

**IMPLEMENTASI *DATA MINING* UNTUK PREDIKSI
DROP-OUT DENGAN MENGGUNAKAN
*RANDOM FOREST METHOD***

**Tesis
untuk memenuhi sebagian persyaratan
mencapai derajat Sarjana S-2 Program Studi
Magister Sistem Informasi**



**Meylani Utari
30000318410029**

**SEKOLAH PASCASARJANA
UNIVERSITAS DIPONEGORO
SEMARANG
2020**

HALAMAN PENGESAHAN

TESIS

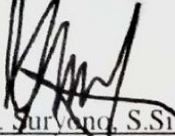
**IMPLEMENTASI DATA MINING UNTUK PREDIKSI DROP-OUT
DENGAN MENGGUNAKAN RANDOM FOREST METHOD**

Oleh:
Meylani Utari
30000318410029

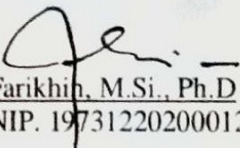
Telah diujikan dan dinyatakan lulus ujian tesis pada tanggal 23 Juni 2020 oleh tim penguji Program Studi Magister Sistem Informasi Sekolah Pascasarjana Universitas Diponegoro.

Semarang, 29 Juni 2020
Mengetahui,

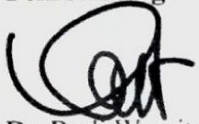
Penguji I


Dr. Suryono, S.Si., M.Si
NIP. 197306301998021001

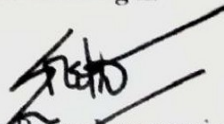
Penguji II


Farikhin, M.Si., Ph.D
NIP. 197312202000121001

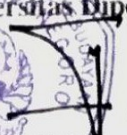
Pembimbing I


Dr. Budi Warsito, S.Si., M.Si
NIP. 197508241999031003

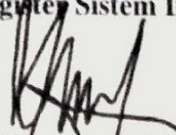
Pembimbing II


Dr. Retno Kusumaningrum, S.Si., M.Kom
NIP. 198104202005012001

Mengetahui :
Dekan Sekolah Pascasarjana
Universitas Diponegoro


Dr. R. B. Sularto, S.H., M.Hum
NIP. 196707011991031005

Ketua Program Studi
Magister Sistem Informasi


Dr. Suryono, S.Si., M.Si
NIP. 197306301998021001

PERNYATAAN

Dengan ini saya menyatakan bahwa dalam tesis ini tidak terdapat karya yang pernah diajukan untuk memperoleh gelar akademik di suatu perguruan tinggi, dan sepanjang pengetahuan saya juga tidak terdapat karya atau pendapat yang pernah ditulis atau diterbitkan oleh orang lain, kecuali secara tertulis diacu dalam naskah ini dan disebutkan dalam daftar pustaka.

Semarang, 26 Juni 2020



Meylani Utari

PERNYATAAN PERSETUJUAN PUBLIKASI TESIS UNTUK KEPENTINGAN AKADEMIS

Sebagai sivitas akademik Universitas Diponegoro, saya yang bertanda tangan di bawah ini :

Nama : Meylani Utari
NIM : 30000318410029
Program Studi : Magister Sistem Informasi
Program : Pascasarjana
Jenis Karya : Tesis

Demi pengembangan ilmu pengetahuan, menyetujui untuk memberikan kepada Universitas Diponegoro Hak Bebas Royalti Noneksklusif atas karya ilmiah saya yang berjudul :

IMPLEMENTASI *DATA MINING* UNTUK PREDIKSI *DROP-OUT* DENGAN MENGGUNAKAN *RANDOM FOREST METHOD*

beserta perangkat yang ada. Dengan Hak Bebas Royalti Noneksklusif ini Program Studi Magister Sistem Informasi Sekolah Pascasarjana Universitas Diponegoro berhak menyimpan, mengalihmedia/formatkan, mengelola dalam bentuk pangkalan data (*database*), merawat, dan mempublikasikan tesis saya selama tetap mencantumkan nama saya sebagai penulis/pencipta dan sebagai pemilik hak cipta.

