

**TEKNIK BALANCING DATASET MENGGUNAKAN  
ALGORITMA ADASYN DAN BORDERLINE-SMOTE**

**SKRIPSI**

Diajukan untuk memenuhi salah satu syarat mencapai derajat Sarjana  
Program Studi S1 Informatika



disusun oleh

**TRIAS HANDAYANI**

**19.11.2733**

Kepada

**FAKULTAS ILMU KOMPUTER  
UNIVERSITAS AMIKOM YOGYAKARTA  
YOGYAKARTA  
2023**

**TEKNIK BALANCING DATASET MENGGUNAKAN  
ALGORITMA ADASYN DAN BORDERLINE-SMOTE**

**SKRIPSI**

untuk memenuhi salah satu syarat mencapai derajat Sarjana  
Program Studi S1 Informatika



disusun oleh

**TRIAS HANDAYANI**

**19.11.2733**

Kepada

**FAKULTAS ILMU KOMPUTER  
UNIVERSITAS AMIKOM YOGYAKARTA  
YOGYAKARTA**

**2023**

**HALAMAN PERSETUJUAN**

**SKRIPSI**

**TEKNIK BALANCING DATASET MENGGUNAKAN ALGORITMA  
ADASYN DAN BORDERLINE-SMOTE**

yang disusun dan diajukan oleh

**Trias Handayani**

**19.11.2733**

telah disetujui oleh Dosen Pembimbing Skripsi  
pada tanggal 21 Februari 2023

**Dosen Pembimbing,**



**Mardhiya Hayaty, S.T., M.Kom**  
**NIK. 190302108**

**HALAMAN PENGESAHAN**  
**SKRIPSI**  
**TEKNIK BALANCING DATASET MENGGUNAKAN ALGORITMA**  
**ADASYN DAN BORDERLINE-SMOTE**

yang disusun dan diajukan oleh

**Trias Handayani**

**19.11.2733**

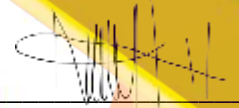
Telah dipertahankan di depan Dewan Penguji  
pada tanggal 21 Februari 2023

**Susunan Dewan Penguji**

**Nama Penguji**

**Tanda Tangan**

**Norhikmah, M.Kom**  
**NIK. 190302245**



**Erni Seniwati, S.Kom, M.Cs**  
**NIK. 190302231**



**Mardhiya Hayaty, S.T., M.Kom.**  
**NIK. 190302108**



Skripsi ini telah diterima sebagai salah satu persyaratan  
untuk memperoleh gelar Sarjana Komputer  
Tanggal 21 Februari 2023

**DEKAN FAKULTAS ILMU KOMPUTER**



**Hanif Al Fatta, S.Kom., M.Kom.**  
**NIK. 190302096**

## HALAMAN PERNYATAAN KEASLIAN SKRIPSI

Yang bertandatangan di bawah ini,

Nama mahasiswa : Trias Handayani  
NIM : 19.11.2733

Menyatakan bahwa Skripsi dengan judul berikut:

**Teknik Balancing Dataset Menggunakan Algoritma ADASYN dan Borderline-SMOTE**

Dosen Pembimbing : Mardhiya Hayaty, S.T., M.Kom

1. Karya tulis ini adalah benar-benar ASLI dan BELUM PERNAH diajukan untuk mendapatkan gelar akademik, baik di Universitas AMIKOM Yogyakarta maupun di Perguruan Tinggi lainnya.
2. Karya tulis ini merupakan gagasan, rumusan dan penelitian SAYA sendiri, tanpa bantuan pihak lain kecuali arahan dari Dosen Pembimbing.
3. Dalam karya tulis ini tidak terdapat karya atau pendapat orang lain, kecuali secara tertulis dengan jelas dicantumkan sebagai acuan dalam naskah dengan disebutkan nama pengarang dan disebutkan dalam Daftar Pustaka pada karya tulis ini.
4. Perangkat lunak yang digunakan dalam penelitian ini sepenuhnya menjadi tanggung jawab SAYA, bukan tanggung jawab Universitas AMIKOM Yogyakarta.
5. Pernyataan ini SAYA buat dengan sesungguhnya, apabila di kemudian hari terdapat penyimpangan dan ketidakbenaran dalam pernyataan ini, maka SAYA bersedia menerima SANKSI AKADEMIK dengan pencabutan gelar yang sudah diperoleh, serta sanksi lainnya sesuai dengan norma yang berlaku di Perguruan Tinggi.

