



**KLASIFIKASI SENTIMEN DATA TIDAK SEIMBANG
MENGGUNAKAN ALGORITMA SMOTE DAN *K*-NEAREST
NEIGHBOR PADA ULASAN PENGGUNA APLIKASI
PEDULILINDUNGI**

SKRIPSI

SHEILA GABRIELA BARUS

1810511072

**UNIVERSITAS PEMBANGUNAN NASIONAL “VETERAN”
JAKARTA**

FAKULTAS ILMU KOMPUTER

PROGRAM STUDI S1 INFORMATIKA

2022



**KLASIFIKASI SENTIMEN DATA TIDAK SEIMBANG
MENGGUNAKAN ALGORITMA SMOTE DAN K-NEAREST
NEIGHBOR PADA ULASAN PENGGUNA APLIKASI
PEDULILINDUNGI**

SKRIPSI

**Diajukan Sebagai Salah Satu Syarat Untuk Memperoleh
Gelar Sarjana Komputer**

SHEILA GABRIELA BARUS

1810511072

**UNIVERSITAS PEMBANGUNAN NASIONAL “VETERAN”
JAKARTA
FAKULTAS ILMU KOMPUTER
PROGRAM STUDI S1 INFORMATIKA**

2022

PERNYATAAN ORISINALITAS

Tugas skripsi ini adalah hasil karya sendiri, dan semua sumber yang dikutip maupun yang dirujuk telah saya nyatakan dengan benar :

Nama : Sheila Gabriela Barus

NIM : 1810511072

Tanggal : 24 Juni 2022

Bilamana dikemudian hari ditemukan ketidaksesuaian dengan pernyataan saya ini, maka saya bersedia dituntut dan diproses sesuai dengan ketentuan yang berlaku.

Jakarta, 24 Juni 2022

Yang Menyatakan,



(Sheila Gabriela Barus)

PERNYATAAN PERSETUJUAN PUBLIKASI

Sebagai civitas akademik Universitas Pembangunan Nasional Veteran Jakarta, saya yang bertandatangan di bawah ini :

Nama : Sheila Gabriela Barus
NIM : 1810511072
Fakultas : Ilmu Komputer
Program Studi : S1 Informatika

Demi pengembangan ilmu pengetahuan, menyetujui untuk memberikan kepada Universitas Pembangunan Nasional Veteran Jakarta Hak Bebas Royalti Non Eksklusif (*Non-Exclusive Royalty Free Right*) atas karya ilmiah saya yang berjudul:

KLASIFIKASI SENTIMEN DATA TIDAK SEIMBANG MENGGUNAKAN ALGORITMA SMOTE DAN K-NEAREST NEIGHBOR PADA ULASAN PENGGUNA APLIKASI PEDULILINDUNGI

Beserta perangkat yang ada (jika diperlukan). Dengan Hak Bebas Royalti ini Universitas Pembangunan Nasional Veteran Jakarta berhak menyimpan, mengalih, media/formatkan, mengelola dalam bentuk pangkalan data (*database*), merawat dan mempublikasikan Skripsi saya selama tetap mencantumkan nama saya sebagai penulis dan sebagai pemilik Hak Cipta. Demikian pernyataan ini saya buat dengan sebenarnya.

Dibuat Di : Jakarta
Pada tanggal : 24 Juni 2022
Yang Menyatakan,



(Sheila Gabriela Barus)

LEMBAR PENGESAHAN

Dengan ini menyatakan bahwa Proposal Skripsi berikut :

Nama : Sheila Gabriela Barus
NIM : 1810511072
Program Studi : S1 Informatika
Judul : Klasifikasi Sentimen Data Tidak Seimbang Menggunakan Algoritma SMOTE dan *K-Nearest Neighbor* pada Ulasan Pengguna Aplikasi PeduliLindungi

Telah berhasil dipertahankan dihadapan Tim Penguji dan diterima sebagai bagian persyaratan yang diperlukan untuk memperoleh gelar Sarjana Komputer di Program Studi S1 Informatika Fakultas Ilmu Komputer Universitas Pembangunan Nasional Veteran Jakarta.

Henki Bayu Setia, S.Kom, MTI

Desta Sandya Prasvita, S.Kom., M.Kom.

