

**EVALUASI DAMPAK RANDOM OVER SAMPLING
TERHADAP KINERJA INDOBERT UNTUK ANALISIS
SENTIMEN BAHASA INDONESIA**

SKRIPSI NON REGULER JALUR SCIENTIST

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



Disusun oleh :

DIMAS RAMADHAN ALFINSYAH

22.11.4742

**FAKULTAS ILMU KOMPUTER
UNIVERSITAS AMIKOM YOGYAKARTA
YOGYAKARTA**

2025

**EVALUASI DAMPAK RANDOM OVER SAMPLING TERHADAP
KINERJA INDOBERT UNTUK ANALISIS
SENTIMEN BAHASA INDONESIA**

SKRIPSI NON REGULER JALUR SCIENTIST

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



Disusun oleh :

DIMAS RAMADHAN ALFINSYAH

22.11.4742

**FAKULTAS ILMU KOMPUTER
UNIVERSITAS AMIKOM YOGYAKARTA
YOGYAKARTA**

2025

HALAMAN PERSETUJUAN

JALUR NON-REGULER

**EVALUASI DAMPAK RANDOM OVER SAMPLING TERHADAP
KINERJA INDOBERT UNTUK ANALISIS
SENTIMEN BAHASA INDONESIA**

yang disusun dan diajukan oleh
DIMAS RAMADHAN ALFINSYAH
22.11.4742

telah disetujui oleh Dosen Pembimbing
pada tanggal 17 Desember 2025

Dosen Pembimbing,



Bambang Pili Hartato, S.Kom., M.Eng.
NIK. 190302707

HALAMAN PENGESAHAN

JALUR NON-REGULER

**EVALUASI DAMPAK RANDOM OVER SAMPLING TERHADAP
KINERJA INDOBERT UNTUK ANALISIS
SENTIMEN BAHASA INDONESIA**

yang disusun dan diajukan oleh

DIMAS RAMADHAN ALFINSYAH

22.11.4742

Telah dipertahankan di depan Dewan Penguji

pada tanggal 17 Desember 2025

Susunan Dewan Penguji

Nama Penguji

Tanda Tangan

Uvoek Anggoro Saputro, S.Kom., M.Kom.

NIK. 190302419



Windha Mega Pradnya Dhuhita, S.Kom., M.Kom.

NIK. 190302185



Bambang Pili Hartato, S.Kom., M.Eng.

NIK. 190302707



Laporan ini telah diterima sebagai salah satu persyaratan

untuk memperoleh gelar Sarjana Komputer

Tanggal 17 Desember 2025

DEKAN FAKULTAS ILMU KOMPUTER



Prof. Dr. Kusriani, M.Kom.

NIK. 190302106

HALAMAN PERNYATAAN KEASLIAN KARYA

Yang bertandatangan di bawah ini,

Nama mahasiswa : DIMAS RAMADHAN ALFINSYAH

NIM : 22.11.4742

Menyatakan bahwa Laporan dengan judul berikut:

Evaluasi Dampak Random Over Sampling Terhadap Kinerja Indobert Untuk Analisis Sentimen Bahasa Indonesia

Dosen Pembimbing : **Bambang Pilo Hartato, S.Kom., M.Eng.**

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 kegiatan 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 ketidak-benaran 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, 17 Desember 2025

Yang Menyatakan,



Dimas Ramadhan Alfinsyah

HALAMAN PERSEMBAHAN

Puji syukur penulis panjatkan ke hadirat **Allah SWT** atas rahmat, taufik, dan hidayah-Nya sehingga laporan ini dapat terselesaikan dengan baik. Penyusunan laporan ini tidak lepas dari bantuan, dukungan, serta doa dari berbagai pihak yang senantiasa memberikan semangat dan motivasi. Dengan penuh rasa hormat, laporan ini penulis persembahkan kepada:

