IMPLEMENTASI NATURAL LANGUAGE PROCESSING (NLP) DASAR PADA ANALISIS SENTIMEN TAGAR #KABURAJADULU

(Skripsi)

Oleh

ZIKWAN ISMAIL NPM 2117051051



FAKULTAS MATEMATIKA DAN ILMU PENGETAHUAN ALAM UNIVERSITAS LAMPUNG BANDAR LAMPUNG

2025

ABSTRAK

IMPLEMENTASI NATURAL LANGUAGE PROCESSING (NLP) DASAR PADA ANALISIS SENTIMEN TAGAR #KABURAJADULU

Oleh

ZIKWAN ISMAIL

Fenomena penggunaan tagar #KaburAjaDulu di media sosial mencerminkan respons publik terhadap kondisi sosial, ekonomi dan politik di Indonesia, khususnya dari masyarakat yang mempertimbangkan untuk menetap di luar negeri. Penelitian ini bertujuan untuk menganalisis opini publik yang terungkap melalui unggahan bertagar tersebut menggunakan pendekatan *Natural Language Processing* (NLP) dasar. Peneliti mengumpulkan 5.198 data unggahan dari platform X (*Twitter*) dan *TikTok* melalui proses *scraping*, kemudian menerapkan tahap preprocessing teks seperti *cleaning, tokenization, stopword removal, case folding* dan *stemming* menggunakan *library Python* seperti *nltk* dan Sastrawi. Data yang telah diproses dilabeli secara otomatis oleh *volunteer* dan digunakan untuk melatih dua model klasifikasi sentimen: *Naïve Bayes* dan *Logistic Regression*. Evaluasi model dilakukan dengan metrik *accuracy, precision, recall* dan *F1-score* melalui pendekatan *cross-validation*.

Hasil penelitian menunjukkan bahwa *Logistic Regression* memberikan performa klasifikasi yang lebih unggul dibandingkan *Naïve Bayes* pada seluruh *dataset*, dengan akurasi mencapai 89% pada data *TikTok*, dibandingkan 65% oleh *Naïve Bayes*. Selain itu, model *Logistic Regression* menunjukkan ketahanan yang lebih baik dalam menangani distribusi data yang kompleks dan tidak seimbang. Secara keseluruhan, penelitian ini menunjukkan efektivitas metode NLP dasar dalam memetakan opini publik secara digital dan menegaskan bahwa *Logistic Regression* lebih adaptif dalam memahami konteks sentimen dari teks sosial media. Temuan ini memberikan kontribusi penting bagi studi lanjutan mengenai analisis opini publik dan pengambilan keputusan berbasis data sosial.

Kata Kunci: Analisis Sentimen, Pemrosesan Bahasa Alami (NLP), *Twitter* (X), TikTok, #kaburajadulu, *Naïve Bayes*, Regresi Logistik.

ABSTRACT

IMPLEMENTATION OF BASIC NATURAL LANGUAGE PROCESSING (NLP) IN THE ANALYSIS OF SENTIMENT HASHTAG #KABURAJADULU

By

ZIKWAN ISMAIL

The phenomenon of using the hashtag #KaburAjaDulu on social media reflects the public's response to social, economic, and political conditions in Indonesia, particularly among those considering settling abroad. This study aims to analyze public opinion expressed through posts using this hashtag using a basic Natural Language Processing (NLP) approach. The researchers collected 5,198 posts from the X (Twitter) and TikTok platforms through scraping, then applied text preprocessing steps, including cleaning, tokenization, stopword removal, case folding, and stemming, using Python libraries such as nltk and Sastrawi. The processed data was manually labeled by volunteers and used to train two sentiment classification models: Naïve Bayes and Logistic Regression. Model evaluation was conducted using accuracy, precision, recall, and F1-score metrics through a crossvalidation approach.

The results showed that Logistic Regression outperformed Naïve Bayes across all datasets, achieving 89% accuracy on TikTok data compared to 65% by Naïve Bayes. Additionally, the Logistic Regression model demonstrated better resilience in handling complex and imbalanced data distributions. Overall, this study highlights the effectiveness of basic NLP methods in mapping public opinion digitally and confirms that Logistic Regression is more adaptive.

Keywords: Sentiment Analysis, NLP, Twitter(X), TikTok, #kaburajadulu, Naïve Bayes, Logistic

IMPLEMENTASI NATURAL LANGUAGE PROCESSING (NLP) DASAR PADA ANALISIS SENTIMEN TAGAR #KABURAJADULU

Oleh

ZIKWAN ISMAIL

Skripsi

Sebagai Salah Satu Syarat untuk Mencapai Gelar SARJANA KOMPUTER

Pada

Jurusan Ilmu Komputer
Fakultas Matematika dan Ilmu Pengetahuan Alam
Universitas Lampung



FAKULTAS MATEMATIKA DAN ILMU PENGETAHUAN ALAM UNIVERSITAS LAMPUNG BANDAR LAMPUNG 2025

Judul Skripsi IMPLEMENTASI NATURAL LANGUAGE PROCESSING (NLP) DASAR PADA ANALISIS SENTIMEN TAGAR **#KABURAJADULU**

Nama Mahasiswa/i Zikwan Ismail

Nomor Pokok Mahasiswa/i: 2117051051

UNIVERSITAS LAMPUNG

UNIVERSITAS LAMITY

UNIVERSITAS LAMPI

UNIVERSITAS LAMPL

UNIVERSITAS LAMPUT

UNIVERSITAS LAMPUN

INIVERSITAS LAMPUNO

UNIVERSITAS LAMPUNU

UNIVERSITAS LAMPUNG

15 LAMPUNG

15 LAMPUNG

15 LAMPUNG

AS LAMPUNG

SITASLAMPUNG

Cres

Program Studi S1 Ilmu Komputer

Fakultas Matematika dan Ilmu Pengetahuan Alam

> USAY ERSITAS LAMPUNC INIVERSILANTAMPUNG ONIVERSITASLASIMING NIVERSITANLAMENTAC Komisi Pembimbing UNIVERSITAS LAMPUNG Ossy Dwi Endah Wulansari, S.Si., M.T INIVERSITAS LAMPUNG INIVERSITAS LAMPLING NIP. 197407132003122002

MENGETAHUI

BMINTERSTERS

CONTRACTOR STATE

Ketua Jurusan Ilmu Komputer

Dwi Sakethi, S.Si., M.Kom.

