

TESIS

**KLASIFIKASI KELULUSAN MAHASISWA MENGGUNAKAN METODE
NAÏVE BAYES DENGAN KOMPARASI PRAPEMROSESAN *RANDOM
OVERSAMPLING* SERTA SELEKSI FITUR *INFORMATION GAIN* DAN
*FORWARD SELECTION***

(Studi Kasus: Fakultas Sains dan Teknologi UIN SUSKA Riau)



Oleh:

Dony Fahrudy

NIM: 20206052006

**STATE ISLAMIC UNIVERSITY
SUNAN KALIJAGA
YOGYAKARTA**

**PROGRAM STUDI INFORMATIKA
PROGRAM MAGISTER
FAKULTAS SAINS DAN TEKNOLOGI
UIN SUNAN KALIJAGA**

YOGYAKARTA

2022

LEMBAR PENGESAHAN



KEMENTERIAN AGAMA
UNIVERSITAS ISLAM NEGERI SUNAN KALIJAGA
FAKULTAS SAINS DAN TEKNOLOGI
Jl. Marsda Adisucipto Telp. (0274) 540971 Fax. (0274) 519739 Yogyakarta 55281

PENGESAHAN TUGAS AKHIR

Nomor : B-1764/Un.02/DST/PP.00.9/08/2022

Tugas Akhir dengan judul : KLASIFIKASI KELULUSAN MAHASISWA MENGGUNAKAN METODE NAIVE BAYES DENGAN KOMPARASI PRAPEMROSESAN RANDOM OVERSAMPLING SERTA SELEKSI FITUR INFORMATION GAIN DAN FORWARD SELECTION (STUDI KASUS; FAKULTAS SAINT DAN TEKNOLOGI UIN SUSKA RIAU)

yang dipersiapkan dan disusun oleh:

Nama : DONY FAHRUDY, S.T.
Nomor Induk Mahasiswa : 20206052006
Telah diujikan pada : Senin, 25 Juli 2022
Nilai ujian Tugas Akhir : A

dinyatakan telah diterima oleh Fakultas Sains dan Teknologi UIN Sunan Kalijaga Yogyakarta

TIM UJIAN TUGAS AKHIR



Ketua Sidang

Dr. Ir. Shofwatul 'Uyun, S.T., M.Kom.
SIGNED

Valid ID: 62fb6981c95a7



Penguji I

Ir. Maria Ulfah Siregar, S.Kom., MIT., Ph.D.
SIGNED

Valid ID: 62fd1349c8fc2



Penguji II

Dr. Ir. Bambang Sugiantoro, S.Si., M.T.
SIGNED

Valid ID: 62eba8cc7e2af



Yogyakarta, 25 Juli 2022
UIN Sunan Kalijaga
Dekan Fakultas Sains dan Teknologi

Dr. Dra. Hj. Khurul Wardati, M.Si.
SIGNED

Valid ID: 62f0e282d1ef

PERNYATAAN KEASLIAN / BEBAS PLAGIASI

SURAT PERNYATAAN KEASLIAN / BEBAS PLAGIASI

Yang bertanda tangan dibawah ini:

Nama Mahasiswa : Dony Fahrudy
NIM : 20206052006
Program Studi : Informatika (S2)
Fakultas : Sains dan Teknologi

Menyatakan dengan sesungguhnya bahwa laporan tesis saya yang berjudul: **“Klasifikasi Kelulusan Mahasiswa Menggunakan Metode Naïve Bayes dengan Komparasi Prapemrosesan Random Oversampling serta Seleksi Fitur Information Gain dan Forward Selection (Studi Kasus: Fakultas Sains dan Teknologi UIN SUSKA Riau)”** adalah hasil karya pribadi yang tidak mengandung plagiarisme dan tidak berisi materi yang dipublikasikan atau ditulis orang lain, kecuali bagian – bagian tertentu yang penulis ambil sebagai acuan dan tata cara yang diberikan secara ilmiah.

Jika terbukti pernyataan ini tidak benar, maka penulis siap mempertanggung jawabkan sesuai hukum yang berlaku.

Yogyakarta, 14 Juli 2022

Yang menyatakan,



METERAI
TEMPER
14BAJX978550212

Dony Fahrudy

NIM. 20206052006

STATE ISLAMIC UNIVERSITY
SUNAN KALIJAGA
YOGYAKARTA

HALAMAN PERSETUJUAN



Universitas Islam Negeri Sunan Kalijaga

SURAT PERSETUJUAN TUGAS AKHIR

Hal : Persetujuan Tugas Akhir
Lamp : -

Kepada Yth.,
Dekan Fakultas Sains dan Teknologi
UIN Sunan Kalijaga Yogyakarta

Assalamualaikum wr. wb.

Setelah membaca, meneliti, memberikan petunjuk dan mengoreksi serta mengadakan perbaikan seperlunya, maka saya selaku pembimbing berpendapat bahwa naskah tesis Saudara:

Nama : Dony Fahrudy
NIM : 20206052006
Judul Tesis : Klasifikasi Kelulusan Mahasiswa Menggunakan Metode *Naïve Bayes* dengan Komparasi Prapemrosesan Random Oversampling serta Seleksi Fitur Information Gain dan Forward Selection (Studi Kasus: Fakultas Sains dan Teknologi UIN SUSKA Riau)

Saya berpendapat bahwa tesis tersebut sudah dapat diajukan kepada Program Studi Magister Informatika Fakultas Sains dan Teknologi UIN Sunan Kalijaga Yogyakarta untuk diujikan dalam rangka memperoleh gelar Magister Informatika.

Walaikumsalam wr. wb.

Yogyakarta, 14 Juli 2022

Pembimbing,

Dr. Ir. Shofwatul 'Uyun, S.T., M.Kom.

NIP. 19820511 200604 2 002

ABSTRAK

Ketidakseimbangan jumlah data pada setiap kelas serta memiliki dimensi atribut tinggi pada *datasets*, seringkali menjadi masalah dalam proses klasifikasi yang dapat mempengaruhi kinerja algoritma dalam proses komputasi, karena terdapat jumlah data pada setiap kelas yang tidak seimbang dan atribut yang tidak relevan harus diproses, sehingga perlunya dilakukan teknik mengatasi ketidakseimbangan kelas untuk menyeimbangkan jumlah data pada setiap kelas dan seleksi fitur untuk mengurangi kompleksitas data dan fitur-fitur yang tidak relevan. Oleh karena itu, penelitian ini menggunakan teknik *Random Oversampling* (ROs) untuk mengatasi ketidakseimbangan kelas, serta dua metode seleksi fitur dengan algoritma *Information Gain* dan *Forward Selection* yang dikomparasikan untuk menentukan metode seleksi fitur yang lebih unggul, lebih efektif, dan lebih cocok digunakan. Hasil dari seleksi fitur digunakan untuk melakukan klasifikasi kelulusan mahasiswa dengan membangun model klasifikasi menggunakan algoritma *Naïve Bayes*. Hasil penelitian menunjukkan peningkatan rata-rata akurasi dari metode *Naïve Bayes* dengan tanpa prapemrosesan ROs dan seleksi fitur, penggunaan ROs, penggunaan *Information Gain* dengan 3 fitur terpilih serta *Forward Selection* dengan 2 fitur terpilih secara berurutan adalah 81.83%; 83.84%; 86.03% dan 86.42%, sehingga terjadi peningkatan akurasi sebesar 4.2% dari tanpa prapemrosesan ke *Information Gain* dan 4.59% dari tanpa prapemrosesan ke *Forward Selection*. Oleh karena itu, metode seleksi fitur terbaik adalah *Forward Selection* dengan 2 fitur terpilih (Ip Semester 8 dan IPK), penggunaan ROs dan kedua seleksi fitur terbukti meningkatkan kinerja metode *Naïve Bayes*.

Kata Kunci: *Forward Selection, Information Gain, Kelulusan Mahasiswa, Naïve Bayes, ROs*

STATE ISLAMIC UNIVERSITY
SUNAN KALIJAGA
YOGYAKARTA

ABSTRACT

The imbalance in the amount of data in each class and having high attribute dimensions in datasets, is often a problem in the classification process that can affect the algorithm's performance in the computational process, because there is an unbalanced amount of data in each class and irrelevant attributes must be processed, so it is necessary to do this. class imbalance techniques to balance the amount of data in each class and feature selection to reduce data complexity and irrelevant features. Therefore, this study uses the Random Oversampling (ROs) technique to overcome class imbalances, as well as two feature selection methods with Information Gain and Forward Selection algorithms which are compared to determine which feature selection method is superior, more effective, and more suitable. used. The results of feature selection are used to classify student graduation by building a classification model using the Naïve Bayes algorithm. The results showed an increase in the average accuracy of the Naïve Bayes method with no ROs preprocessing and feature selection, the use of ROs, the use of Information Gain with 3 selected features and Forward Selection with 2 selected features sequentially was 81.83%; 83.84%; 86.03% and 86.42%, so that there is an increase in accuracy of 4.2% from no pre-processing to Information Gain and 4.59% from no pre-processing to Forward Selection. Therefore, the best feature selection method is Forward Selection with 2 selected features (Ip Semester 8 and GPA), the use of ROs and both feature selections are proven to improve the performance of the Naïve Bayes method.