1. **Kedua orang tua tercinta**, Bapak Marjadi dan Ibu Ika Kumarasari , serta kakak Desta, yang selalu memberikan doa, kasih sayang, dan dukungan tanpa henti. Terima kasih atas segala pengorbanan dan kerja keras yang telah diberikan agar penulis dapat menempuh pendidikan hingga tahap ini.
2. **Bapak/Ibu dosen pembimbing dan penguji**, yang telah membimbing dan memberikan arahan dengan penuh kesabaran serta memberikan masukan berharga demi kesempurnaan laporan ini.
3. **Teman-teman seperjuangan** di kelas S1-Informatika 03, yang selalu memberikan semangat, kebersamaan, dan dukungan selama masa perkuliahan hingga penyusunan laporan ini. Terima kasih atas tawa, bantuan, dan kebersamaan yang telah menjadi bagian dari perjalanan ini.
4. **Semua pihak** yang telah membantu secara langsung maupun tidak langsung. Terima kasih atas doa dan dukungannya hingga laporan ini dapat terselesaikan dengan baik.

Akhir kata, semoga laporan ini dapat memberikan manfaat, menambah wawasan, serta menjadi kontribusi kecil bagi pengembangan ilmu pengetahuan di lingkungan akademik AMIKOM.

KATA PENGANTAR

Puji dan syukur penulis panjatkan ke hadirat **Tuhan Yang Maha Esa** atas segala rahmat, karunia, dan petunjuk-Nya sehingga penulis dapat menyelesaikan **Laporan Skripsi Non Reguler Jalur Scientist** yang berjudul “Evaluasi Dampak Random Over Sampling Terhadap Kinerja Indobert Untuk Analisis Sentimen Bahasa Indonesia”.

Penyusunan laporan ini merupakan salah satu syarat untuk menyelesaikan program studi **S1 Informatika** di **Universitas Amikom Yogyakarta**. Laporan ini disusun sebagai hasil dari penelitian dan pembelajaran yang telah penulis jalani selama masa perkuliahan, khususnya dalam bidang keilmuan yang berkaitan dengan topik penelitian ini.

Dalam proses penyusunan laporan ini, penulis banyak menerima bantuan, dukungan, dan bimbingan dari berbagai pihak. Oleh karena itu, pada kesempatan ini penulis ingin menyampaikan ucapan terima kasih yang sebesar-besarnya kepada:

1. **Allah SWT**, atas limpahan rahmat dan karunia-Nya sehingga laporan ini dapat terselesaikan dengan baik.
2. **Bapak Prof. Dr. M. Suyanto, M.M.**, selaku Rektor Universitas Amikom Yogyakarta.
3. **Ibu Eli Pujastuti, M.Kom.**, selaku Ketua Program Studi S1 Informatika Universitas Amikom Yogyakarta.
4. **Bapak Bambang Pilu Hartato, S.Kom., M.Eng.**, selaku dosen pembimbing yang telah memberikan arahan, bimbingan, serta motivasi selama proses penyusunan laporan ini.
5. **Seluruh dosen dan staf Universitas Amikom Yogyakarta**, yang telah memberikan ilmu dan bantuan administrasi selama penulis menempuh pendidikan.

6. **Kedua orang tua dan keluarga tercinta**, yang selalu memberikan doa, dukungan moral, dan semangat tanpa henti.
7. **Rekan-rekan dan sahabat seperjuangan**, yang senantiasa memberikan semangat, bantuan, serta kebersamaan selama proses penyusunan laporan ini.

Penulis menyadari bahwa laporan ini masih jauh dari sempurna karena keterbatasan pengetahuan dan pengalaman. Oleh karena itu, penulis sangat mengharapkan kritik dan saran yang membangun demi perbaikan di masa mendatang.

Akhir kata, semoga laporan ini dapat memberikan manfaat bagi penulis sendiri maupun bagi pembaca yang tertarik dengan topik yang dibahas.

