



**IMPLEMENTASI KOMBINASI METODE *RESAMPLING* PADA
KLASIFIKASI PENYAKIT STROKE DENGAN ALGORITMA *K-NEAREST NEIGHBOR* DAN SELEKSI FITUR *INFORMATION GAIN***

SKRIPSI

**MUHAMMAD FATHURRAHMAN
NIM. 1910511058**

**PROGRAM STUDI INFORMATIKA
FAKULTAS ILMU KOMPUTER
UNIVERSITAS PEMBANGUNAN NASIONAL VETERAN JAKARTA
2023**



**IMPLEMENTASI KOMBINASI METODE *RESAMPLING* PADA
KLASIFIKASI PENYAKIT STROKE DENGAN ALGORITMA *K-NEAREST NEIGHBOR* DAN SELEKSI FITUR *INFORMATION GAIN***

SKRIPSI

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

**MUHAMMAD FATHURRAHMAN
NIM. 1910511058**

**PROGRAM STUDI INFORMATIKA
FAKULTAS ILMU KOMPUTER
UNIVERSITAS PEMBANGUNAN NASIONAL VETERAN JAKARTA
2023**

PERNYATAAN ORISINALITAS

PERNYATAAN ORISINALITAS

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

Nama : Muhammad Fathurrahman

NRP : 1910511058

Tanggal : 16 juli 2023

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

Jakarta, 16 Juli 2023

Yang Menyatakan



(Muhammad Fathurrahman)

PERNYATAAN PERSETUJUAN PUBLIKASI SKRIPSI

UNTUK KEPENTINGAN AKADEMIS

PERNYATAAN PERSETUJUAN PUBLIKASI SKRIPSI

UNTUK KEPENTINGAN AKADEMIS

Sebagai civitas akademis Universitas Pembangunan Nasional Veteran Jakarta.
Saya bertanda tangan di bawah ini:

Nama : Muhammad Fathurrahman

NRP : 1910511058

Fakultas : Ilmu Komputer

Program Studi : SI Informatika

Demi pengembangan ilmu pengetahuan, menyetujui untuk memberikan kepada
Universitas Pembangunan Nasional Veteran Jakarta Hak Bebas Royalti Non-
Ekslusif (*Non-Exclusive Royalty Free Right*) atas karya Ilmiah saya yang berjudul:

IMPLEMENTASI KOMBINASI METODE RESAMPLING PADA KLASIFIKASI PENYAKIT STROKE DENGAN ALGORITMA K- NEAREST NEIGHBOR DAN SELEKSI FITUR INFORMATION GAIN

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 Skripsi saya selama tetap mencantumkan nama
saya sebagai penulis/pencipta dan sebagai pemilik Hak Cipta.

Demikian pernyataan ini saya buat dengan sebenarnya.

Dibuat di : Jakarta

Pada Tanggal : 16 Juli 2023

Yang Menyatakan,



(Muhammad Fathurrahman)

LEMBAR PENGESAHAN

LEMBAR PENGESAHAN

Tugas Akhir ini diajukan oleh:

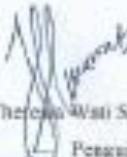
Nama : Muhammad Fathurrahman
NIM : 1910511058
Program Studi : SI Informatika
Judul Tugas Akhir : Implementasi Kombinasi Teknik *Resampling* Pada Klasifikasi Penyakit Stroke menggunakan Algoritma K-Nearest Neighbor dan Seleksi Fitur Information Gain.

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



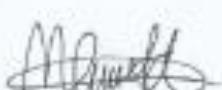
Bayu Hananto, S.Kom., M.Kom.

Pengaji I



Therawati S.Kom., M.TI.

Pengaji II



Nurhafizah Matondang, S.Kom., MM

Pemimpin



Dr. Emanuella, M.Kom.

Dekan



Dr. Widya Chisti, M.I.T.

Kepala Program Studi

Ditetapkan di : Jakarta
Tanggal Ujian : 06 Juli 2023



IMPLEMENTASI KOMBINASI METODE RESAMPLING PADA KLASIFIKASI PENYAKIT STROKE DENGAN ALGORITMA K-NEAREST NEIGHBOR DAN SELEKSI FITUR INFORMATION GAIN