Keywords : *Forward Selection, Information Gain, Student Graduation, Naïve Bayes, ROs*

STATE ISLAMIC UNIVERSITY
SUNAN KALIJAGA
YOGYAKARTA

MOTTO

“Sesungguhnya Allah tidak akan mengubah keadaan suatu kaum, sebelum mereka mengubah keadaan diri mereka sendiri.”

(QS Ar Rad 11)

“Rahasia kesuksesan adalah mengetahui yang orang lain belum ketahui.”

(Aristotle Onassis)

“Pengetahuan yang baik adalah yang memberikan manfaat, bukan hanya diingat.”

(Imam Syafi'i)

“Hari esok harus lebih baik dari hari ini.”

(Dony Fahrudy)

STATE ISLAMIC UNIVERSITY
SUNAN KALIJAGA
YOGYAKARTA

KATA PENGANTAR

Puji dan Syukur Kehadirat Allah SWT atas segala rahmat dan karunia-Nya, akhirnya penulis dapat menyelesaikan penyusunan tesis yang berjudul “Klasifikasi Kelulusan Mahasiswa Menggunakan Metode *Naïve Bayes* dengan Komparasi Prapemrosesan *Random Oversampling* serta Seleksi Fitur *Information Gain* dan *Forward Selection* (Studi Kasus : Fakultas Sains dan Teknologi UIN SUSKA Riau)”.

Tesis ini ditulis dalam rangka memenuhi persyaratan untuk memperoleh gelar Magister di UIN Sunan Kalijaga Yogyakarta. Penulis menyadari bahwa tesis dapat diselesaikan berkat dukungan dan bantuan dari berbagai pihak, oleh karena itu penulis berterima kasih kepada semua pihak yang secara langsung maupun tidak langsung memberikan kontribusi dalam menyelesaikan Tesis ini. Selanjutnya ucapan terima kasih penulis sampaikan kepada :

1. Bapak Dr. Phil. Sahiron, M.A. sebagai Plt. Rektor UIN Sunan Kalijaga Yogyakarta.
2. Ibu Dr. Khurul Wardati, M,Si, selaku dekan fakultas SAINTEK UIN Sunan Kalijaga.
3. Bapak Dr. Bambang Sugiantoro, S.Si., M.T. sebagai Kepala Program Studi Magister Informatika Fakultas Sains dan Teknologi UIN Sunan Kalijaga Yogyakarta sekaligus penguji kedua, yang telah memberikan izin dan kemudahan, sehingga penulis dapat menyelesaikan studi di Program Studi Program Magister Informatika Fakultas Sains dan Teknologi UIN Sunan Kalijaga Yogyakarta.
4. Ibu Dr. Ir. Shofwatul 'Uyun, S.T., M.Kom. sebagai pembimbing akademik sekaligus pembimbing tugas akhir, dengan kesibukannya masih tetap meluangkan waktunya untuk memberikan bimbingan, petunjuk, dan mendorong semangat penulis untuk menyelesaikan penulisan Tesis ini.
5. Ibu Ir. Maria Ulfah Siregar, S.Kom., MIT., Ph.D, M.Kom. sebagai penguji pertama, dengan kesibukannya masih tetap meluangkan waktunya untuk memberikan petunjuk, saran, masukan dan mendorong semangat penulis untuk

menyelesaikan penulisan Tesis ini.

6. Seluruh dosen dan staf administrasi pada Program Studi Informatika Program Magister Fakultas Sains dan Teknologi UIN Sunan Kalijaga Yogyakarta, yang secara langsung atau tidak langsung telah memberi bantuan kepada penulis dalam menyelesaikan penulisan tesis.
7. Segenap staf UPT. Pusat Teknologi Informasi dan Pangkalan Data UIN Sunan Kalijaga Yogyakarta yang secara langsung atau tidak langsung telah memberi bantuan kepada penulis dalam menyelesaikan penulisan tesis.
8. Orang Tua tercinta yang mendidik dengan penuh rasa kasih sayang dan senantiasa memberi semangat dan dorongan kepada penulis.
9. Keluarga tercinta yang telah memberikan do'a dan dukungan kepada penulis.
10. Seluruh rekan-rekan mahasiswa di Magister Informatika yang telah saling mendukung untuk melalui perjuangan bersama-sama.

Akhirnya penulis berharap semoga Tesis ini dapat bermanfaat dan permintaan maaf yang tulus jika seandainya dalam penulisan ini terdapat kekurangan dan kekeliruan, penulis juga menerima kritik dan saran yang bersifat membangun demi menyempurnakan penulisan Tesis ini.

Yogyakarta, 25 Juli 2022

Penulis

Dony Fahrudy

NIM. 20206052006

DAFTAR ISI

LEMBAR PENGESAHAN	ii
PERNYATAAN KEASLIAN / BEBAS PLAGIASI.....	iii
HALAMAN PERSETUJUAN	iv
ABSTRAK	v
ABSTRACT	vi
MOTTO	vii
KATA PENGANTAR.....	viii
DAFTAR ISI.....	x
DAFTAR TABEL	xv
DAFTAR GAMBAR.....	xx
DAFTAR LAMPIRAN.....	xxiii
BAB I PENDAHULUAN.....	1
A. Latar Belakang	1
B. Rumusan Masalah	6
C. Batasan Masalah.....	6
D. Tujuan Penelitian	7
E. Manfaat Penelitian	7
F. Sistematika Penulisan.....	7
BAB II KAJIAN PUSTAKA DAN LANDASAN TEORI.....	9
A. Kajian Pustaka.....	9
B. Landasan Teori	14
1. Kelulusan Mahasiswa	14
2. KDD dan <i>Data Mining</i>	15

3. Pengelompokan Data Mining	16
4. Tahapan Data Mining dalam KDD	16
a. Pembersihan Data (<i>Data Cleaning</i>).....	17
b. <i>Resampling</i>	17
c. Seleksi Data (<i>Data Selection</i>).....	18
1). Seleksi Fitur (<i>Feature Selection</i>).....	18
2). Algoritma <i>Information Gain</i>	19
3). Algoritma <i>Forward Selection</i>	20
d. Transformasi Data (<i>Data Transformation</i>).....	22
1). Konversi Data (<i>Data Convert</i>).....	22
2). Normalisasi Data (<i>Data Normalization</i>).....	22
e. <i>Data Mining</i>	23
1). Klasifikasi.....	23
2). Algoritma <i>Naïve Bayes</i>	24
f. Evaluasi Pola (<i>Pattern Evaluation</i>).....	26
5. K-Fold Cross Validation	27
6. Pengujian Akurasi	28
BAB III METODE PENELITIAN	30
A. Perumusan Masalah	31
B. Pengumpulan Data	31
1. Studi Pustaka.....	31
2. Pengambilan Data.....	31
C. Analisa Permasalahan	31
1. Analisa Data.....	32
a. Data Sekunder.....	32

b.	Atribut Data	32
2.	Analisa Proses KDD	33
a.	Pembersihan Data	33
b.	<i>Resampling</i>	34
c.	Seleksi Fitur	34
1).	<i>Information gain</i>	35
2).	<i>Forward Selection</i>	36
d.	Transformasi Data.....	37
1).	Konversi Data	37
2).	Normalisasi Data.....	37
e.	Pembagian data	37
f.	<i>Data Mining</i>	38
D.	Perancangan	39
1.	Perancangan Program	39
2.	Perancangan Basis Data (<i>Database</i>).....	40
E.	Implementasi	40
F.	Evaluasi Kinerja	40
G.	Kesimpulan dan Saran	41
BAB IV	HASIL DAN PEMBAHASAN	42
A.	Analisa Permasalahan	43
1.	Analisa Data.....	43
a.	Data Sekunder	43
b.	Atribut Data	43
2.	Analisa Proses KDD	43
a.	Pembersihan Data	44

b.	<i>Resampling</i>	45
c.	Seleksi Fitur	48
1).	<i>Information Gain</i>	48
2).	<i>Forward Selection</i>	56
d.	Transformasi Data.....	68
1).	Konversi Data	68
2).	Normalisasi Data.....	69
e.	Pembagian Data	71
f.	<i>Data Mining</i>	75
1).	<i>Naïve Bayes</i> berbasis <i>Information Gain</i>	76
2).	<i>Naïve Bayes</i> berbasis <i>Forward Selection</i>	83
B.	Perancangan	89
1.	Perancangan Program	89
2.	Perancangan <i>Database</i> (Basis Data).....	89
C.	Implementasi	96
1.	Batasan Implementasi	96
2.	Implementasi Program.....	97
3.	Output Program.....	97
a.	Output Pemrosesan <i>Datasets</i>	97
b.	Output Pemrosesan Pembersihan Data	98
c.	Output Pemrosesan <i>Random Oversampling</i>	98
d.	Output Pemrosesan Seleksi Fitur	100
e.	Output Pemrosesan Transformasi Data	102
f.	Output Pemrosesan Pembagian Data	103
g.	Output Pemrosesan <i>Data Mining</i>	105