UNIVERSITAS LAMPUNG

PATRICE

NIP. 196806111998021001

LINIVERSITAS LAMPUNG

UNIVERSITAS LAMPUSICI

INIVEROR

Ketua Program Studi S1 Ilmu Komputer

UNIVERSITAS LAMPING

ONIVERSITAS LAMPUNC

UNIVERSITAS LAMPUNC

INIVERSETAS CAMPUNO

UNIN LESITAS LAMPUNG

THIVERSHAS LAMPUNE

ONLY ERSTLANDING

UNIVERSITAS LAMPING

UNIVERSITAS LANIPUNG

UNIVERSITAS LAMPLING

UNIVERSITAS LAMPUNG

UNIVERSITAS LAMPUNG

UNIVERSITAS LAMPUNG

UNIVERSITAS LAAPUNG

UNIVERSITAS LAMPUNG

UNIVERSITAS LAMPUNG

UNIVERSITAS LAMPUNG

UNIVERSITAS LAMPUNG

UNIVERSITAS LAMPUNG

UNIVERSITAS LAMPUNG

UNIVERSITAS LAMPLING

UNIVERSITAN LAMPLING

INIVERSETAS LAMPUNG

Tristiyanto, M.I.S., Ph.D. NIP. 198104142005011001

LANGE OF THE LANGE OF THE PARTY OF THE PARTY

UNIVERSITES LAMBURY CONTVERSITES LAMPENG LINIVERSITES LAMPENG LINIVERSITES LAMPENG LINIVERSITES

MALAMPERN UNIVERSITAS LAMPUNG UNIVERSITAS LAMPUNI UNIVERSIFAS LAMPUN UNIVERSITAS LAMPUNG UNIVERSITAS LAMPUNG UNIVERSITAS LAMBUNG UNIVERSIDISLAMPUNU UNIVERSITAS LANDING UNIVERSITASLAMPUNG UNIVERSITAS LAMPUNO INIVERSITAS LAMPUNG UNIVERSITAS LAMPUNG UNIVERSITASLAMPING UNIVERSITAS LAMPUNC UNIVERSITASLAMPUNG UNIVERSITAS LAMPUNG UNIVERSITAS LAMPUNG UNIVERSITASLAMPUNG UNIVERSITAS LAMPUNG UNIVERSITAS LAMPLING MENGESAHKAN NIVERSITAS LAMPUNG UNIVERSITAS LAMPUNC UNIVERSITAS LANDUNG NIN 1. RS Tim Penguji NIVERSITAS LAMPUNG UNIVERSITAS LAMPUNG UNIVERSITAS LAMPUNG UNIVERSITAS LAMPUNG Ketua LAMPUNO : Ossy Dwi Endah Wulansari, S.Si., M.T UNIVERSITAS LAMPUNC UNIVERSITASIAMPUNG UNIVERSITAS LAMPUNG UNIVERSITAS LAMPUNG UNIVERSITAS LAMPUNG UNIVERSITAS LAMPUNG LIMPUNG Rizky Prabowo, M.Kom. Sekretaris Penguji Utama Dewi Asiah Shofiana, S.Komp., M.Kom. INIVERSITAS LAMPUN UNIVERSITAGILAMPUNU UNIVERSITAS LAMPUNG UNIVERSITAS LAMPUNG UNIVERSITAS LABOUNG Dekan Fakultar Matematika dan Ilmu Pengetahun Alam UNIVERSITAS LAMPUNG UNIVERSITAS LAURUNG UNIVERSITAS LAMPUNG UNIVERSITAS LAMPUNG INIVERSITAS LAMPUNG. ON TOTAS MATENANT ST UNIVERSITAS LAMPUNO UNIVERSITAS LAMPONG UNIVERSITAS LAMPUNG CNAVERSITAS LANGUNG Dr. Eng. Heri Satria, S.Si., M.Si. UNIVERSITAS LAMPUNG NIVERSITASLAMPUNG UNIVERSITASLANDUNG UNIVERSIDAS LAMPENG NIP. 197110012005011002 UNIVERSITAS LAMPUNG UNIVERSITAS LAMPUNG UNIVERSITAS LAWRUNG UNIVERSITAS LAMPUNG UNIVERSITAS LAMPUNG UNIVERSITASIANDUNG UNIVERSITAS LAMPUNG UNIVERSITIS LAMPUNO UNIVERSITAS LAMPLING UNIVERSITAS LAMPUNG UNIVERSITAS LAMPUNC UNIVERSITASCANDUNG UNIVERSITAS LAMPUNG ONIVERSITAS LAMPUNG STAMPUNG NIVERSITAS LAMPUNG UNIVERSITASIAMPUNG UNIVERSITAS LAMPUNG UNIVERSITAS LAMPUNG UNIVERSITANLAMPUNG UNIVERSITAS LAMPUNG UNIVERSITANLAMPLING UNIVERSITAS LAMPUNG NIVERSITAS LANDUNG UNIVERSITAS LANDUNG UNIVERSITAS LAMPUNG UNIVERSITAS LAMPUNC Tanggal Lulus Ujian Skripsi: 29 Juli 2025 UNIVERSITISLAMEUNG UNIVERSITAS LAMPUNG UNIVERSITAS LAMPUNG STAMPLING UNIVERSITASLAMPUNG ONIVERSITAS LAMPUNE UNIVERSITASLAMPUNG DRIVERSTEAS CAMPUNG SI AMPUNG INIVERSITAS LANDUNG ONIVERSITASIAMPUNO LINIVERSITAS LAMPUNG UNIVERSITAS LAMPUNG UNIVERSITAS LAMPUNG STAMPUNG ANIVERSITAS LANDENG UNIVERSITAS LAMPUNO UNIGERSITAS LAMPIDIG NIVERSITALS LAMPUNG UNIVERSITAS LAMPUNG CINIVERSITAN LAMPUNO UNIVERSITAS LAMPUNG LAMPING UNIVERSITAS LANDRING UNIVERSITAS CAMPUNO LAMPUNG NIVERS LAS LAMPLING UNIVERSITAS LAMPUNG LAMPUNG UNIVERSITAS LAMPUNG SI AMPLISA CHILDERS/TAS LAMPUNG UNIVERSITAS LAMPUNG UNIVERSITANLAMPUNG SLAMPLING UNIVERSITASTAMPLING UNIVERSITAS LAMPENG USINERSITAS LAMPATNO UNIVERSITASIAMPUNK UNIVERSITAS LAMPUNG IS LAMPUNG SOUTOCLE AND DEPOSIT STREET, THE PERSON

MINIVI

UNIV

PERNYATAAN

Saya bertanda tangan di bawah ini:

Nama : Zikwan Ismail

Npm : 2117051051

Program Studi/ Jurusan : S1 Ilmu Komputer

Fakultas : Matematika dan Ilmu Pengetahuan Alam

menyatakan bahwa skripsi saya yang berjudul "Implementasi Natural Language Processing (NLP) Dasar Pada Analisis Sentimen Tagar #Kaburajadulu" merupakan hasil karya saya sendiri dan bukan karya orang lain. Seluruh tulisan yang tertuang di skripsi ini telah mengikuti kaidah penulisan karya ilmiah Universitas Lampung. Apabila di kemudian hari terbukti skripsi saya merupakan hasil penjiplakan atau dibuat orang lain, maka saya bersedia menerima sanksi berupa pencabutan gelar yang telah saya terima.

Bandar Lampung, 11 Agustus 2025

Penulis



Zikwan Ismail

NPM. 2117051051

RIWAYAT HIDUP



Penulis dilahirkan pada tanggal 22 Oktober 2000 di Tulang Bawang, Kabupaten Tulang Bawang, Provinsi Lampung dari pasangan Bapak Suparmin dan Ibu Rumini. Penulis menyelesaikan pendidikan formal pertama kali di Taman Kanak-kanak Dharma Wanita Bumi Dipasena Mulya 2006, lalu pendidikan Sekolah Dasar (SD) diselesaikan di Sekolah Dasar Negeri 01 Bumi Dipasena Mulya Pada Tahun 2013.

Kemudian penulis menamatkan pendidikan Sekolah Menengah Pertama (SMP) di Sekolah Menengah Pertama Negeri 1 Rawajitu Timur pada Tahun 2016, lalu menyelesaikan pendidikan Sekolah Menengah Kejuruan (SMK) di SMK Al-Hikmah Kalirejo Lampung Tengah pada Tahun 2019.

Pada Tahun 2021 penulis terdaftar sebagai mahasiswa Jurusan Ilmu Komputer Fakultas Matematika dan Ilmu Pengetahuan Alam Universitas Lampung melalui jalur SBMPTN. Selama menjadi mahasiswa, penulis melakukan beberapa kegiatan antara lain.

- Menjadi Pimpinan Dinas Advokasi Kesejahteraan Mahasiswa (ADKESMA)
 Badan Eksekutif Mahasiswa (BEM) Fakultas Matematika dan Ilmu
 Pengetahuan pada Tahun 2024.
- 2. Menjadi Peserta Pada Kompetisi Olahraga Jurusan Ilmu Komputer (KOKOM) dan Menjadi Juara 2 Pada Cabang Olahraga *PUBG Mobile* pada Tahun 2021.
- 3. Menjadi Asisten Dosen Mata Kuliah Pemrograman Web pada Tahun 2024.
- 4. Menjadi panitia anggota Divisi Acara pada acara Pekan Raya Jurusan Ilmu Komputer Tahun 2022.

- Melaksanakan Kerja Praktik di PLN ULP Natar periode 2023/2024 dengan program kerja pengembangan Aplikasi Penjadwalan Dan Pemetaan Pelanggan Untuk Tugas Operasi Petugas Ganti Meter Dengan Appsheet Dan Web Programing.
- Menjadi Pesrta Dalam Kegiatan Kuliah Umum: Iot and Artificial Intelegence in Business yang diselenggarakan oleh Jurusan Ilmu Komputer pada Tahun 2023.
- 7. Mengikuti *Short Class UI/UX Design* yang diselenggarakan oleh *MySkill* pada Tahun 2023.
- 8. Mengikuti *Short Class Python Introduction* yang diselenggarakan oleh *MySkill* pada Tahun 2023.
- 9. Menjadi Peserta Dalam Kegiatan *Social Engineering: The Human Element in Cybersecurity Breaches*, yang diselenggarakan oleh Jurusan Ilmu Komputer pada Tahun 2024

MOTTO

- 1. "Nilai seseorang tergantung pada apa yang ia anggap berharga." (Ali Bin Abu Thalib RA)
- 2. "The Beautiful Thing About Learning is That no One Can Take it Away From You."

