

**KLASIFIKASI PENYAKIT DIABETES *MELLITUS*
MENGUNAKAN METODE *SYNTHETIC MINORITY OVER-
SAMPLING TECHNIQUE (SMOTE) RANDOM FOREST***

SKRIPSI

**Diajukan untuk memenuhi salah satu syarat kelulusan
di Program Studi Sains Data**



Disusun Oleh:

NOFA AULIYATUL MAULIDIYYAH

20083010029

**UNIVERSITAS PEMBANGUNAN NASIONAL "VETERAN" JAWA TIMUR
FAKULTAS ILMU KOMPUTER
PROGRAM STUDI SAINS DATA
SURABAYA
2024**

LEMBAR PENGESAHAN

**KLASIFIKASI PENYAKIT *DIABETES MELLITUS* MENGGUNAKAN
METODE *SYNTHETIC MINORITY OVER-SAMPLING TECHNIQUE*
(SMOTE) *RANDOM FOREST***

SKRIPSI

Diajukan untuk memenuhi salah satu syarat memperoleh gelar Sarjana Sains Data
pada : Senin, 15 juli 2024

**Program Studi S-1 Sains Data
Fakultas Ilmu Komputer**

**Universitas Pembangunan Nasional Veteran Jawa Timur
Surabaya**

Oleh :

NOFA AULIYATUL MAULIDIYYAH

NPM. 20083010029

Disetujui oleh Tim Penguji Skripsi :

Penguji 1

Penguji 2

Tresna Maulana Fahrudin, S.ST., M.T.

NIP. 1993050120220301007

Pembimbing 1

Amri Muhaimin, S.Stat., M.Stat., M.S.

NIP. 21119950723270

Pembimbing 2

Trimono, S.Si., M.Si

NIP. 199509082022031003

Fakultas Ilmu Komputer
Dekan,

Avioli Terza Damaliana, S.Si., M.Stat

NIP. 199408022022032015

Mengetahui,

Program Studi Sains Data
Fakultas Ilmu Komputer
Koordinator,

Prof. Dr. Ir. Novrina Hendrasarie, MT

NIP. 196811261994032001

Dr. Eng. Ir. Dwi Arman Prasetya, ST, MT, IPU.

NIP. 198012052005011002

Surabaya, Juli, 2024

SURAT PERNYATAAN

Saya yang bertanda tangan di bawah ini:

Nama : Nofa Auliyatul Maulidiyyah
NPM : 20083010029
Program Studi : Sains Data

Menyatakan bahwa judul Skripsi / Tugas Akhir sebagai berikut:

Klasifikasi Penyakit *Diabetes Mellitus* Menggunakan Metode *Synthetic Minority Over-Sampling Technique (SMOTE)* Random Forest

Bukan merupakan plagiat dari Skripsi/ Tugas Akhir/ Penelitian orang lain dan juga bukan merupakan produk/ *software*/ hasil karya yang saya beli dari orang lain

Saya juga menyatakan bahwa Skripsi/ Tugas Akhir ini adalah pekerjaan saya sendiri, kecuali yang dinyatakan dalam Daftar Pustaka, dan tidak pernah diajukan untuk syarat memperoleh gelar di Universitas Pembangunan Nasional "Veteran" Jawa Timur maupun di institusi pendidikan lain.

Jika ternyata dikemudian hari pernyataan ini terbukti tidak benar, maka Saya bertanggung jawab penuh dan siap menerima segala konsekuensi, termasuk pembatalan ijazah dikemudian hari

Surabaya, 30 April 2024

Hormat Saya



Nofa Auliyatul Maulidiyyah
NPM.