h. Output Pemrosesan Evaluasi	109
D. Evaluasi Kinerja	114
1. Rencana Pengujian.....	114
2. Data Pengujian	114
3. <i>K-Fold Cross Validation</i>	115
a. Pengujian Kinerja <i>Naïve Bayes</i> tanpa Prapemrosesan <i>Random Oversampling</i> dan Seleksi Fitur.....	115
b. Pengujian Kinerja <i>Naïve Bayes</i> dengan Prapemrosesan <i>Random Oversampling</i>	117
c. Pengujian Kinerja <i>Naïve Bayes</i> dengan Prapemrosesan <i>Information Gain</i>	118
d. Pengujian Kinerja <i>Naïve Bayes</i> dengan Prapemrosesan <i>Forward Selection</i>	120
e. Perbandingan Kinerja <i>Naïve Bayes</i>	122
4. Grafik Evaluasi Kinerja	122
5. Analisis Hasil Pengujian.....	123
BAB V PENUTUP.....	125
A. Kesimpulan.....	125
B. Saran	126
DAFTAR PUSTAKA.....	xxiv
LAMPIRAN.....	xxxii
DAFTAR RIWAYAT HIDUP	cxxxvi

DAFTAR TABEL

Tabel	Halaman
Tabel 2.1 Skema <i>10-Fold Cross Validation</i>	27
Tabel 2.2 <i>Confusion Matrix</i> untuk Klasifikasi <i>Class Biner</i>	28
Tabel 3.1 Atribut Data.....	32
Tabel 4.1 <i>Datasets</i> Kelulusan Mahasiswa	44
Tabel 4.2 Data Hasil <i>Cleaning</i>	44
Tabel 4.3 Data Hasil Duplikasi Pertama Kelas Minoritas	46
Tabel 4.4 Data Hasil Duplikasi Kedua Kelas Minoritas	46
Tabel 4.5 Data Hasil <i>Random Oversampling</i>	47
Tabel 4.6 Jumlah Data Tiap Kelas	47
Tabel 4.7 Jumlah Data Tiap Kelas – <i>Random Oversampling</i>	47
Tabel 4.8 Jumlah Data Setiap Kelas Pada Jenis Kelamin	48
Tabel 4.9 Jumlah Data Setiap Kelas Pada Tempat Lahir.....	48
Tabel 4.10 Jumlah Data Setiap Kelas Pada Tempat Tinggal	48
Tabel 4.11 Jumlah Data Setiap Kelas Pada Alat Transportasi	49
Tabel 4.12 Jumlah Data Setiap Kelas Pada Tahun Angkatan	49
Tabel 4.13 Nilai <i>Entropy</i> Jenis Kelamin.....	50
Tabel 4.14 Nilai <i>Entropy</i> Tempat Lahir.....	50
Tabel 4.15 Nilai <i>Entropy</i> Tempat Tinggal.....	50
Tabel 4.16 Nilai <i>Entropy</i> Alat Transportasi.....	51
Tabel 4.17 Nilai <i>Entropy</i> Tahun Angkatan.....	51
Tabel 4.18 Nilai <i>Information Gain</i> Jenis Kelamin.....	52
Tabel 4.19 Nilai <i>Information Gain</i> Tempat Lahir	52

Tabel 4.20 Nilai <i>Information Gain</i> Tempat Tinggal.....	53
Tabel 4.21 Nilai <i>Information Gain</i> Alat Transportasi.....	53
Tabel 4.22 Nilai <i>Information Gain</i> Tahun Angkatan.....	54
Tabel 4.23 Nilai <i>Information Gain</i>	54
Tabel 4.24 Atribut <i>Ranking</i>	55
Tabel 4.25 Atribut Terpilih.....	56
Tabel 4.26 Parameter Hasil Seleksi <i>Information Gain</i>	56
Tabel 4.27 Data Hasil Seleksi <i>Information Gain</i>	56
Tabel 4.28 Data Uji Pertama.....	58
Tabel 4.29 Hasil Pengujian Data Uji Variabel Jenis Kelamin.....	58
Tabel 4.30 Akurasi Pengujian Variabel.....	59
Tabel 4.31 Variabel Optimal.....	60
Tabel 4.32 Data Uji Pertama.....	61
Tabel 4.33 Hasil Pengujian Data Uji Variabel Ip Semester 8 dan Jenis Kelamin	61
Tabel 4.34 Akurasi Pengujian Kombinasi Variabel.....	62
Tabel 4.35 Variabel Optimal.....	63
Tabel 4.36 Data Uji Pertama.....	64
Tabel 4.37 Hasil Pengujian Data Uji Variabel Ip Semester 8, Ipk, dan Jenis Kelamin	65
Tabel 4.38 Akurasi Pengujian Kombinasi Variabel.....	66
Tabel 4.41 Parameter Hasil Seleksi <i>Forward Selection</i>	68
Tabel 4.42 Data Hasil Seleksi <i>Forward Selection</i>	68
Tabel 4.43 Data Hasil Konversi – <i>Information Gain</i>	68
Tabel 4.44 Data Hasil Konversi – <i>Forward Selection</i>	69

Tabel 4.45 Data Hasil Normalisasi – <i>Information Gain</i>	70
Tabel 4.46 Data Hasil Normalisasi – <i>Forward Selection</i>	71
Tabel 4.47 Pembagian Data <i>Training</i> dan Data <i>Testing</i>	72
Tabel 4.48 Data Pelatihan Pembangunan Model <i>Fold-1 - Information Gain</i>	73
Tabel 4.49 Data Pengujian Klasifikasi <i>Fold-1 - Information Gain</i>	73
Tabel 4.50 Data Pelatihan Pembangunan Model <i>Fold-1 - Forward Selection</i>	74
Tabel 4.51 Data Pengujian Klasifikasi <i>Fold-1 - Forward Selection</i>	75
Tabel 4.52 Nilai Mean (Rata-rata) – <i>Information Gain</i>	77
Tabel 4.53 Nilai Standar Deviasi – <i>Information Gain</i>	78
Tabel 4.54 Hasil Perhitungan Gaussian Atribut – <i>Information Gain</i>	81
Tabel 4.56 Hasil Perhitungan Probabilitas Kelas – <i>Information Gain</i>	82
Tabel 4.57 Hasil Perhitungan Gaussian Akhir – <i>Information Gain</i>	82
Tabel 4.58 Hasil Prediksi – <i>Information Gain</i>	82
Tabel 4.59 Nilai Mean (Rata-rata) – <i>Forward Selection</i>	84
Tabel 4.60 Nilai Standar Deviasi – <i>Forward Selection</i>	85
Tabel 4.61 Hasil Perhitungan Gaussian Atribut – <i>Forward Selection</i>	87
Tabel 4.62 Hasil Perhitungan Probabilitas Kelas – <i>Forward Selection</i>	87
Tabel 4.63 Hasil Perhitungan Gaussian Akhir – <i>Forward Selection</i>	88
Tabel 4.64 Hasil Prediksi – <i>Forward Selection</i>	88
Tabel 4.65 Akurasi	89
Tabel 4.66 <i>Confusion Matrix</i>	89
Tabel 4.67 <i>Datasets</i>	90
Tabel 4.68 <i>Datasets Cleaning</i>	90
Tabel 4.69 <i>Datasets Convert</i>	91
Tabel 4.70 <i>Datasets Imbalanced</i>	91

Tabel 4.71 <i>Datasets Normalize</i>	92
Tabel 4.72 <i>Datasets Selection - Information Gain</i>	92
Tabel 4.73 <i>Datasets Selection - Forward Selection</i>	93
Tabel 4.74 <i>Data Hasil K-Fold – Information Gain</i>	93
Tabel 4.75 <i>Data Hasil K-Fold Seluruh Atribut</i>	93
Tabel 4.76 <i>Data Hasil K-Fold – Forward Selection</i>	93
Tabel 4.77 <i>Gaussian Testing Seluruh Atribut</i>	94
Tabel 4.78 <i>Gaussian Testing – Information Gain</i>	95
Tabel 4.79 <i>Gaussian Testing – Forward Selection</i>	95
Tabel 4.80 <i>Hitung Forward Selection</i>	95
Tabel 4.81 <i>Hitung Information Gain</i>	95
Tabel 4.82 <i>Hitung Naïve Bayes</i>	95
Tabel 4.83 <i>Prediksi Naïve Bayes</i>	96
Tabel 4.84 <i>Hasil Pengujian (k = 10) Fold-1</i>	115
Tabel 4.85 <i>Pengujian Pada Confusion Matrix Fold-1</i>	116
Tabel 4.86 <i>Rata-Rata Akurasi Naïve Bayes Nilai k = 10</i>	116
Tabel 4.87 <i>Hasil Pengujian (k = 10) Fold-1</i>	117
Tabel 4.88 <i>Pengujian Pada Confusion Matrix Fold-1</i>	117
Tabel 4.89 <i>Rata-Rata Akurasi Naïve Bayes - ROs Nilai k = 10</i>	118
Tabel 4.90 <i>Hasil Pengujian (k = 10) Fold-1</i>	118
Tabel 4.91 <i>Pengujian Pada Confusion Matrix Fold-1</i>	119
Tabel 4.92 <i>Rata-Rata Akurasi Naïve Bayes - Information Gain Nilai k = 10</i>	119
Tabel 4.93 <i>Hasil Pengujian (k = 10) Fold-1</i>	120
Tabel 4.94 <i>Pengujian Pada Confusion Matrix Fold-1</i>	121