Dibuat di : Semarang
Pada Tanggal : 26 Juni 2020
Yang menyatakan



Meylani Utari
NIM. 30000318410029

KATA PENGANTAR

Puji dan syukur kepada Tuhan yang Maha Esa, atas segala berkat, rahmat, dan karunia yang dicurahkan, sehingga tesis dengan judul Implementasi *Data Mining* untuk Prediksi *Drop-Out* dengan Menggunakan *Random Forest Method* ini dapat diselesaikan. Tesis ini disusun untuk memenuhi salah satu persyaratan memperoleh gelar Magister Komputer (M.Kom) pada Program Studi Magister Sistem Informasi Universitas Diponegoro. Pada kesempatan ini penulis menyampaikan terima kasih yang sebesar – besarnya kepada :

1. Dr. Budi Warsito, S.Si., M.Si selaku Pembimbing I dan Dr. Retno Kusumaningrum, S.Si, M.Kom selaku Pembimbing II yang telah memberikan waktu, ilmu, saran, kritik, semangat, dan nasihat selama penulisan tesis ini.
2. Dr. Suryono, S.Si, M.Si selaku Ketua Program Studi Magister Sistem Informasi Sekolah Pascasarjana Universitas Diponegoro Semarang.
3. Program Beasiswa PasTi Tahun 2018 yang diselenggarakan oleh Kemeristekdikti dan Kemendikbud yang telah memberikan kesempatan kepada Penulis sehingga dapat melanjutkan pendidikan Magister di Universitas Diponegoro Semarang dan menyelesaikan penelitian ini.
4. Universitas Sriwijaya dan Fakultas Ilmu Komputer sebagai tempat Penulis bekerja yang telah memberikan dukungan kepada Penulis dengan memberikan izin untuk melanjutkan Pendidikan Magister di Universitas Diponegoro Semarang.
5. Suami, kedua orangtua, sahabat, teman-teman, serta pihak-pihak lain yang tidak dapat disebutkan satu per-satu, yang telah memberikan kontribusi baik moral maupun material sehingga tesis ini dapat terselesaikan dengan baik.

Penulis menyadari bahwa dalam penyusunan tesis ini masih jauh dari sempurna. Oleh karena itu, saran dan kritik yang sifatnya membangun sangat diharapkan. Akhirnya, penulis berharap semoga tulisan ini bermanfaat. Amin.

Semarang, 26 Juni 2020

Penulis,

Meylani Utari

DAFTAR ISI

HALAMAN PENGESAHAN	ii
PERNYATAAN	iii
PERNYATAAN PERSETUJUAN	iv
KATA PENGANTAR	v
DAFTAR ISI	vi
DAFTAR GAMBAR	ix
DAFTAR TABEL	x
ABSTRAK	xii
ABSTRACT	xiii
BAB I PENDAHULUAN	1
1.1 Latar Belakang	1
1.2 Tujuan Penelitian	3
1.3 Manfaat Penelitian	4
BAB II TINJAUAN PUSTAKA DAN DASAR TEORI	5
2.1 Tinjauan Pustaka	5
2.2 Dasar Teori	7
2.2.1 <i>Drop-Out</i>	7
2.2.2 <i>Business Intelligence (BI)</i>	8
2.2.3 <i>Data Mining</i>	9
2.2.4 Tahapan <i>Data Mining</i>	11
2.2.5 Prediksi	12
2.2.6 Klasifikasi	12
2.2.7 Synthetic Minority Over-Sampling Technique (SMOTE)	13
2.2.8 <i>Decision Tree (DT)</i>	16
2.2.9 <i>Random Forest (RF)</i>	17
2.2.10 <i>K-Fold Cross Validation</i>	19
2.2.11 <i>Confusion Matrix (Pengukuran Kinerja)</i>	20
BAB III METODE PENELITIAN	23
3.1 Bahan dan Alat Penelitian	23
3.2 Prosedur Penelitian	24
3.2.1 <i>Business Understanding (Pemahaman Bisnis)</i>	24
3.2.2 <i>Data Understanding (Pemahaman Data)</i>	25
3.2.3 <i>Data Preparation (Persiapan Data)</i>	27
3.2.4 <i>Modeling (Pemodelan)</i>	31
3.2.5 <i>Evaluation (Evaluasi)</i>	33
3.2.6 <i>Deployment (Penyebaran)</i>	33
3.3 Kerangka Sistem Informasi	34
3.4 Diagram Alur Kerja Sistem Prediksi	36
3.5 Desain Antarmuka	38
BAB IV HASIL PENELITIAN DAN PEMBAHASAN	40
4.1 Hasil Penelitian	40
4.2 Pembahasan	54
4.3 <i>Deployment</i>	80
4.3.1 Model Prediksi <i>Drop-Out</i>	80
4.3.2 Implementasi Antarmuka	83

