

**PERBANDINGAN ANALISIS SENTIMEN TERHADAP ULASAN
PELANGGAN AMAZON ANTARA ALGORITMA LOGISTIC
REGRESSION DAN MULTINOMIAL NAIVE BAYES**



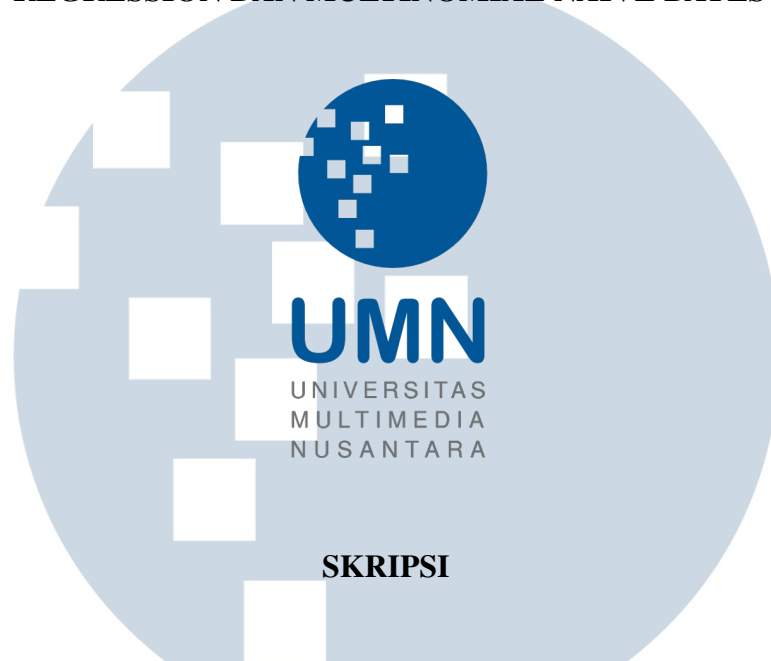
SKRIPSI

Diajukan sebagai salah satu syarat untuk memperoleh
Gelar Sarjana Komputer (S.Kom.)

Sarah Nabila Khairunisa
00000042223

**PROGRAM STUDI INFORMATIKA
FAKULTAS TEKNIK DAN INFORMATIKA
UNIVERSITAS MULTIMEDIA NUSANTARA
TANGERANG
2024**

**PERBANDINGAN ANALISIS SENTIMEN TERHADAP ULASAN
PELANGGAN AMAZON ANTARA ALGORITMA LOGISTIC
REGRESSION DAN MULTINOMIAL NAIVE BAYES**



SKRIPSI

Diajukan sebagai salah satu syarat untuk memperoleh
Gelar Sarjana Komputer (S.Kom.)

Sarah Nabila Khairunisa

00000042223

UMN

UNIVERSITAS

MULTIMEDIA

NUSANTARA

**PROGRAM STUDI INFORMATIKA
FAKULTAS TEKNIK DAN INFORMATIKA
UNIVERSITAS MULTIMEDIA NUSANTARA**

TANGERANG

2024

HALAMAN PERNYATAAN TIDAK PLAGIAT

Dengan ini saya,

Nama : Sarah Nabila Khairunisa
Nomor Induk Mahasiswa : 00000042223
Program Studi : Informatika

Skripsi dengan judul:

Perbandingan Analisis Sentimen terhadap Ulasan Pelanggan Amazon antara Algoritma Logistic Regression dan Multinomial Naive Bayes

merupakan hasil karya saya sendiri bukan plagiat dari karya ilmiah yang ditulis oleh orang lain, dan semua sumber baik yang dikutip maupun dirujuk telah saya nyatakan dengan benar serta dicantumkan di Daftar Pustaka.

Jika di kemudian hari terbukti ditemukan kecurangan/ penyimpangan, baik dalam pelaksanaan Skripsi maupun dalam penulisan laporan Skripsi, saya bersedia menerima konsekuensi dinyatakan TIDAK LULUS untuk Tugas akhir yang telah saya tempuh.