Yogyakarta, 6 Desember 2025

Penulis

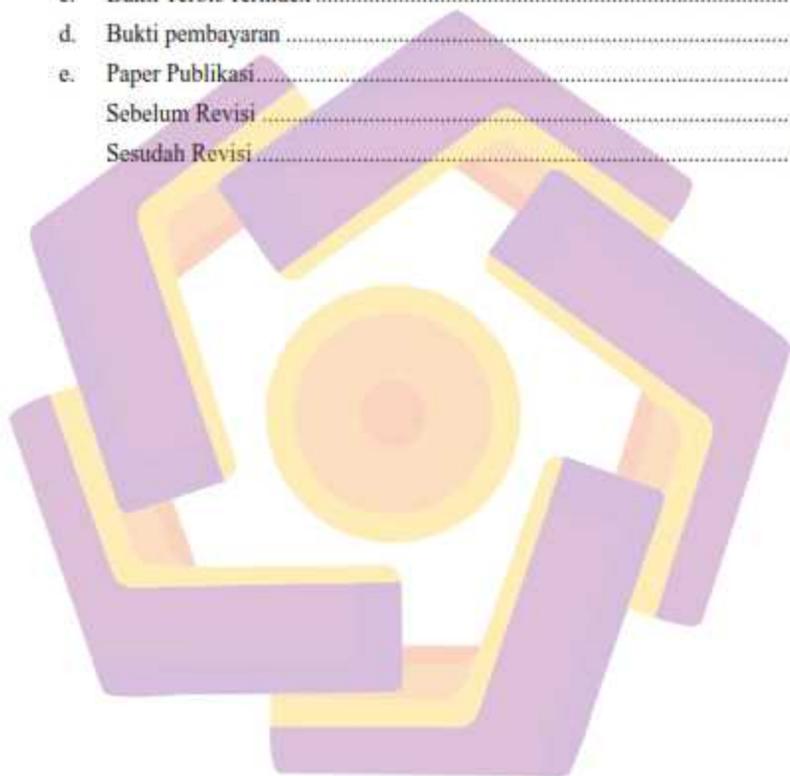


DAFTAR ISI

Halaman Judul.....	i
Halaman Persetujuan.....	ii
Halaman Pengesahan.....	iii
Halaman Pernyataan Keaslian Karya.....	iv
Halaman Persembahan.....	v
Kata Pengantar.....	vi
Daftar Isi.....	viii
Daftar Tabel.....	xi
Daftar Gambar.....	xii
Daftar Lampiran.....	xiii
Daftar Lambang dan Singkatan.....	xiv
Daftar Istilah.....	xvi
Intisari.....	xix
<i>Abstract</i>	xx
Bab I Pendahuluan.....	1
1.1 Gambaran Umum.....	1
1.2 Rumusan Masalah.....	5
1.3 Batasan Masalah.....	6
1.4 Tujuan Penelitian.....	7
Bab II Tinjauan Pustaka.....	8
2.1 Studi Literatur.....	8
2.2 Landasan Teori.....	11
2.2.1 Natural Language Processing (NLP).....	11
2.2.2 IndoBERT.....	12
2.2.3 Imbalance Data pada Klasifikasi Sentimen.....	12
2.2.4 Random Over Sampler.....	13
2.2.5 Preprocessing Data untuk Analisis Sentimen.....	14
2.2.6 Evaluasi Model.....	15
BAB III Metode Penelitian.....	16

3.1	Studi Literatur	17
3.2	Merumuskan Masalah	17
3.3	Merencanakan Eksperimen	18
3.4	Pengumpulan Data	18
3.5	Labeling.....	19
3.6	Preprocessing Data.....	20
3.6.1	Punctuation Removal dan Case Folding	20
3.6.2	Stopword Removal	21
3.6.3	Tokenizing	22
3.6.4	Normalisasi dengan Kamus Colloquial Indonesian Lexicon	22
3.6.5	Stemming	23
3.6.6	Mengembalikan Token menjadi kalimat utuh	24
3.7	Percabangan Alur Eksperimen	25
3.7.1.	Tanpa Random Over Sampler (ROS)	25
3.7.2.	Dengan Random Over Sampler (ROS)	26
3.8	Persiapan Data untuk Training	27
3.9	Modeling	28
3.10	Evaluasi dan Hasil	29
3.11	Analisis Hasil	31
3.12	Dokumentasi	31
3.13	Publikasi	32
BAB IV	Pembahasan	33
4.1	Performa Training dan Validation Model.....	33
4.2	Analisis Grafik Loss.....	34
4.3	Analisis Grafik Accuracy	35
4.4	Analisis Hasil Confusion Matrix.....	36
4.5	Hasil Performa IndoBERT	38
4.6	Pengujian Model dengan Dataset Baru	39
4.7	Publikasi.....	40
BAB V	Kesimpulan	42
5.1	Kesimpulan	42
5.2	Saran.....	42

