dari karya 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 Skripsi maupun dalam penulisan laporan Skripsi, saya bersedia menerima konsekuensi dinyatakan **TIDAK LULUS** untuk Tugas akhir yang telah saya tempuh.

Tangerang, 15 Mei 2024



(Nehemia Gueldi)

UNIVERSITAS
MULTIMEDIA
NUSANTARA

HALAMAN PENGESAHAN

Skripsi dengan judul

IMPLEMENTASI ALGORITMA RANDOM FOREST UNTUK KLASIFIKASI TINGKAT RISIKO MENGIDAP KANKER PARU-PARU

oleh

Nama : Nehemia Gueldi
NIM : 00000043202
Program Studi : Informatika
Fakultas : Fakultas Teknik dan Informatika

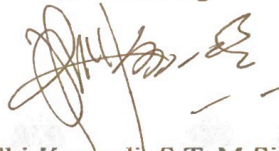
Telah diujikan pada hari Jumat, 07 Juni 2024

Pukul 10.00 s/d 12.00 dan dinyatakan

LULUS

Dengan susunan penguji sebagai berikut

Ketua Sidang



(Adhi Kushadi, S.T, M.Si.)

NIDN: 0303037304

Penguji

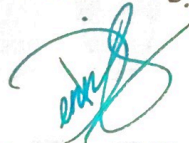


(Dr. Maria Irmina Prasetyowati, S.Kom.,

M.T.)

NIDN: 0725057201

Pembimbing



(Dennis Gunawan, S.Kom., M.Sc.)

NIDN: 0320059001

PPS Ketua Program Studi Informatika,



Dr. Eng. Niki Prastomo, S.T., M.Sc.

NIDN: 0419128203

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

Yang bertanda tangan di bawah ini:

Nama : Nehemia Gueldi
NIM : 00000043202
Program Studi : Informatika
Jenjang : S1
Jenis Karya : Skripsi

Menyatakan dengan sesungguhnya bahwa:

- Saya bersedia memberikan izin sepenuhnya kepada Universitas Multimedia Nusantara untuk mempublikasikan hasil karya ilmiah saya di repositori Knowledge Center, sehingga dapat diakses oleh Civitas Akademika/Publik. Saya menyatakan bahwa karya ilmiah yang saya buat tidak mengandung data yang bersifat konfidensial dan saya juga tidak akan mencabut kembali izin yang telah saya berikan dengan alasan apapun.
- Saya tidak bersedia karena dalam proses pengajuan untuk diterbitkan ke jurnal/konferensi nasional/internasional (dibuktikan dengan *letter of acceptance*)**.

Tangerang, 3 Juni 2024

Yang menyatakan



Nehemia Gueldi

UMMN
UNIVERSITAS
MULTIMEDIA
NUSANTARA

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

Halaman Persembahan / Motto

"A good name is to be more desired than great wealth, Favor is better than silver and gold."

Proverbs 22:1 (NASB)



UMMN
UNIVERSITAS
MULTIMEDIA
NUSANTARA

KATA PENGANTAR

Puji Syukur atas berkat dan rahmat kepada Tuhan Yang Maha Esa, atas selesainya penulisan laporan Skripsi ini dengan judul: Implementasi Algoritma Random Forest Untuk Klasifikasi Tingkat Risiko Mengidap Kanker Paru-Paru dilakukan untuk memenuhi salah satu syarat untuk mencapai gelar Sarjana Komputer Jurusan Informatika Pada Fakultas Teknik dan Informatika Universitas Multimedia Nusantara. Saya menyadari bahwa, tanpa bantuan dan bimbingan dari berbagai pihak, dari masa perkuliahan sampai pada penyusunan skripsi ini, sangatlah sulit bagi saya untuk menyelesaikan skripsi ini. Oleh karena itu, saya mengucapkan terima kasih kepada:

1. Bapak Dr. Ninok Leksono, 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 Dr. Eng. Niki Prastomo, S.T., M.Sc., selaku Pjs. Ketua Program Studi Informatika Universitas Multimedia Nusantara.
4. Bapak Dennis Gunawan, S.Kom., M.Sc., sebagai Pembimbing pertama yang telah banyak meluangkan waktu untuk memberikan bimbingan, arahan dan motivasi atas terselesainya skripsi ini.
5. Orang Tua, dan keluarga saya yang telah memberikan bantuan dukungan material dan moral, sehingga penulis dapat menyelesaikan tesis ini.

Semoga skripsi ini bermanfaat, baik sebagai sumber informasi maupun sumber inspirasi, bagi para pembaca.