Tangerang, 3 April 2024



(Sarah Nabila Khairunisa)

UNIVERSITAS
MULTIMEDIA
NUSANTARA

HALAMAN PENGESAHAN

Skripsi dengan judul

**PERBANDINGAN ANALISIS SENTIMEN TERHADAP ULASAN PELANGGAN
AMAZON ANTARA ALGORITMA LOGISTIC REGRESSION DAN
MULTINOMIAL NAIVE BAYES**

oleh

Nama : Sarah Nabila Khairunisa
NIM : 00000042223
Program Studi : Informatika
Fakultas : Fakultas Teknik dan Informatika

Telah diujikan pada hari Selasa, 4 Juni 2024

Pukul 15.00 s/d 17.00 dan dinyatakan

LULUS

Dengan susunan penguji sebagai berikut

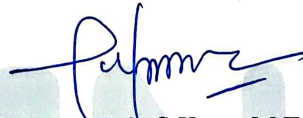
Ketua Sidang

Penguji



(Dr. Ir. Winarno, M.Kom.)

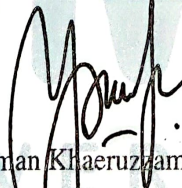
NIDN: 330106002



(Sy. Yuliani Yakub, S.Kom., M.T., Ph.D.)

NIDN: 0411037904

Pembimbing



(Yaman Khaeruzaman, M.Sc.)

NIDN: 0413057104

P. Ketua Program Studi Informatika,



(Dr. Eng. Niki Prastomo, S.T., M.Sc.)

NIDN: 0419128203

HALAMAN PERSETUJUAN PUBLIKASI KARYA ILMIAH

Yang bertanda tangan di bawah ini:

Nama : Sarah Nabila Khairunisa

NIM : 00000042223

Program Studi : Informatika

Jenjang : S1


Jenis Karya : Skripsi

Menyatakan dengan sesungguhnya bahwa:

- Saya bersedia memberikan izin sepenuhnya kepada Universitas Multimedia Nusantara untuk mempublikasikan hasil karya ilmiah saya di repositori Knowledge Center, sehingga dapat diakses oleh Civitas Akademika/Publik. Saya menyatakan bahwa karya ilmiah yang saya buat tidak mengandung data yang bersifat konfidensial dan saya juga tidak akan mencabut kembali izin yang telah saya berikan dengan alasan apapun.
- Saya tidak bersedia karena dalam proses pengajuan untuk diterbitkan ke jurnal/konferensi nasional/internasional (dibuktikan dengan *letter of acceptance*)**.

Tangerang, 3 April 2024

Yang menyatakan



Sarah Nabila Khairunisa

UMN
UNIVERSITAS
MULTIMEDIA
NUSANTARA

** Jika tidak bisa membuktikan LoA jurnal/HKI selama enam bulan ke depan, saya bersedia mengizinkan penuh karya ilmiah saya untuk diunggah ke KC UMN dan menjadi hak institusi UMN.

Halaman Persembahan / Motto

"The best way to not being hopeless is to get up and do something."

Barrack Obama



UMMN
UNIVERSITAS
MULTIMEDIA
NUSANTARA

KATA PENGANTAR

Puji Syukur atas berkat dan rahmat kepada Tuhan Yang Maha Esa, atas selesainya penulisan laporan Skripsi ini dengan judul: Perbandingan Analisis Sentimen terhadap Ulasan Pelanggan Amazon antara Algoritma Logistic Regression dan Multinomial Naive Bayes dilakukan untuk memenuhi salah satu syarat untuk mencapai gelar Sarjana Komputer Jurusan Informatika Pada Fakultas Teknik dan Informatika Universitas Multimedia Nusantara. Saya menyadari bahwa, tanpa bantuan dan bimbingan dari berbagai pihak, dari masa perkuliahan sampai pada penyusunan skripsi ini, sangatlah sulit bagi saya untuk menyelesaikan skripsi ini. Oleh karena itu, saya mengucapkan terima kasih kepada:

1. Bapak Dr. Ninok Leksono, selaku Rektor Universitas Multimedia Nusantara.
2. Bapak Dr. Eng. Niki Prastomo, S.T., M.Sc., selaku Dekan Fakultas Teknik dan Informatika dan Ketua Program Studi Informatika Universitas Multimedia Nusantara.
3. Bapak Yaman Khaeruzzaman, M.Sc. sebagai dosen pembimbing saya yang telah memberikan bimbingan untuk menyelesaikan skripsi saya.
4. Keluarga dan teman dekat saya yang telah memberikan semangat sehingga saya dapat menyelesaikan skripsi ini.

