



**PERBANDINGAN METODE NAIVE BAYES DAN K-
NEAREST NEIGHBOR PADA KLASIFIKASI MORFOLOGI
GEN SEL DARAH PUTIH**

SKRIPSI

MUHAMMAD NUR'ADLI HASBI GUMAY

1810511103

UNIVERSITAS PEMBANGUNAN NASIONAL "VETERAN" JAKARTA

FAKULTAS ILMU KOMPUTER

PROGRAM STUDI INFORMATIKA

2022



**PERBANDINGAN METODE NAIVE BAYES DAN K-
NEAREST NEIGHBOR PADA KLASIFIKASI MORFOLOGI
GEN SEL DARAH PUTIH**

SKRIPSI

**Diajukan Sebagai Salah Satu Syarat untuk Memperoleh
Gelar Sarjana Komputer**

MUHAMMAD NUR'ADLI HASBI GUMAY

1810511103

UNIVERSITAS PEMBANGUNAN NASIONAL "VETERAN" JAKARTA

FAKULTAS ILMU KOMPUTER

PROGRAM STUDI INFORMATIKA

2022

LEMBAR PENGESAHAN

LEMBAR PENGESAHAN

Dengan ini dinyatakan bahwa Skripsi berikut:

Nama : Muhammad Nur'adli Hasbi Gumay
Nim : 1810511103
Program Studi : S1 Informatika
Judul : PERBANDINGAN METODE NAIVE BAYES DAN
K-NEAREST NEIGHBOR PADA KLASIFIKASI
MORFOLOGI GEN SEL DARAH PUTIH

Telah berhasil dipertahankan di hadapan Tim penguji dan diterima sebagai bagian persyaratan yang diperlukan untuk memperoleh gelar Sarjana Komputer pada Program Studi S1 Informatika, Fakultas Ilmu Komputer, Universitas Pembangunan Nasional Veteran Jakarta.



Iin Ernawati, S.Kom., M.Si.
Penguji 1



Noor Falih, S.Kom., M.T.
Penguji 2



Yuni Widiastiti, S.Kom., M.Si.
Pembimbing 1



Mayanda Mega Santoni, M.Kom.
Pembimbing 2



Dr. Ernawati, M.Kom.
Dekan



Desta Sandya Prasvita, S.Kom.,
M.Kom.
Ketua Program Studi

Ditetapkan di : Jakarta
Tanggal Pengesahan : 27 Juli 2022



PERNYATAAN ORISINALITAS

Tugas Skripsi ini adalah hasil karya sendiri, dan semua sumber yang dikutip maupun yang dirujuk telah saya nyatakan dengan benar.

Nama : Muhammad Nur'adli Hasbi Gumay

NIM : 18105011103

Tanggal : Kamis, 23 Juni 2022

Bilamana di kemudian hari ditemukan ketidaksesuaian dengan pernyataan saya ini, maka saya bersedia dituntut dan diproses sesuai dengan ketentuan yang berlaku.

Jakarta, 23 Juni 2022

Yang Menyatakan,



(Muhammad Nur'adli
Hasbi Gumay)

PERNYATAAN PERSETUJUAN PUBLIKASI SKRIPSI UNTUK KEPENTINGAN AKADEMIS

Sebagai civitas akademik Universitas Pembangunan Nasional “Veteran” Jakarta, saya yang bertanda tangan di bawah ini:

Nama : Muhammad Nur’adli Hasbi Gumay
NIM : 1810511103
Fakultas : Ilmu Komputer
Program Studi : Informatika

Demi pengembangan ilmu pengetahuan, menyetujui untuk memberikan kepada Universitas Pembangunan Nasional “Veteran” Jakarta Hak Bebas Royalti Non eksklusif (*Non-exclusive Royalty Free Right*) atas karya ilmiah saya yang berjudul:

PERBANDINGAN METODE NAIVE BAYES DAN K-NEAREST NEIGHBOR PADA KLASIFIKASI MORFOLOGI GEN SEL DARAH PUTIH

Beserta perangkat yang ada (jika diperlukan). Dengan Hak Bebas Royalti ini Universitas Pembangunan Nasional “Veteran” Jakarta berhak menyimpan, mengalih media/formatkan, mengelola dalam bentuk pangkalan data (*database*), merawat, dan mempublikasikan Tugas Skripsi saya selama tetap mencantumkan nama saya sebagai penulis/pencipta dan sebagai pemilih Hak Cipta. Demikian pernyataan ini saya buat dengan sebenarnya.