Tabel 4.95 Rata-Rata Akurasi *Naïve Bayes* - *Forward Selection* Nilai $k = 10$

.....121



DAFTAR GAMBAR

Gambar	Halaman
Gambar 3.1 Metodologi Penelitian.....	30
Gambar 3.2 Tahapan Teknik <i>Random Oversampling</i>	34
Gambar 3.3 Tahapan Algoritma <i>Information Gain</i>	35
Gambar 3.4 Tahapan Algoritma <i>Forward Selection - Naïve Bayes</i>	36
Gambar 3.5 Tahapan <i>Training</i> Algoritma <i>Naïve Bayes</i>	39
Gambar 3.6 Tahapan <i>Testing</i> Algoritma <i>Naïve Bayes</i>	39
Gambar 4.1 Tahapan Penelitian.....	42
Gambar 4.2 <i>Datasets</i>	97
Gambar 4.3 Hasil Pembersihan Data.....	98
Gambar 4.4 Visualisasi Histogram Jumlah Data Sebelum Teknik ROs	98
Gambar 4.5 Jumlah Data Sebelum dan Setelah Teknik ROs.....	99
Gambar 4.6 Hasil Data Sebelum Teknik ROs.....	99
Gambar 4.7 Hasil Data Setelah Teknik ROs.....	99
Gambar 4.8 Data Perhitungan <i>Information Gain</i>	100
Gambar 4.9 Proses <i>Ranking Information Gain</i>	100
Gambar 4.10 <i>Datasets</i> Hasil Seleksi <i>Information Gain</i>	101
Gambar 4.11 Hasil Penambahan Atribut Optimal <i>Forward Selection</i>	101
Gambar 4.12 Hasil Atribut Optimal dan Kinerja <i>Forward Selection</i>	101
Gambar 4.13 <i>Datasets</i> Hasil Seleksi <i>Forward Selection</i>	102
Gambar 4.14 Hasil Konversi Data	102
Gambar 4.15 Hasil Normalisasi Data.....	103
Gambar 4.16 Data Pembagian <i>Fold - Information Gain</i>	103

Gambar 4.17 Data Pembagian <i>Fold - Forward Selection</i>	104
Gambar 4.18 Data Pembagian <i>Fold – Without Ros</i>	104
Gambar 4.19 Data Pembagian <i>Fold –ROs</i>	104
Gambar 4.20 Hasil Perhitungan <i>Learning NB - Information Gain</i>	105
Gambar 4.22 Hasil Perhitungan Prediksi <i>NB - Information Gain</i>	105
Gambar 4.23 Hasil Prediksi <i>NB - Information Gain</i>	106
Gambar 4.24 Hasil Perhitungan <i>Learning NB – Forward Selection</i>	106
Gambar 4.25 Hasil Perhitungan Prediksi <i>NB - Forward Selection</i>	107
Gambar 4.26 Hasil Prediksi <i>NB - Forward Selection</i>	107
Gambar 4.27 Hasil Perhitungan <i>Learning NB – Without ROs</i>	108
Gambar 4.28 Hasil Perhitungan Prediksi <i>NB - Without ROs</i>	108
Gambar 4.29 Hasil Prediksi <i>NB - Without ROs</i>	108
Gambar 4.30 Hasil Perhitungan <i>Learning NB – ROs</i>	109
Gambar 4.31 Hasil Perhitungan Prediksi <i>NB - ROs</i>	109
Gambar 4.32 Hasil Prediksi <i>NB - ROs</i>	109
Gambar 4.33 Pengujian <i>K-Fold Cross Validation – Information Gain</i>	110
Gambar 4.34 Hasil <i>Confusion Matrix Fold – Information Gain</i>	110
Gambar 4.35 Rata-rata Pengujian <i>K-Fold Cross Validation – Information Gain</i>	110
Gambar 4.36 Rata-rata Akurasi <i>K-Fold Cross Validation – Information Gain</i>	110
Gambar 4.37 Pengujian <i>K-Fold Cross Validation – Forward Selection</i>	111
Gambar 4.38 Hasil <i>Confusion Matrix Fold – Forward Selection</i>	111
Gambar 4.39 Rata-rata Pengujian <i>K-Fold Cross Validation – Forward Selection</i> 111	

Gambar 4.40 Rata-rata Akurasi K-Fold Cross Validation – Forward Selection	112
Gambar 4.41 Pengujian K-Fold Cross Validation – Without ROs	112
Gambar 4.42 Hasil Confusion Matrix Fold – Without ROs	112
Gambar 4.43 Rata-rata Pengujian K-Fold Cross Validation – Without ROs	112
Gambar 4.44 Rata-rata Akurasi K-Fold Cross Validation – Without ROs ...	113
Gambar 4.45 Pengujian K-Fold Cross Validation – ROs	113
Gambar 4.46 Hasil Confusion Matrix Fold – ROs	113
Gambar 4.47 Rata-rata Pengujian K-Fold Cross Validation – ROs	114
Gambar 4.48 Rata-rata Akurasi K-Fold Cross Validation – ROs	114
Gambar 4.49 Grafik Perbandingan Kinerja Naïve Bayes	123

DAFTAR LAMPIRAN

Lampiran	Halaman
Lampiran 1. Perancangan Program	xxxii
Lampiran 2. Implementasi Program.....	lxxi



BAB I

PENDAHULUAN

A. Latar Belakang

Kelulusan mahasiswa merupakan faktor keberhasilan yang menjadi salah satu penilaian akreditasi pada perguruan tinggi. Apabila mahasiswa lulus tepat waktu dalam menyelesaikan studinya, maka dapat mendukung penilaian akreditasi tersebut (ACCJC/WASC, 2012). Salah satu kriteria keberhasilan mahasiswa dalam memperoleh gelar sarjana pada perguruan tinggi adalah lulus tepat waktu, sehingga mahasiswa dapat lulus tepat waktu apabila telah menyelesaikan studinya selama kurang dari atau sama dengan empat tahun. Namun dalam pelaksanaannya mahasiswa tidak selalu dapat menyelesaikan pendidikan sarjana selama kurang dari atau sama dengan empat tahun (Nuffic, 2017).

Berdasarkan *datasets* yang diperoleh pada perguruan tinggi, mahasiswa tidak selalu dapat menyelesaikan pendidikan sarjana secara tepat waktu, sehingga diperlukan peran perguruan tinggi untuk mengantisipasi hal tersebut. Otoritas pendidikan, administrasi akademik, dan orang tua prihatin dengan penurunan tingkat kelulusan mahasiswa di lembaga pendidikan tinggi. Pihak perguruan tinggi mencoba untuk meningkatkan kelulusan mahasiswa. Kemampuan klasifikasi untuk mengantisipasi kelulusan mahasiswa secara akurat sangat penting karena memungkinkan lembaga pendidikan untuk menetapkan program strategis untuk membantu dan meningkatkan kinerja mahasiswa menuju kelulusan tepat waktu (Bassi et al., 2019). Oleh karena itu, perlunya peningkatan kualitas akademik oleh institusi dan mengoptimalkan sumber daya agar dapat membantu mahasiswa menyelesaikan studi mereka secara tepat waktu (Lei & Li, 2015). Oleh karena itu, diperlukan sebuah model untuk mengklasifikasi kelulusan mahasiswa yang merupakan langkah terpenting untuk mengetahui kualitas mahasiswa (Bassi et al., 2019). Ketersediaan data latih dan uji untuk setiap kelas merupakan salah satu kriteria yang menentukan keberhasilan suatu model. Ada beberapa isu yang seringkali ditemukan dalam komputasi, yaitu: jumlah data antar kelas yang tidak

seimbang dan data dengan dimensi atribut tinggi.

Masalah ketidakseimbangan kelas dapat terjadi ketika *instance* dari satu kelas melebihi jumlah *instance* (kelas mayoritas) dari kelas lain (kelas minoritas), sehingga menyebabkan kesalahan klasifikasi pada kelas minoritas yang berdampak bias pada hasil klasifikasi kelas mayoritas. Secara umum, ada tiga teknik *resampling* yaitu *random undersampling*, *random oversampling* dan *hybrid methods* (Q. Wang et al., 2017). Untuk mengatasi masalah ketidakseimbangan pada kelas minoritas dapat menggunakan teknik *random oversampling* dengan mereplikasi *instance* pada kelas minoritas secara acak (Naganjaneyulu & Kuppa, 2013).