Referensi	43
Curriculum Vitae.....	47
Lampiran dan Bukti Pendukung.....	48
a. Letter of Acceptance (LOA).....	48
b. Lembar Review	49
c. Bukti Terbit/Terindex	59
d. Bukti pembayaran	60
e. Paper Publikasi.....	61
Sebelum Revisi	61
Sesudah Revisi.....	73

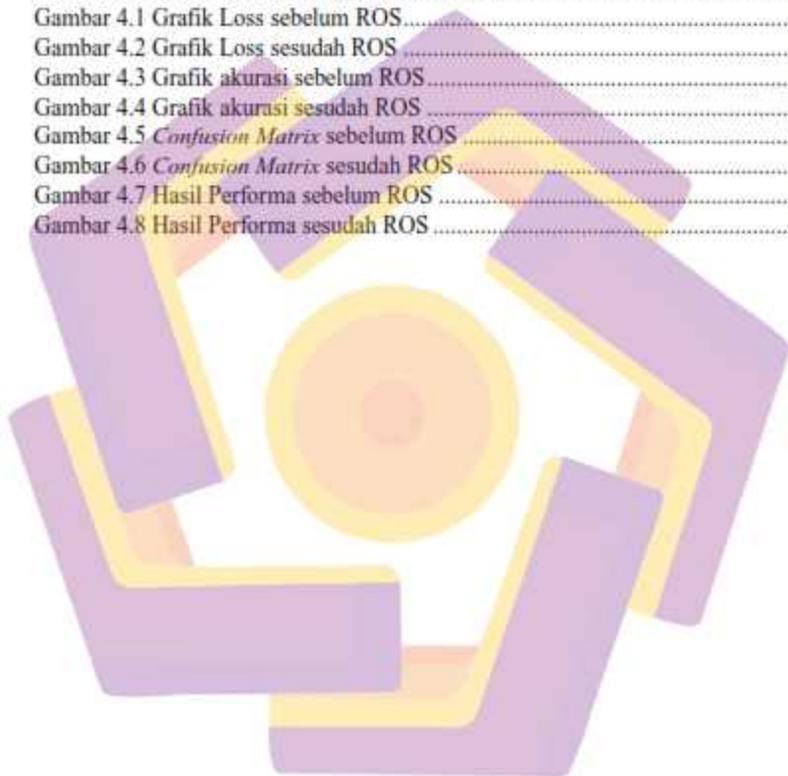


DAFTAR TABEL

Tabel 1.1 Distribusi Ulasan 15 Aplikasi pada Play Store.....	1
Tabel 3.1 Hasil Scraping Data.....	19
Tabel 3.2 Hasil Labeling	20
Tabel 3.3 Punctuation Removal dan Case Folding	21
Tabel 3.4 Stopword Removal.....	22
Tabel 3.5 Tokenizing	22
Tabel 3.6 Hasil Normalisasi	23
Tabel 3.7 Hasil Stemming	24
Tabel 3.8 Hasil mengembalikan Token menjadi kalimat utuh	24
Tabel 3.9 Hyperparameter.....	29
Tabel 4.1 Hasil Training dan Validation sebelum ROS.....	33
Tabel 4.2 Hasil Training dan Validation sesudah ROS	33
Tabel 4.3 Perbandingan Sebelum dan Sesudah ROS	39
Tabel 4.4 Pengujian Dengan Dataset Baru.....	40