(B.B King)

3. "Sometimes love isn't about who can complete us, but who makes us feel whole."

(Ariana Grande)

4. "A Difficult Struggle Will Get A Satisfactory Result"

PERSEMBAHAN

Alhamdulillahirabbil'alamiin

Puji syukur kehadirat Allah Subhanahu Wa Ta'ala atas segala rahmat dan karunia-Nya sehingga skripsi ini dapat diselesaikan dengan sebaik-baiknya. Shalawat serta salam senantiasa tercurahkan kepada suri teladan Nabi Muhammad Shallallahu 'Alaihi Wasallam.

Kupersembahkan karya ini kepada:

Bapak dan Mamaku Tercinta

Untuk Bapak dan Mama tercinta, yang namanya selalu terucap dalam setiap doaku. Terima kasih atas kasih sayang yang tak pernah lekang, doa yang tak pernah henti, serta pengorbanan yang tak terukur nilainya. Setiap langkah yang kutempuh adalah jejak dari bimbingan dan keteladanan kalian. Karya ini hadir sebagai wujud kecil dari cinta dan rasa terima kasihku, serta sebagai bukti bahwa doa dan restu kalian adalah kekuatan terbesar dalam hidupku.

Seluruh Keluarga Besar Ilmu Komputer 2021

Yang telah memberikan ilmu, bantuan, dukungan, juga suka duka selama masa perkuliahan.

Almamater Tercinta, Universitas Lampung dan Jurusan Ilmu Komputer

Tempat di mana mimpi-mimpi mulai dirangkai, pengetahuan ditempa dan persahabatan terjalin. Terima kasih telah menjadi ruang untuk bertumbuh, belajar dan membekali diri demi menapaki masa depan

SANWACANA

Puji dan Syukur dipanjatkan kehadirat Allah Subhanahu Wa Ta'ala atas limpahan nikmat, rahmat dan karunia-Nya. Shalawat serta salam senantiasa tercurahkan kepada junjungan Nabi Muhammad Shallallahu 'Alaihi Wasallam, yang syafa'at nya sangat diharapkan di yaumil akhir kelak. Skripsi berjudul "Implementasi Natural Language Processing (NLP) Dasar Pada Analisis Sentimen Tagar #Kaburajadulu" telah disusun dengan sebaik-baiknya dan sebagai salah satu syarat untuk mendapatkan gelar sarjana ilmu komputer di Universitas Lampung.

Terima kasih penulis ucapkan kepada pihak-pihak yang telah memberikan arahan, dukungan dan bimbingan selama penulis menyelesaikan penyusunan skripsi ini. Ucapan terima kasih sebesar-besarnya ditujukan kepada:

- 1. Allah SWT yang telah memberikan nikmat, petunjuk dan *ridho*-Nya sehingga karya ini dapat selesai dengan tepat waktu dan sebaik-baiknya.
- 2. Ayah dan Ibu Tercinta, Bapak Suparmin dan Mama Rumini yang selalu menjadi rumah terhangat untuk hati ini, tempat pulang yang penuh doa, kasih dan ketulusan. Terima kasih atas setiap dukungan, motivasi dan ruang diskusi yang membentuk prinsip hidupku. Atas pengorbanan yang tak pernah terhitung kesabaran yang tak pernah surut, serta kepercayaan yang kalian berikan hingga penulis dapat menyelesaikan perjalanan ini. Karya ini adalah buah dari cinta, doa dan perjuangan kalian. Semoga Allah memanjangkan usia kalian dalam kebaikan dan memberi kesempatan bagiku untuk terus membalas kebaikan itu, meski aku tahu, kasih sayang kalian takkan pernah terbalas seluruhnya.

- 3. Bapak Dr. Eng. Heri Satria, S.Si., M.Si. selaku Dekan FMIPA Universitas Lampung.
- 4. Bapak Dwi Sakethi, S.Si., M.Kom. selaku Ketua Jurusan Ilmu Komputer FMIPA Universitas Lampung.
- 5. Ibu Yunda Heningtyas, M. Kom. selaku Sekretaris Jurusan Ilmu Komputer FMIPA Universitas Lampung tahun 2025.
- 6. Bapak Tristiyanto, S. Kom., M.I.S., Ph.D. selaku Ketua Program Studi Ilmu Komputer FMIPA Universitas Lampung tahun 2025.
- Bapak Dr. Aristoteles, S.Si., M.Si. selaku Pembimbing Akademik yang telah banyak memberikan arahan dan bimbingannya selama penulis melaksanakan kegiatan perkuliahan.
- 8. Ibu Ossy Dwi Endah Wulansari, S.Si., M.T. selaku Pembimbing Utama yang telah banyak memberikan arahan, ide, kritik, serta saran hingga penyelesaian penelitian dan karya tulis ini.
- 9. Ibu Dewi Asiah Shofiana, S.Komp., M.Kom. selaku Pembahas Pertama yang selalu menjadi tempat diskusi dan selalu meluangkan waktunya untuk membimbing, memberikan ide dan bantuan kepada penulis selama menjalankan penelitian ini.
- 10. Bapak Rizky Prabowo, M.Kom selaku Pembahas Kedua yang telah memberikan masukan dan saran perbaikan untuk penelitian dan skripsi ini.
- 11. Bapak dan Ibu Dosen Jurusan Ilmu Komputer FMIPA Universitas Lampung yang telah memberikan ilmu, motivasi dan pengalaman hidup selama penulis menempuh pendidikan di Jurusan Ilmu Komputer Universitas Lampung.
- 12. Seluruh Staf dan Karyawan Jurusan Ilmu Komputer, Ibu Ade Nora Maela, Bang Zainuddin, Mas Syam dan Mas Nofal yang telah melayani segala urusan administrasi, laboratorium dan pinjaman ruangan kepada penulis selama masa perkuliahan di Jurusan Ilmu Komputer.

- 13. Abel Ivani Fitri Salsabila, S.H, pasangan hidup selama menjadi mahasiswa di Universitas Lampung yang selalu memberikan motivasi dan arahan dari awal hingga akhir perkuliahan. Terima kasih sudah selalu mendukung, memberikan waktunya tanpa ragu, menjadi tempat cerita, berkeluh kesah, diskusi berat hingga sepele. Serta canda tawa humor yang mengisi hubungan kita. Terima kasih juga untuk *support* nya dalam setiap kegiatan yang penulis lakukan. Terima kasih sudah memahami dan mengerti penulis. Penulis harap hubungan ini tetap terjalin dengan baik hingga kita tua.
- 14. Teman-teman Pimpinan Badan Eksekutif Mahasiswa FMIPA Unila 2024 (Abdurachman dan Erwin Kesuma) terima kasih sudah menjadi *partner* kerja yang mengesankan selama masa jabatan. Dari pengalaman itu banyak pembelajaran berharga yang membentuk penulis hingga saat ini. Terima kasih sudah saling percaya dan bertanggung jawab untuk segala sesuatunya.
- 15. Teman-teman Ilmu Komputer Angkatan 2021, terima kasih telah menjadi rekan kelompok, rekan diskusi dan rekan berjuang selama menjalankan studi di Jurusan Ilmu Komputer Universitas Lampung.
- 16. Teman-teman KKN Desa Hurun (Abdurachman, Amel, dll) yang telah membersamai serta memberikan semangat dan dukungan kepada penulis dari semasa KKN hingga masa penelitian ini selesai.
- 17. Kedua Kakak saya, Eka Julianto dan Muhammad Dwi Susilo, Terimakasih telah senantiasa mendukung, memberikan arahan serta saran yang berharga kepada penulis. Do'a, perhatian dan teladan yang kalian berikan menjadi dorongan besar bagi penulis untuk tumbuh menjadi pribadi yang lebih baik.
- 18. Serta semua pihak yang tidak tersebut disini, terima kasih atas dukungan dan apresiasi yang telah diberikan secara langsung maupun tidak langsung kepada penulis selama perkuliahan hingga penyelesaian karya tulis ini.

19. Diri saya sendiri, Zikwan Ismail, Terima kasih telah bertahan dan mengusahakan segala sesuatunya dengan maksimal, juga mampu mengendalikan diri di tengah berbagai tekanan dan peristiwa sulit selama masa perkuliahan. Terima kasih sudah berani dan bertanggung jawab menyelesaikan apa yang telah dimulai. Tetap semangat untuk menjalankan rencana-rencana berikutnya!