Semoga skripsi ini bermanfaat, baik sebagai sumber informasi maupun sumber inspirasi, bagi para pembaca.

Tangerang, 3 April 2024


Sarah Nabila Khairunisa

**PERBANDINGAN ANALISIS SENTIMEN TERHADAP ULASAN
PELANGGAN AMAZON ANTARA ALGORITMA LOGISTIC
REGRESSION DAN MULTINOMIAL NAIVE BAYES**

Sarah Nabila Khairunisa

ABSTRAK

Amazon adalah salah satu penyedia jasa penjualan daring terbesar di dunia. Amazon memiliki banyak pembeli dan berbagai macam produk dan layanan yang dijual. Terkadang para pembeli amazon bingung saat ingin membeli produk yang cocok untuk mereka. Oleh karena itu pembeli harus bisa mengetahui ulasan pelanggan sebelumnya sebelum membeli produk. Pembeli harus bisa membedakan ulasan positif dan negatif agar bisa mengetahui apakah produk yang ingin dibeli bagus atau tidak. Pada penelitian ini, dibuat model analisis sentimen untuk mengetahui apakah sebuah ulasan bersifat positif atau negatif dengan membandingkan 2 algoritma yaitu *Logistic Regression* dan *Multinomial Naive Bayes*. Skenario pengujian dilakukan dengan rasio data latih dan data uji sebesar 80:20 dan 70:30 untuk kedua algoritma. Setelah itu, pada masing-masing model diterapkan metode *Synthetic Minority Oversampling Technique (SMOTE)* untuk menyeimbangkan antara data-data dengan kelas positif dan negatif pada kumpulan data latih. Hasil klasifikasi menunjukkan bahwa model dengan algoritma *Logistic Regression* dengan rasio data latih dan data uji sebesar 80:20 dan telah diterapkan SMOTE memiliki hasil akurasi yang paling tinggi yaitu sebesar 86.47%. Hasil presisi pada model tersebut adalah 56.31%, *recall* sebesar 72.96% dan *f1-score* sebesar 63.56%.

Kata kunci: Amazon, Analisis Sentimen, *Logistic Regression*, *Multinomial Naive Bayes*, Python, *Synthetic Minority Oversampling Technique*

U N I V E R S I T A S
M U L T I M E D I A
N U S A N T A R A

THE COMPARISON BETWEEN LOGISTIC REGRESSION AND MULTINOMIAL NAIVE BAYES IN THE SENTIMENT ANALYSIS OF AMAZON CUSTOMER REVIEWS

Sarah Nabila Khairunisa

ABSTRACT

Amazon is one of the biggest e-commerce platforms in the world. It has a lot of customers and offers a lot of products and services for sale. Sometimes Amazon customers get confused while picking the right product for them. Therefore, customers have to know previous customer reviews before buying a product. Customers should be able to tell the difference between positive and negative reviews to determine if the product they want to buy is good or bad. In this research, a sentiment analysis model was made to determine if a review is positive or negative by comparing 2 different algorithms which are Logistic Regression and Multinomial Naive Bayes. There were a few different scenarios, first with 80:20 train-test split and second with 70:30 train-test split for both algorithms. Then, for each model, *Synthetic Minority Oversampling Technique* (SMOTE) was implemented to balance the amount of positive-labelled reviews and negative-labelled reviews in the train data. The classification result showed that the model with Logistic Regression and 80:20 train-test data ratio that had been implemented SMOTE had the biggest accuracy which was 86.47%. The precision was 56.31%, the recall was 72.96% and the f1-score was 63.56%.