BAB V KESIMPULAN DAN SARAN	94
5.1 Kesimpulan.....	94
5.2 Saran.....	95



DAFTAR GAMBAR

Gambar 2.1 Tahapan <i>Data Mining</i> (Fayyad, 1996)	11
Gambar 2.2. Model Klasifikasi (Sumarlin, 2015).....	13
Gambar 2.3 Prosedur 3-Fold Cross-Validation (Refaeilzadeh dkk., 2009).....	20
Gambar 3.1 Siklus CRISP-DM (Gutierrez dkk., 2018).....	24
Gambar 3.2 Ilustrasi Data <i>Preprocessing</i>	28
Gambar 3.3. Ilustrasi <i>10-fold cross validation</i>	32
Gambar 3.4. Diagram Alur Sistem Prediksi <i>Drop-Out</i>	37
Gambar 3.5. Halaman Utama.....	38
Gambar 3.6. Halaman Proses <i>Training Data</i>	38
Gambar 3.7. Halaman Implementasi Model	39
Gambar 4.1. Arsitektur Sistem.....	41
Gambar 4.2. <i>Chart</i> Sebaran Data Sebelum <i>Preprocessing</i>	42
Gambar 4.3. <i>Chart</i> Sebaran Data Setelah <i>Preprocessing</i>	42
Gambar 4.4. <i>Chart</i> Sebaran Data Setelah Proses SMOTE.....	42
Gambar 4.5. Pembentukan Pohon Keputusan – <i>Decision Tree</i>	48
Gambar 4.6. Pembentukan Pohon Keputusan – <i>Random Forest</i>	48
Gambar 4.7. Proses Pelatihan Data dengan dan RF secara keseluruhan	49
Gambar 4.8. Proses Pelatihan Data dengan RF tanpa SMOTE	50
Gambar 4.9. Lanjut Proses Pelatihan Data dengan RF tanpa SMOTE.....	50
Gambar 4.10. Proses Pelatihan Data dengan RF dengan SMOTE	51
Gambar 4.11. Lanjut Proses Pelatihan Data dengan RF dengan SMOTE	51
Gambar 4.12. Proses Pelatihan Data dengan DT tanpa SMOTE.....	52
Gambar 4.13. Lanjut Proses Pelatihan Data dengan DT dengan SMOTE	52
Gambar 4.14. Proses Pelatihan Data dengan DT dengan SMOTE.....	53
Gambar 4.15. Lanjut Proses Pelatihan Data dengan DT dengan SMOTE	53
Gambar 4.16. Grafik <i>Performance Random Forest</i> tanpa SMOTE	80
Gambar 4.17. Grafik <i>Performance Random Forest</i> dengan SMOTE	81
Gambar 4.18. Grafik <i>Performance Decision Tree</i> tanpa SMOTE.....	82
Gambar 4.19. Grafik <i>Performance Decision Tree</i> dengan SMOTE.....	82
Gambar 4.20. Halaman Utama.....	84
Gambar 4.21. Halaman Proses <i>Training Data</i>	85
Gambar 4.22. Lanjut Halaman Proses <i>Training Data</i>	85
Gambar 4.23. Lanjut Halaman Proses <i>Training Data</i>	86
Gambar 4.24. Lanjut Halaman Proses <i>Training Data</i>	87
Gambar 4.25. Lanjut Halaman Proses <i>Training Data</i>	87
Gambar 4.26. Lanjut Halaman Proses <i>Training Data</i>	88
Gambar 4.27. Lanjut Halaman Proses <i>Training Data</i>	89
Gambar 4.28. Lanjut Halaman Proses <i>Training Data</i>	89
Gambar 4.29. Halaman Implementasi Model RF (Sistem Prediksi DO).....	90
Gambar 4.30. Halaman Implementasi Model DT (Sistem Prediksi DO)	90
Gambar 4.31. Halaman Implementasi Model RF (Sistem Prediksi DO).....	91
Gambar 4.32. Halaman Implementasi Model DT (Sistem Prediksi DO)	92