ABSTRAK

KLASIFIKASI PENYAKIT *DIABETES MELLITUS* MENGGUNAKAN METODE *SYNTHETIC MINORITY OVER-SAMPLING TECHNIQUE* (SMOTE) *RANDOM FOREST*

Nama Mahasiswa / NPM : Nofa Auliyatul Maulidiyyah / 20083010029
Program Studi : Sains Data, FASILKOM, UPN Veteran Jatim
Dosen Pembimbing 1 : Trimono, S.Si., M.Si
Dosen Pembimbing 2 : Aviolla Terza Damaliana, S.Si., M.Stat

Abstrak

Indonesia menempati peringkat kelima di dunia, dengan 19,47 juta orang diperkirakan menderita diabetes pada tahun 2021, dan akan meningkat menjadi 23,32 pada tahun 2030. Penyakit diabetes dapat dicegah dengan melakukan deteksi dini dan melakukan pendekatan machine learning untuk memprediksi penyakit tersebut. Namun sering kali prediksi menjadi kurang akurat karena data yang tersedia tidak seimbang dalam distribusi kelas. Penelitian ini bertujuan untuk menangani masalah ketidakseimbangan data serta membandingkan performa model dalam memprediksi penyakit diabetes. Metode yang digunakan yaitu metode *Synthetic Minority Oversampling Technique* (SMOTE) untuk mengatasi ketidakseimbangan data, *Random Forest* untuk mengklasifikasikan penyakit diabetes dan *Confusion Matrix* untuk menghitung evaluasi kinerja model. Hasil penelitian menunjukkan bahwa model yang dilatih dengan data yang melalui proses SMOTE menunjukkan akurasi lebih tinggi yaitu 98% dibandingkan model tanpa SMOTE yang menghasilkan akurasi 93%. Penelitian menunjukkan bahwa penggunaan SMOTE dalam proses pelatihan model *Random Forest* secara signifikan mengurangi kesalahan prediksi dan meningkatkan akurasi klasifikasi.

Kata kunci: *Diabetes Mellitus, Machine Learning, SMOTE, Random Forest, Confusin Matrix*

ABSTRACT

CLASSIFICATION OF DIABETES MELLITUS USING THE SYNTHETIC MINORITY OVER-SAMPLING TECHNIQUE (SMOTE) RANDOM FOREST METHOD

Student Name / NPM : Nofa Auliyatul Maulidiyyah / 20083010029
Study Program : Sains Data, FASILKOM,UPN Veteran Jatim
Advisor 1 : Trimono, S.Si., M.Si
Advisor 2 : Aviolla Terza Damaliana, S.Si., M.Stat

Abstract

Indonesia ranks fifth globally, with an estimated 19.47 million people suffering from diabetes in 2021, projected to rise to 23.32 million by 2030. Diabetes can be prevented by carrying out early detection and using a machine learning approach to predict the disease. However, predictions often become less accurate because the available data is not balanced in the class distribution. This research aims to address the problem of data imbalance and compare the performance of models in predicting diabetes. The methods used are the Synthetic Minority Oversampling Technique (SMOTE) method to overcome data imbalances, Random Forest to classify diabetes, and Confusion Matrix to calculate model performance evaluations. The research results show that the model trained with data that went through the SMOTE process showed higher accuracy, namely 98%, compared to the model without SMOTE, which produced an accuracy of 93%. Research shows that the use of SMOTE in the Random Forest model training process significantly reduces prediction errors and increases classification accuracy.

Keywords : *Diabetes Mellitus, Machine Learning, SMOTE, Random Forest, Confusin Matrix*

KATA PENGANTAR

Segala puja dan puji syukur atas kehadiran Allah SWT, yang telah melimpahkan rahmat dan hidayahnya kepada kita semua, sehingga penulis dapat menyelesaikan Tugas Akhir yang merupakan persyaratan dalam menyelesaikan mata kuliah Skripsi pada Program Studi S1 Sains Data di Universitas Pembangunan Nasional “Veteran” Jawa Timur.

Shalawat dan salam semoga tetap tercurahkan kepada Nabi besar kita Nabi Muhammad SAW yang telah membimbing kita dari jalan kegelapan menuju jalan terang benderang yakni addinul islam ahlussunna waljamaah.