Dibuat Di : Jakarta
Pada Tanggal : 23 Juni 2022

Yang Menyatakan,



(Muhammad Nur’adli Hasbi
Gumay)

LEMBAR PERSETUJUAN

Dengan ini menyatakan bahwa laporan tugas akhir sebagai berikut:

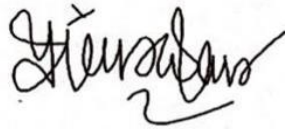
Nama : Muhammad Nur'adli Hasbi Gumay
NIM : 1810511103
Program Studi : Informatika
Judul : Perbandingan Metode Naive Bayes dan K-Nearest Neighbor Pada Klasifikasi Morfologi Gen Sel Darah Putih.

Sebagai bagian persyaratan yang diperlukan untuk mengikuti ujian Sidang Tugas Akhir/Skripsi pada Program Studi Informatika, Fakultas Ilmu Komputer, Universitas Pembangunan Nasional Veteran Jakarta.

Menyetujui,

Dosen Pembimbing 1

Dosen Pembimbing 2

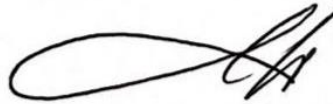


Yuni Widiastiwi, S.Kom., M.Si.

Mayanda Mega Santoni, M.Kom.

Mengetahui,

Ketua Program Studi



Desta Sandya Prasvita, S. Komp., M.Kom.

Ditetapkan : Jakarta

Tanggal Persetujuan : Kamis, 23 Juni 2022

PERBANDINGAN METODE NAIVE BAYES DAN K-NEAREST NEIGHBOR PADA KLASIFIKASI MORFOLOGI GEN SEL DARAH PUTIH

Muhammad Nur'adli Hasbi Gumay

ABSTRAK

Dibidang kesehatan, mendiagnosis penyakit leukemia merupakan hal yang sulit karena masih didiagnosis secara manual dengan bantuan dokter. Diagnosis manual tersebut dapat mengalami kesalahan yang disebabkan oleh kelalaian manusia. Dari permasalahan tersebut, maka dibutuhkan diagnosis jenis penyakit leukemia menggunakan kecanggihan teknologi yaitu *machine learning* untuk mengatasi permasalahan tersebut. Dalam penelitian ini, *machine learning* tersebut mengolah data yang berasal dari jenis leukemia yaitu *Acute Myeloid Leukemia* (AML) dan *Acute Lymphoblastic Leukemia* (ALL) berdasarkan ciri morfologi gen sel darah putih tersebut. Metode pengklasifikasian data yang digunakan untuk penelitian ini yaitu K-Nearest Neighbor (K-NN) dan Naïve Bayes yang kemudian kedua metode klasifikasi tersebut dibandingkan untuk melihat metode klasifikasi yang terbaik. Penelitian ini menggunakan praproses *data cleaning*, seleksi fitur, dan *scaling* untuk meningkatkan nilai akurasi. Hasil dari penelitian ini adalah metode klasifikasi K-Nearest Neighbors (K-NN) merupakan klasifikasi yang terbaik dengan nilai akurasi yang menggunakan kurva ROC/AUC bernilai 0.952 jika dibandingkan dengan metode klasifikasi Naïve Bayes yaitu 0.912.

Kata Kunci: Perbandingan, Naïve Bayes, K-Nearest Neighbors, Leukemia, *Machine Learning*

COMPARISON OF NAIVE BAYES AND K-NEAREST NEIGHBOR METHODS IN WHITE BLOOD CELL GENE MORPHOLOGICAL CLASSIFICATION

Muhammad Nur'adli Hasbi Gumay

ABSTRACT

In the health sector, the diagnosis of leukemia is a difficult thing because it is still diagnosed manually with the help of a doctor. The diagnosis manual may suffer from errors caused by human negligence. From these problems, it is necessary to diagnose the type of leukemia using advanced technology, namely Machine learning to overcome these problems. In this study, the machine learning processes data from the types of leukemia, namely Acute Myeloid Leukemia (AML) and Acute Lymphoblastic Leukemia (ALL) based on the morphological characteristics of the white blood cell genes. The data classification methods used for this research are K-Nearest Neighbor (K-NN) and Naïve Bayes, then the two classification methods are compared to see the best classification method. This study uses preprocessing of data cleaning, feature selection, and scaling to increase the accuracy value. The results of this study are the K-Nearest Neighbors (K-NN) classification method is the best classification with an accuracy value using the ROC/AUC curve worth 0.952 when compared to the Naïve Bayes classification method, which is 0.912.