Yogyakarta, 21 Februari 2023

Yang Menyatakan,

  
Trias Handayani

## HALAMAN PERSEMBAHAN

Puji syukur penulis ucapkan kehadirat Allah SWT yang melimpahkan karunia-Nya sehingga penulis mampu menyelesaikan penyusunan skripsi dengan judul “**Teknik Balancing Dataset Menggunakan Algoritma ADASYN dan Borderline-SMOTE**”. Skripsi ini penulis persembahkan kepada:

1. Orang tua dan keluarga yang selalu memberikan semangat dan doa tanpa henti.
2. Ibu Mardhiya Hayaty, S.T., M.Kom selaku dosen pembimbing yang telah membantu dalam penyusunan skripsi.
3. Teman-teman kelas 19-IF-03 yang telah menemani dari awal hingga akhir semester yang telah membantu saya dalam menyelesaikan studi.
4. Teman-teman yang selalu mendengarkan keluh kesah dan memberikan motivasi untuk menyelesaikan skripsi.

## KATA PENGANTAR

Puji syukur penulis ucapkan kehadirat Allah SWT yang melimpahkan karunia-Nya sehingga penulis mampu menyelesaikan penyusunan skripsi dengan judul **“Teknik Balancing Dataset Menggunakan Algoritma ADASYN dan Borderline-SMOTE”** sebagai syarat untuk menyelesaikan studi sarjana di Universitas Amikom Yogyakarta.

Dalam proses penyusunan skripsi ini penulis banyak mengalami kendala. Namun, berkat bantuan dari berbagai pihak skripsi ini dapat diselesaikan. Oleh karena itu, Penulis mengucapkan terimakasih kepada :

1. Bapak Prof. Dr. M. Suyanto, M.M. selaku rektor Universitas Amikom Yogyakarta.
2. Ibu Mardhiya Hayaty, S.T., M.Kom selaku dosen pembimbing yang telah membimbing dalam penyusunan skripsi.
3. Orang tua dan keluarga yang telah mendoakan dan memberikan motivasi.
4. Seluruh dosen penguji yang telah memberikan saran sehingga skripsi ini menjadi lebih baik.

Penulis berharap, skripsi ini dapat memberikan manfaat bagi para pembaca. Penulis menyadari bahwa skripsi ini terdapat banyak kekurangan sehingga saran dan kritik dibutuhkan untuk menyempurnakan skripsi ini.