Penulis menyadari bahwa penyusunan skripsi ini masih jauh dari kata sempurna. Namun penulis sangat mengharapkan skripsi ini dapat bermanfaat bagi para civitas akademik Universitas Lampung, khususnya mahasiswa Ilmu Komputer.

Bandar Lampung, 11 Agustus 2025

\bigver_\tag{\frac{1}{2}}

Zikwan Ismail

NPM. 2117051051

DAFTAR ISI

	Halaman
DAFTAR TABEL	viii
DAFTAR GAMBAR	ix
I.PENDAHULUAN	1
1.1. Latar Belakang	1
1.2. Rumusan Masalah	2
1.4. Tujuan	3
1.5. Manfaat	3
II. TINJAUAN PUSTAKA	4
2.1 Penelitian Terkait	4
2.2 Media Sosial	5
2.3 Natural Language Processing (NLP)	5
2.4 Naive Bayes	6
2.5 Logistic Regression	7
2.7 Text Processing	9
2.8 Pembobotan TF	11
2.9 Tagar (<i>Hashtag</i>)	11
2.10 Confusion Matrix	12
III. METODOLOGI PENELITIAN	17
3.1 Waktu dan Tempat Penelitian	17
3.2 Perangkat Penelitian	17
3.2.1 Perangkat Keras (<i>Hardware</i>)	17
3.2.2 Perangkat Lunak (Software)	17
3.3 Tahapan Penelitian	19
3.3.1 Analisis Permasalahan	20
3.3.2 Studi Literatur	21
3.3.3 Pengambilan Data	21

	3.3.4 Preprocessing Datasets	22
	3.3.5 Pelabelan <i>Dataset</i>	23
	3.3.6 Klasifikasi Sentimen	24
	3.3.7 Evaluasi	24
IV	. HASIL DAN PEMBAHASAN	25
4	4.1 Deskripsi Data	25
	4.1.1 Scraping Data Platform X	25
	4.1.2 Scraping Data Platform Tiktok	26
4	4.2 Hasil <i>Preprocessing</i> Data	26
	4.2.1 Tahap Cleaning	27
	4.2.2 Tahap Tokenization	28
	4.2.3 Tahap Stopword Removal	28
	4.2.4 Tahap Case Folding	29
	4.2.5 Tahap Stemming/Lemmatization	30
	4.2.6 Tahap Pembobotan Term Frequency (TF)	30
	4.2.7 Pembuangan Data Kosong	30
	4.2.8 Tahap Pelabelan Data	30
	4.2.8.1 Evaluasi Pelabelan Sentimen: Peneliti vs Volunteer	31
	4.2.8.2 Kesimpulan Evaluasi Pelabelan	33
2	4.3 Hasil Pelabelan Sentimen	33
	4.3.1 Klasifikasi Hasil Sentimen Pada Naive Bayes	34
	4.3.1.1 Datasets X(Twitter)	34
	4.3.1.2 Tiktok 946 Datasets	35
	4.3.1.3 Tiktok 2940 Datasets	36
	4.3.2 Klasifikasi Hasil Sentimen Pada Logistic Regression	37
	4.3.2.1 Datasets X (Twitter)	38
	4.3.2.2 Tiktok 946 Datasets	39
	4.3.2.3 Tiktok 2940 Datasets	40
4	4.4 Hasil Pelatihan Model	41
4	4.5 Evaluasi dan Perbandingan Model	42
	4.5.1 Evaluasi Cross-Validation	45
	4.5.2 Visualisasi Word Cloud Sentimen	47
	4.5.3 Analisis Ketidakseimbangan Data	52

4.5.4 Analisis Kesalahan Model	53
4.6 Interpretasi Hasil	54
4.7 Visualisasi Perbandingan Accuracy Model	54
V. SIMPULAN DAN SARAN	56
5.1 Kesimpulan	56
5.2 Saran	57
DAFTAR PUSTAKA	58

DAFTAR TABEL

Tabel	Halaman
1. Penelitian Terkait.	4
2. Evaluasi Confusion Matrix Multi-Class.	12
3. Spesifikasi Lingkungan Google Colab.	18
4. Perangkat Lunak dan <i>Library</i> yang Digunakan	19
5. Hasil Cleaning Datasets	27
6. Hasil Stopword Removal.	28
7. Hasil Penggunaan Case Folding.	29
8. Kesimpulan Evaluasi Pelabelan.	33
9. Perbandingan Model Naïve Bayes dan Logistic Regression Datasets X	<i></i> 42
10. Perbandingan Model Naïve Bayes dan Logistic Regression Tiktok 94	6
Datasets	43
11. Perbandingan Model Naïve Bayes dan Logistic Regression Tiktok 29	940
Datasets	44

DAFTAR GAMBAR

Gambar Halar	nan
1.Tahapan Text Processing	10
2. Tahapan Penelitian.	20
3. Instalasi <i>Library Node.Js</i> .	25
4. Tahapan Scrapping Data X.	26
5. Tahapan Scrapping Data Tiktok.	26
6. Tahap Cleaning Text.	27
7. Implementasi <i>Library nltk</i>	28
8. Implementasi <i>Library</i> Sastrawi.	29
9. Implementasi Case Folding.	29
10. Implementasi Stemmer Factory.	30
11. Confusion Matrix Datasets X(Twitter) Naive Bayes	35
12. Confusion Matrix 946 Datasets Tiktok Naive Bayes	36
13. Confusion Matrix 2940 Datasets Tiktok Naive Bayes	37
14. Confusion Matrix Datasets X(Twitter) Logistic Regression	38
15. Confusion Matrix 946 Datasets Tiktok Logistic Regression	39
16. Confusion Matrix 2940 Datasets Tiktok Logistic Regression	40
17. Grafik Accuracy Perbandingan Iterasi Tiap Datasets.	41
18. Grafik Cross Validation Model Data X.	45
19. Grafik Cross Validation Model 946 Data Tiktok.	46
20 Cross Validation Per Fold	47

21. Visualisasi Word Cloud Sentimen Positif 946 Data Tiktok	48
22. Visualisasi Word Cloud Sentimen Positif 2940 Data Tiktok	48
23. Visualisasi Word Cloud Sentimen Positif Data X.	49
24. Visualisasi Word Cloud Sentimen Netral 946 Data Tiktok	49
25. Visualisasi Word Cloud Sentimen Netral 2940 Data Tiktok	50
26. Visualisasi Word Cloud Sentimen Netral Data X.	50
27. Visualisasi Word Cloud Sentimen Negatif 946 Data Tiktok	51
28. Visualisasi Word Cloud Sentimen Negatif 2940 Data Tiktok	51
29. Visualisasi Word Cloud Sentimen Negatif Data X	51
30. Grafik Perbandingan Accuracy Model Data X.	54
31. Grafik Perbandingan <i>Accuracy</i> Model 946 Data Tiktok	55
32. Grafik Perbandingan Accuracy Model 2940 Data Tiktok	55

I. PENDAHULUAN

1.1. Latar Belakang

Di era digital, media sosial berperan signifikan dalam membentuk opini publik. Salah satu bentuk interaksi yang umum adalah penggunaan hashtag atau tagar sebagai sarana menyampaikan pandangan, menggerakkan percakapan, dan membentuk tren diskusi. Tagar menjadi media komunikasi yang efektif untuk menyuarakan dukungan, kritik, maupun aspirasi terhadap suatu isu.

Salah satu menarik perhatian publik adalah tagar yang sempat #KABURAJADULU, yang banyak digunakan oleh masyarakat Indonesia, khususnya yang bekerja atau tinggal di luar negeri. Unggahan dengan tagar ini umumnya menggambarkan pengalaman dan pandangan terkait kehidupan di luar negeri yang dianggap lebih baik dibandingkan di Indonesia. Fenomena ini memunculkan perbedaan sudut pandang antara masyarakat yang berada di luar negeri dan pihak-pihak di dalam negeri, termasuk pemerintah, sehingga memicu diskusi pro dan kontra di media sosial (Fransiska Vina Sari & Arief Wibowo, 2019). Analisis sentimen terhadap tagar ini penting dilakukan untuk memahami persepsi publik dan dinamika diskusi yang terjadi. Namun, proses klasifikasi opini di media sosial memiliki tantangan tersendiri, seperti penggunaan bahasa informal yang tidak baku, konteks kalimat yang ambigu, serta distribusi sentimen yang tidak seimbang. Tantangan tersebut menuntut pemilihan metode pemrosesan bahasa alami (Natural Language Processing/NLP) dan algoritma klasifikasi yang tepat agar hasil analisis akurat dan representatif.