Keywords: *Comparasion, Naïve Bayes, K-Nearest Neighbors, Leukemia, Machine Learning*

KATA PENGANTAR

Puji dan syukur peneliti panjatkan kehadirat Allah SWT karena berkat rahmat dan anugerahNya sehingga peneliti dapat menyelesaikan Laporan Tugas Akhir ini dengan baik dan lancar. Penulisan laporan ini adalah untuk memenuhi dan sebagai syarat untuk menempuh proses lebih lanjut pelaksanaan Tugas Akhir di Program Studi Informatika, Fakultas Ilmu Komputer, Universitas Pembangunan Nasional “Veteran” Jakarta. Judul yang dipilih dalam penelitian ini adalah “*Perbandingan Metode Naive Bayes dan K-Nearest Neighbor Pada Klasifikasi Morfologi Gen Sel Darah Putih*”.

Peneliti menyadari bahwa tanpa bantuan dan bimbingan dari semua pihak penyusunan laporan Tugas Akhir ini tidak dapat berjalan dengan baik. Oleh karena itu peneliti mengucapkan terima kasih kepada:

1. Kepada kedua orang tua peneliti yang telah memberikan doa dan semangat selama proses pembuatan Proposal Tugas Akhir.
2. **Ibu Dr. Ermatita, M. Kom.** selaku Dekan Fakultas Ilmu Komputer Universitas Pembangunan Nasional “Veteran” Jakarta.
3. **Pak Desta Sandya Prasvita, S. Komp., M. Kom.** selaku kepala program studi Informatika Universitas Pembangunan Nasional “Veteran” Jakarta
4. **Ibu Yuni Widiastiwi, S. Kom., M.Si.** selaku dosen Pembimbing 1
5. **Ibu Mayanda Mega Santoni, S. Kom., M.Kom.** selaku Dosen Pembimbing 2.
6. Kepada Lovya Neysa serta teman-teman mahasiswa Informatika angkatan 2018 yang selalu memberikan semangat, bantuan, serta dukungan untuk peneliti dalam menyelesaikan Proposal Tugas Akhir.

Jakarta, 27 Mei 2022

Penulis



Muhammad Nur'adli Hasbi Gumay

DAFTAR ISI

LEMBAR PENGESAHAN	iii
PERNYATAAN ORISINALITAS	iv
PERNYATAAN PERSETUJUAN PUBLIKASI SKRIPSI UNTUK KEPENTINGAN AKADEMIS	v
LEMBAR PERSETUJUAN.....	vi
ABSTRAK	vii
<i>ABSTRACT</i>	viii
KATA PENGANTAR	ix
DAFTAR ISI.....	x
DAFTAR TABEL.....	xiii
DAFTAR GAMBAR	xiv
DAFTAR LAMPIRAN.....	xv
DAFTAR SIMBOL.....	xvi
BAB 1 PENDAHULUAN	1
1.1 Latar Belakang Masalah	1
1.2 Rumusan Masalah	3
1.3 Tujuan Penelitian.....	3
1.4 Manfaat Penelitian.....	3
1.5 Ruang Lingkup	3
1.6 Sistematika Penulisan.....	4
BAB 2 TINJAUAN PUSTAKA	6
2.1 Sel Kanker dan Sel Non-Kanker	6
2.2 Sel Darah Putih.....	7
2.2.1 <i>Acute Myeloid Leukemia (AML)</i>	7

2.2.2	<i>Lymphoblastic Leukemia (ALL)</i>	8
2.3	<i>Data Mining</i>	8
2.4	Pra Proses	9
2.4.1	<i>Data Cleaning</i>	9
2.4.2	Seleksi Fitur	10
2.4.3	<i>Scaling</i>	10
2.5	<i>Principal Component Analysis (PCA)</i>	11
2.6	Naïve Bayes.....	12
2.7	<i>K-Nearest Neighbor (KNN)</i>	13
2.8	Evaluasi	15
2.8.1	Confusion Matrix	15
2.8.2	Kurva <i>Receiver Operating Characteristic (ROC)</i>	15
2.9	Penelitian Terdahulu.....	18
BAB 3	METODOLOGI PENELITIAN	22
3.1	Tahapan Penelitian	22
3.2	Identifikasi Masalah	23
3.3	Studi Literatur.....	23
3.4	Pengumpulan <i>Dataset</i>	23
3.5	<i>Data Cleaning</i>	24
3.6	Pembagian Data.....	24
3.7	Praproses Data	25
3.7.1	<i>Scaling</i>	25
3.7.2	Seleksi Fitur	25
3.8	<i>Principal Component Analysis (PCA)</i>	25
3.9	Klasifikasi.....	26
3.10	Evaluasi.....	26