DAFTAR TABEL

Tabel 2.1 Penelitian terkait	5
Tabel 2.2 Algoritma <i>Random Forest</i> (Breiman, 2001).....	18
Tabel 2.3 Algoritma <i>Out-Of-Bag</i> pada <i>Random Forest</i> (Breiman, 2001)	19
Tabel 2.4 Confusion Matrix	21
Tabel 3.1 Atribut Dataset	26
Tabel 3.2 Seleksi Atribut	29
Tabel 3.3. Kerangka Sistem Informasi.....	34
Tabel 4.1. Data Kelas Minoritas (<i>drop-out</i>).....	44
Tabel 4.2. Lanjutan Data Kelas Minoritas (<i>drop-out</i>)	45
Tabel 4.3. Pembentukan Data Sintetis	46
Tabel 4.4. Lanjutan Pembentukan Data Sintetis	47
Tabel 4.5. Hasil Pelatihan Data dengan RF tanpa SMOTE.....	50
Tabel 4.6. Hasil Pelatihan Data Algoritma RF dengan SMOTE	51
Tabel 4.7. Hasil Pelatihan Data dengan <i>Decision Tree</i> tanpa SMOTE	52
Tabel 4.8. Hasil Pelatihan Data Algoritma <i>Decision Tree</i> dengan SMOTE.....	53
Tabel 4.9. <i>Confusion Matrix</i> Skenario 1	55
Tabel 4.10. <i>Confusion Matrix</i> Skenario 2	56
Tabel 4.11. <i>Confusion Matrix</i> Skenario 3	57
Tabel 4.12. <i>Confusion Matrix</i> Skenario 3	58
Tabel 4.13. <i>Confusion Matrix</i> Skenario 5	59
Tabel 4.14. <i>Confusion Matrix</i> Skenario 6	61
Tabel 4.15. <i>Confusion Matrix</i> Skenario 7	62
Tabel 4.16. <i>Confusion Matrix</i> Skenario 8	63
Tabel 4.17. <i>Confusion Matrix</i> Skenario 9	64
Tabel 4.18. <i>Confusion Matrix</i> Skenario 10	65
Tabel 4.19. <i>Confusion Matrix</i> Skenario 11	66
Tabel 4.20. <i>Confusion Matrix</i> Skenario 12	67
Tabel 4.21. <i>Confusion Matrix</i> Skenario 13	69
Tabel 4.22. <i>Confusion Matrix</i> Skenario 14	70
Tabel 4.23. <i>Confusion Matrix</i> Skenario 15	71
Tabel 4.24. <i>Confusion Matrix</i> Skenario 16	72
Tabel 4.25. <i>Confusion Matrix</i> Skenario 17	73
Tabel 4.26. <i>Confusion Matrix</i> Skenario 18	74
Tabel 4.27. <i>Confusion Matrix</i> Skenario 19	76
Tabel 4.28. <i>Confusion Matrix</i> Skenario 20	77
Tabel 4.29. <i>Confusion Matrix</i> Skenario 21	78
Tabel 4.30. <i>Confusion Matrix</i> Skenario 22	79
Tabel 4.31. <i>Performance Random Forest</i> tanpa SMOTE.....	80
Tabel 4.32. Perbandingan <i>Performance Random Forest</i> + SMOTE	81
Tabel 4.33. <i>Performance Decision Tree</i> tanpa SMOTE	81
Tabel 4.34. Perbandingan <i>Performance Decision Tree</i> + SMOTE	82