Yogyakarta, 21 Februari 2023

Penulis

## DAFTAR ISI

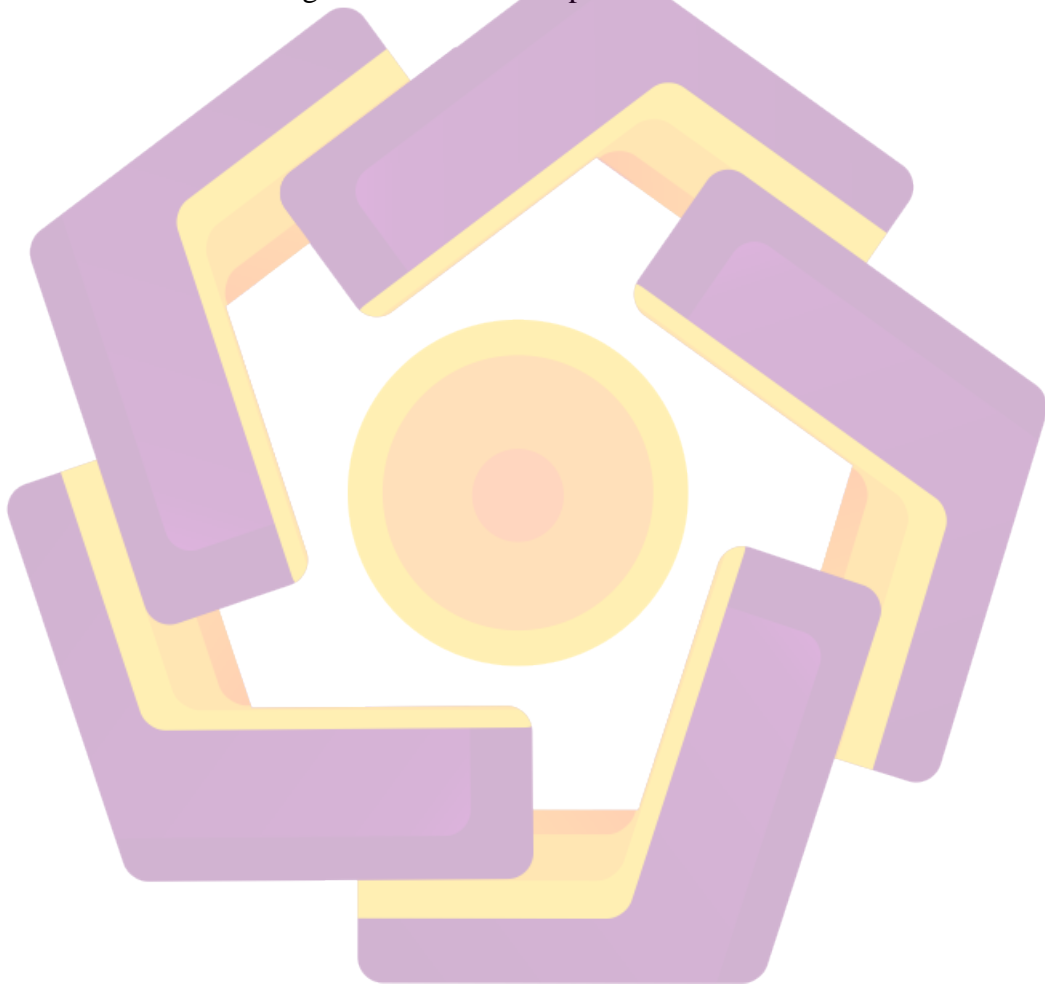
HALAMAN JUDUL .....	i
HALAMAN PERSETUJUAN.....	ii
HALAMAN PENGESAHAN .....	iii
HALAMAN PERNYATAAN KEASLIAN SKRIPSI .....	iv
HALAMAN PERSEMBAHAN .....	v
KATA PENGANTAR .....	vi
DAFTAR ISI.....	vii
DAFTAR TABEL.....	ix
DAFTAR GAMBAR.....	x
DAFTAR ISTILAH .....	xi
INTISARI .....	xii
ABSTRACT.....	xiii
BAB I PENDAHULUAN.....	1
1.1 Latar Belakang .....	1
1.2 Rumusan Masalah.....	2
1.3 Batasan Masalah .....	2
1.4 Tujuan Penelitian .....	3
1.5 Manfaat Penelitian .....	3
1.6 Sistematika Penulisan .....	3
BAB II TINJAUAN PUSTAKA .....	5
2.1 Studi Literatur .....	5
2.2 Dasar Teori .....	10
2.2.1 Klasifikasi .....	10
2.2.2 Imbalance Dataset.....	11
2.2.3 Preprocessing .....	11
2.2.4 Teknik Oversampling.....	12
2.2.5 Adaptive Synthetic (ADASYN) .....	12
2.2.6 Borderline-SMOTE.....	14
2.2.7 Random Forest .....	16
2.2.8 Stratified K-Fold Cross Validation .....	17



2.2.9	Evaluasi.....	18
<b>BAB III METODE PENELITIAN .....</b>		<b>21</b>
3.1	Alur Penelitian .....	21
3.2	Preprocessing .....	22
3.3	Balancing Dataset .....	22
3.4	Klasifikasi .....	22
3.5	Cross Validasion.....	23
3.6	Evaluasi.....	23
3.7	Alat dan Bahan.....	23
3.7.1	Alat.....	23
3.7.2	Bahan .....	24
<b>BAB IV HASIL DAN PEMBAHASAN .....</b>		<b>27</b>
4.1	Preproccesing.....	27
4.2	Balancing Dataset.....	29
4.2.1	Balancing ADASYN.....	29
4.4.2	Borderline-SMOTE.....	31
4.3	Klasifikasi dan Evaluasi.....	32
<b>BAB V PENUTUP .....</b>		<b>35</b>
5.1	Kesimpulan.....	35
5.2	Saran .....	35
<b>REFERENSI .....</b>		<b>36</b>