3.11	Perbandingan	30
3.12	Bahan dan Alat Pendukung yang Digunakan	30
3.12.1	<i>Hardware</i>	30
3.12.2	<i>Software</i>	31
3.13	Jadwal Penelitian	31
BAB 4	HASIL DAN PEMBAHASAN	33
4.1	Pengumpulan Data	33
4.2	Data <i>Cleaning</i>	36
4.3	Pembagian Data.....	38
4.4	Praproses Data	39
4.4.1	<i>Scaling</i>	40
4.4.2	Seleksi Fitur	41
4.5	<i>Principal Component Analysis (PCA)</i>	42
4.6	Klasifikasi.....	43
4.6.1	Naïve Bayes	43
4.6.2	KNN	45
4.6.3	Perbandingan Metode Klasifikasi Naïve Bayes dan KNN	47
4.7	Evaluasi	48
BAB 5	PENUTUP	53
5.1	Kesimpulan	53
5.2	Saran.....	53
	DAFTAR PUSTAKA	54
	RIWAYAT HIDUP.....	58
	LAMPIRAN.....	59

DAFTAR TABEL

Tabel 2.1 Klasifikasi Metode Nearest Neighbour.....	13
Tabel 3.1 Confusion Matrix	28
Tabel 3.2. Kriteria Nilai AUC.....	29
Tabel 3.3 Jadwal Tahapan Penelitian.....	31
Tabel 4.1 Ukuran Dataset Train dan Dataset Test	35
Tabel 4.2 Hasil Pembagian Dataset	38
Tabel 4.3 Dataset train_target	38
Tabel 4.4 Dataset test_target.....	39
Tabel 4.5 Jumlah Dataset Train Hasil Logistic Regresion.....	42
Tabel 4.6 Nilai Confusion Matrix metode Naïve Bayes	44
Tabel 4.7 Hasil Model Klasifikasi menggunakan Naïve Bayes.....	45
Tabel 4.8 Nilai Confusion Matrix metode KNN.....	46
Tabel 4.9 Hasil untuk Metode Klasifikasi K-Nearest Neighbors	46


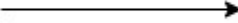



DAFTAR GAMBAR

Gambar 2.1 Komponen pada citra sel darah	7
Gambar 2.2 Arsitektur Sistem Data Mining	9
Gambar 3.1 Flowchart Metode Penelitian	22
Gambar 4.1 Bagan Pembagian Jumlah Dataset Target	34
Gambar 4.2 Data Train Kanker	34
Gambar 4.3 Data Test Kanker.....	35
Gambar 4.4 Ukuran Korelasi dari ALL dan AML.....	35
Gambar 4.5 Dataset yang telah di Cleaning	38
Gambar 4.6 Hasil sebelum proses Scaling	41
Gambar 4.7 Hasil setelah proses Scaling	41
Gambar 4.8 Persebaran Dataset setelah menggunakan PCA	43
Gambar 4.9 Perbandingan Nilai KNN dan Naïve Bayes	47
Gambar 4.10 AUROC/PR CURVES False Positive dan True Positive Rate	49
Gambar 4.11 AUROC/PR CURVES precision dan recall.....	49
Gambar 4.12 Hasil pembagian dataset 70% dan 30%	51

DAFTAR LAMPIRAN

Lampiran 1 Dataset Target.....	59
Lampiran 2 Dataset Train.....	62
Lampiran 3 Dataset Test	64
Lampiran 4 Gabungan Dataset Train dan Test Sesudah Transpose.....	66
Lampiran 5 Gabungan Dataset Train dan Test setelah Cleaning.....	67
Lampiran 6 Pembagian Dataset	68
Lampiran 7 Hasil Scaling Dataset Train	70
Lampiran 8 Dataset Train setelah Proses PCA (Principal Component Analysis)	71
Lampiran 9 Hasil Uji Turnitin	72

DAFTAR SIMBOL

No.	Gambar Simbol	Nama	Keterangan
1.		<i>Start / End</i>	Simbol oval yang menggambarkan dimulainya atau berhenti suatu diagram alir.
2.		<i>Flow</i>	Simbol garis dengan tanda panah yang menggambarkan aliran data atau prosedur dari satu tahap ke tahap yang lain.
3.		<i>Process</i>	Simbol berbentuk persegi panjang dan menggambarkan proses atau kegiatan apa yang sedang terjadi.
4.		<i>Input / Output</i>	Menyatakan proses <i>input</i> atau <i>output</i>
5.		<i>Document</i>	Menyatakan bahwa <i>input</i> berasal dari dokumen dalam bentuk fisik atau <i>output</i> yang perlu dicetak