DAFTAR LAMPIRAN

Lampiran 1. Dataset Penelitian	100
Lampiran 2. Dataset Penelitian	111



ABSTRAK

Drop-out merupakan salah satu isu yang paling umum pada dunia pendidikan, tingkat *drop-out* yang tinggi memberikan dampak yang buruk terhadap perguruan tinggi seperti reputasi dan akreditasi yang kurang baik. Prediksi awal mahasiswa yang berisiko *drop-out* sangat penting dilakukan untuk membantu mengurangi persentase tingkat *drop-out*. Metode klasifikasi statistik telah banyak diimplementasikan untuk mengetahui persentase tingkat *drop-out*, namun hasil akhir pengklasifikasian yang diberikan tidak signifikan serta tingkat ketepatan dugaan yang lebih rendah dibandingkan dengan penggunaan metode seperti pohon *ensemble*. Pada penelitian ini, saya menerapkan teknik *data mining* dengan menggunakan algoritma *random forest* dan *decision tree* yang disertai dengan teknik *imbalance dataset* yaitu *synthetic minority oversampling technique* (SMOTE) untuk melakukan prediksi *drop-out* terhadap mahasiswa sarjana dan diploma pada Fakultas ABC di Perguruan Tinggi XYZ yang mendaftar pada tahun pendaftaran 2008 hingga 2012 dengan konsep *business intelligence*. Dalam penelitian ini, saya membandingkan kinerja dari algoritma *random forest* dan *decision tree* dengan SMOTE dan tanpa SMOTE. Dari hasil percobaan pada penelitian ini, *random forest* tanpa SMOTE menghasilkan nilai akurasi sebesar 92.27%, *random forest* dengan SMOTE menghasilkan nilai akurasi terbaik pada saat nilai $k = 2$ sebesar 94.72%, *decision tree* tanpa SMOTE menghasilkan nilai akurasi sebesar 91.31% dan *decision tree* dengan SMOTE menghasilkan nilai akurasi terbaik pada saat nilai $k = 10$. Hasil percobaan menunjukkan bahwa algoritma *random forest* dengan SMOTE yang memiliki nilai akurasi terbaik dibandingkan dengan ketiga algoritma lainnya. Dimasa depan, hasil penelitian ini dapat digunakan sebagai acuan untuk membangun model prediksi dini terhadap mahasiswa yang berpotensi tinggi untuk *drop-out* serta memungkinkan untuk mengetahui variabel-variabel yang menyebabkan mahasiswa tersebut dapat *drop-out*.

Kata Kunci : *drop out*; *random forest*; *synthetic minority over sampling*; SMOTE; data pendidikan; *data mining*; *classification*; *prediction*; *imbalance dataset*; *business intelligence*.

ABSTRACT

Drop-out is one of the most common issues in education, high drop-out rates have a bad impact on universities such as bad reputation and low accreditation. Prediction of students who are at risk of drop-out is very important to be done in order to help reduce the percentage of students drop-out rates. Statistical classification methods have been widely implemented to determine the percentage drop-out rate, but the final results of the classification given are insignificant and the level of accuracy of prediction is lower compared to the ensemble trees method. In this research, the author applies data mining technique using a random forest algorithm and decision tree, accompanied by an imbalance dataset technique, namely synthetic minority oversampling technique (SMOTE) to predict drop-out of undergraduate and diploma students at the ABC Faculty at XYZ University, which registered between 2008 to 2012 along with business intelligence concept. In this research, the author also compared the performance of the random forest algorithm and decision tree with SMOTE and without SMOTE. From the experimental results in this research, random forest without SMOTE produces an accuracy value of 92.27%, random forest with SMOTE produces the best accuracy value when the $k = 2$ value is 94.72%, decision tree without SMOTE produces an accuracy value of 91.31% and decision tree with SMOTE produces the best accuracy value when the value $k = 10$. The experimental results show that the random forest algorithm with SMOTE has the best accuracy compared to the other three algorithms. In the future, the results of this research can be used as a reference to build early prediction models of students with high potential for drop-out and make it possible to find out the variables that cause students to drop out.

Keywords: drop out; random forest; synthetic minority over-sampling; SMOTE; educational data; data mining; classification; prediction; imbalance dataset; business intelligence.