Penguji I

Penguji II

Dr. Didit Widjianto, S.Kom, M.Si

Pembimbing I

Mayanda Mega Santoni, S.Kom., M.Kom.

Pembimbing II



Dr. Ermawita, M.Kom.

Dekan

Desta Sandya Prasvita, S.Kom., M.Kom.

Ketua Program Studi

Ditetapkan : Jakarta

Tanggal Persetujuan : 24 Juni 2022



**KLASIFIKASI SENTIMEN DATA TIDAK SEIMBANG MENGGUNAKAN
ALGORITMA SMOTE DAN *K-NEAREST NEIGHBOR* PADA ULASAN
PENGGUNA APLIKASI PEDULILINDUNG**

SHEILA GABRIELA BARUS

ABSTRAK

Salah satu penanganan pemerintah dalam mengatasi penyebaran Covid-19 yang terjadi di Indonesia yaitu dengan membuat sebuah aplikasi yaitu aplikasi PeduliLindungi. Aplikasi ini berfungsi dalam melacak dan memantau penyebaran Covid-19, oleh karena itu banyak masyarakat Indonesia yang harus mempunyai aplikasi ini. Banyak juga ulasan yang diberikan pada aplikasi ini, dari komentar yang positif hingga komentar negatif. Ulasan tersebut yang menjadi data dalam penelitian ini untuk mengetahui hasil sentimen masyarakat dan menguji klasifikasi algoritma *K-Nearest Neighbor*. Pengumpulan data dilakukan dengan *scraping* di google play menggunakan bahasa pemrograman *Python*, dimana data yang diperoleh mendapatkan 750 label negatif dan 250 label positif. Sehingga data yang tidak seimbang ini harus diseimbangkan dengan teknik *undersampling* dan *oversampling* SMOTE. Oleh karena itu, penelitian ini dilakukan tiga data yang berbeda jumlah yaitu dari data yang tidak seimbang, data yang sudah di *undersampling* dan data yang sudah di *oversampling* dengan SMOTE. Hasil dari ketiga percobaan tersebut diperoleh nilai terbaik menggunakan teknik SMOTE pada K = 1 dengan nilai akurasi sebesar 0.9766, nilai presisi sebesar 0.9691, nilai *F1 score* 0.9781, nilai spesifisitas sebesar 0.9645, dan nilai sensitivitas sebesar 0.9874.

Kata Kunci : Sentimen, *K-Nearest Neighbor*, SMOTE, PeduliLindungi

***CLASSIFICATION OF UNBALANCED DATA SENTIMENTS USING
SMOTE ALGORITHM AND K-NEAREST NEIGHBOR ON USER REVIEWS
OF CARE APPLICATIONS***

SHEILA GABRIELA BARUS

ABSTRACT

One of the government's methods in dealing with the spread of Covid-19 that occurred in Indonesia is to create an application, namely the PeduliLindungi application. This application functions in tracking and monitoring the spread of Covid-19, therefore many Indonesian people must have this application. Many reviews are also given on this application, from positive comments to negative comments. These reviews are used as data in this study to determine the results of community sentiment and to test the classification of the K-Nearest Neighbor algorithm. Data collection was done by scraping on google play using the Python programming language, where the data obtained got 750 negative labels and 250 positive labels. So this unbalanced data must be balanced with SMOTE undersampling and oversampling techniques. Therefore, this study carried out three experiments, namely from unbalanced data, data that had been undersampled and data that had been oversampled with SMOTE. The results of the three experiments obtained the best value using the SMOTE technique at $K = 1$ with an accuracy value of 0.9766, a precision value of 0.9691, an F1 score of 0.9781, a specificity value of 0.9645, and a sensitivity value of 0.9874.