Berbagai penelitian terdahulu telah memanfaatkan metode NLP untuk analisis sentimen menggunakan beragam algoritma, seperti *Naïve Bayes, Support Vector Machine* (SVM), dan *Logistic Regression*. Akan tetapi, penelitian yang secara khusus membandingkan performa *Naïve Bayes* dan *Logistic Regression* pada konteks bahasa informal Indonesia, khususnya pada data dari media sosial seperti

Twitter(X) dan TikTok dengan isu spesifik seperti #KABURAJADULU, masih terbatas (Syahputra & Wibowo, 2023). Natural Language Processing (NLP) merupakan bidang penelitian yang berfokus pada bagaimana komputer dapat memahami dan memproses bahasa manusia secara efektif (Astiningrum et al., 2018). NLP banyak diterapkan dalam berbagai bidang, termasuk analisis sentimen, terjemahan otomatis, serta pencarian informasi. Pada penelitian sebelumnya yang dilakukan oleh (Prasetya et al., 2024) dengan judul "Implementasi NLP (Natural Language Processing) Dasar pada Analisis Sentimen Review Spotify", telah dilakukan analisis sentimen terhadap ulasan aplikasi Spotify di Google Play Store menggunakan metode NLP untuk memahami pola opini pengguna terhadap aplikasi tersebut.

Berdasarkan kondisi tersebut, penelitian ini dilakukan untuk menerapkan NLP dasar dalam analisis sentimen terhadap unggahan yang menggunakan tagar #KABURAJADULU di media sosial, sekaligus membandingkan kinerja *Naïve Bayes* dan *Logistic Regression* dalam mengklasifikasikan sentimen positif, negatif, dan netral. Hasil penelitian diharapkan dapat memberikan gambaran yang lebih jelas mengenai opini publik terhadap isu ini serta menawarkan pendekatan analisis yang dapat diterapkan pada isu-isu serupa di masa depan.

1.2. Rumusan Masalah

Berdasarkan latar belakang, rumusan masalah pada penelitian ini adalah sebagai berikut:

- Apa saja tantangan dalam mengidentifikasi dan mengklasifikasikan opini publik pada unggahan media sosial bertagar #kaburajadulu menggunakan metode NLP?
- 2. Bagaimana perbandingan akurasi dan performa algoritma Naive Bayes dan Logistic Regression dalam klasifikasi sentimen unggahan dengan tagar #kaburajadulu?

1.3. Batasan Masalah

Adapun batasan masalah pada penelitian ini antara lain adalah sebagai berikut:

- Dataset yang digunakan dalam penelitian ini adalah berupa unggahan pengguna yang terbatas pada unggahan media sosial yang menggunakan tagar #kaburajadulu.
- Penelitian ini akan fokus pada implementasi Natural Language Processing
 (NLP) dasar. Teknik Natural Language Processing (NLP) yang lebih
 canggih, seperti deep learning atau model transformer, mungkin tidak
 dieksplorasi secara mendalam.

1.4. Tujuan

Adapun tujuan dari penelitian ini adalah sebagai berikut:

- 1. Menganalisis opini publik terhadap fenomena sosial yang tercermin melalui unggahan bertagar #KABURAJADULU menggunakan pendekatan *Natural Language Processing* (NLP) dasar.
- Mengidentifikasi dan mengklasifikasikan sentimen pengguna (positif, negatif, netral) terhadap berbagai aspek yang berkaitan dengan tagar #KABURAJADULU, serta memahami pola opini publik terhadap isu yang diangkat.

1.5. Manfaat

Adapun manfaat dari penelitian ini antara lain sebagai berikut:

- Memberikan pemahaman tentang kinerja metode Natural Language Processing (NLP) dasar dalam analisis sentimen ulasan tagar #kaburajadulu.
- 2. Memberikan informasi kepada pemangku kepentingan atau pihak terkait mengenai metode NLP dasar yang lebih efektif untuk menganalisis sentimen pengguna terhadap tagar #KABURAJADULU, sehingga dapat digunakan sebagai dasar dalam memahami opini publik dan tren yang berkembang di media sosial.
- 3. Menyediakan dasar untuk penelitian lebih lanjut dalam pengembangan dan penerapan teknik *machine learning* pada analisis sentimen, khususnya untuk pro kontra dan juga tagar yang sedang naik daun di media sosial saat ini.

II. TINJAUAN PUSTAKA

2.1 Penelitian Terkait

Penelitian sebelumnya mengenai implementasi *Natural Language Processing* (NLP) untuk analisis sentimen telah banyak dilakukan dengan menggunakan berbagai metode dan data yang berbeda. Sebagai landasan teoritis, tabel 1 berikut adalah beberapa penelitian sebelumnya yang relevan dengan topik analisis sentimen berbasis *Natural Language Processing* (NLP), yang digunakan untuk membandingkan metode, pendekatan, serta efektivitas algoritma dalam konteks data berbeda.

Tabel 1. Penelitian Terkait.

Peneliti	Judul	Metode	Hasil
Marsha et al. (2024)	Implementasi NLP	Naïve Bayes	Logistic Regression
	Dasar pada Analisis	dan <i>Logistic</i>	memberikan akurasi
	Sentimen Review	Regression	tertinggi sebesar 79% ,
	Spotify		sementara <i>Naïve Bayes</i>
			hanya 76% , karena
			Logistic Regression
			menghitung bobot
			berdasarkan probabilitas fitur.
Astiningrum et al.	Implementasi NLP	Levenshtein	Presisi mencapai 94%
(2018)	dengan Konversi	Distance, TF-	dan <i>recall</i> 85% dalam
	Kata pada Sistem	IDF, dan	pengujian 60
	Chatbot Konsultasi	Cosine	pertanyaan. <i>User</i>
	Laktasi	Similarity	Acceptance Test
			menunjukkan 77%
			pengguna menyatakan
			setuju atau sangat setuju
			terhadap kelayakan
			sistem.
Jovanica et al.	Analisis Pengaruh	Social Network	Ditemukan 12 aktor dan
(2022)	Aktor pada Tagar	Analysis	40 interaksi (ties) dalam
	#roketchina	(SNA), Ucinet,	penggunaan tagar
	menggunakan	dan teori	#roketchina. Data
	Social Network	graph.	dianalisis untuk
	Analysis (SNA)		memetakan struktur
			jaringan dan hubungan
			antar pengguna.

Tabel 2. Penelitian Terkait.

Lidinillah et al.
(2023)

Analisis Sentimen di
Twitter terhadap
Layanan Steam

SVM

Regression dan
SVM

lebih tinggi dibanding
Logistic Regression
dalam klasifikasi
sentimen pengguna
terhadap layanan digital.

2.2 Media Sosial

Media sosial telah menjadi ruang publik digital yang memungkinkan penggunanya berinteraksi, membentuk jaringan sosial, serta menyuarakan opini secara bebas dan cepat (Kim et al., 2013). Tidak hanya sebagai alat komunikasi, media sosial kini berperan sebagai wadah untuk menyampaikan pendapat terhadap isu-isu sosial dan politik, serta memobilisasi opini publik melalui berbagai fitur, salah satunya adalah tagar (hashtag). Tagar #kaburajadulu merupakan contoh fenomena digital yang digunakan masyarakat, khususnya anak muda, untuk menyuarakan ketidakpuasan terhadap kondisi pemerintahan, pendidikan, dan budaya korupsi di Indonesia. Namun, media sosial juga dapat menjadi tempat berkembangnya ekspresi negatif seperti ujaran kebencian atau konflik SARA (Yulianto et al., 2018). Oleh karena itu, penting dilakukan analisis terhadap sentimen di balik penggunaan tagar tersebut untuk memahami pola opini publik yang berkembang secara daring.

