



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

SKRIPSI

ALESSANDRO SAMUEL TAMADA

1810511070

PROGRAM STUDI S1 INFORMATIKA

FAKULTAS ILMU KOMPUTER

UNIVERSITAS PEMBANGUNAN NASIONAL VETERAN JAKARTA

2024

PERNYATAAN ORISINALITAS

Skripsi ini adalah hasil karya sendiri, dan semua sumber yang dikutip maupun yang dirujuk telah saya nyatakan dengan benar.

Nama : Alesandro Samuel Tamada

NIM : 1810511070

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

Jakarta, 10 Juli 2024

Yang Menyatakan,



Alesandro Samuel Tamada

PERNYATAAN PERSETUJUAN PUBLIKASI

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

Nama : Alesandro Samuel Tamada
NIM : 1810511070
Fakultas : Ilmu Komputer
Program Studi : S1 Informatika
Jenis Karya : Skripsi

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 & K-Nearest Neighbor Pada Ulasan Pengguna Aplikasi Maxim

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/pencipta dan sebagai pemilik Hak Cipta.

Demikian pernyataan ini saya buat dengan sebenarnya.

Dibuat Di : Jakarta
Pada tanggal : 19 Juli 2024
Yang Menyatakan,



Alesandro Samuel Tamada

LEMBAR PENGESAHAN SKRIPSI

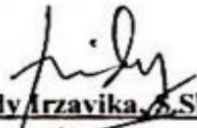
Dengan ini dinyatakan bahwa Tugas Akhir/Skripsi berikut :

Nama : Alessandro Samuel Tamada
NIM : 1810511070
Program Studi : SI Informatika
Judul Tugas Akhir : Klasifikasi Sentimen Data Tidak Seimbang Menggunakan Algoritma *SMOTE* & *K-NEAREST NEIGHBOR* Pada Ulasan Pengguna Aplikasi MAXIM

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


(Dr. Widya Cholli, M.I.T)

Penguji I


(Nindy Arzavika, S.SI, M.T.)

Penguji II


(Dr. Ermatita, M.Kom)

Pembimbing I


(Neny Rosniawati, M.Kom)

Pembimbing II


(Prof. Dr. H. Subandiyanto, ST., M.Sc., IPM)

Dekan Fakultas Ilmu Komputer


(Dr. Widya Cholli, M.I.T)

Ketua Program Studi

Disahkan di : Jakarta

Tanggal : 30 Juli 2024

**KLASIFIKASI SENTIMEN DATA TIDAK SEIMBANG
MENGUNAKAN ALGORITMA SMOTE & K-NEAREST NEIGHBOR PADA
ULASAN PENGGUNA APLIKASI MAXIM**

Alessandro Samuel Tamada

Abstrak

Maxim adalah aplikasi ojek *online* yang dipakai banyak masyarakat Indonesia, sebagai pendatang baru namun mampu bersaing dengan ojek *online* lainnya yang terlebih dahulu, dengan memberi banyak diskon dan harga yang lebih murah, pada awalnya banyak kontroversi mengenai aplikasi ini karena pelayanannya yang kurang baik, namun dari penelitian ini didapatkan bahwa banyak penilaian yang baik dari pada yang buruk. Penelitian ini mengambil data dari ulasan aplikasi Maxim dan diolah sehingga dapat diklasifikasi dengan algoritma KNN. Namun dari data yang didapatkan juga terdapat selisih yang cukup jauh antara data berlabel negatif dan data berlabel positif, karena itu harus menggunakan algoritma SMOTE agar mendapatkan label data yang seimbang. Hasil dari penelitian ini menggunakan data yang tidak seimbang dengan data positif sebesar 842 dan data negatif 158, didapatkan nilai tertinggi pada nilai $K = 2$, yaitu akurasi sebesar 0.875, presisi 0.877, *f1-score* 0.932, spesifisitas 0.111, dan sensitivitas 0.994. Sedangkan penelitian kedua dilakukan *oversampling* dengan teknik SMOTE agar data dapat seimbang yang kemudian data positif menjadi 842 dan data negatif juga 842, hasilnya tertinggi didapatkan pada nilai $K = 1$ dan 2 dengan nilai yang sama yaitu akurasi 0.985, presisi 1, *f1-score* 0.986, spesifisitas 1 dan sensitivitas 0.973. Penelitian dengan data yang sudah melalui proses SMOTE ini mendapatkan hasil yang sangat baik dapat dilihat dari nilai spesifisitas dan sensitivitas yang memiliki selisih tipis, menunjukkan bahwa data penelitian yang digunakan seimbang. Sehingga dapat disimpulkan bahwa data yang baik untuk penelitian adalah data yang seimbang.

Kata kunci : Maxim, KNN, SMOTE.