Keywords : Sentiment, K-Nearest Neighbor, SMOTE, PeduliLindungi

KATA PENGANTAR

Puji dan syukur penulis panjatkan kehadiran Tuhan Yang Maha Esa atas segala anugerah-Nya sehingga proposal ini dapat terselesaikan. Judul dari penelitian ini yang dilaksanakan sejak bulan November 2021 ini adalah “Klasifikasi Sentimen Data Tidak Seimbang Menggunakan Algoritma SMOTE dan *K-Nearest Neighbor* Pada Ulasan Pengguna Aplikasi PeduliLindungi” berhasil diselesaikan. Tak lupa penulis ingin mengucapkan banyak terima kasih kepada:

1. Orang tua, keluarga yang selalu memberikan dukungan kepada penulis sehingga dapat menyelesaikan skripsi ini.
2. Bapak Dr. Didit Widiyanto, S.Kom, M.Si dan Ibu Mayanda Mega Santoni, S.Kom., M.Kom. selaku dosen pembimbing yang telah memberikan saran yang bermanfaat selama proses pembuatan Proposal hingga menyelesaikan skripsi.
3. Bapak Henki Bayu Seta, S.Kom, MTI. selaku dosen pembimbing akademik.
4. Bapak Desta Sandya Prasvita, M.Kom. selaku Kaprodi Informatika yang telah memberikan informasi mengenai tugas akhir.
5. Ibu Dr. Ermatita, M.Kom. selaku dekan Fakultas Ilmu Komputer.
6. Teman-teman Informatika 2018 yang telah berjuang bersama dalam setiap proses perkuliahan serta saling memberikan semangat untuk dapat menyelesaikan Skripsi.

Dan semua pihak yang telah membantu penulisan dalam menyelesaikan skripsi ini.

Jakarta, 24 Juni 2022

Sheila Gabriela Barus

DAFTAR ISI

PERNYATAAN ORISINALITAS	iii
PERNYATAAN PERSETUJUAN PUBLIKASI	iv
ABSTRAK	vi
KATA PENGANTAR	viii
DAFTAR ISI.....	ix
DAFTAR GAMBAR	xii
DAFTAR TABEL.....	xiii
DAFTAR LAMPIRAN	xiv
BAB I PENDAHULUAN	1
1.1. Latar Belakang Masalah	1
1.2. Rumusan Masalah	2
1.3. Tujuan Penelitian.....	2
1.4. Ruang Lingkup	2
1.5. Luaran Penelitian.....	3
1.6. Manfaat Penelitian.....	3
1.7. Sistematika Penulisan.....	3
BAB II TINJAUAN PUSTAKA.....	5
2.1. Aplikasi PeduliLindungi.....	5
2.2. <i>Text Mining</i>	6
2.3. Analisis Sentimen.....	6
2.4. <i>Text Preprocessing</i>	7
2.4.1. <i>Data Cleaning</i>	7
2.4.2. <i>Case Folding</i>	8
2.4.3. <i>Spelling Normalization</i>	8
2.4.4. <i>Tokenizing</i>	8
2.4.5. <i>Filtering</i>	8
2.4.6. <i>Stemming</i>	8

2.5.	<i>Term Frequency-Inverse Document Frequency (TF-IDF)</i>	9
2.6.	<i>Machine Learning</i>	10
2.7.	Klasifikasi KNN	10
2.8.	SMOTE (<i>Synthetic Minority Oversampling Technique</i>).....	11
2.9.	Penelitian Terkait	12
	BAB III METODOLOGI PENELITIAN	15
3.1.	Kerangka Pikir.....	15
3.2.	Identifikasi dan Perumusan Masalah.....	15
3.3.	Studi Literatur.....	16
3.4.	Scraping Data	16
3.5.	Pelabelan kelas sentimen.....	16
3.6.	Praproses teks	17
3.7.	Pembobotan	18
3.8.	SMOTE (<i>Synthetic Minority Oversampling Technique</i>).....	18
3.9.	Pembagian Data.....	19
3.10.	Proses Klasifikasi.....	19
3.11.	Evaluasi.....	19
3.12.	Alat Pendukung.....	20
3.13.	Jadwal	21
	BAB IV HASIL DAN PEMBAHASAN	22
4.1.	Data	22
4.2.	Pelabelan data.....	22
4.3.	Praproses data.....	24
4.3.1.	<i>Data cleaning</i>	25
4.3.2.	<i>Case Folding</i>	26
4.3.3.	<i>Spelling normalization</i>	27
4.3.4.	<i>Tokenizing</i>	28
4.3.5.	<i>Filtering</i>	29
4.3.6.	<i>Stemming</i>	30
4.4.	Pembobotan data dengan TF IDF.....	31
4.5.	<i>Word Cloud</i> data fitur	33
4.5.1.	<i>Word Cloud</i> data positif	33