2.3 Natural Language Processing (NLP)

Natural Language Processing (NLP) adalah bidang penelitian dan aplikasi yang mengeksplorasi bagaimana komputer dapat digunakan untuk memahami dan memanipulasi teks bahasa alami. Penelitian Natural Language Processing (NLP) bertujuan untuk mengumpulkan pengetahuan tentang bagaimana manusia memahami dan menggunakan bahasa sehingga peralatan dan teknik pemasangan dapat dikembangkan untuk membuat sistem komputer memahami dan memanipulasi bahasa alami untuk melakukan tugas yang disukai (Syahputra & Wibowo, 2023).

Sebagai contoh nyata penerapan NLP dalam analisis sentimen, (Hindarto, 2023) melakukan studi komparatif antara metode tradisional seperti *Naive Bayes* dan *Support Vector Machine* (SVM) *versus model neural network Multi-Layer Perceptron* (MLP) dalam mengklasifikasi sentimen ulasan layanan pelanggan.

Hasilnya menunjukkan model MLP mencapai akurasi penuh (100%), sedangkan *Naive Bayes* hanya 68,8% dan SVM 87,5%, menunjukkan bahwa jaringan saraf lebih unggul dalam menangani konten yang kompleks.

2.4 Naive Bayes

Metode *Naive Bayes* merupakan metode klasifikasi dengan nilai yang ditentukan oleh probabilitas data (Prasetya et al., 2024). Klasifikasi *Naive Bayes* merupakan metode klasifikasi berbasis *supervised learning*. Metode ini memiliki keunggulan dalam waktu klasifikasi yang relatif singkat, sehingga efisien digunakan dalam pemrosesan sistem analisis sentimen. Persamaan dasar dari algoritma *Naive Bayes* dituliskan pada persamaan 1 berikut:

$$P(A|x) \propto P(A) \prod_{f \in x} P(f|A)$$
 (1)

Persamaan dijelaskan dengan:

- 1. P(A|x) adalah probabilitas teks x termasuk dalam kelas A (positif, negatif, atau netral).
- 2. P(A) adalah probabilitas awal dari kelas A (disebut *prior*).
- 3. P(f|A) adalah probabilitas suatu fitur (kata) f muncul dalam teks yang termasuk kelas A.
- Produk ∏ dihitung karena dalam *naive bayes* diasumsikan bahwa semua fitur saling independen (meskipun ini tidak selalu terjadi dalam kenyataan, asumsi ini justru membuat perhitungan lebih sederhana).

Persamaan dasar dari algoritma *Naive Bayes* secara umum tidak digunakan secara eksplisit dalam penelitian ini. Namun demikian, beberapa varian dan bentuk formulasi dari metode ini tetap relevan dan dapat dipertimbangkan sebagai alternatif pendekatan dalam proses klasifikasi teks. Penelitian ini secara khusus memilih pendekatan dalam bentuk logaritmik karena dinilai lebih efisien dan stabil dalam menangani nilai probabilitas yang sangat kecil (Ceri et al., 2013).

Adapun beberapa formulasi varian *Naive Bayes* yang sering digunakan dan relevan dalam konteks pengolahan teks adalah sebagai berikut:

- a. Rumus Probabilistik Penuh
- b. Bentuk Logaritmik

- c. Laplace Smoorthing
- d. Kombinasi dengan distribusi multinomial dan gaussian

Dalam penelitian ini, pendekatan yang digunakan adalah varian logaritmik dari algoritma *Naive Bayes*, karena metode ini dinilai paling sesuai untuk klasifikasi teks berbasis frekuensi kata. Bentuk logaritmik dari algoritma tersebut direpresentasikan melalui persamaan 2 berikut:

$$\log P(A|x) = \log P(A) + \sum_{f \in x} \log P(f|A)$$
 (2)

Rumus *logaritmik* pada algoritma *Naive Bayes* digunakan dalam penelitian ini karena memiliki keunggulan dalam hal kestabilan numerik, khususnya untuk menghindari terjadinya *underflow* akibat perkalian banyak nilai probabilitas kecil yang umum terjadi dalam data teks berdimensi tinggi. Selain itu, bentuk logaritmik juga memberikan efisiensi yang lebih baik dalam implementasi karena telah tersedia dalam pustaka (*library*) *scikit-learn*, yang mengintegrasikan algoritma ini secara langsung dan efisien (Barupal & Fiehn, 2019).

Rumus-rumus lanjutan seperti *Laplace Smoothing* memang tidak diterapkan secara eksplisit dalam skripsi ini, namun telah secara otomatis diimplementasikan oleh *MultinomialNB*, yaitu varian dari *Naive Bayes* yang digunakan. Oleh karena itu, tidak dilakukan perhitungan manual terhadap probabilitas fitur, karena data yang digunakan memiliki ukuran yang cukup besar dan kompleks. Melakukan perhitungan manual dalam konteks ini dinilai tidak efisien serta berpotensi meningkatkan risiko kesalahan perhitungan. Fokus utama dari penelitian ini adalah mengevaluasi dan membandingkan performa dua model klasifikasi, *Naive Bayes* dan *Logistic Regression*, dalam konteks klasifikasi sentimen berbasis teks media sosial, bukan pada aspek formulasi probabilistik secara matematis. Oleh karena itu, penggunaan fungsi-fungsi bawaan dari pustaka *machine learning* yang telah terverifikasi dianggap lebih tepat dan relevan dengan tujuan penelitian.

2.5 Logistic Regression

Logistic Regression merupakan sebuah model statistik yang digunakan untuk menentukan apakah independent variable mempunyai pengaruh terhadap sebuah binary dependent variable (Prasetya et al., 2024). Logistic Regression memiliki sejumlah keunggulan yang membuatnya menjadi pilihan yang baik dalam berbagai

tugas klasifikasi. Salah satu keunggulan utamanya adalah interpretabilitas model, variable atau hasilnya mudah diinterpretasikan, memungkinkan pemahaman yang jelas tentang hubungan antara fitur dan probabilitas kelas yang dihasilkan. Metode ini juga termasuk ke dalam algoritma klasifikasi. Dalam penelitian ini, bentuk dasar *logistic regression* digunakan sebagaimana dinyatakan dalam persamaan sebagai berikut:

$$P(Y) = \frac{1}{(1 + e - \beta_0 + \beta_1 X_1 + \dots + \beta_k X_k)}$$
(3)

Keterangan:

P(Y) =: Probabilitas variabel dependen Y sama dengan 1

e : Basis logaritma natural

 $X_1, X_2, \dots X_k$: Variable-variable predictor

 β_0 : *Intercept* (konstanta)

 $\beta_1,\beta_2\cdots\cdots\beta_k$: Koefisien yang menunjukkan seberapa besar perubahan dalam *variable independent*.

Tujuan utama dari model regresi logistik adalah mengestimasi parameter $\beta_1, \beta_2 \cdots \beta_k$ yang dilakukan menggunakan pendekatan *Maximum Likelihood Estimation* (MLE). Metode ini mencari nilai parameter yang memaksimalkan probabilitas kemunculan data yang diamati. Setiap koefisien β_i menunjukkan besarnya perubahan pada logit ($log \ odds$) akibat perubahan satu unit pada variabel prediktor X_i , dengan asumsi variabel lain konstan. Dengan demikian, model ini sangat berguna dalam mengukur derajat hubungan antara masing-masing variabel independen terhadap hasil (outcome) yang bersifat biner (Boateng & Abaye, 2019).

2.6 Analisis Sentimen

Analisis sentimen merupakan salah satu bidang dalam pemrosesan bahasa alami *Natural Language Processing* (NLP) yang bertujuan untuk mengidentifikasi opini, sentimen, dan emosi yang terkandung dalam bentuk teks (Ratnawati, 2018). Proses ini digunakan untuk mengkategorikan pendapat atau ulasan terhadap suatu topik menjadi tiga label utama, yaitu positif, negatif, dan netral.

Dalam pelaksanaannya, terdapat dua pendekatan utama dalam proses pelabelan data. Pertama, *manual labeling*, yaitu metode pelabelan yang dilakukan secara

manual menggunakan bantuan kamus kata (*lexicon*) berisi daftar kata positif, netral, dan negatif. Setiap kalimat dianalisis dengan mencocokkan kata-kata dalam teks terhadap daftar tersebut untuk menentukan kategori sentimen yang sesuai. (Mäntylä et al., 2018). Kedua, pelabelan otomatis (*library python*), yaitu proses pemberian label sentimen secara langsung pada setiap kalimat berdasarkan makna dan konteksnya dengan bantuan algoritma bawaan dari *library* yang digunakan. Pendekatan pelabelan otomatis dalam penelitian ini melibatkan dua *volunteer* serta peneliti, yang masing-masing menggunakan bantuan *library Python* untuk menentukan label sentimen. Proses pelabelan dilakukan dengan pendekatan berbasis leksikal menggunakan Sastrawi untuk *preprocessing* dan *VADER* untuk analisis sentimen pada teks berbahasa Indonesia, serta pendekatan berbasis transformer menggunakan model *Sentiment-analysis pipeline* dari *Hugging Face Transformers*. Hasil pelabelan dari masing-masing pihak kemudian dibandingkan guna mengevaluasi konsistensi klasifikasi sentimen, terutama dalam menangani konteks emosional dan ironi dalam teks media sosial (Amal & Jayanta, 2023).

