

**KLASIFIKASI STADIUM KANKER PROSTAT
MENGGUNAKAN LR DAN RF DENGAN
SELEKSI FITUR LASSO-RFE BERBASIS
DATA GEN DAN MIRNA**



SKRIPSI

**IVANDY WIJAYA
00000061844**

**PROGRAM STUDI INFORMATIKA
FAKULTAS TEKNIK DAN INFORMATIKA
UNIVERSITAS MULTIMEDIA NUSANTARA
TANGERANG
2025**

**KLASIFIKASI STADIUM KANKER PROSTAT
MENGGUNAKAN LR DAN RF DENGAN
SELEKSI FITUR LASSO-RFE BERBASIS
DATA GEN DAN MIRNA**



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

**IVANDY WIJAYA
00000061844**

UMN
UNIVERSITAS
MULTIMEDIA
PROGRAM STUDI INFORMATIKA
FAKULTAS TEKNIK DAN INFORMATIKA
UNIVERSITAS MULTIMEDIA NUSANTARA
TANGERANG
2025

HALAMAN PERNYATAAN TIDAK PLAGIAT

Dengan ini saya,

Nama : Ivandy Wijaya
Nomor Induk Mahasiswa : 00000061844
Program Studi : Informatika

Skripsi dengan judul:

Klasifikasi Stadium Kanker Prostat Menggunakan LR-RF dengan Seleksi Fitur Lasso-RFE Berbasis Data Gen dan Mirna

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, 26 Juni 2025



UNIVERSITAS
MULTIMEDIA
NUSANTARA

HALAMAN PENGESAHAN



Ketua Sidang

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

NIDN: 0320059001

Pembimbing I

Penguji

(Arya Wicaksana, S.Kom., M.Eng.Sc.,
OCA)

NIDN: 0315109103

Pembimbing II

(Moeljono Widjaja, B.Sc., M.Sc., Ph.D.)

NIDN: 0311106903

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

NIDN: 0525088601

Ketua Program Studi Informatika,

(Arya Wicaksana, S.Kom., M.Eng.Sc., OCA)

NIDN: 0315109103

HALAMAN PERSETUJUAN PUBLIKASI KARYA ILMIAH UNTUK KEPENTINGAN AKADEMIS

Yang bertanda tangan di bawah ini:

Nama : Ivandy Wijaya
NIM : 00000061844
Program Studi : Informatika
Jenjang : S1
Judul Karya Ilmiah : Klasifikasi Stadium Kanker Prostat
Menggunakan LR dan RF dengan
Seleksi Fitur Lasso-RFE Berbasis
Data Gen dan Mirna

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, 26 Juni 2025
Yang menyatakan

UNIVERSITAS
MULTIMEDIA
NUSANTARA

Ivandy Wijaya

**Jika tidak bisa membuktikan LoA jurnal/HKI, saya bersedia mengizinkan penuh karya ilmiah saya untuk dipublikasikan 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)

UMN
UNIVERSITAS
MULTIMEDIA
NUSANTARA

KATA PENGANTAR

Puji Syukur atas berkat dan rahmat kepada Tuhan Yang Maha Esa, atas selesainya penulisan laporan Skripsi ini sebagai 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 laporan magang ini, sangatlah sulit bagi saya untuk menyelesaikan laporan magang ini. Oleh karena itu, saya mengucapkan terima kasih kepada:

1. Bapak Dr. Ir. Andrey Andoko, M.Sc., 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, selaku Ketua Program Studi Informatika Universitas Multimedia Nusantara.
4. Bapak Moeljono Widjaja, B.Sc., M.Sc., Ph.D., sebagai Pembimbing pertama yang telah memberikan bimbingan, arahan, dan motivasi atas terselesainya tugas akhir ini.
5. Bapak David Agustriawan, S.Kom., M.Sc., Ph.D, sebagai Pembimbing kedua yang turut memberikan bimbingan, arahan, dan motivasi atas terselesainya tugas akhir ini.
6. Keluarga saya yang telah memberikan bantuan dukungan material dan moral, sehingga penulis dapat menyelesaikan tugas akhir ini.