4.5.2. <i>Word Cloud</i> data negatif	34
4.6. SMOTE (<i>Synthetic Minority Oversampling Technique</i>).....	35
4.7. Klasifikasi.....	37
4.8. Evaluasi	38
4.8.1. Klasifikasi dengan data tidak seimbang.....	38
4.8.2. Klasifikasi dengan <i>Undersampling</i>	40
4.8.3. Klasifikasi dengan <i>Oversampling</i>	41
4.8.4. Perbandingan rata-rata	43
BAB V PENUTUP.....	44
5.1. Kesimpulan.....	44
5.2. Saran	45
DAFTAR PUSTAKA	46

DAFTAR GAMBAR

Gambar 2. 1. Proses <i>Text Mining</i>	6
Gambar 2. 2. Tahap Praproses Teks.....	7
Gambar 2. 3. Proses <i>Stemming</i>	9
Gambar 3. 1. Kerangka Pikir.....	15
Gambar 4. 1. Label negatif dan positif.....	24
Gambar 4. 2. Diagram alir <i>data cleaning</i>	25
Gambar 4. 3. Diagram alir <i>Case Folding</i>	26
Gambar 4. 4. Diagram alir <i>Spelling normalization</i>	27
Gambar 4. 5. Diagram alir <i>Tokenizing</i>	28
Gambar 4. 6. Diagram alir <i>Filtering</i>	29
Gambar 4. 7. Diagram alir <i>Stemming</i>	30
Gambar 4. 8. <i>Word Cloud</i> Label Positif.....	33
Gambar 4. 9. Jumlah <i>Word Cloud</i> Label Positif	34
Gambar 4. 10. <i>Word Cloud</i> Label Negatif	34
Gambar 4. 11. Jumlah <i>Word Cloud</i> Label Negatif.....	35
Gambar 4. 12. Grafik perbandingan rata-rata	43

DAFTAR TABEL

Tabel 2. 1. Matriks Penelitian Terdahulu.....	13
Tabel 3. 1. Tingkat Nilai Kappa.....	17
Tabel 3. 2. Jadwal Kegiatan Penelitian	21
Tabel 4. 1. Sampel Hasil Pengambilan Data.....	22
Tabel 4. 2. Hasil Pelabelan Data	23
Tabel 4. 3. Sebelum dan Setelah <i>Data Cleaning</i>	25
Tabel 4. 4. Sebelum dan Sesudah <i>Case Folding</i>	26
Tabel 4. 5. Sebelum dan Sesudah <i>Spelling Normalization</i>	27
Tabel 4. 6. Sebelum dan Sesudah <i>Tokenizing</i>	28
Tabel 4. 7. Sebelum dan Sesudah <i>Filtering</i>	29
Tabel 4. 8. Sebelum dan Sesudah <i>Stemming</i>	30
Tabel 4. 9. Sampel Ulasan.....	31
Tabel 4. 10. Nilai TFIDF.....	32
Tabel 4. 11. Sampel Data untuk SMOTE	36
Tabel 4. 12. Hasil Setelah <i>Oversampling</i> SMOTE	37
Tabel 4. 13. Hasil evaluasi dengan data tidak seimbang.....	39
Tabel 4. 14. Hasil evaluasi dengan <i>undersampling</i>	40
Tabel 4. 15. Hasil Evaluasi dengan <i>Oversampling</i>	42

DAFTAR LAMPIRAN

Lampiran 1 Pelabelan Data Ulasan	50
Lampiran 2 Daftar Stopword Sastrawi.....	115
Lampiran 3 Daftar Kata-Kata Normalization	123
Lampiran 4 Perhitungan Evaluasi Data Tidak Seimbang	155
Lampiran 5 Perhitungan Evaluasi Data <i>Undersampling</i>	156
Lampiran 6 Perhitungan Evaluasi Data <i>Oversampling</i>	157
Lampiran 7 Aplikasi PeduliLindungi.....	159
Lampiran 8 Visualisasi data.....	160