MUHAMMAD FATHURRAHMAN

ABSTRAK

Salah satu masalah utama dalam dunia medis adalah stroke. Stroke menjadi penyebab kematian kedua di dunia. Berdasarkan hasil Riset Kesehatan Dasar (Riskesdar) tahun 2018 prevalensi stroke di Indonesia sebanyak 713,783 orang yang menderita stroke setiap tahunnya. Namun, dalam mendiagnosa stroke diperlukan waktu yang cukup panjang. Mengingat tiap menit ada sel yang mati karena penyumbatan aliran pada otak. Perlu dilakukan diagnosis sedini mungkin untuk mengurangi resiko dari pasien. *Data mining* dapat digunakan sebagai prediksi penyakit. Dalam pembuatan model *data mining*, *imbalanced data* merupakan masalah karena dapat berdampak buruk pada hasil klasifikasi dimana model *machine learning* akan lebih memperhatikan kelas mayoritas dan mengabaikan kelas minoritas. Pada penelitian telah dilakukan prediksi penyakit stroke menggunakan algoritma K-Nearest Neighbor dengan mengkombinasikan teknik *resampling* seperti SMOTE, *Tomek Links* dan ENN. Serta penelitian dilakukan untuk mengetahui pengaruh seleksi fitur *information gain* terhadap model. Melalui proses *10 fold cross validation* diketahui model *machine learning* K-NN dengan SMOTE dan *Tomek Links* mampu memprediksi stroke dengan akurasi 83,5%, *f1-score* 12,5%, dan recall 24,7%. Kemudian untuk K-NN dengan SMOTE dan ENN diperoleh akurasi 78%, *f1-score* 16,8%, dan *recall* 45%. Ketika dilakukan seleksi fitur *information gain* terdapat peningkatan performa pada kedua metode tersebut. SMOTE dan *Tomek Links* menghasilkan akurasi 79,9%, *f1-score* 18,3%, dan recall 46,6% serta kombinasi SMOTE dan ENN diperoleh akurasi 76%, *f1-score* 20%, dan *recall* 59%. Setelah dilakukan pengujian diketahui bahwa teknik *resampling* dapat meningkatkan performa model pada kasus data yang tidak seimbang dari nilai *recall* dan *f1-score* sebesar 54% dan 7%.

Kata kunci : *Synthetic Minority Over-sampling*, *K-Nearest-Neighbor*, *Stroke*, *Tomek Links*, *Edited Nearest Neighbor*, *Information gain*.

IMPLEMENTATION OF A COMBINATION OF RESAMPLING METHODS IN STROKE CLASSIFICATION USING K-NEAREST NEIGHBOR ALGORITHM AND INFORMATION GAIN FEATURE SELECTION

MUHAMMAD FATHURRAHMAN

ABSTRACT

One of the main problems in the medical world is stroke. Stroke is the second cause of death in the world. Based on the results of Basic Health Research (Rskesdar) in 2018, the prevalence of stroke in Indonesia is 713,783 people who suffer from stroke every year. However, diagnosing a stroke takes quite a long time. Considering that every minute there are cells that die due to blockage of flow in the brain. Data mining can be used as a prediction of disease. In making data mining models, data imbalance is a problem because it can have a negative impact on the classification results where the machine learning model will pay more attention to the majority class and ignore the minority class. In this study, stroke prediction was carried out using the K-Nearest Neighbor algorithm by combining resampling techniques such as SMOTE, Tomek Links, and ENN. As well as research conducted to determine the effect of the search feature information obtained on the model. Through a 10 fold cross validation process, it is known that the K-NN machine learning model with SMOTE and Tomek Links is able to predict stroke with an accuracy of 83.5%, an f1-score of 12.5%, and a recall of 24.7%. Then K-NN with SMOTE and ENN obtained 78% accuracy, f1 score 16.8%, and recall 45%. When the selection of information gain features is carried out, there is an increase in performance in both methods. SMOTE and Tomek Links produce 79.9% accuracy, 18.3% f1-score, and 46.6% recall and the combination of SMOTE and ENN obtains 76% accuracy, 20% f1-score, and 59% recall. After the experiments, it is known that the resampling technique can improve the performance of the model in the case of imbalanced data from the recall and f1-score values by 54% and 7%.

Keyword: *Synthetic Minority Over-sampling, K-Nearest-Neighbor, Stroke, Tomek Links, Edited Nearest Neighbor, Information gain.*

KATA PENGANTAR

Segala puji dan syukur kehadirat Allah SWT atas hadirat dan nikmatnya, rahmat yang selalu dilimpahkan, sehingga penulis dapat menyelesaikan skripsi yang berjudul “Implementasi Kombinasi Metode *Resampling* Pada Klasifikasi Penyakit Stroke Dengan Algoritma *K-Nearest Neighbor* dan Seleksi Fitur *Information Gain*” sebagai salah satu syarat dalam menyelesaikan Program S1 Informatika Universitas Pembangunan Nasional Veteran Jakarta.