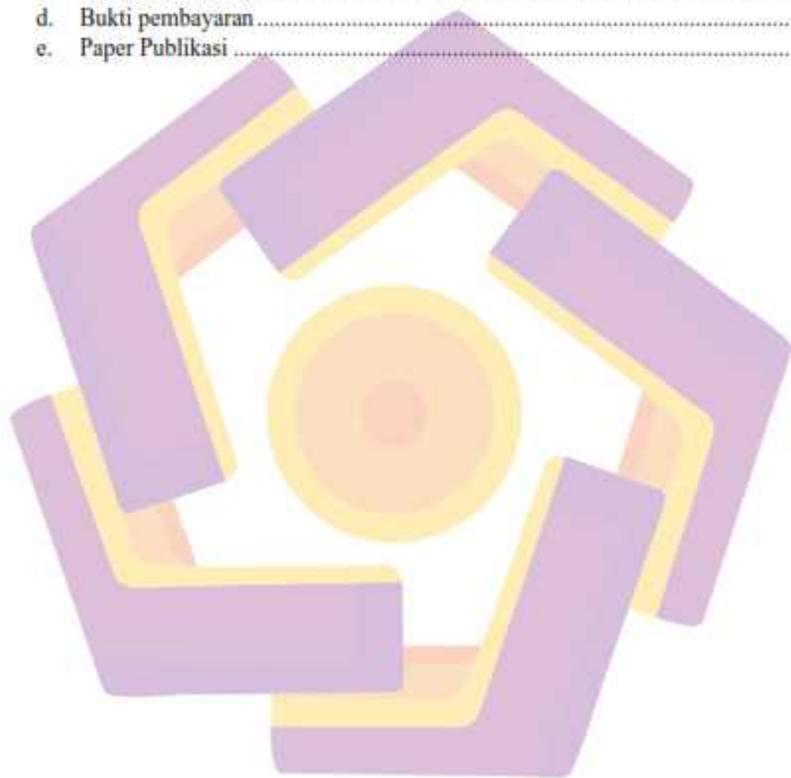
DAFTAR GAMBAR

Gambar 3.1 Alur Eksperimen.....	16
Gambar 3.2 Hasil Persebaran Label sebelum ROS.....	26
Gambar 3.3 Hasil Persebaran Label sesudah ROS.....	27
Gambar 4.1 Grafik Loss sebelum ROS.....	34
Gambar 4.2 Grafik Loss sesudah ROS.....	35
Gambar 4.3 Grafik akurasi sebelum ROS.....	35
Gambar 4.4 Grafik akurasi sesudah ROS.....	36
Gambar 4.5 <i>Confusion Matrix</i> sebelum ROS.....	37
Gambar 4.6 <i>Confusion Matrix</i> sesudah ROS.....	38
Gambar 4.7 Hasil Performa sebelum ROS.....	38
Gambar 4.8 Hasil Performa sesudah ROS.....	39



DAFTAR LAMPIRAN

a. Letter of Acceptance (LOA)	48
b. Lembar Review	49
c. Bukti Terbit/Terindex	59
d. Bukti pembayaran	60
e. Paper Publikasi	61