Tangerang, 3 Juni 2024



Nehemia Gueldi

IMPLEMENTASI ALGORITMA RANDOM FOREST UNTUK KLASIFIKASI TINGKAT RISIKO MENGIDAP KANKER PARU-PARU

Nehemia Gueldi

ABSTRAK

Kanker paru-paru merupakan salah satu penyakit yang menjadi fokus utama dalam bidang kesehatan global karena memiliki jumlah kasus kematian yang tinggi. Meskipun jumlah kasus kanker payudara pada wanita melampaui kanker paru-paru, kanker paru-paru tetap menjadi penyebab utama kematian. Faktor risiko utama kanker paru-paru adalah merokok dan paparan polusi udara. Dengan membangun sejumlah pohon keputusan secara terpisah dan menggabungkan hasilnya, algoritma *Random Forest* dapat melakukan klasifikasi yang akurat dan mencegah *overfitting*. Pengujian dilakukan menggunakan *dataset Cancer Patient Datasets* untuk mengukur *accuracy*, *precision*, *recall*, dan *F1 score* yang diperoleh. Hasil yang didapatkan dalam penelitian ini adalah algoritma *Random Forest* berhasil diimplementasikan untuk klasifikasi tingkat risiko mengidap kanker paru-paru dengan nilai *accuracy* sebesar 97%, *precision* sebesar 97.26%, *recall* sebesar 97%, dan *f1-score* sebesar 96.98% berdasarkan hasil evaluasi metrik model yang terbentuk.

Kata kunci: Merokok, Kanker Paru-Paru, Kesehatan, Klasifikasi, *Random Forest*



Implementation Random Forest Algorithm to Classify Risk Level of Lung Cancer

Nehemia Gueldi

ABSTRACT

Lung cancer is one of the diseases that has become a major focus in global health due to its high mortality rate. Although the number of breast cancer cases in women surpasses lung cancer, lung cancer remains a leading cause of death. The main risk factors for lung cancer are smoking and exposure to air pollution. By constructing a number of decision trees separately and combining their results, the Random Forest algorithm can accurately classify and prevent overfitting. Testing was conducted using the Cancer Patient Datasets to measure the accuracy, precision, recall, and F1 score obtained. The results obtained in this study indicate that the Random Forest algorithm has been successfully implemented for classifying the risk level of developing lung cancer with an accuracy of 97%, precision of 97.26%, recall of 97%, and F1-score of 96.98% based on the evaluation metrics of the formed model.

Keywords: *Classification, Health, Lung Cancer, Random Forest, Smoking*



DAFTAR ISI

HALAMAN JUDUL	i
PERNYATAAN TIDAK MELAKUKAN PLAGIAT	ii
HALAMAN PENGESAHAN	iii
HALAMAN PERSETUJUAN PUBLIKASI ILMIAH	iv
HALAMAN PERSEMBAHAN/MOTO	v
KATA PENGANTAR	vi
ABSTRAK	vii
ABSTRACT	viii
DAFTAR ISI	ix
DAFTAR GAMBAR	x
DAFTAR TABEL	xi
DAFTAR KODE	xii
DAFTAR LAMPIRAN	xiii
BAB 1 PENDAHULUAN	1
1.1 Latar Belakang Masalah	1
1.2 Rumusan Masalah	2
1.3 Batasan Permasalahan	3
1.4 Tujuan Penelitian	3
1.5 Manfaat Penelitian	3
1.6 Sistematika Penulisan	4
BAB 2 LANDASAN TEORI	5
2.1 Kanker Paru-Paru	5
2.2 Ensemble Learning	6
2.3 Decision Tree	7
2.4 Random Forest	8
2.5 Confusion Matrix	9
BAB 3 METODOLOGI PENELITIAN	14
3.1 Gambaran Umum	14
3.1.1 Spesifikasi Perangkat	15
3.2 Studi Literatur	16
3.3 Pengumpulan Data	16
3.4 Pembangunan Model	17
3.5 Evaluasi	18
3.6 Penulisan Laporan dan Konsultasi	19
BAB 4 HASIL DAN DISKUSI	20
4.1 Source Code Implementasi Algoritma	20
4.1.1 Read Dataset	20
4.1.2 Pembersihan Data	20
4.1.3 Mapping	21
4.1.4 Pembagian Data	21
4.1.5 Grid Search	22
4.2 Skenario Uji Coba	23
4.3 Hasil Uji Coba	24
4.4 Diskusi	28
BAB 5 SIMPULAN DAN SARAN	31
5.1 Simpulan	31
5.2 Saran	31
DAFTAR PUSTAKA	32