## DAFTAR TABEL

Tabel 2.1 Keaslian Penelitian .....	8
Tabel 2.2 Kategori Nilai AUC .....	20
Tabel 3.1 Atribut Dataset .....	24
Tabel 3.2 Dataset .....	26
Tabel 4.1 Distribusi Data Sebelum Balancing VS ADASYN .....	30
Tabel 4. 2 Distribusi Data Sebelum Balancing VS Borderline-SMOTE.....	31
Tabel 4. 3 Perbandingan Nilai AUC Setiap Skenario.....	32

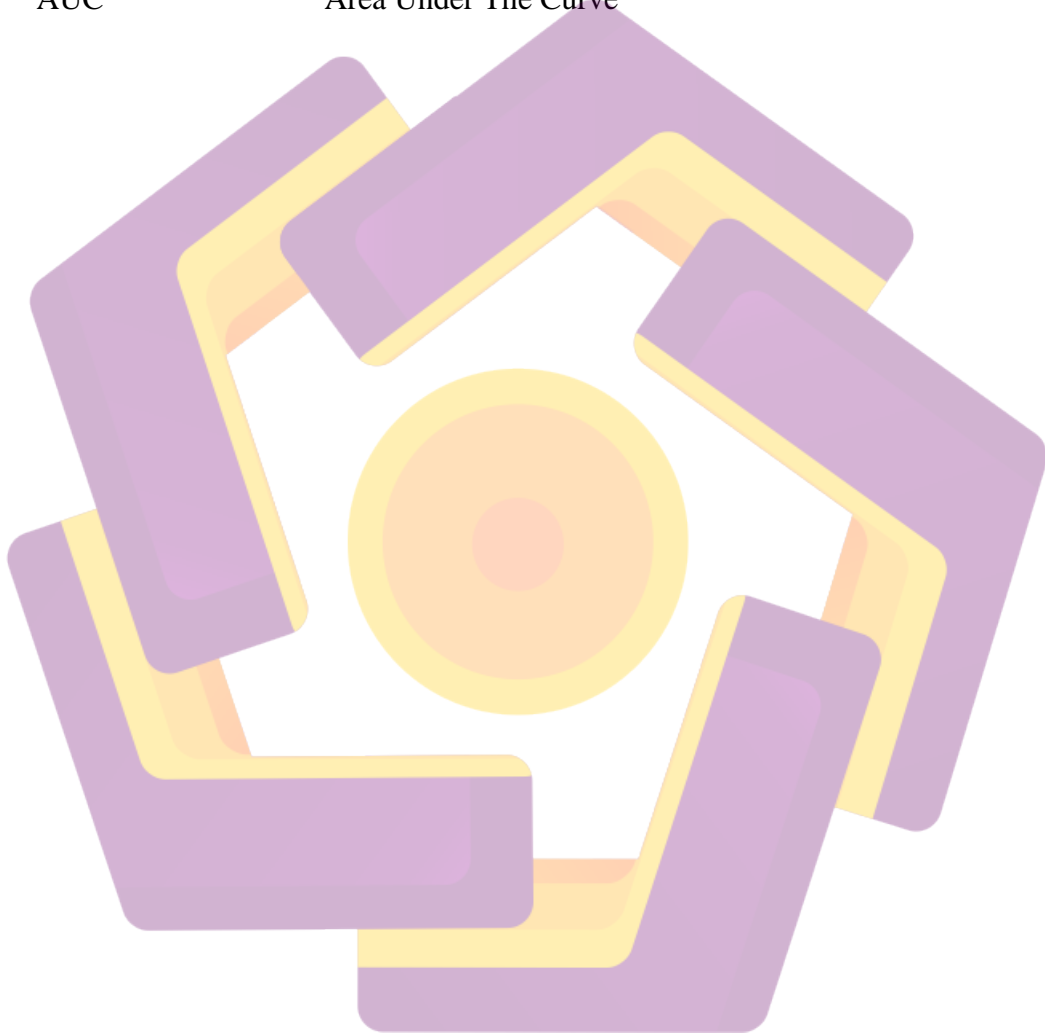


## DAFTAR GAMBAR

Gambar 2.1 Proses Klasifikasi.....	10
Gambar 2. 2 Dataset sebelum dan sesudah penerapan ADASYN.....	12
Gambar 2. 3 Ilustrasi kategori Borderline-SMOTE.....	15
Gambar 2. 4 Ilustrasi Random Forest. ....	17
Gambar 2. 5 Ilustrasi Stratified K-Fold(Scikit-learn.org, 2023).....	18
Gambar 2. 6 Kurva ROC Berbagai Nilai. ....	19
Gambar 4. 1 Data Cleaning Pada Data duplikat. ....	27
Gambar 4. 2 Atribut Bertipe Kategori. ....	28
Gambar 4. 3 Hasil Transformasi Data Kategori ke Numerik. ....	28
Gambar 4. 4 Countplot Perbandingan Data Antar Kelas.....	29
Gambar 4. 5 Distribusi Data Sebelum dan Sesudah ADASYN.....	30
Gambar 4. 6 Distribusi Sebelum dan Sesudah Borderline-SMOTE.....	31
Gambar 4. 7 Distribusi ADASYN VS Borderline-SMOTE.....	32
Gambar 4. 8 ROC Curve Sebelum Balancing.....	33
Gambar 4. 9 ROC Curve ADASYN.....	33
Gambar 4. 10 ROC Curve Borderline-SMOTE.....	34

## DAFTAR ISTILAH

ADASYN	Adaptive Synthetic
IR	Imbalance Ratio
ROC Curve	Receiver Operating Characteristic Curve
AUC	Area Under The Curve



## INTISARI

Ketidakseimbangan dataset merupakan permasalahan umum yang sering ditemukan pada klasifikasi. Ketidakseimbangan dataset merupakan kondisi dimana sebuah dataset memiliki perbedaan jumlah data antara kelas satu dengan lainnya.

Masalah ketidakseimbangan ini mengakibatkan salah klasifikasi yang dapat menurunkan kinerja dari model klasifikasi. Penurunan kinerja model terjadi karena data lebih mengenali data kelas mayoritas yang lebih dominan sehingga kelas minoritas sering salah diklasifikasikan. Untuk mengatasi permasalahan klasifikasi ini, maka dataset diseimbangkan dengan meningkatkan data di kelas minoritas. Pada penelitian ini, ADASYN dan Borderline-SMOTE digunakan untuk menyeimbangkan dataset.