**UNBALANCED DATA SENTIMENT CLASSIFICATION USING SMOTE
& K-NEAREST NEIGHBOR ALGORITHM ON USER REVIEWS OF THE
MAXIM APPLICATION**

Alessandro Samuel Tamada

Abstract

Maxim is an online motorbike taxi application that is used by many Indonesian people, as a newcomer but able to compete with other online motorbike taxis that were earlier, by providing lots of discounts and cheaper prices, at first there was a lot of controversy regarding this application because the service was not good, however From this research it was found that there were more good assessments than bad. This research takes data from reviews of the Maxim application and processes it so that it can be classified using the KNN algorithm. However, from the data obtained there is also quite a large difference between data labeled as negative and data labeled as positive, therefore you have to use the SMOTE algorithm to get balanced data labels. The results of this study used unbalanced data with positive data of 842 and negative data of 158. The highest value was obtained at $K = 2$, namely accuracy of 0.875, precision of 0.877, f1-score of 0.932, specificity of 0.111, and sensitivity of 0.994. while the second research carried out oversampling using the SMOTE technique so that the data could be balanced, then positive data became 842 and negative data also 842, the highest results were obtained at values $K = 1$ and 2 with the same values, namely accuracy 0.985, precision 1, f1-score 0.986, specificity 1 and sensitivity 0.973. Research with data that has gone through the SMOTE process has obtained very good results which can be seen from the specificity and sensitivity values which have a slight difference, indicating that the research data used is balanced. So it can be concluded that good data for research is balanced data.

Keywords : Maxim, KNN, SMOTE.

KATA PENGANTAR

Puji dan syukur penulis panjatkan kehadirat Tuhan Yang Maha Esa atas segala anugerah-Nya, sehingga skripsi 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 Maxim” 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 Prof. Dr. Ir. Supriyanto, ST., M.Sc., IPM. selaku dekan Fakultas Ilmu Komputer.
3. Ibu Dr. Widya Cholil, M.I.T. selaku Kaprodi Informatika yang telah memberikan informasi mengenai tugas akhir.
4. Bapak Henki Bayu Seta, S.Kom, MTI. selaku dosen pembimbing akademik.
5. Ibu Dr. Ermatita, M.Kom . selaku dosen pembimbing yang telah memberikan saran yang bermanfaat selama proses pembuatan Proposal hingga menyelesaikan skripsi.
6. Ibu Neny Rosmawarni, M.Kom. selaku dosen pembimbing yang telah memberikan saran yang bermanfaat selama proses pembuatan Proposal hingga menyelesaikan skripsi.
7. 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, sehingga dapat terselesaikan dengan baik.

Jakarta, 30 Juli 2024

Penulis

DAFTAR ISI

PERNYATAAN ORISINALITAS	ii
PERNYATAAN PERSETUJUAN PUBLIKASI	iii
LEMBAR PENGESAHAN SKRIPSI	iv
ABSTRAK	v
<i>ABSTRACT</i>	vi
KATA PENGANTAR.....	vii
DAFTAR ISI	viii
DAFTAR TABEL	xi
DAFTAR GAMBAR.....	xii
DAFTAR SIMBOL.....	xiii
DAFTAR LAMPIRAN.....	xiv
BAB I.....	1
PENDAHULUAN	1
1.1 Latar Belakang.....	1
1.2 Rumusan Masalah	2
1.3 Batasan Masalah	3
1.4 Tujuan Penelitian.....	3
1.5 Manfaat Penelitian.....	3
1.6 Luaran yang Diharapkan.....	3
1.7 Sistematika Penulisan.....	4
BAB II.....	5
TINJAUAN PUSTAKA.....	5
2.1.Aplikasi Maxim	5
2.2. <i>Text Mining</i>	5
2.3.Analisis Sentimen.....	6
2.4. <i>Text Preprocessing</i>	6
2.4.1. <i>Data Cleaning</i>	6

2.4.2. <i>Case Folding</i>	6
2.4.3. <i>Spelling Normalization</i>	7
2.4.4. <i>Tokenizing</i>	7
2.4.5. <i>Filtering</i>	7
2.4.6. <i>Stemming</i>	7
2.5. <i>Term Frequency-Inverse Document Frequency (TFIDF)</i>	8
2.6. <i>Machine Learning</i>	9
2.7. Klasifikasi KNN	9
2.8. <i>SMOTE (Synthetic Minority Oversampling Technique)</i>	10
2.9. Penelitian Terkait	10
BAB III	15
METODOLOGI PENELITIAN	15
3.1. Alur Penelitian	15
3.1.1. Identifikasi dan Perumusan Masalah	16
3.1.2. Studi Literatur	16
3.1.3. Pengumpulan Data	16
3.1.4. Pelabelan kelas sentimen	17
3.1.5. Praproses teks	18
3.1.6. Pembobotan	19
3.1.7. Pembagian Data	19
3.1.8. Proses Klasifikasi	19
3.1.9. Evaluasi	19
3.1.10. ... Hasil	21
3.2. Alat Pendukung	21
3.3. Jadwal	22
BAB IV	24
HASIL DAN PEMBAHASAN	24
4.1. Pengambilan Data	24
4.2. Pelabelan data	25
4.3. Praproses data	28