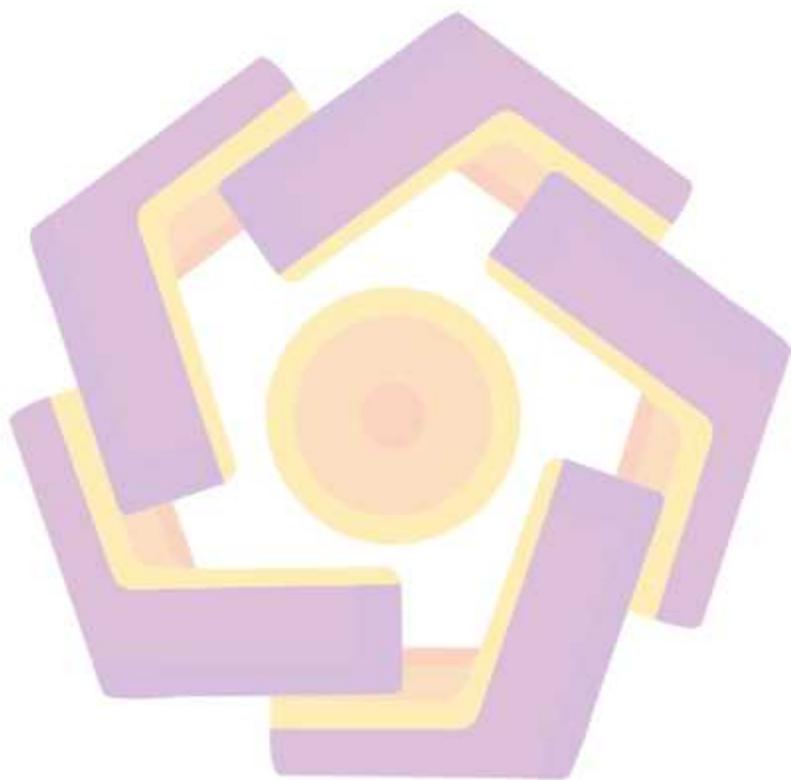


DAFTAR LAMBANG DAN SINGKATAN

ACC	Accuracy
ADASYN	Adaptive Synthetic Sampling
AdamW	Adaptive Moment Estimation with Weight Decay
API	Application Programming Interface
BERT	Bidirectional Encoder Representations from Transformers
BERTForSequenceClassification	Arsitektur model klasifikasi berbasis BERT
CLS	Classification token (token awal pada model BERT)
COLAB	Google Colaboratory
CSV	Comma Separated Values
F1	F1-Score
FN	False Negative
FP	False Positive
GPU	Graphics Processing Unit
ID	Identifier
IDN App	Aplikasi berita digital "IDN" (konteks penelitian)
IndoBERT	Indonesian Bidirectional Encoder Representations from Transformers
LDA	Latent Dirichlet Allocation
MBTI	Myers-Briggs Type Indicator
NLP	Natural Language Processing
GPU NVIDIA T4	Jenis Graphics Processing Unit dari NVIDIA tipe T4
PREC	Precision
REC	Recall
ROS	Random Over Sampler / Random Over Sampling
RUS	Random Under Sampling
SEP	Separator token (penanda akhir urutan pada model BERT)
SMOTE	Synthetic Minority Over-sampling Technique
TN	True Negative