Semoga karya ilmiah ini dapat memberikan manfaat bagi pengembangan ilmu pengetahuan, serta menjadi referensi bagi studi lanjutan di masa mendatang.

Tangerang, 26 Juni 2025



Ivandy Wijaya

**KLASIFIKASI STADIUM KANKER PROSTAT MENGGUNAKAN LR
DAN RF DENGAN SELEKSI FITUR LASSO-RFE BERBASIS DATA GEN
DAN MIRNA**

Ivandy Wijaya

ABSTRAK

Kanker prostat merupakan salah satu jenis kanker yang paling umum terjadi pada pria dan menghadirkan tantangan dalam proses penentuan stadium kanker yang dini dan akurat. Penelitian ini bertujuan untuk menggunakan pendekatan berbasis pembelajaran mesin dalam mengklasifikasikan kanker prostat stadium II dan III menggunakan data transkriptomik (RNA-seq dan miRNA) dari *The Cancer Genome Atlas Prostate Adenocarcinoma* (TCGA-PRAD). Dua pendekatan seleksi fitur digunakan, yaitu analisis gen yang diekspresikan secara berbeda (Differentially Expressed Genes/DEG) dan penyaringan statistik (LASSO + RFE). Delapan skenario eksperimen dikembangkan dengan menggabungkan berbagai jenis data, metode seleksi fitur, dan algoritma klasifikasi (Logistic Regression dan Random Forest). Performa terbaik diperoleh dari kombinasi data ekspresi gen dengan metode seleksi fitur LASSO + RFE serta algoritma Logistic Regression. Model yang dihasilkan mencapai akurasi sebesar 99,17%, precision sebesar 99,05%, recall sebesar 90%, dan F1-score sebesar 90%. Sebanyak 25 gen terpilih digunakan dalam model optimal akhir. Analisis lanjutan menunjukkan bahwa meskipun kemampuan diskriminatif tiap gen secara individu terbatas, penggunaan gabungan gen-gen tersebut memberikan performa klasifikasi yang unggul, yang menyoroti pentingnya interaksi antar gen.

Kata kunci: *Ekspresi Gen, Kanker prostat, miRNA, Pembelajaran mesin, Seleksi fitur*



PROSTATE CANCER STAGE CLASSIFICATION USING LR AND RF WITH LASSO-RFE FEATURE SELECTION BASED ON GENE AND MIRNA DATA

Ivandy Wijaya

ABSTRACT

Prostate cancer is one of the most prevalent cancers among men and presents significant challenges in early and accurate staging. This study proposes a machine learning-based approach to classify stage II and III prostate cancer using transcriptomic data (RNA-seq and miRNA) from the The Cancer Genome Atlas Prostate Adenocarcinoma (TCGA-PRAD) cohort. Two different feature selection approaches were employed, involving Differentially Expressed Genes (DEG) analysis and statistical filtering (LASSO + RFE). Eight experimental scenarios were developed by combining different data types, feature selection methods, and classification algorithms (Logistic Regression and Random Forest). The best performance was observed in the combination of gene expression data with the LASSO + RFE feature selection approach and the Logistic Regression algorithm. The resulting model achieved 99.17% accuracy, 99.05% precision, 90% recall, and a 90% F1-score. A set of 25 selected genes was used in the final optimal model. Further analysis revealed limited discriminative power individually, but collective use of selected genes provided superior classification performance, highlighting the importance of gene interactions.