Dalam penyusunan Skripsi ini, penulis tidak ada kata yang dapat penulis ucapkan, kecuali berterimakasih yang sebesar-besarnya kepada semua pihak yang telah membantu dalam penulisan ini baik secara langsung maupun tidak langsung. Maka dalam kesempatan ini, penulis mengucapkan terimakasih kepada :

1. Dua orang paling berjasa dalam hidup penulis, ayah dan ibu. Terimakasih atas kepercayaan yang telah diberikan atas izin merantau dari kalian, serta cinta, do'a, motivasi, semangat, dan nasihat-nasihat nya. Semoga Allah SWT selalu menjaga kalian dalam kebaikan dan kemudahan.
2. Prof. Dr. Ir. Akhmad Fauzi, M.MT., IPU selaku Rektor Universitas Pembangunan Nasional “Veteran” Jawa Timur.
3. Ibu Prof. Dr. Ir. Novirina Hendrasarie, M.T.selaku Dekan Fakultas Ilmu Komputer Universitas Pembangunan Nasional “Veteran” Jawa Timur
4. Bapak Dr. Eng. Ir. Dwi Arman Prasetya, ST., MT., IPU selaku Koordinator Program Studi Sains Data Universitas Pembangunan Nasional “Veteran” Jawa Timur.
5. Bapak Trimono, S.Si., M.Si dan Ibu Aviolla Terza Damaliana, S.Si., M.Stat selaku Dosen Pembimbing 1 dan 2.
6. Bapak dan Ibu dosen Program Studi Sains Data UPN “Veteran” Jawa Timur yang sudah berkenan untuk memberikan waktu untuk berkontribusi dalam penelitian ini.
7. Inisial AR yang memaksa untuk namanya dimasukkan ke dalam buku skripsi, terimakasih atas bantuan terhadap jalannya pembuatan skripsi ini.

8. Terakhir, kepada diri saya sendiri, terimakasih sudah berjuang dan bertahan sejauh ini. Apresiasi sebesar-besarnya karena bertanggung jawab untuk menyelesaikan apa yang telah dimulai. Terimakasih untuk tidak menyerah dalam hal sesulit apapun dalam proses penyusunan skripsi ini. Tetap bersyukur dan rendah hati.

Penulis menyadari bahwa masih terdapat banyak kekurangan dalam penyusunan Skripsi ini, namun penulis berharap semoga Skripsi ini dapat memberikan kontribusi terhadap perkembangan ilmu pengetahuan, khususnya dalam bidang ilmu sains data.

Surabaya, 05 Juli 2024

Nofa Auliyatul Maulidiyyah

DAFTAR ISI

LEMBAR PENGESAHAN	ii
SURAT PERNYATAAN.....	iii
ABSTRAK.....	iv
<i>ABSTRACT</i>	v
KATA PENGANTAR	vi
DAFTAR ISI.....	viii
DAFTAR GAMBAR	x
DAFTAR TABEL.....	xi
DAFTAR LAMPIRAN.....	xii
BAB I PENDAHULUAN.....	1
1.1 Latar Belakang	1
1.2 Rumusan Masalah	5
1.3 Batasan Masalah.....	5
1.4 Tujuan Penelitian	5
1.5 Manfaat Penelitian	6
BAB II TINJAUAN PUSTAKA.....	8
2.1 Dasar Teori.....	8
2.1.1 Diabetes <i>Mellitus</i>	8
2.1.2 <i>Machine Learning</i>	12
2.1.3 <i>Imbalance Data</i>	14
2.1.4 SMOTE	16
2.1.5 <i>Random Forest</i>	19
2.1.6 <i>Confusion Matrix</i>	23
2.2 Penelitian Terdahulu	25
BAB III METODOLOGI PENELITIAN.....	30
3.1 Variabel Penelitian dan Sumber Data	30
3.2 Langkah Analisis.....	32
3.3 Diagram Alir Penelitian	33
3.4 Jadwal Penelitian.....	34
BAB IV HASIL DAN PEMBAHASAN	35
4.1 Hasil Penelitian	35