Banyak tantangan yang penulis temui dalam penyusunan skripsi ini dan penulis berterima kasih kepada semua pihak yang telah membantu penulis selama masa percobaan ini. Selama menyusun skripsi ini telah banyak hambatan yang penyusun lewati dan tanpa bantuan banyak pihak tentu akan sulit untuk penyusun menyelesaikan skripsi ini, untuk itu penulis mengucapkan terima kasih kepada :

1. Dr. Ermatita, M.Kom. Selaku Dekan Fakultas Ilmu Komputer Universitas Pembangun Nasional Veteran Jakarta.
2. Dr. Widya Cholil, M.I.T. selaku Ketua Program Studi Sarjana Jurusan Informatika.
3. Nur Hafifah Matondang, S.Kom., M.M., M.T.I selaku Pembimbing Akademik yang memberikan bimbingan dan mengarahkan penyusun dalam menyelesaikan skripsi ini.
4. Kedua Orang tua yang telah mendukung saya dalam hal menuntut ilmu sedari dulu.
5. Teman seperjuangan yang telah menyelesaikan skripsi bersama saya.

Akhir kata penulis ucapan banyak terima kasih kepada semua pihak yang membantu serta semoga Allah SWT selalu melimpahkan karunianya dalam semua amal kebaikan kita serta diberikan balasan yang baik pula. Amin.

Jakarta, Juni 2023

Penulis

DAFTAR ISI

Halaman

LEMBAR SAMPUL

LEMBAR JUDUL	i
PERNYATAAN ORISINALITAS.....	iii
PERNYATAAN PERSETUJUAN PUBLIKASI SKRIPSI UNTUK KEPENTINGAN AKADEMIS	iv
LEMBAR PENGESAHAN	v
ABSTRAK	vi
ABSTRACT	vii
KATA PENGANTAR.....	viii
DAFTAR ISI.....	ix
DAFTAR SIMBOL	xii
DAFTAR GAMBAR.....	xiii
DAFTAR TABEL	xv
BAB I.....	1
PENDAHULUAN.....	1
1.1. Latar Belakang	1
1.2. Rumusan Masalah	3
1.3. Tujuan Penelitian	4
1.4. Ruang Lingkup	4
1.5. Manfaat Penelitian	4
1.6. Luaran yang Diharapkan	5
1.7. Sistematika Penulisan	5
BAB II	7
TINJAUAN PUSTAKA	7
2.1. Stroke.....	7
2.2. Data mining	8
2.3. Klasifikasi.....	10
2.4. K-Nearest Neighbor (K-NN)	11
2.5. Min-Max Normalization	12
2.6. Boxplot	13

2.7. <i>Imbalanced Data</i>	14
2.8. <i>Resampling Data</i>	15
2.9. <i>Synthetic Minority Oversampling Technique (SMOTE)</i>	16
2.10. <i>Tomek Links</i>	17
2.11. <i>Edited Nearest Neighbors (ENN)</i>	19
2.12. <i>Information Gain</i>	19
2.13. <i>10-Fold Cross Validation</i>	20
2.14. <i>Classification Metrics</i>	21
2.15. Penelitian Terdahulu	25
BAB III.....	28
METODOLOGI PENELITIAN	28
3.1. Kerangka Berpikir	28
3.1.1. Identifikasi Masalah.....	28
3.1.2. Studi Literatur	28
3.1.3. Pengumpulan Data	29
3.1.4. Praproses Data.....	29
3.1.5. Seleksi Fitur	30
3.1.6. <i>10 Fold Cross Validation</i>.....	30
3.1.7. Evaluasi Model	30
3.1.8. Pelatihan Model.....	31
3.1.9. Implementasi Sistem	31
3.2. Alat Bantu Penelitian	31
3.2.1. Perangkat Keras	31
3.2.2. Perangkat Lunak.....	31
3.3. Jadwal Kegiatan Penelitian	33
BAB IV	34
HASIL DAN PEMBAHASAN	34
4.1. <i>Dataset</i>.....	34
4.2. Praproses Data.....	37
4.2.1. <i>Data Cleansing</i>	37
4.2.2. <i>Label Encoding</i>.....	42
4.2.3. Normalisasi.....	43
4.3. Seleksi Fitur	44

4.4. 10 Fold Cross Validation.....	46
4.5. Train Model	60
4.6. Rancangan Sistem	62
4.7. Implementasi Sistem	63
4.7.1 Halaman Input Data.....	63
BAB V.....	66
PENUTUP.....	66
5.1. Kesimpulan	66
5.2. Saran	67
DAFTAR PUSTAKA.....	68
LAMPIRAN.....	76
 RIWAYAT HIDUP	 89