Keywords: Feature selection, Gene Expression, Machine learning, miRNA, Prostate cancer



DAFTAR ISI

HALAMAN JUDUL	i
PERNYATAAN TIDAK MELAKUKAN PLAGIAT	ii
HALAMAN PENGESAHAN	iii
HALAMAN PERSETUJUAN PUBLIKASI KARYA ILMIAH	iv
HALAMAN PERSEMBAHAN/MOTO	v
KATA PENGANTAR	vi
ABSTRAK	vii
ABSTRACT	viii
DAFTAR ISI	ix
DAFTAR TABEL	xi
DAFTAR GAMBAR	xii
DAFTAR KODE	xiii
DAFTAR RUMUS	xiv
DAFTAR LAMPIRAN	xv
BAB 1 PENDAHULUAN	1
1.1 Latar Belakang Masalah	1
1.2 Rumusan Masalah	4
1.3 Batasan Permasalahan	4
1.4 Tujuan Penelitian	5
1.5 Manfaat Penelitian	5
1.6 Sistematika Penulisan	6
BAB 2 LANDASAN TEORI	7
2.1 Kanker Prostat	7
2.2 Ekspresi Gen (RNA-seq)	8
2.3 Ekspresi mikroRNA	9
2.4 Differentially Expressed Genes (Limma)	9
2.5 Recursive Feature Elimination	10
2.6 Logistic Regression	12
2.7 Random Forest	15
2.8 Metrik Evaluasi	17
2.8.1 Matriks Kebingungan (Confusion Matrix)	18
2.8.2 Akurasi (Accuracy)	18
2.8.3 Precision	19
2.8.4 Recall	19
2.8.5 F1-Score	20
2.8.6 Receiver Operating Characteristic (ROC)	20
2.8.7 Area Under Curve (AUC)	20
BAB 3 METODOLOGI PENELITIAN	22
3.1 Alur Kerja Penelitian	22
3.2 Spesifikasi Perangkat	23
3.2.1 Perangkat Keras	23
3.2.2 Perangkat Lunak	23
3.3 Studi Literatur	23
3.4 Pengumpulan Data	24
3.5 Praproses Data (Data Preprocessing)	25
3.6 Seleksi Fitur	27
3.6.1 Seleksi Fitur Berbasis DEG	28
3.6.2 Seleksi Fitur Berbasis Metode Statistik	28

3.7	Pembangunan Model Klasifikasi	29
3.8	Evaluasi	30
3.9	Validasi Eksperimen dengan Random Seed Berulang	30
3.10	Skenario Eksperimen	31
3.11	Analisis Biomarker	31
BAB 4	HASIL DAN DISKUSI	32
4.1	Pengumpulan Data	32
4.2	Pra-pemrosesan Data	32
4.2.1	Penanganan Missing Values	33
4.2.2	Filtering Race	33
4.2.3	Labeling dan Filter Stage	34
4.2.4	Penggabungan Dataset	35
4.2.5	Normalisasi Data	37
4.3	Seleksi Fitur	38
4.3.1	Seleksi Fitur Pendekatan DEG	38
4.3.2	Seleksi Fitur Pendekatan Statistik	41
4.4	Pembangunan dan Evaluasi Model	43
4.5	Hasil Skenario	46
4.5.1	Skenario 1	46
4.5.2	Skenario 2	49
4.5.3	Skenario 3	50
4.5.4	Skenario 4	52
4.5.5	Skenario 5	53
4.5.6	Skenario 6	55
4.5.7	Skenario 7	56
4.5.8	Skenario 8	59
4.6	Perbandingan Penggunaan Data Gen dan miRNA	61
4.7	Perbandingan Pendekatan Seleksi Fitur	64
4.8	Perbandingan Penggunaan Algoritma Klasifikasi Logistic Regression dan Random Forest	67
4.9	Hasil Terbaik	70
4.10	Analisis Kandidat Biomarker	71
BAB 5	SIMPULAN DAN SARAN	76
5.1	Simpulan	76
5.2	Saran	77
DAFTAR PUSTAKA		78