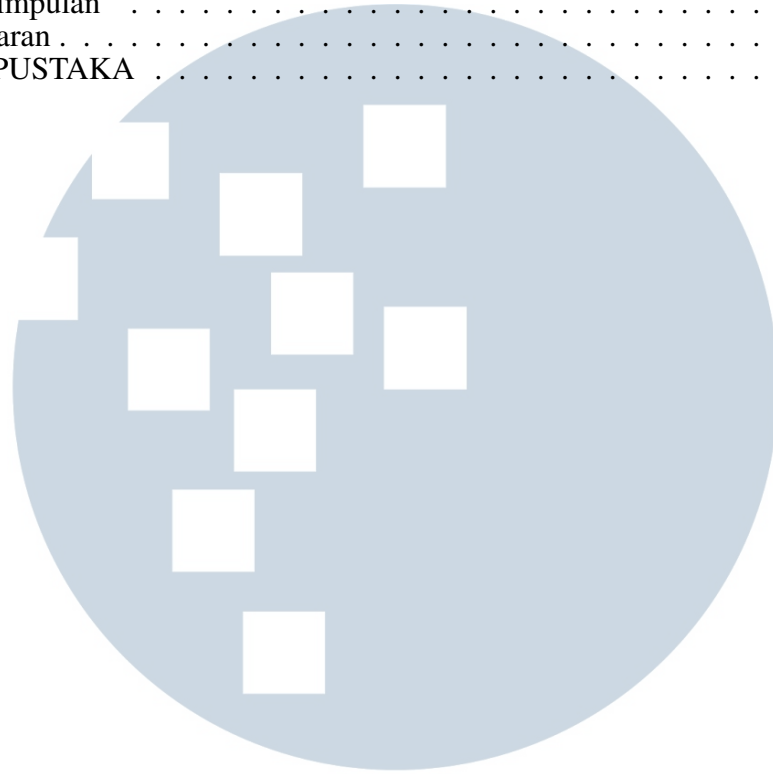
Keywords: Amazon, *Logistic Regression*, *Multinomial Naive Bayes*, Python, *Synthetic Minority Oversampling Technique*

U M N
U N I V E R S I T A S
M U L T I M E D I A
N U S A N T A R A

DAFTAR ISI

HALAMAN JUDUL	i
PERNYATAAN TIDAK MELAKUKAN PLAGIAT	ii
HALAMAN PERSETUJUAN	iii
HALAMAN PENGESAHAN	iv
HALAMAN PERSETUJUAN PUBLIKASI ILMIAH	v
HALAMAN PERSEMBAHAN/MOTO	vi
KATA PENGANTAR	vii
ABSTRAK	viii
ABSTRACT	ix
DAFTAR ISI	x
DAFTAR GAMBAR	xii
DAFTAR TABEL	xiii
DAFTAR KODE	xiv
DAFTAR LAMPIRAN	xv
BAB 1 PENDAHULUAN	1
1.1 Latar Belakang Masalah	1
1.2 Rumusan Masalah	3
1.3 Batasan Permasalahan	3
1.4 Tujuan Penelitian	4
1.5 Manfaat Penelitian	4
1.6 Sistematika Penulisan	4
BAB 2 LANDASAN TEORI	6
2.1 Amazon	6
2.2 Machine Learning	6
2.3 Natural Language Processing	6
2.4 Preprocessing	7
2.5 Term Frequency - Inverse Document Frequency	8
2.5.1 Term Frequency - Inverse Document Frequency	9
2.6 Synthetic Minority Oversampling Technique (SMOTE)	10
2.7 Naive Bayes	10
2.8 Logistic Regression	11
2.9 Confusion Matrix	12
BAB 3 METODOLOGI PENELITIAN	15
3.1 Gambaran Umum Penelitian	15
3.2 Perancangan Sistem	16
3.2.1 Preprocessing	16
3.2.2 Labelling	18
3.2.3 TF-IDF	19
3.2.4 Modelling dan Evaluasi	20
BAB 4 HASIL DAN DISKUSI	21
4.1 Spesifikasi Sistem	21
4.2 Hasil Implementasi	21
4.2.1 Import Libraries	21
4.2.2 Import Data	22
4.2.3 Preprocessing	23
4.2.4 Pembuatan Tabel-tabel Baru	27
4.2.5 WordCloud	28
4.2.6 Modelling	29