4.1.1	Dataset Diabetes <i>Mellitus</i>	35
4.1.2	<i>Preprocessing</i>	36
4.1.3	<i>Split Data</i>	41
4.1.4	SMOTE	43
4.1.5	Klasifikasi <i>Random Forest</i>	43
4.1.6	Evaluasi Model.....	44
4.2	Pembahasan.....	50
4.3	<i>Deployment</i> Aplikasi.....	52
BAB V PENUTUP.....		55
5.1	Kesimpulan	55
5.2	Saran.....	56
DAFTAR PUSTAKA		57
LAMPIRAN.....		64
BIODATA PENULIS		65

DAFTAR GAMBAR

Gambar 2. 1 Teknik SMOTE	16
Gambar 2. 2 Ilustrasi Proses dari Algoritma <i>Random Forest</i>	21
Gambar 2. 3 Confusion Matrix	23
Gambar 4. 1 Diagram Batang Diabetes <i>Mellitus</i>	35
Gambar 4. 2 Diagram Batang Rata-Rata Fitur per Diagnosa.....	39
Gambar 4. 3 Korelasi Antara Fitur dengan Diagnosa	40
Gambar 4. 4 Fitur yang Paling Berpengaruh terhadap Diagnosa.....	41
Gambar 4. 5 Uji Proporsi data tanpa SMOTE	42
Gambar 4. 6 Uji Proporsi data dengan SMOTE	42
Gambar 4. 7 Jumlah Data Setelah SMOTE	43
Gambar 4. 8 <i>Confusion Matrix Random Forest</i> tanpa SMOTE Rasio 80:20	45
Gambar 4. 9 <i>Confusion Matrix Random Forest</i> dengan SMOTE Rasio 80:20	45
Gambar 4. 10 <i>Confusion Matrix Random Forest</i> tanpa SMOTE Rasio 70:30	46
Gambar 4. 11 <i>Confusion Matrix Random Forest</i> dengan SMOTE Rasio 70:30 ..	46
Gambar 4. 12 <i>Confusion Matrix Random Forest</i> tanpa SMOTE Rasio 90:10	47
Gambar 4. 13 <i>Confusion Matrix Random Forest</i> dengan SMOTE Rasio 90:10 ..	47
Gambar 4. 14 Bobot Fitur Model <i>Random Forest</i> tanpa SMOTE.....	50
Gambar 4. 15 Bobot Fitur <i>Random Forest</i> dengan SMOTE	50
Gambar 4. 16 Hasil Performa <i>Matrix</i> Rasio 80:20 dan 90:10.....	51
Gambar 4. 17 Hasil Performa <i>Matrix</i> Rasio 70:30	51
Gambar 4. 18 Menu Home.....	52
Gambar 4. 19 <i>Click Button Predict</i>	52
Gambar 4. 20 Menu Diabetes <i>Mellitus Disease</i>	53
Gambar 4. 21 Contoh Input dalam Kumpulan Data	53
Gambar 4. 22 Contoh Data yang sudah di Prediksi	54

DAFTAR TABEL

Tabel 2. 1 Kategori ketidakseimbangan data	15
Tabel 2. 2 <i>Studi Literatur</i>	25
Tabel 3. 1 Variabel Penelitian	31
Tabel 3. 2 Jadwal Penelitian.....	34
Tabel 4. 1 Contoh Data Pasien Diabetes	36
Tabel 4. 2 Data Setelah Tahapan <i>Preprocessing</i>	37
Tabel 4. 3 Deskripsi Statistik Data.....	38
Tabel 4. 4 Pembagian Data Tanpa SMOTE dan Dengan SMOTE	41
Tabel 4. 5 Perbandingan Performa Model ' <i>random_state</i> ' 42.....	48
Tabel 4. 6 Hasil Prediksi Model Pada Sampel Data	49

DAFTAR LAMPIRAN

Lampiran 1. 1 Hasil Uji Plagiasi.....	64
Lampiran 1. 2 Data Penelitian.....	64
Lampiran 1. 3 <i>Source Code</i>	64