Pada umumnya algoritma klasifikasi menggunakan semua fitur pada data untuk membangun sebuah model, namun tidak semua fitur tersebut relevan untuk proses klasifikasi. Jika hal tersebut terdapat pada data yang memiliki ukuran dan dimensi yang sangat besar, maka dapat membuat kinerja algoritma menjadi tidak efektif akibat banyak fitur yang tidak relevan harus diproses. Salah satu solusi yang dapat mengatasi masalah tersebut adalah dengan menggunakan seleksi fitur. Seleksi fitur merupakan salah satu tahap dalam *preprocessing* pada klasifikasi yang dilakukan dengan cara memilih atribut yang relevan terhadap data yang mempengaruhi hasil klasifikasi. Seleksi fitur juga digunakan untuk meningkatkan efektifitas dan efisiensi kinerja dari algoritma klasifikasi, serta dapat mengurangi dimensi data dan fitur yang tidak relevan (Vanaja & Ramesh Kumar, 2014). Tujuan dari seleksi fitur adalah untuk mengurangi tingkat kompleksitas dari sebuah algoritma klasifikasi, meningkatkan akurasi dari algoritma klasifikasi, dan dapat mengetahui atribut yang mempengaruhi kinerja algoritma (C.Deisy et al., 2007). Secara umum, dasar ide algoritma *feature selection* yaitu mencari semua kemungkinan kombinasi dari atribut dalam data untuk menemukan subset dari feature yang terbaik untuk klasifikasi (Ting et al., 2011).

Terdapat tiga pendekatan dalam seleksi fitur, yaitu *filter*, *wrapper* dan *embedded* (Vanaja & Ramesh Kumar, 2014). Metode *filter* merupakan metode seleksi fitur yang memilih variabel berdasarkan peringkat dan menghilangkan variabel yang berada di bawah ambang batas (Chandrashekar & Sahin, 2014).

Metode ini menggunakan kriteria teknik *ranking* yang sederhana dan menghasilkan fitur yang relevan. Seleksi fitur pada pendekatan *filter* menggunakan teknik evaluasi yang terpisah dengan algoritma pembelajaran (Visalakshi & Radha, 2015). Kelebihan dari metode *filter* adalah lebih cepat, terukur untuk kumpulan data yang besar, dan independen, selain itu teknik *filter* mudah diterapkan (Mwadulo, 2016). Kekurangan dari metode *filter* adalah dapat mengabaikan ketergantungan fitur dan kurangnya interaksi dengan *classifier* karena dilakukan secara terpisah (Ang et al., 2016).

Subset fitur yang dihasilkan metode *wrapper* menggunakan sekumpulan kombinasi fitur dari teknik pencarian kemudian melatih model prediktif pada himpunan bagian dari fitur dan mengevaluasi *subset* tersebut menggunakan algoritma pembelajaran terawasi untuk mengetahui akurasi kinerja. Setiap kombinasi fitur dibandingkan satu dengan yang lainnya dan digunakan model algoritma untuk mengevaluasi kumpulan kombinasi yang ada (Gnana et al., 2016). Metode *wrapper* mengatasi masalah dalam metode *filter*. Metode *wrapper* dapat berinteraksi dengan model pembelajaran *classifier* dan memperhatikan keterkaitan antar atribut (Mwadulo, 2016). Metode *wrapper* dioptimalkan bersamaan dengan menggunakan algoritma pembelajaran klasifikasi. Oleh karena itu, secara umum, metode ini menghasilkan akurasi kinerja yang lebih baik daripada metode *filter*. Namun, metode ini secara komputasi mahal dan lebih kompleks daripada metode *filter* dan cenderung memerlukan waktu pemrosesan lebih lama (Jain & Singh, 2018).

Metode *embedded* mencari subset fitur yang optimal untuk algoritma klasifikasi tertentu ketika membangun sebuah *classifier* (Khaira & Dhanalakshmi, 2019). Metode ini membuat keputusan tergantung pada pengklasifikasi, sehingga pilihan fitur dipengaruhi oleh hipotesis pengklasifikasi dan tidak cocok dengan beberapa pengklasifikasi lainnya (Jain & Singh, 2018). Oleh karena itu, penelitian ini mengkomparasikan metode seleksi fitur dengan pendekatan *filter* dan *wrapper* dengan algoritma *information gain* dan *forward selection*.

Information gain adalah pendekatan pembelajaran mesin yang umum digunakan sebagai kriteria pemilihan atribut. *Information gain* secara khusus

memberikan *ranking* pada setiap *feature* (atribut) yang ada dan mengesampingkan *feature* (atribut) yang tidak memenuhi standar tertentu yang disusun berdasarkan nilai tertinggi ke terendah (Khemphila & Boonjing, 2011). Fitur dengan nilai *information gain* yang tinggi lebih bagus daripada fitur lainnya, yang mencerminkan bahwa semakin banyak informasi atribut yang berhubungan dengan kelas. *Information gain* tidak menghapus fitur yang berlebihan (Wah et al., 2018).

Forward selection merupakan algoritma pencarian paling sederhana. *Forward selection* adalah salah satu teknik untuk mereduksi dimensi *dataset* dengan menghilangkan atribut-atribut yang tidak relevan atau redundan (Panthong & Srivihok, 2015). Metode *forward selection* merupakan pemodelan yang diawali dengan nol *variable* (*empty model*) atau tidak ada *variable* dalam model, selanjutnya *variable* dimasukkan satu persatu. Untuk setiap *variable* yang ditambahkan, kinerja akan dievaluasi. Hanya atribut yang memiliki kinerja tertinggi ditambahkan ke seleksi untuk fungsi objek hingga kriteria tertentu dipenuhi (Kalyani & Karnan, 2010). Hasil dari seleksi fitur digunakan untuk membangun model klasifikasi menggunakan metode *data mining* dengan algoritma *naïve bayes* yang bertujuan untuk mengetahui teknik pemilihan fitur terbaik, lebih efisien dan lebih tepat untuk digunakan.

Menurut (Vanaja & Ramesh Kumar, 2014) *Data mining* merupakan kegiatan yang dilakukan berupa mengumpulkan data berukuran besar kemudian data tersebut diekstraksi menjadi sebuah pengetahuan (*knowledge*) yang dapat digunakan. Dalam *data mining*, terdapat beberapa metode yang dapat digunakan salah satunya adalah klasifikasi. Klasifikasi merupakan metode *data mining* yang digunakan untuk mendapatkan kelas yang belum diketahui yang dapat digunakan untuk memprediksi suatu kelas atau label tertentu.

Klasifikasi adalah suatu tipe data yang dianalisis yang dapat membantu orang menentukan kelas label dari sampel data yang ingin di klasifikasi dan mencari hubungan antara fitur masukan (*input*) dan fitur target (kelas) (Wah et al., 2018). Klasifikasi banyak digunakan untuk memprediksi suatu kelas atau label tertentu, yaitu dengan mengklasifikasi data dengan cara membangun model *berdasarkan data training* (data latih) dan memprediksi kelas atau label baru yang tidak

diketahui dari dataset yang mempunyai kelas menggunakan model dari klasifikasi untuk memprediksi data yang baru (Seh, 2019). Penelitian yang telah dilakukan oleh (Ashari et al., 2013) tentang perbandingan kinerja *Decision Tree*, *Naïve Bayes*, dan *K-Nearest Neighbor* adalah tiga teknik klasifikasi untuk mencari alternatif desain dalam alat replikasi energi. Hasil penelitian menunjukkan bahwa algoritma klasifikasi *Naïve Bayes* memiliki kinerja terbaik berdasarkan *Precision*, *Recall*, *Fmeasure*, *Accuracy*, dan AUC. *Naïve Bayes* mengungguli *Decision Tree* dan *K-Nearest Neighbor* pada semua parameter kecuali presisi. Oleh karena itu, algoritma klasifikasi yang digunakan untuk klasifikasi adalah algoritma *Naïve Bayes*.

Algoritma *Naïve Bayes* merupakan salah satu metode machine learning yang menggunakan perhitungan probabilitas dan statistik untuk memprediksi probabilitas atau peluang keanggotaan suatu class (Artaye, 2015). *Thomas Bayes* merupakan seorang ilmuwan Inggris yang menciptakan metode *Naïve Bayes*. Teorema Bayes menyatakan bahwa *Naïve Bayes* dapat memprediksi probabilitas masa depan berdasarkan data historis (Ali et al., 2021). Algoritma *Naïve Bayes* memiliki akurasi dan kecepatan yang tinggi saat diaplikasikan ke dalam database dengan data yang besar (Artaye, 2015).