4.2.7	Evaluasi	33
BAB 5	SIMPULAN DAN SARAN	42
5.1	Simpulan	42
5.2	Saran	42
DAFTAR PUSTAKA	43



UMMN
 UNIVERSITAS
 MULTIMEDIA
 NUSANTARA

DAFTAR GAMBAR

Gambar 2.1	Tabel Confusion Matrix	12
Gambar 3.1	Gambaran Umum Penelitian	15
Gambar 3.2	<i>Labelling</i>	19
Gambar 4.1	Data yang Telah Diimpor	22
Gambar 4.2	Jumlah Kolom dalam Tabel Setelah Diimpor	23
Gambar 4.3	Jumlah Baris Setelah Fungsi <code>drop_duplicates()</code> Dijalankan .	23
Gambar 4.4	Jumlah Baris Setelah Baris-baris Data dengan Nilai Kosong Dihapus	24
Gambar 4.5	Total Label Positif dan Negatif dalam Tabel	27
Gambar 4.6	WordCloud Kata-kata Berlabel Positif	28
Gambar 4.7	WordCloud Kata-kata Berlabel Negatif	29
Gambar 4.8	<i>Confusion Matrix</i> untuk Model <i>Logistic Regression</i> (80:20)	33
Gambar 4.9	<i>Confusion Matrix</i> untuk Model <i>Logistic Regression</i> (70:30)	34
Gambar 4.10	<i>Confusion Matrix</i> untuk Model <i>Multinomial Naive Bayes</i> (80:20)	35
Gambar 4.11	<i>Confusion Matrix</i> untuk Model <i>Multinomial Naive Bayes</i> (70:30)	36
Gambar 4.12	<i>Confusion Matrix</i> untuk Model <i>Logistic Regression</i> (80:20) ditambah SMOTE	38
Gambar 4.13	<i>Confusion Matrix</i> untuk Model <i>Logistic Regression</i> (70:30) ditambah SMOTE	39
Gambar 4.14	<i>Confusion Matrix</i> untuk Model <i>Multinomial Naive Bayes</i> (80:20) ditambah SMOTE	40
Gambar 4.15	<i>Confusion Matrix</i> untuk Model <i>Multinomial Naive Bayes</i> (70:30) ditambah SMOTE	41



DAFTAR TABEL

Tabel 3.1	Perbandingan Data Original dan yang Sudah di <i>Casefolding</i>	17
Tabel 3.2	Perbandingan Data Original dan yang Sudah di <i>Cleaning</i>	17
Tabel 3.3	Perbandingan Data Original dan yang Sudah Dilakukan <i>Tokenizing</i>	17
Tabel 3.4	Perbandingan Data Original dan yang Sudah Dilakukan <i>Stopword Removal</i>	18
Tabel 3.5	Perbandingan Data Original dan yang Sudah Dilakukan <i>Lemmatization</i>	18
Tabel 4.1	Jumlah Data <i>Train</i> dan Data <i>Test</i>	29
Tabel 4.2	Hasil Pemodelan Tanpa SMOTE untuk Kedua Algoritma	33
Tabel 4.3	Hasil Pemodelan dengan SMOTE untuk Kedua Algoritma	37