UNIVERSITAS
MULTIMEDIA
NUSANTARA

DAFTAR TABEL

Tabel 2.1	Staging Kanker Prostat Berdasarkan TNM, PSA, dan Grade Group	8
Tabel 2.2	Confusion Matrix	18
Tabel 2.3	Interpretasi Nilai AUC	21
Tabel 3.1	Dataset Penelitian	24
Tabel 3.2	Tuning Hyperparameter untuk LR	29
Tabel 3.3	Tuning Hyperparameter untuk RF	30
Tabel 3.4	Daftar Skenario Seleksi Fitur dan Klasifikasi	31
Tabel 4.1	Dataset Penelitian	32
Tabel 4.2	Distribusi ras setelah penghapusan missing values	33
Tabel 4.3	Distribusi label setelah <i>filtering</i> ras	34
Tabel 4.4	Jumlah Sampel Data Pasca-pemrosesan	35
Tabel 4.5	Jumlah Sampel Data Gabungan	37
Tabel 4.6	Lima baris pertama hasil analisis ekspresi gen menggunakan limma	40
Tabel 4.7	Lima baris pertama hasil analisis ekspresi miRNA menggunakan limma	40
Tabel 4.8	Daftar Gen Berdasarkan Nilai logFC Tertinggi	47
Tabel 4.9	Hasil Evaluasi Data Testing pada Skenario 1 (Menggunakan Macro Average)	48
Tabel 4.10	Daftar miRNA Berdasarkan Nilai logFC Tertinggi	49
Tabel 4.11	Hasil Evaluasi Data Testing pada Skenario 2 (Menggunakan Macro Average)	50
Tabel 4.12	Hasil Evaluasi Data Testing pada Skenario 3 (Menggunakan Macro Average)	51
Tabel 4.13	Hasil Evaluasi Data Testing pada Skenario 4 (Menggunakan Macro Average)	52
Tabel 4.14	Hasil Evaluasi Data Testing pada Skenario 5 (Menggunakan Macro Average)	54
Tabel 4.15	Hasil Evaluasi Data Testing pada Skenario 6 (Menggunakan Macro Average)	56
Tabel 4.16	Hasil Evaluasi Data Testing pada Skenario 7 (Menggunakan Macro Average)	58
Tabel 4.17	Hasil Evaluasi Data Testing pada Skenario 8 (Menggunakan Macro Average)	60
Tabel 4.18	Hasil Terbaik pada masing-masing skenario	70
Tabel 4.19	Daftar Ensemble ID dan Simbol Gen terpilih	72
Tabel 4.20	Daftar Gen dan Temuan Referensi Kanker Prostat	73
Tabel 4.21	Daftar Gen dan Temuan Referensi Kanker	74

N U S A N T A R A

DAFTAR GAMBAR

Gambar 3.1	Diagram alur kerja penelitian	22
Gambar 3.2	Diagram alur kerja pra-proses data	25
Gambar 3.3	Diagram alur kerja seleksi fitur	27
Gambar 4.1	Top 25 Genes with highest Feature Importance	53
Gambar 4.2	Top 25 miRNAs with highest Feature Importance	55
Gambar 4.3	Top 25 Genes with highest Feature Importance	57
Gambar 4.4	Top 25 miRNAs with highest Feature Importance	59
Gambar 4.5	Plot perbandingan akurasi antara skenario 1 dan 2	61
Gambar 4.6	Plot perbandingan akurasi antara skenario 3 dan 4	62
Gambar 4.7	Plot perbandingan akurasi antara skenario 5 dan 6	62
Gambar 4.8	Plot perbandingan akurasi antara kenario 7 dan 8	63
Gambar 4.9	Plot perbandingan akurasi antara skenario 1 dan 5	64
Gambar 4.10	Plot perbandingan akurasi antara skenario 2 dan 6	65
Gambar 4.11	Plot perbandingan akurasi antara skenario 3 dan 7	65
Gambar 4.12	Plot perbandingan akurasi antara kenario 4 dan 8	66
Gambar 4.13	Plot perbandingan akurasi antara skenario 1 dan 3	67
Gambar 4.14	Plot perbandingan akurasi antara skenario 2 dan 4	68
Gambar 4.15	Plot perbandingan akurasi antara skenario 5 dan 7	68
Gambar 4.16	Plot perbandingan akurasi antara kenario 6 dan 8	69
Gambar 4.17	Plot perbandingan akurasi data training dan testing pada skenario 5	71
Gambar 4.18	Kurva ROC dari 25 gen kandidat secara individual.	74