Fakultas Sains dan Teknologi merupakan salah satu fakultas yang ada di Universitas Islam Negeri Sultan Syarif Kasim Riau di Pekanbaru. *Datasets* diperoleh terdiri dari 4 jurusan yaitu Teknik Industri, Teknik Elektro, Sistem Informasi dan Matematika menunjukkan bahwa banyak mahasiswa yang menyelesaikan studinya dengan status tidak tepat waktu dapat berpengaruh pada penilaian akreditasi, sehingga perlu diperhatikan oleh perguruan tinggi dalam menentukan kualitas perguruan tinggi yang menjadi penilaian akreditasi. Pihak fakultas berusaha untuk membantu tenaga pengajar dalam menganalisa cara belajar, mendeteksi mahasiswa yang membutuhkan dukungan lebih. Hal tersebut perlu dilakukan untuk mencegah kegagalan akademik mahasiswa. Berdasarkan permasalahan tersebut, maka dalam hal ini peneliti akan melakukan penelitian tugas akhir tesis yang berjudul “Klasifikasi Kelulusan Mahasiswa Menggunakan Metode *Naïve Bayes* dengan Komparasi Prapemrosesan *Random Oversampling* serta Seleksi Fitur *Information Gain* dan *Forward Selection* (Studi Kasus : Fakultas Sains dan

Teknologi UIN SUSKA Riau)”. Oleh karena itu, kontribusi penting dari penelitian ini adalah mengklasifikasi kelulusan mahasiswa untuk mengetahui kinerja akademik mahasiswa, menganalisis berbagai fitur untuk mengetahui fitur-fitur relevan untuk analisis data klasifikasi, selanjutnya dapat mengetahui perbandingan seleksi fitur mana yang terbaik dalam melakukan klasifikasi kelulusan mahasiswa serta mengetahui perbandingan kinerja algoritma *naïve bayes* sebelum dan setelah diterapkan prapemrosesan *random oversampling* dan seleksi fitur, selanjutnya diharapkan dapat membantu perguruan tinggi dan fakultas untuk mengetahui variabel yang berpengaruh terhadap tingkat kelulusan mahasiswa.

B. Rumusan Masalah

Berdasarkan permasalahan yang dijelaskan pada latar belakang, maka dapat dirumuskan masalah yang akan dijelaskan pada laporan Tugas Akhir Tesis ini adalah sebagai berikut :

1. Bagaimana menerapkan teknik *random oversampling* serta seleksi fitur *information gain* dan *forward selection* pada algoritma *naïve bayes* untuk klasifikasi kelulusan mahasiswa di Fakultas Sains dan Teknologi UIN SUSKA Riau ?
2. Bagaimana mengevaluasi perbandingan kinerja akurasi dalam penerapan teknik *random oversampling* pada algoritma *naïve bayes*, seleksi fitur *information gain* pada algoritma *naïve bayes*, seleksi fitur *forward selection* pada algoritma *naïve bayes* dan tanpa prapemrosesan *random oversampling* dan seleksi fitur untuk klasifikasi kelulusan mahasiswa di Fakultas Sains dan Teknologi UIN SUSKA Riau ?

C. Batasan Masalah

Ada beberapa batasan masalah dalam penelitian ini yaitu sebagai berikut :

1. Data yang dibutuhkan sebagai *datasets* utama adalah data akademik mahasiswa Fakultas Sains dan Teknologi UIN SUSKA Riau dari tahun 2016-2019 yang sudah dinyatakan lulus sebanyak 1420 data.

2. Kelas atau label yang digunakan sebagai hasil klasifikasi yaitu tepat waktu dan tidak tepat waktu.

D. Tujuan Penelitian

Tujuan dalam penelitian ini adalah sebagai berikut :

1. Menerapkan teknik *random oversampling* serta seleksi fitur *information gain* dan *forward selection* pada algoritma *naïve bayes* untuk klasifikasi kelulusan mahasiswa di Fakultas Sains dan Teknologi UIN SUSKA Riau.
2. Mengevaluasi perbandingan kinerja akurasi dalam penerapan teknik *random oversampling* pada algoritma *naïve bayes*, seleksi fitur *information gain* pada algoritma *naïve bayes*, seleksi fitur *forward selection* pada algoritma *naïve bayes* dan tanpa prapemrosesan *random oversampling* dan seleksi fitur untuk klasifikasi kelulusan mahasiswa di Fakultas Sains dan Teknologi UIN SUSKA Riau.

E. Manfaat Penelitian

Adapun manfaat penelitian adalah :

1. Dengan menganalisis berbagai fitur dapat mengetahui fitur-fitur relevan, terutama untuk analisis data klasifikasi.
2. Dengan adanya penelitian komparasi teknik *random oversampling* serta seleksi fitur *information gain* dan *forward selection* pada algoritma *naïve bayes* untuk klasifikasi kelulusan mahasiswa, dapat mengetahui perbandingan seleksi fitur mana yang terbaik dalam melakukan klasifikasi kelulusan mahasiswa serta mengetahui perbandingan kinerja algoritma *naïve bayes* sebelum dan setelah diterapkan prapemrosesan *random oversampling* dan seleksi fitur.
3. Dapat membantu perguruan tinggi dan fakultas untuk mengetahui variabel yang berpengaruh terhadap kelulusan mahasiswa.

F. Sistematika Penulisan

Laporan penelitian ini ditulis dengan sistematika penulisan secara sistematis yaitu sebagai berikut :

BAB I PENDAHULUAN

Bab ini menjelaskan tentang latar belakang masalah, rumusan masalah, batasan masalah, tujuan penelitian, manfaat penelitian, serta sistematika penulisan.

BAB II KAJIAN PUSTAKA DAN LANDASAN TEORI

Bab ini menjelaskan tentang kajian hasil penelitian yang telah dilakukan para peneliti lain pada kajian pustaka, sedangkan landasan teori merupakan kajian dari teori-teori para ahli yang dijadikan landasan dalam penelitian.

BAB III METODE PENELITIAN

Bab ini berisi mengenai objek penelitian dan metode yang digunakan dalam penelitian, serta tahapan-tahapan penelitian yang dilakukan.

BAB IV HASIL DAN PEMBAHASAN

Bab ini menjelaskan proses analisa terhadap permasalahan serta solusi yang diberikan, perancangan program berupa perancangan database, implementasi hasil perancangan ke dalam sebuah program serta menjelaskan tentang hasil pengujian dari program yang sudah dibuat.

BAB V PENUTUP

Bab ini berisi kesimpulan yang diperoleh dari pembahasan dalam penelitian serta saran untuk penelitian selanjutnya.

BAB V

PENUTUP

A. Kesimpulan

Berdasarkan penelitian yang telah dilakukan, maka kesimpulan yang diperoleh dari penelitian ini adalah sebagai berikut.

1. Metode klasifikasi *Naïve Bayes* dapat digunakan untuk mengklasifikasikan kelulusan mahasiswa. Pengujian algoritma *Naïve Bayes* dilakukan untuk mengetahui perbandingan kinerja klasifikasi sebelum dan setelah diterapkan prapemrosesan *Random Oversampling* dan seleksi fitur *Information Gain* dan *Forward Selection*. Berdasarkan hasil pengujian kinerja, diperoleh hasil bahwa kinerja algoritma klasifikasi memiliki kinerja yang baik terhadap *datasets* kelulusan mahasiswa. Algoritma *Naïve Bayes* mengalami peningkatan kinerja setelah dilakukan prapemrosesan *Random Oversampling* dan seleksi fitur. Metode seleksi fitur *Forward Selection* memiliki pengaruh kinerja lebih baik dalam mengklasifikasi data, sehingga dapat mengurangi kompleksitas data dan fitur-fitur yang tidak relevan serta meningkatkan kinerja algoritma klasifikasi *Naïve Bayes*.
2. Akurasi diperoleh pada algoritma *Naïve Bayes* tanpa prapemrosesan *Random Oversampling* dan seleksi fitur dengan metode evaluasi *k-fold cross validation* adalah 81.83%, dengan prapemrosesan *Random Oversampling* menghasilkan akurasi sebesar 83.84%, dengan seleksi fitur *Information Gain* diperoleh akurasi 86.03% dengan 3 fitur terpilih yaitu ip semester 8, ip semester 7, dan ipk, dengan seleksi fitur *Forward Selection* diperoleh akurasi 86.42% dengan 2 fitur terpilih yaitu ip semester 8 dan ipk, sehingga terjadi peningkatan akurasi sebesar 4.2% dari tanpa prapemrosesan ke *Information Gain* dan 4.59% dari tanpa prapemrosesan ke *Forward Selection*. Oleh karena itu, dapat diketahui bahwa tahapan prapemrosesan *Random Oversampling* dan seleksi fitur berpengaruh pada algoritma klasifikasi *Naïve Bayes* dan seleksi fitur *Forward Selection* memiliki pengaruh kinerja lebih baik dengan 2 atribut terpilih dalam

mengklasifikasi data secara tepat, sehingga dapat meningkatkan kinerja algoritma klasifikasi *Naïve Bayes*.

B. Saran

Saran yang diberikan untuk penelitian selanjutnya adalah melakukan prapemrosesan lain seperti *outlier detection* dan prapemrosesan seleksi fitur dengan pendekatan lain selain pendekatan *filter* dan *wrapper* seperti *embedded* untuk mengetahui seleksi fitur yang paling berpengaruh pada klasifikasi. Selain itu diharapkan juga untuk menguji metode klasifikasi selain *Naïve Bayes* untuk mengetahui metode klasifikasi mana yang paling baik kinerjanya.