TP

True Positive

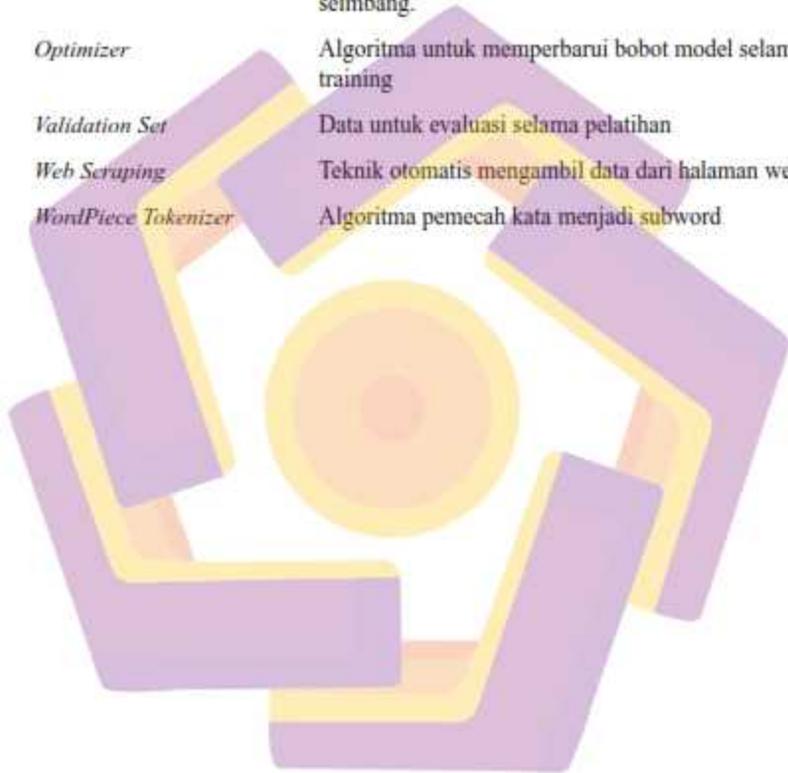


DAFTAR ISTILAH

<i>Accuracy</i>	Persentase prediksi benar
Analisis Sentimen	Proses menentukan sentimen positif atau negatif dalam teks
<i>Attention Mechanism</i>	Mekanisme fokus model pada token penting
<i>Attention Mask</i>	Penanda token mana yang diproses oleh model
<i>Binary Classification</i>	Klasifikasi dua kelas
<i>Case Folding</i>	Mengubah huruf menjadi bentuk kecil (seragam)
<i>Class Weighting</i>	Teknik pembobotan agar kelas minoritas berpengaruh lebih besar saat pelatihan
<i>Colloquial Lexicon</i>	Kamus untuk normalisasi kata tidak baku ke bentuk formal
<i>Confusion Matrix</i>	Tabel evaluasi hasil klasifikasi
<i>Cost Sensitive Learning</i>	Metode yang memberi penalti lebih besar pada kesalahan di kelas minoritas.
<i>Cross Entropy Loss</i>	Mengukur perbedaan antara distribusi prediksi model dan label sebenarnya
<i>Data Augmentation</i>	Penambahan variasi data
<i>DataFrame</i>	Struktur data tabel di Pandas
<i>DataLoader</i>	Pemuat data dalam batch pada PyTorch
<i>Dataset</i>	Kumpulan data untuk pelatihan model
<i>Early Stopping</i>	Menghentikan pelatihan sebelum overfitting
<i>Epoch</i>	Satu siklus penuh pelatihan model pada seluruh dataset
<i>F1-Score</i>	Rata-rata harmonik precision dan recall
<i>Fine-tuning</i>	Penyesuaian model pre-trained
<i>Inference</i>	Proses penggunaan model yang telah dilatih untuk memprediksi data baru
Ketidakseimbangan Data	Jumlah data antar kelas tidak seimbang

<i>Labeling</i>	Pemberian label kelas pada data
<i>Learning Curve</i>	Grafik yang menunjukkan perkembangan performa model selama pelatihan
<i>Learning Rate</i>	Kecepatan perubahan bobot model saat pelatihan
<i>Library</i>	Kumpulan fungsi atau modul siap pakai dalam pemrograman
<i>Noise</i>	Data acak atau tidak relevan yang dapat mengganggu proses pelatihan model
Normalisasi	Mengubah kata tidak baku menjadi baku
<i>Overfitting</i>	Model terlalu cocok pada data latih
<i>Over sampling</i>	Menambah sampel pada kelas minoritas untuk menyeimbangkan data.
<i>Padding</i>	Menambahkan token kosong agar panjang urutan seragam
<i>Pipeline</i>	Rangkaian proses atau alur kerja otomatis dalam sistem pemrosesan data
<i>Precision</i>	Ketepatan prediksi positif
<i>Preprocessing</i>	Tahap pembersihan dan persiapan data
<i>Pseudocode</i>	Representasi langkah algoritma dalam bentuk mirip kode
<i>Punctuation Removal</i>	Menghapus tanda baca yang tidak perlu
<i>Recall</i>	Kemampuan menangkap data positif
<i>Regularisasi</i>	Pencegahan overfitting pada model
<i>Resampling</i>	Menyeimbangkan jumlah data antar kelas
<i>Sastrawi Stemmer</i>	Algoritma stemming untuk Bahasa Indonesia
<i>Scheduler</i>	Pengatur perubahan learning rate selama training
<i>Stemming</i>	Mengubah kata berimbuhan ke bentuk dasar
<i>Stopword</i>	Kata umum yang diabaikan dalam analisis
<i>Tensor</i>	Representasi data multidimensi
<i>Test Set</i>	Data untuk pengujian akhir model

<i>Tokenisasi</i>	Memecah teks menjadi kata atau token
<i>Training Set</i>	Data yang digunakan untuk melatih model
<i>Transformer</i>	Arsitektur model berbasis attention
<i>Truncating</i>	Memotong urutan teks yang terlalu panjang
<i>Under sampling</i>	Mengurangi sampel pada kelas mayoritas agar data seimbang.
<i>Optimizer</i>	Algoritma untuk memperbarui bobot model selama training
<i>Validation Set</i>	Data untuk evaluasi selama pelatihan
<i>Web Scraping</i>	Teknik otomatis mengambil data dari halaman web
<i>WordPiece Tokenizer</i>	Algoritma pemecah kata menjadi subword