Dataset yang digunakan adalah dataset lung cancer yang tidak seimbang. Random Forest diterapkan sebagai algoritma klasifikasi dengan nilai  $n\_estimators=100$  dan divalidasi menggunakan 5 stratified k-fold cross validation. Hasil penelitian menunjukkan adanya peningkatan kinerja berdasarkan nilai AUC setelah dilakukan penyeimbangan baik menggunakan ADASYN maupun Borderline-SMOTE. Peningkatan kinerja berdasarkan nilai AUC yang terjadi antara 1%-13%. Pada beberapa k-fold ADASYN dan Borderline-SMOTE menghasilkan nilai AUC sebesar 100% yang menandakan bahwa model mampu mengklasifikasikan seluruh data dengan sempurna. Dari k-fold 2,4, dan 5 hasil nilai AUC menunjukkan bahwa ADASYN lebih unggul dari Borderline-SMOTE.

**Kata kunci:** Data Tidak Seimbang, Oversampling, ADASYN, Borderline SMOTE, Random Forest.

## ABSTRACT

*Dataset imbalance is a common problem that is often found in classification. Dataset imbalance is a condition where a dataset has a difference in the amount of data between one class and another.*

*This imbalance problem results in misclassification which can reduce the performance of the classification model. The decline in model performance occurs because the data better recognizes the dominant majority class data so that the minority class is often misclassified. To overcome this classification problem, the dataset is balanced by increasing the data in the minority class. In this study, ADASYN and Borderline-SMOTE are used to balance the dataset.*

*The dataset used is an unbalanced lung cancer dataset. Random Forest is applied as a classification algorithm with a value of  $n\_estimators=100$  and validated using 5 stratified k-fold cross validation. The results showed that there was an increase in performance based on the AUC value after balancing using both ADASYN and Borderline-SMOTE. Performance improvement based on the AUC value that occurs between 1% -13%. In several k-folds ADASYN and Borderline-SMOTE produce an AUC value of 100% which indicates that the model is able to classify all data perfectly. From k-fold 2, 4, and 5 the results of the AUC value show that ADASYN is superior to Borderline-SMOTE.*

**Keyword:** Imbalance Data, Oversampling, ADASYN, Borderline SMOTE, Random Forest.