DAFTAR PUSTAKA

- ACCJC/WASC. (2012). *Guide To Evaluating Institutions*. 56. www.g-fras.org
- Albarak, M., Alrazgan, M., & Bahsoon, R. (2017). *Identifying and Managing Technical Debt in Database Normalization Using Machine Learning and Trade-off Analysis*. <http://arxiv.org/abs/1711.06109>
- Ali, A., Khairan, A., Tempola, F., & Fuad, A. (2021). Application Of Naïve Bayes to Predict the Potential of Rain in Ternate City. *E3S Web of Conferences*, 328, 04011. <https://doi.org/10.1051/e3sconf/202132804011>
- Ang, J. C., Mirzal, A., Haron, H., & Hamed, H. N. A. (2016). Supervised, unsupervised, and semi-supervised feature selection: A review on gene selection. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 13(5), 971–989. <https://doi.org/10.1109/TCBB.2015.2478454>
- Artaye, K. (2015). *International Conference On Information Technology And Business ISSN 2460-7223 IMPLEMENTATION OF NAÏVE BAYES CLASSIFICATION METHOD TO PREDICT GRADUATION TIME OF IBI DARMAJAYA SCHOLAR Z. A . Pagar Alam Street No . 93 Bandar Lampung. August*, 284–290.
- Ashari, A., Paryudi, I., & Tjoa, A. M. (2013). Performance Comparison between Naïve Bayes, Decision Tree and k-Nearest Neighbor in Searching Alternative Design in an Energy Simulation Tool. (*IJACSA*) *International Journal of Advanced Computer Science and Applications*, 4(11), 33–39.
- Astuti, F. D. (2018). Seleksi Fitur Forward Selection pada Algoritma Naive Bayes untuk Klasifikasi Benih Gandum. *Jurnal Informasi Interaktif*, 3(1), 161–166.
- BAN-PT. (2015). Akreditasi Institusi Perguruan Tinggi. *Pedoman Penyusunan Borang*.
- Bassi, J. S., Dada, E. G., Hamidu, A. A., & Elijah, M. D. (2019). Students Graduation on Time Prediction Model Using Artificial Neural Network. *IOSR Journal of Computer Engineering (IOSR-JCE)*, 21(3), 28–35. <https://doi.org/10.9790/0661-2103012835>

- Bermejo, P., Gámez, J. A., & Puerta, J. M. (2014). Speeding up incremental wrapper feature subset selection with Naive Bayes classifier. *Knowledge-Based Systems*, 55, 140–147. <https://doi.org/10.1016/j.knosys.2013.10.016>
- Berrar, D. (2018). Cross-validation. *Encyclopedia of Bioinformatics and Computational Biology: ABC of Bioinformatics*, 1–3(January 2018), 542–545. <https://doi.org/10.1016/B978-0-12-809633-8.20349-X>
- Bimantoro, D. A., & ‘Uyun, S. (2017). *PENGARUH PENGGUNAAN INFORMATION GAINUNTUK SELEKSI FITUR CITRA TANAH DALAM RANGKA MENILAI KESESUAIAN LAHAN PADA TANAMAN CENGKEH*. 2(1), 42–52.
- C.Deisy, B.Subbulakshmi, S.Baskar, D., & Dr.N.Ramaraj. (2007). Efficient Dimensionality Reduction Approaches for Feature Selection. *Proceedings - International Conference on Computational Intelligence and Multimedia Applications, ICCIMA 2007*, 4, 270–272. <https://doi.org/10.1109/ICCIMA.2007.288>
- Chandrashekar, G., & Sahin, F. (2014). A survey on feature selection methods. *Computers and Electrical Engineering*, 40(1), 16–28. <https://doi.org/10.1016/j.compeleceng.2013.11.024>
- Fali Oklilas, A., Tasmi, Desy Siswanti, S., Afrina, M., & Setiawan, H. (2019). Attribute selection using information gain and naïve bayes for traffic classification. *Journal of Physics: Conference Series*, 1196(1). <https://doi.org/10.1088/1742-6596/1196/1/012021>
- Fanani, M. R. (2020). Penggabungan Forward Selection untuk Pemilihan Fitur pada Prediksi Bimbingan Konseling Siswa dengan Menggunakan Algoritma Naive Bayes. *Smart Comp :Jurnalnya Orang Pintar Komputer*, 9(2), 85–88. <https://doi.org/10.30591/smartcomp.v9i2.1924>
- Fayyad, U., Piatetsky-Shapiro, G., & Smyth, P. (1997). Knowledge Discovery and Data Mining: Towards a Unifying Framework. *Journal of Japan Society for Fuzzy Theory and Systems*, 9(6), 851–860. https://doi.org/10.3156/jfuzzy.9.6_851
- Gnana, D. A. A., Singh, Balamurugan, S. A., & Leavline, E. J. (2016). Literature

- Review on Feature Selection Methods for High-Dimensional Data. *International Journal of Computer Applications*, 136(1), 9–17.
- Grzymala-Busse, J. W., & Grzymala-Busse, W. J. (2010). Handling Missing Attribute Values. *Data Mining and Knowledge Discovery Handbook*, 33–51. https://doi.org/10.1007/978-0-387-09823-4_3
- Han, J., Kamber, M., & Pei, J. (2011). Data Mining: Concepts and Techniques. In *Third Edition. 2011. British Library Cataloguing-in-Publication. Morgan Kaufmann*. <https://doi.org/10.1109/ICMIRA.2013.45>
- Harahap, F., Harahap, A. Y. N., Ekadiansyah, E., Sari, R. N., Adawiyah, R., & Harahap, C. B. (2019). Implementation of Naïve Bayes Classification Method for Predicting Purchase. *2018 6th International Conference on Cyber and IT Service Management, CITSM 2018, Citsm*, 1–5. <https://doi.org/10.1109/CITSM.2018.8674324>
- Hashemi, H. Z., Parvasideh, P., Larijani, Z. H., & Moradi, F. (2018). Analyze students performance of a national exam using feature selection methods. *2018 8th International Conference on Computer and Knowledge Engineering, ICCKE 2018, Iccke*, 7–11. <https://doi.org/10.1109/ICCKE.2018.8566671>
- Henderi, H. (2021). Comparison of Min-Max normalization and Z-Score Normalization in the K-nearest neighbor (kNN) Algorithm to Test the Accuracy of Types of Breast Cancer. *IJIIS: International Journal of Informatics and Information Systems*, 4(1), 13–20. <https://doi.org/10.47738/ijiis.v4i1.73>
- Heranova, O. (2021). *Synthetic Minority Oversampling Technique pada Averaged One Dependence Estimators untuk Klasifikasi Credit Scoring*. 1(10), 10–12.
- Jain, D., & Singh, V. (2018). Feature selection and classification systems for chronic disease prediction: A review. *Egyptian Informatics Journal*, 19(3), 179–189. <https://doi.org/10.1016/j.eij.2018.03.002>
- Jurafsky, D., & Martin, J. (2019). Naive bayes and sentiment classification. *Speech and Language Processing*, 1024. <https://web.stanford.edu/~jurafsky/slp3/6.pdf>
- Kalyani, P., & Karnan, D. M. (2010). Attribute Reduction using Forward Selection

- and Relative Reduct Algorithm. *International Journal of Computer Applications*, 11(3), 8–12. <https://doi.org/10.5120/1564-1499>
- Khaire, U. M., & Dhanalakshmi, R. (2019). Stability of feature selection algorithm: A review. *Journal of King Saud University - Computer and Information Sciences*, xxx. <https://doi.org/10.1016/j.jksuci.2019.06.012>
- Khemphila, A., & Boonjing, V. (2011). Heart disease classification using neural network and feature selection. *Proceedings - ICSEng 2011: International Conference on Systems Engineering*, 2007, 406–409. <https://doi.org/10.1109/ICSEng.2011.80>
- Lasarudin, A., & Purwanto. (2018). KLASIFIKASI PENGADUAN MASYARAKAT MENGGUNAKAN NAIVE BAYES BERBASIS SELEKSI ATRIBUT INFORMATION GAIN. *Jurnal Teknologi Informasi*, 14(2), 63. <http://research.pps.dinus.ac.id/index.php/Cyberku/article/view/65>
- Lei, C., & Li, K. F. (2015). Academic Performance Predictors. *Proceedings - IEEE 29th International Conference on Advanced Information Networking and Applications Workshops*, WAINA 2015, 577–581. <https://doi.org/10.1109/WAINA.2015.114>
- Mani, K., & Kalpana, P. (2016). An Efficient Feature Selection based on Bayes Theorem, Self Information and Sequential Forward Selection. *International Journal of Information Engineering and Electronic Business*, 8(6), 46–54. <https://doi.org/10.5815/ijieeb.2016.06.06>
- Markoulidakis, I., Rallis, I., Georgoulas, I., Kopsiaftis, G., Doulamis, A., & Doulamis, N. (2021). Multiclass Confusion Matrix Reduction Method and Its Application on Net Promoter Score Classification Problem. *Technologies*, 9(4), 81. <https://doi.org/10.3390/technologies9040081>
- Muin, A. A. (2016). Penerapan seleksi atribut weights by information gain dan select by weights pada algoritma Naive Bayes untuk prediksi kolektibilitas pembiayaan usaha kecil dan menengah. *Technologia*, 7(4), 245–249.
- Mulyadi, C., & Sukron, . (2020). *Prediction of Timeliness of Graduating with Naive Bayes Algorithm*. *Icri* 2018, 3043–3050. <https://doi.org/10.5220/0009946430433050>