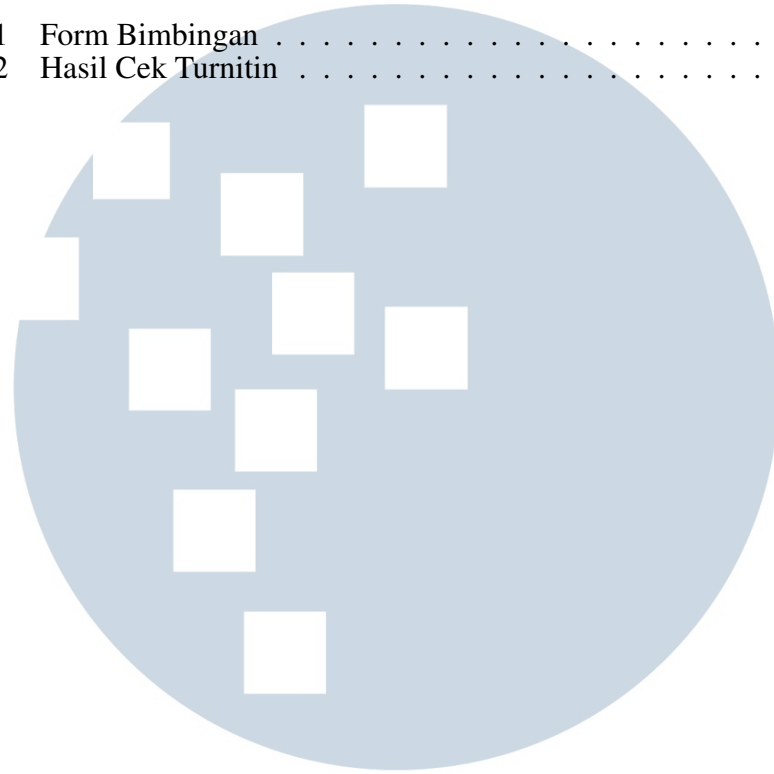
DAFTAR KODE

4.1	Semua <i>Library</i> yang Dibutuhkan	21
4.2	Kode untuk Mengimpor Data	22
4.3	Kode untuk Menghitung Baris dalam Tabel	23
4.4	Kode untuk Menghilangkan Data Duplikat	23
4.5	Kode untuk Menghilangkan Data dengan Nilai Kosong	24
4.6	Kode untuk Melakukan <i>Casefolding</i>	24
4.7	Kode untuk <i>Data Cleaning</i>	24
4.8	Kode untuk Melakukan <i>Tokenization</i>	25
4.9	Kode untuk <i>Stopword Removal</i>	25
4.10	Kode untuk <i>Lemmatization</i>	26
4.11	Kode untuk <i>Labelling</i>	26
4.12	Kode untuk Membuat dan Menyimpan Tabel <i>preprocessed_data</i>	28
4.13	Kode untuk Membuat Tabel <i>final_df</i>	28
4.14	Kode untuk Melakukan <i>Train-test Split</i> dengan Rasio 80:20	30
4.15	Kode untuk Melakukan <i>Train-test Split</i> dengan Rasio 70:30	30
4.16	Kode untuk Melakukan Pemodelan dengan <i>Logistic Regression</i> dengan Rasio <i>Train-test Split</i> Sebesar 80:20	30
4.17	Kode untuk Melakukan Pemodelan dengan <i>Logistic Regression</i> dengan Rasio <i>Train-test Split</i> Sebesar 70:30	30
4.18	Kode untuk Melakukan Pemodelan dengan <i>Logistic Regression</i> dengan Rasio <i>Train-test Split</i> Sebesar 80:20 Ditambah SMOTE	31
4.19	Kode untuk Melakukan Pemodelan dengan <i>Logistic Regression</i> dengan Rasio <i>Train-test Split</i> Sebesar 70:30 Ditambah SMOTE	31
4.20	Kode untuk Melakukan Pemodelan dengan <i>Multinomial Naive</i> <i>Bayes</i> dengan Rasio <i>Train-test Split</i> Sebesar 80:20	31
4.21	Kode untuk Melakukan Pemodelan dengan <i>Multinomial Naive</i> <i>Bayes</i> dengan Rasio <i>Train-test Split</i> Sebesar 70:30	32
4.22	Kode untuk Melakukan Pemodelan dengan <i>Multinomial Naive</i> <i>Bayes</i> dengan Rasio <i>Train-test Split</i> Sebesar 80:20 Ditambah SMOTE	32
4.23	Kode untuk Melakukan Pemodelan dengan <i>Multinomial Naive</i> <i>Bayes</i> dengan Rasio <i>Train-test Split</i> Sebesar 70:30 Ditambah SMOTE	32

UNIVERSITAS
MULTIMEDIA
NUSANTARA

DAFTAR LAMPIRAN

Lampiran 1	Form Bimbingan	45
Lampiran 2	Hasil Cek Turnitin	46



UMN
UNIVERSITAS
MULTIMEDIA
NUSANTARA