DAFTAR SIMBOL

Simbol <i>Flowchart</i>			
No	Gambar	Nama	Keterangan
1.		<i>Terminal (start, end)</i>	Menggambarkan bagaimana kegiatan dimulai atau kegiatan berakhir
2.		<i>Flow Direction</i>	Menggambarkan hubungan antar simbol yang menyatakan suatu jalannya proses dalam sistem
3.		<i>Process</i>	Menggambarkan deskripsi dari proses yang dijalankan

DAFTAR GAMBAR

	Halaman
Gambar 2.1. Tahapan Data Mining (Sumber: Ha et al., 2011)	9
Gambar 2.2. Proses Klasifikasi (Sumber: Ha et al., 2011)	11
Gambar 2.3. <i>Plotbox</i> (Sumber: MIT, 2016).....	14
Gambar 2.4. Contoh Penerapan SMOTE $k = 5$ (Sumber: Beckmann et al., 2015)	17
Gambar 2.5. <i>Tomek Links</i> (Sumber: Pereira dkk., 2020)	18
Gambar 2.6. <i>10-fold Cross Validation</i> (Sumber: Berrar, 2018).....	21
Gambar 2.7. <i>Receiver Operating Characteristic</i> (Sumber: Kulkarni dkk., 2020)24	24
Gambar 3.1. Kerangka Berpikir	28
Gambar 4.1. Proporsi Data.....	36
Gambar 4.2. Indikator <i>Missing Value</i>	38
Gambar 4.3. Informasi <i>Dataset</i>	39
Gambar 4.4. Boxplot Pada Atribut bmi	39
Gambar 4.5. Setelah Penghapusan Nilai <i>Outlier</i>	41
Gambar 4.6. Proporsi Data Setelah Dilakukan Praproses Data	41
Gambar 4.7. <i>10-fold Cross Validation</i> (Sumber: (Andrade dkk., 2020).....	46
Gambar 4.8. <i>Oversampling k-Fold Cross Validation</i>	47
Gambar 4.9. Sebaran Data Latih	48
Gambar 4.10. <i>Confusion matrix</i> dari model KNN dengan data yang <i>Imbalanced</i>	50
Gambar 4.11. <i>Confusion Matrix</i> dari model KNN dengan data yang <i>Balance</i> dengan metode <i>Resampling</i> SMOTE.....	51
Gambar 4.12. <i>Confusion Matrix</i> dari model KNN dengan data yang <i>Balance</i> dengan metode <i>Resampling</i> SMOTE dan <i>Tomek Links</i>	52
Gambar 4.13. Confusion Matrix dari model KNN dengan data yang <i>Balance</i> dengan metode <i>Resampling</i> SMOTE dan ENN	53
Gambar 4.14. Hasil Evaluasi K-NN <i>Imbalance</i>	54
Gambar 4.15. Hasil Evaluasi SMOTE	55
Gambar 4.16. Hasil Evaluasi SMOTE Tomek.....	56
Gambar 4.17. Hasil Evaluasi SMOTE ENN	57

Gambar 4.18. Perbandingan Evaluasi	58
Gambar 4.19. Kurva ROC dari SMOTE,SMOTETomek, dan SMOTEENN	59
Gambar 4.20. Rancangan Sistem	63
Gambar 4.21. Halaman <i>Input</i> Data (kosong)	63
Gambar 4.22. Halaman <i>Input</i> Data	64
Gambar 4.23. Tampilan Hasil Prediksi	64

DAFTAR TABEL

	Halaman
Tabel 3.1. Jadwal Kegiatan Penelitian	33
Tabel 4.1. Informasi Nilai Pada Dataset	34
Tabel 4.2. Sampel Data Missing Value.....	37
Tabel 4.3. Pemetaan Atribut	42
Tabel 4.4. Data sebelum Normalisasi	43
Tabel 4.5. Data Setelah Normalisasi.....	44
Tabel 4.6. Hasil <i>Information Gain</i>	45
Tabel 4.7. Jumlah Data latih	47
Tabel 4.8. Hasil Evaluasi <i>10-fold Cross Validation</i>	49
Tabel 4.9 Sampel data latih.....	60
Tabel 4.10 Sampel data uji.....	60
Tabel 4.11 Sampel data latih setelah dihitung jarak	61
Tabel 4.12 Data uji setelah diprediksi.....	61