- Muqorobin, M., Kusrini, K., & Luthfi, E. T. (2019). Optimasi Metode Naive Bayes Dengan Feature Selection Information Gain Untuk Prediksi Keterlambatan Pembayaran Spp Sekolah. *Jurnal Ilmiah SINUS*, 17(1), 1. <https://doi.org/10.30646/sinus.v17i1.378>
- Mwadulo, M. W. (2016). A Review on Feature Selection Methods For Classification Tasks. *International Journal of Computer Applications Technology and Research*, 5(6), 395–402. <https://doi.org/10.1109/ICISC.2017.8068746>
- Naganjaneyulu, S., & Kuppa, M. R. (2013). A novel framework for class imbalance learning using intelligent under-sampling. *Progress in Artificial Intelligence*, 2(1), 73–84. <https://doi.org/10.1007/s13748-012-0038-2>
- Normawati, D., & Ismi, D. P. (2019). K-Fold Cross Validation for Selection of Cardiovascular Disease Diagnosis Features by Applying Rule-Based Datamining. *Signal and Image Processing Letters*, 1(2), 23–35. <https://doi.org/10.31763/simple.v1i2.3>
- Nuffic. (2017). *Education system Indonesia*. www.nuffic.nl/en/home/copyright
- Nugroho, M. F., & Wibowo, S. (2017). Fitur Seleksi Forward Selection Untuk Menentukan Atribut Yang Berpengaruh Pada Klasifikasi Kelulusan Mahasiswa Fakultas Ilmu Komputer UNAKI Semarang Menggunakan Algoritma Naive Bayes. *Jurnal Informatika Upgris*, 3(1), 63–70. <https://doi.org/10.26877/jiu.v3i1.1669>
- Nurjoko, H., & Kurniawan. (2016). *Aplikasi datamining untuk memprediksi tingkat kelulusan mahasiswa menggunakan algoritma apriori di ibi darmajaya bandar lampung*. 02(01), 79–93.
- Panthong, R., & Srivihok, A. (2015). Wrapper Feature Subset Selection for Dimension Reduction Based on Ensemble Learning Algorithm. *Procedia Computer Science*, 72, 162–169. <https://doi.org/10.1016/j.procs.2015.12.117>
- Patil, T. R., & Sherekar, M. S. S. (2019). Performance Analysis of ANN and Naive Bayes Classification Algorithm for Data Classification. *International Journal of Intelligent Systems and Applications in Engineering*, 7(2), 88–91. <https://doi.org/10.18201/ijisae.2019252786>

- Prasetyowati, M. I., Maulidevi, N. U., & Surendro, K. (2021). Determining threshold value on information gain feature selection to increase speed and prediction accuracy of random forest. *Journal of Big Data*, 8(1). <https://doi.org/10.1186/s40537-021-00472-4>
- Purnana, P., & Supriyanto, D. C. (2013). Deteksi Penyakit Diabetes Type Ii Dengan Naive Bayes Berbasis Particle Swarm Optimization. *Jurnal Teknologi Informasi*, 9(2), 1414–9999. <http://research.pps.dinus.ac.id>
- Putri, L. R., Mubarak, M. S., & Adiwijaya. (2017). *KLASIFIKASI SENTIMEN ULASAN BUKU BERBAHASA INGGRIS MENGGUNAKAN INFORMATION GAIN DAN NAÏVE BAYES*. 4(3), 4659–4666.
- Ridwan, M., Suyono, H., & Sarosa, M. (2013). Penerapan Data Mining Untuk Evaluasi Kinerja Akademik Mahasiswa Menggunakan Algoritma Naive Bayes Classifier. *Eeccis*, 7(1), 59–64. <https://doi.org/10.1038/hdy.2009.180>
- Rozzaqi, A. R. (2015). *Naïve Bayes dan Filtering Feature Selection Information Gain untuk Prediksi Ketepatan Kelulusan Mahasiswa*. 30–41.
- Saifudin, A., Ekawati, Yulianti, & Desyani, T. (2020). Forward Selection Technique to Choose the Best Features in Prediction of Student Academic Performance Based on Naïve Bayes. *Journal of Physics: Conference Series*, 1477(2). <https://doi.org/10.1088/1742-6596/1477/3/032007>
- Seh, A. H. (2019). *A Review on Heart Disease Prediction Using Machine Learning Techniques A Review on Heart Disease Prediction Using Machine Learning Techniques*. 9(April), 208–224.
- Shivali, Joni Birla, G. (2015). Knowledge Discovery in Data-Mining. *International Journal of Engineering Research & Technology (IJERT)*, 3(10), 1–5. <https://www.ijert.org/research/knowledge-discovery-in-data-mining-IJERTCONV3IS10051.pdf>
- Singh, M., & Gupta, S. (2020). Sentiment Analysis using Naive Bayes Classifier and Information Gain Feature Selection over Twitter. *International Journal of Computer Trends and Technology*, 68(5), 84–91. <https://doi.org/10.14445/22312803/ijctt-v68i5p117>
- Sugriyono, S., & Siregar, M. U. (2020). Preprocessing kNN algorithm classification

- using K-means and distance matrix with students' academic performance dataset. *Jurnal Teknologi Dan Sistem Komputer*, 8(4), 311–316. <https://doi.org/10.14710/jtsiskom.2020.13874>
- Tangirala, S. (2020). Evaluating the impact of GINI index and information gain on classification using decision tree classifier algorithm. *International Journal of Advanced Computer Science and Applications*, 11(2), 612–619. <https://doi.org/10.14569/ijacsa.2020.0110277>
- Ting, S. L., Ip, W. H., & Tsang, A. H. C. (2011). Is Naïve bayes a good classifier for document classification? *International Journal of Software Engineering and Its Applications*, 5(3), 37–46.
- Uyun, S., & Choridah, L. (2018). Feature selection mammogram based on breast cancer mining. *International Journal of Electrical and Computer Engineering*, 8(1), 60–69. <https://doi.org/10.11591/ijece.v8i1.pp60-69>
- Uyun, S., & Sulistyowati, E. (2020). Feature selection for multiple water quality status: Integrated bootstrapping and SMOTE approach in imbalance classes. *International Journal of Electrical and Computer Engineering*, 10(4), 4331–4339. <https://doi.org/10.11591/ijece.v10i4.pp4331-4339>
- van Lith, J. W., & Vanschoren, J. (2021). *From Strings to Data Science: a Practical Framework for Automated String Handling*. 1–19. <https://arxiv.org/abs/2111.01868v1>
- Vanaja, S., & Ramesh Kumar, K. (2014). Analysis of Feature Selection Algorithms on Classification: A Survey. *International Journal of Computer Applications*, 96(17), 29–35. <https://doi.org/10.5120/16888-6910>
- Visalakshi, S., & Radha, V. (2015). A literature review of feature selection techniques and applications: Review of feature selection in data mining. *2014 IEEE International Conference on Computational Intelligence and Computing Research, IEEE ICCIC 2014*, 1997. <https://doi.org/10.1109/ICCIC.2014.7238499>
- Wah, Y. B., Ibrahim, N., Hamid, H. A., Abdul-Rahman, S., & Fong, S. (2018). Feature selection methods: Case of filter and wrapper approaches for maximising classification accuracy. *Pertanika Journal of Science and*

Technology, 26(1), 329–340.

Wang, Q., Luo, Z., Huang, J., Feng, Y., & Liu, Z. (2017). A Novel Ensemble Method for Imbalanced Data Learning. *Computational Intelligence and Neuroscience*, 2017, 1–11.

Wang, Y., Li, Y., Song, Y., Rong, X., & Zhang, S. (2017). Improvement of ID3 algorithm based on simplified information entropy and coordination degree. *Algorithms*, 10(4), 1–18. <https://doi.org/10.3390/a10040124>

Yunianita, S., Setiani, N., dan Mulyati, S. (2018). Prediksi Ketepatan Masa Studi Mahasiswa dengan Algoritma Pohon Keputusan C45. *SNATi*, 23–29.

Zeniarja, J., Widia, K., & Sani, R. R. (2020). Penerapan Algoritma Naive Bayes dan Forward Selection dalam Pengklasifikasian Status Gizi Stunting pada Puskesmas Pandanaran Semarang. *JOINS (Journal of Information System)*, 5(1), 1–9. <https://doi.org/10.33633/joins.v5i1.2745>