2.7 Text Processing

Metode *Preprocessing* memainkan peran yang sangat penting dalam teknik dan aplikasi *text mining*. Ini adalah langkah pertama dalam proses penambangan teks. Dalam makalah ini, penelitian ini membahas tiga langkah utama dari *preprocessing*, yaitu *cleaning*, *case folding*, *stemming*, dan *term frequency* (TF) (Astiningrum et al., 2018).

Major Malanat According Total Policy Thermaly Total Control (Control (C

Gambar 1 berikut merupakan komponen text preprocessing:

Gambar 1. Tahapan Text Processing

- 1. *Normalization* umumnya mengacu pada serangkaian tugas terkait yang dimaksudkan untuk mengkonversi semua teks.
- 2. *Tokenization* adalah langkah yang membagi string teks yang lebih panjang menjadi potongan-potongan yang lebih kecil, atau token.
- 3. *Case folding* yaitu proses ini akan mengubah seluruh teks pada dokumen menjadi bentuk standar yaitu huruf kecil atau *lowercase*.
- 4. *Stemming* adalah teknik prapemrosesan teks dalam *Natural Language Processing* (NLP). Secara khusus, ini adalah proses mengurangi bentuk infleksi dari sebuah kata menjadi satu yang disebut "*stem*", atau bentuk dasar, yang juga dikenal sebagai "*lema*" dalam linguistik.
- 5. Perhitungan *Term Frequency* (tf) menggunakan persamaan = ij Dengan tf adalah *term frequency*, dan adalah banyaknya kemunculan term dalam dokumen, *Term frequency* (tf) dihitung dengan menghitung banyaknya kemunculan term dalam dokumen.

2.8 Pembobotan TF

Metode *Term Frequency* (TF) adalah cara pemberian bobot hubungan suatu kata (*term*) terhadap dokumen (Amrizal, 2018).

- a. Perhitungan *Term Frequency* (tf) menggunakan persamaan = ij Dengan tf adalah *term frequency*, dan adalah banyaknya kemunculan *term* dalam dokumen, *Term frequency* (tf) dihitung dengan menghitung banyaknya kemunculan term dalam dokumen.
- b. Perhitungan *term frequency Inverse Document Frequency* (tf idf), menggunakan persamaan 3 berikut:

$$w_{ij} = t f_{ij} \times \log \frac{N}{n} \cdots \cdots \cdots (3)$$

Keterangan:

wij = bobot kata/term tj terhadap dokumen di

tfij = jumlah kemunculan kata/term tj dalam di

N = jumlah semua dokumen yang ada dalam *database*

n = jumlah dokumen yang mengandung kata/ $term\ tj$ (minimal ada satu kata yaitu $term\ tj$).

2.9 Tagar (Hashtag)

Tagar (hashtag) adalah simbol "#" yang digunakan dalam media sosial untuk mengelompokkan atau mengidentifikasi suatu topik tertentu. Penggunaan tagar memungkinkan pengguna untuk mencari dan mengikuti percakapan mengenai topik yang sama, menjadikannya alat yang ampuh dalam menyebarkan informasi, menggerakkan opini publik, serta meningkatkan interaksi dalam media sosial (Fahrezi et al., 2022). Sementara itu, penelitian yang dilakukan oleh (Jovanica et al., 2022) menunjukkan bahwa tagar juga berperan dalam membentuk opini publik. Dalam studinya tentang penggunaan tagar #roketchina di X, mereka menemukan bahwa tagar dapat menciptakan pola interaksi yang memperlihatkan bagaimana individu dan kelompok berpartisipasi dalam diskusi publik serta bagaimana informasi tersebar di media sosial.

2.10 Confusion Matrix

Confusion matrix merupakan salah satu alat evaluasi yang paling umum dan efektif untuk menilai kinerja algoritma klasifikasi dalam tugas klasifikasi tunggal (single-label classification), di mana setiap data hanya diklasifikasikan ke dalam satu kelas. Matrix ini menyajikan secara jelas jumlah prediksi yang benar maupun salah dalam bentuk True Positive (TP), False Positive (FP), True Negative (TN), dan False Negative (FN), yang kemudian digunakan untuk menghitung metrik evaluasi seperti akurasi, presisi, recall, dan F1-score (Heydarian et al., 2022). Berikut tabel 2 merupakan evaluasi kinerja model klasifikasi multi-kelas:

Keterangan:

- a. Baris = Label aktual
- b. Kolom = Prediksi model

Tabel 3. Evaluasi Confusion Matrix Multi-Class.

	Pred: A	Pred: B	Pred: C
Actual: A	TP	FN	FN
Actual: B	FP	TP	FN
Actual: C	FP	FP	TP

Berdasarkan tabel 2, terlihat bahwa nilai *True Negative* (TN) tidak ditampilkan secara eksplisit, karena dalam klasifikasi multi-kelas tidak terdapat dikotomi yang jelas antara "positif" dan "negatif" seperti pada klasifikasi biner. Meskipun demikian, nilai TN tetap dapat dihitung menggunakan pendekatan *one-vs-rest*, yaitu dengan menganggap satu kelas sebagai "positif", dan dua kelas lainnya sebagai "negatif". Rumus untuk menghitung TN untuk kelas ke-*k* tertulis pada persamaan 4 berikut:

$$TNk = Total \ data - TPk - FPk - FNk \tag{4}$$

Dengan diketahuinya nilai TP, FP, FN, dan TN untuk masing-masing kelas, maka metrik-metrik evaluasi seperti *precision, recall*, dan *F1-score* dapat dihitung secara menyeluruh untuk mengukur kinerja model klasifikasi. Berdasarkan penelitian

yang dilakukan oleh (Rahmad et al., 2020) evaluasi performa model klasifikasi dilakukan dengan menggunakan beberapa metrik utama

Berikut ini merupakan penjelasan masing-masing metrik beserta rumus yang digunakan:

1) Precision:

Presisi mengukur ketepatan model dalam mengklasifikasikan data sebagai positif. Dalam konteks multi-kelas, presisi dihitung untuk setiap kelas dengan cara membandingkan jumlah prediksi positif yang benar dengan jumlah seluruh prediksi positif (baik benar maupun salah). Untuk mengukur presisi dibutuhkan rumus pada persamaan 5 sebagai berikut:

$$Precision = \left(\frac{TP}{TP + FP}\right) \times 100\% \tag{5}$$

2) Recall:

Recall mengukur seberapa banyak data positif yang berhasil diidentifikasi oleh model dari seluruh data positif yang ada. Pada klasifikasi multi-kelas, *recall* dihitung per kelas dengan membandingkan jumlah data positif yang benar dengan jumlah seluruh data positif aktual. Untuk mengukur *recall* dibutuhkan rumus pada persamaan 6 sebagai berikut:

$$Recall = \left(\frac{TP}{TP + FN}\right) \times 100\% \tag{6}$$

3) Accuracy:

Akurasi menunjukkan seberapa banyak prediksi model yang benar dibandingkan dengan total data yang ada. Akurasi dihitung dengan membandingkan jumlah *True Positive* (TP) dan *True Negative* (TN) terhadap seluruh data yang diklasifikasikan. Untuk mengukur akurasi dibutuhkan rumus pada persamaan 7 sebagai berikut:

$$Accuracy = \left(\frac{TP + TN}{TP + TN + FP + FN}\right) \times 100\% (7)$$

4) F1-score:

F1-score adalah metrik yang merupakan *harmonic mean* dari presisi dan *recall. F1-score* memberikan gambaran seimbang antara kedua metrik ini, sangat berguna ketika ada ketidakseimbangan antara presisi dan *recall.* Untuk mengukur *F1-Score* dibutuhkan rumus pada persamaan 8 sebagai berikut:

$$F1 - Score = 2 \times \left(\frac{Precision \times Recall}{Precision + Recall}\right)$$
 (8)

Untuk menghitung metrik evaluasi pada multi-kelas, biasanya digunakan dua pendekatan:

- a) *Macro Average*: Menghitung presisi, dan F1 untuk setiap kelas, lalu dirataratakan. Cocok untuk kasus di mana distribusi kelas seimbang.
- b) *Micro Average*: Menggabungkan seluruh TP, FP, dan FN terlebih dahulu, lalu dihitung metrik evaluasinya. Cocok untuk kasus di mana distribusi kelas tidak seimbang.