INTISARI

Analisis sentimen merupakan salah satu bidang penelitian penting dalam *Natural Language Processing (NLP)*. Untuk bahasa Indonesia, *IndoBERT* telah muncul sebagai model unggulan berkat kinerjanya yang kompetitif. Namun, efektivitasnya sangat dipengaruhi oleh distribusi kelas yang seimbang. Tantangan umum muncul karena ulasan pengguna pada platform digital, seperti *Google Play Store*, sering kali menunjukkan ketidakseimbangan kelas. Penelitian ini mengkaji efektivitas teknik *Random Over Sampler (ROS)* dalam meningkatkan kinerja *IndoBERT* pada kondisi data yang tidak seimbang. Dataset yang digunakan terdiri dari 13.821 ulasan pengguna aplikasi *IDN App* yang dikumpulkan dari *Google Play Store* sejak tahun 2015 hingga Juli 2025. Sebelum proses pemodelan, dilakukan tahap pra-proses data yang meliputi penghapusan tanda baca, *case folding*, penghapusan *stopword*, *tokenizing*, normalisasi, dan *stemming* untuk memastikan konsistensi teks. Ulasan kemudian dikategorikan ke dalam dua kelas sentimen, yaitu positif (rating 3–5 bintang) dan negatif (rating 1–2 bintang). Dua skenario eksperimen dilakukan, yaitu (1) *IndoBERT* tanpa *ROS* dan (2) *IndoBERT* dengan dataset seimbang menggunakan *ROS*. Kinerja model dievaluasi menggunakan metrik *accuracy*, *precision*, *recall*, dan *F1-score*, dengan pembagian data sebesar 70% untuk pelatihan, 20% untuk validasi, dan 10% untuk pengujian. Hasil penelitian menunjukkan adanya peningkatan signifikan setelah penerapan *ROS*, yaitu *accuracy* sebesar 94,55%, *precision* 94,61%, *recall* 94,53%, dan *F1-score* 94,54%. Analisis *confusion matrix* menunjukkan peningkatan dalam klasifikasi kelas minoritas dengan penurunan tingkat kesalahan hingga 49%. Namun, analisis *learning curve* mengindikasikan adanya potensi *overfitting* akibat penggunaan *ROS*. Oleh karena itu, penelitian lanjutan diperlukan untuk mengoptimalkan strategi *ROS* agar dapat mencapai kinerja dan kemampuan generalisasi yang lebih baik.

Kata kunci: Analisis Sentimen, *IndoBERT*, *Random Over Sampler*, Data Tidak Seimbang, Evaluasi Model.

ABSTRACT

Sentiment analysis is a prominent research area in natural language processing (NLP). For the Indonesian language, IndoBERT has emerged as a leading model due to its competitive performance. However, its effectiveness is strongly influenced by balanced class distribution. A common challenge arises because user reviews on digital platforms, such as the Google Play Store, often exhibit imbalanced classes. This study investigates the effectiveness of the Random Over Sampler (ROS) technique in improving IndoBERT's performance under imbalanced data conditions. The dataset consists of 13,821 user reviews of the IDN App collected from the Google Play Store between 2015 and July 2025. Prior to modeling, data preprocessing was performed, including punctuation removal, case folding, stopword removal, tokenizing, normalization, and stemming to ensure textual consistency. Reviews were categorized into two sentiment classes: positive (3–5 stars) and negative (1–2 stars). Two experimental scenarios were conducted: (1) IndoBERT without ROS and (2) IndoBERT with a balanced dataset using ROS. Model performance was evaluated using accuracy, precision, recall, and F1-score, with data split into 70% training, 20% validation, and 10% testing. Results showed a significant improvement after ROS implementation: 94.55% accuracy, 94.61% precision, 94.53% recall, and 94.54% F1-score. Confusion matrix analysis indicated improved classification of the minority class, reducing the error rate by 49%. However, learning curve analysis revealed potential overfitting due to ROS. Further research is needed to optimize ROS strategies for better performance and generalization.

Keywords: *Sentiment Analysis, IndoBERT, Random Over Sampler, Imbalanced Data, Model Evaluation.*