DAFTAR KODE

Kode 2.1	Contoh Fungsi RFE	11
Kode 2.2	Contoh Fungsi LR	14
Kode 2.3	Contoh Fungsi RF	17
Kode 4.1	Pemuatan Data	32
Kode 4.2	Penghapusan Missing Values	33
Kode 4.3	Penyaringan Ras Sampel	34
Kode 4.4	Penyaringan Stadium Kanker	35
Kode 4.5	Persiapan Merging Data	36
Kode 4.6	Penggabungan Data	36
Kode 4.7	Normalisasi Data	37
Kode 4.8	Analisis Limma	39
Kode 4.9	Fungsi Lasso	41
Kode 4.10	Fungsi RFE	41
Kode 4.11	Fungsi Perulangan Rentang Fitur	43
Kode 4.12	Fungsi klasifikasi Logistic Regression	44
Kode 4.13	Fungsi klasifikasi Random Forest	45



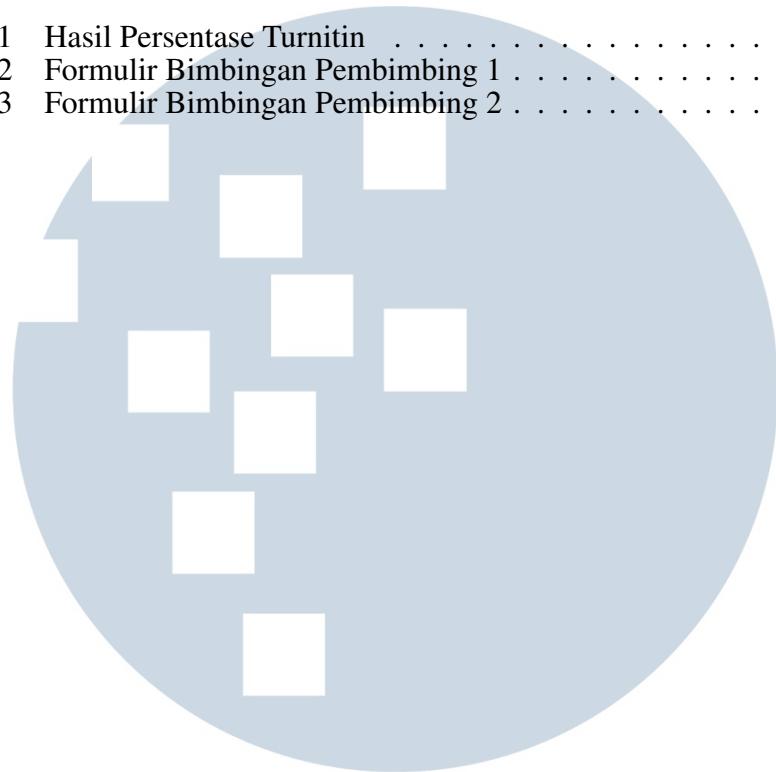
DAFTAR RUMUS

Rumus 2.1	Linear Regression	12
Rumus 2.2	Logistic Function	12
Rumus 2.6	Binary Cross-Entropy	13
Rumus 2.7	L1 Regularization Loss Function	14
Rumus 2.8	L2 Regularization Loss Function	14
Rumus 2.9	Random Forest Majority Voting	15
Rumus 2.10	Gini Index	16
Rumus 2.11	Accuracy	19
Rumus 2.12	Precision	19
Rumus 2.13	Recall	19
Rumus 2.14	F1-Score	20
Rumus 2.15	True Positive Rate	20
Rumus 2.16	False Positive Rate	20



DAFTAR LAMPIRAN

Lampiran 1	Hasil Persentase Turnitin	84
Lampiran 2	Formulir Bimbingan Pembimbing 1	95
Lampiran 3	Formulir Bimbingan Pembimbing 2	98



UMN
UNIVERSITAS
MULTIMEDIA
NUSANTARA