DAFTAR GAMBAR

Gambar 2.1	<i>Decision Tree</i>	8
Gambar 2.2	Visualisasi Algoritma <i>Random Forest</i>	9
Gambar 3.1	<i>Flowchart</i> Gambaran Umum Penelitian	14
Gambar 3.2	Alur Kerja Utama Penelitian	15
Gambar 3.3	<i>Cancer Patient Datasets I</i>	16
Gambar 3.4	<i>Cancer Patient Datasets II</i>	16
Gambar 3.5	<i>Flowchart</i> Pembangunan Model	18
Gambar 3.6	<i>Flowchart</i> Evaluasi Model	19
Gambar 4.1	Grafik Batang Evaluasi Metrik Skenario 1 (<i>Max Depth</i>)	25
Gambar 4.2	Grafik Batang Evaluasi Metrik Skenario 2 (<i>Min Samples Leaf</i>)	25
Gambar 4.3	Grafik Batang Evaluasi Metrik Skenario 3 (<i>Max Leaf Nodes</i>)	26
Gambar 4.4	<i>Confusion Matrix</i> Kombinasi Nilai <i>Hyperparameter</i> Terbaik	27
Gambar 4.5	Contoh salah satu <i>Decision Tree</i> yang terbuat	28
Gambar 4.6	<i>Loss Learning Curve</i> berdasarkan jumlah data <i>training</i> dengan kombinasi <i>hyperparameter</i> terbaik	30



DAFTAR TABEL

Tabel 2.1	Bentuk Umum <i>Confusion Matrix</i> 3x3	10
Tabel 2.2	Bentuk <i>Confusion Matrix</i> 2x2 Kelas A dari <i>Confusion Matrix</i> 3x3	11
Tabel 2.3	Gambaran <i>Confusion Matrix</i> 3x3 Kelas A dari <i>Confusion Matrix</i> 3x3	11
Tabel 2.4	Bentuk <i>Confusion Matrix</i> 2x2 Kelas B dari <i>Confusion Matrix</i> 3x3	11
Tabel 2.5	Gambaran <i>Confusion Matrix</i> 3x3 Kelas B dari <i>Confusion Matrix</i> 3x3	11
Tabel 2.6	Bentuk <i>Confusion Matrix</i> 2x2 Kelas C dari <i>Confusion Matrix</i> 3x3	11
Tabel 2.7	Gambaran <i>Confusion Matrix</i> 3x3 Kelas C dari <i>Confusion Matrix</i> 3x3	12
Tabel 4.1	<i>Default Value Hyperparameter</i> yang digunakan	24
Tabel 4.2	<i>Hyperparameter Tuning Values</i>	24
Tabel 4.3	Kombinasi Nilai <i>Hyperparameter</i> Terbaik setelah proses <i>Grid Search</i>	26



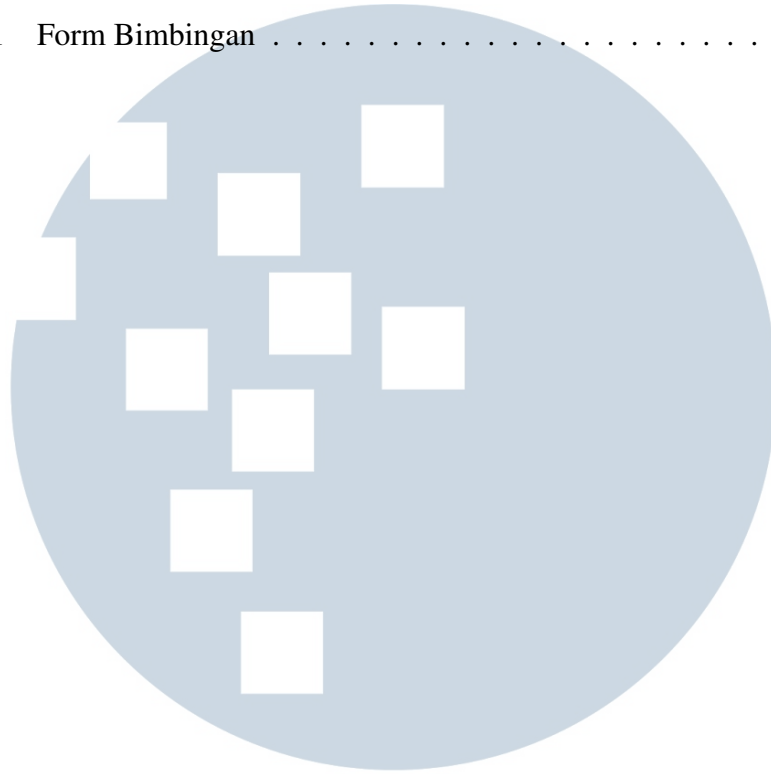
DAFTAR KODE

4.1	Kode <i>Read Dataset</i>	20
4.2	Kode <i>Preprocessing</i>	21
4.3	Kode <i>Mapping</i>	21
4.4	Kode Pembagian <i>Train Data</i> dan <i>Test Data</i>	22
4.5	Kode <i>Grid Search</i>	23



DAFTAR LAMPIRAN

Lampiran 1 Form Bimbingan 40



UMMN
UNIVERSITAS
MULTIMEDIA
NUSANTARA