4.3.1. <i>Data cleaning</i>	28
4.3.2. <i>Case Folding</i>	29
4.3.3. <i>Spelling normalization</i>	31
4.3.4. <i>Tokenizing</i>	32
4.3.5. <i>Filtering</i>	33
4.3.6. <i>Stemming</i>	34
4.4. Pembobotan data dengan TF IDF	36
4.5. <i>Word Cloud</i>	38
4.5.1. <i>Word Cloud data positif</i>	38
4.5.2. <i>Word Cloud data negatif</i>	40
4.6. SMOTE (<i>Synthetic Minority Oversampling Technique</i>)	42
4.7. Klasifikasi	42
4.8. Evaluasi	43
4.8.1. Klasifikasi dengan data tidak seimbang	43
4.8.2. Evaluasi dengan data seimbang (<i>Oversampling</i>)	45
BAB V	47
PENUTUP	47
5.1. Kesimpulan	47
5.2. Saran	48
DAFTAR PUSTAKA	49
RIWAYAT HIDUP	51
LAMPIRAN	52

\



DAFTAR TABEL

Tabel 2.1 Penelitian Terdahulu	12
Tabel 3.1 Tingkat Nilai Kappa.....	17
Tabel 3.2 Jadwal Kegiatan Penelitian	21
Tabel 4.1 Tabel Anotasi	23
Tabel 4.2 Tabel Proses <i>Data Cleaning</i>	26
Tabel 4.3 Proses <i>Case Folding</i>	28
Tabel 4.4 Proses <i>Spelling Normalization</i>	29
Tabel 4.5 Proses <i>Tokenizing</i>	31
Tabel 4.6 Proses <i>Filtering</i>	32
Tabel 4.7 Proses <i>Stemming</i>	34
Tabel 4.8 Sampel Data Komentar	34
Tabel 4.9 Hasil Perhitungan TF-IDF	36
Tabel 4.10 Hasil Perhitungan Menggunakan <i>Confusion Matrix</i>	41
Tabel 4.11 Hasil Perhitungan Menggunakan SMOTE	42

DAFTAR GAMBAR

Gambar 3.1 Alur Penelitian	15
Gambar 4.1 Contoh Data Ulasan	22
Gambar 4.2 Label Negatif Dan Positif.....	25
Gambar 4.3 Proses <i>Data Cleaning</i>	26
Gambar 4.4 <i>Code</i> Proses <i>Data Cleaning</i>	27
Gambar 4.5 Proses <i>Case Folding</i>	27
Gambar 4.6 <i>Code</i> Proses <i>Case Folding</i>	28
Gambar 4.7 Proses <i>Spelling Normalization</i>	29
Gambar 4.8 <i>Code</i> Proses <i>Spelling Normalization</i>	30
Gambar 4.9 Proses <i>Tokenizing</i>	30
Gambar 4.10 <i>Code</i> Proses <i>Tokenizing</i>	31
Gambar 4.11 Proses <i>Filtering</i>	32
Gambar 4.12 <i>Code</i> Proses <i>Filtering</i>	33
Gambar 4.13 Proses <i>Stemming</i>	33
Gambar 4.14 <i>Code</i> Proses <i>Stemming</i>	34
Gambar 4.15 <i>World Cloud</i> Data Positif	37
Gambar 4.16 Jumlah <i>World Cloud</i> Data Positif	37
Gambar 4.17 <i>World Cloud</i> Data Negatif	38
Gambar 4.18 Jumlah <i>World Cloud</i> Data Negatif	38

DAFTAR SIMBOL

No	Simbol	Nama	Keterangan
1.		Terminal (<i>start, end</i>)	Simbol ini menggambarkan proses dimulai atau proses berakhir
2.		<i>Process</i>	Simbol ini menggambarkan penjelasan dari suatu proses yang akan dijalankan
3.		<i>Flow Direction</i>	Simbol ini menggambarkan hubungan antar simbol yang menyatakan berjalannya suatu sistem
4.		<i>Data</i>	Simbol ini menggambarkan suatu proses yang dilakukan memiliki data masukan dan keluaran

DAFTAR LAMPIRAN

Lampiran 1 Daftar <i>Stopword</i> Sastrawi.....	49
Lampiran 2 Daftar Kata <i>Normalization</i>	58