Perhitungan dimulai dari pembentukan *confusion matrix* multi-kelas yang merepresentasikan distribusi hasil prediksi model terhadap label sebenarnya. Pada kasus ini, dataset memiliki tiga kelas sentimen, yaitu positif (A), netral (B), dan negatif (C). *Confusion matrix* disusun dalam bentuk tabel dua dimensi di mana baris merepresentasikan kelas aktual, sedangkan kolom menunjukkan kelas hasil prediksi. Setiap sel dalam matriks menunjukkan jumlah data yang sesuai dengan kombinasi kelas aktual dan kelas prediksi tersebut. Tahap selanjutnya adalah melakukan perhitungan metrik untuk masing-masing kelas secara terpisah dengan pendekatan *one-vs-rest*. Pada pendekatan ini, setiap kelas diperlakukan sebagai kelas positif, sementara gabungan dua kelas lainnya dianggap sebagai kelas negatif.

1. Kelas A (Positif)

Untuk menghitung 4 komponen dalam kelas A dibutuhkan nilai nilai sebagai berikut:

- a. TP A: Jumlah data yang benar-benar kelas A dan diprediksi sebagai A.
- b. FP A: Jumlah data dari kelas B atau C yang salah diprediksi sebagai A.
- c. FN A: Jumlah data dari kelas A yang salah diprediksi sebagai B atau C.
- d. TN A: Semua data selain kelas A yang tidak diprediksi sebagai A.

2. Kelas B (Netral)

Untuk menghitung 4 komponen dalam kelas B dibutuhkan nilai-nilai sebagai berikut:

a. TP B:Jumlah data yang benar-benar kelas B dan diprediksi sebagai
 B.

- b. FP B: Jumlah data dari kelas A atau C yang salah diprediksi sebagai B.
- c. FN B: Jumlah data dari kelas B yang salah diprediksi sebagai A atau C.
- d. TN B: Semua data selain kelas B yang tidak diprediksi sebagai B.

3. Kelas C (Negatif)

Untuk menghitung 4 komponen dalam kelas C dibutuhkan nilai-nilai sebagai berikut:

- a. TP C: Jumlah data yang benar-benar kelas C dan diprediksi sebagai C.
- b. FP C: Jumlah data dari kelas A atau B yang salah diprediksi sebagai C.
- c. FN C: Jumlah data dari kelas C yang salah diprediksi sebagai A atau B.
- d. TN C: Semua data selain kelas C yang tidak diprediksi sebagai C.

Setelah diperoleh nilai *accuracy*, *precision*, *recall*, dan *F1-score* untuk masing-masing kelas (kelas A, kelas B, dan kelas C), tahap selanjutnya adalah menghitung nilai rata-rata makro (*macro average*) untuk mendapatkan gambaran kinerja model secara keseluruhan tanpa memperhitungkan proporsi jumlah data di tiap kelas. Perhitungan *macro average* dilakukan dengan cara merata-ratakan nilai metrik dari seluruh kelas menggunakan rumus persamaan 9, 10, 11, dan 12 berikut:

$$Macro \ Precision = \frac{PrecisionA + PrecisionB + PrecisionC}{3}$$
 (9)
$$Macro \ Recall = \frac{RecallA + RecallB + RecallC}{3}$$
 (10)
$$Macro \ F1 - Score = \frac{F1A + F1B + F1C}{3}$$
 (11)
$$Akurasi = \frac{Total \ Prediksi \ Benar \ ABC}{Total \ Data \ ABC}$$
 (12)

Pendekatan ini memberikan bobot yang sama pada setiap kelas, sehingga sangat cocok digunakan ketika distribusi data pada setiap kelas relatif seimbang. Dengan kata lain, setiap kelas memiliki pengaruh yang sama terhadap nilai rata-rata, meskipun jumlah data antar kelas berbeda.

Hasil dari *macro average* ini merepresentasikan kinerja model secara umum pada semua kelas, sehingga memudahkan peneliti untuk melihat seberapa baik model dapat mengklasifikasikan data dari berbagai kategori secara merata

Penelitian lain oleh (Xu et al., 2020) memperkenalkan pendekatan *Three-Way Decisions* (3WD) untuk menangani ketidakpastian dalam klasifikasi. Pendekatan ini mengintegrasikan konsep *confusion matrix* dengan semantik tiga arah, yaitu:

- 1) Positif (data diterima sebagai kelas tertentu),
- 2) Negatif (data ditolak sebagai kelas tertentu),
- Ambigu (data tidak diklasifikasikan secara langsung karena ketidakpastian tinggi).

Penelitian tersebut mengembangkan tujuh mode pengukuran berbasis semantik tiga arah dan membentuk wilayah keputusan berdasarkan fungsi objektif yang fleksibel, disesuaikan dengan preferensi pemangku kepentingan. Hasilnya menunjukkan bahwa model 3WD lebih unggul dibandingkan metode berbasis *Gini Coefficient* dan *Shannon Entropy*, terutama dalam konteks data yang tidak pasti atau kompleks.

Baik pendekatan evaluasi klasifikasi konvensional menggunakan *confusion matrix* maupun pendekatan berbasis *Three-Way Decisions* (3WD) menunjukkan pentingnya penggunaan struktur evaluasi yang tepat untuk memahami dan menilai kinerja model. Kedua pendekatan tersebut juga menekankan perlunya mempertimbangkan ketidakpastian klasifikasi, terutama dalam konteks pengambilan keputusan yang kompleks dan berorientasi pada kepentingan sosial. Evaluasi yang akurat dan menyeluruh menjadi kunci dalam merancang model klasifikasi yang andal dan responsif terhadap kebutuhan nyata di lapangan.

III. METODOLOGI PENELITIAN

3.1 Waktu dan Tempat Penelitian

1. Waktu Penelitian

Penelitian ini dilaksanakan sejak tanggal dikeluarkan nya izin penelitian dalam kurun waktu kurang lebih 4 (empat) bulan, 2 bulan pertama pengumpulan data dan 2 bulan pengolahan data yang meliputi penyajian dalam bentuk skripsi dan proses bimbingan berlangsung.

2. Tempat Penelitian

Tempat pelaksanaan penelitian ini adalah di Jurusan Ilmu Komputer Fakultas Matematika dan Ilmu Pengetahuan Alam, Universitas Lampung.

3.2 Perangkat Penelitian

Penelitian ini dilakukan dengan menggunakan alat dan bahan sebagai berikut:

3.2.1 Perangkat Keras (*Hardware*)

Perangkat keras yang digunakan dalam penelitian ini adalah satu unit laptop *Lenovo* dengan spesifikasi:

a. *Processor* : AMD Ryzen 3 5300U

b. RAM : 8,00 GB

c. GPU : AMD Radeon(TM)

d. Storage : 256GB

3.2.2 Perangkat Lunak (Software)

Perangkat lunak yang digunakan dalam penelitian analisis sentimen ini melibatkan beberapa komponen utama, yaitu sistem operasi Windows 11 sebagai platform lokal, peramban *Google Chrome* untuk mengakses layanan daring, serta *Google Colab* sebagai lingkungan eksekusi kode berbasis *cloud*. Lingkungan *Google Colab* menyediakan sarana eksekusi berbasis

Jupyter Notebook yang telah dilengkapi dengan spesifikasi perangkat keras dan perangkat lunak yang cukup memadai untuk keperluan komputasi ringan hingga menengah.

Tabel 3 berikut menunjukkan spesifikasi teknis dari lingkungan *Google Colab* yang digunakan dalam penelitian ini:

Tabel 4. Spesifikasi Lingkungan Google Colab.

Kategori	Spesifikasi
Runtime	Google Colab (cloud-based Jupyter
Environment	Notebook)
Jenis Instance	Free-tier (standar)
Prosesor (CPU)	2-core Intel Xeon (virtual)
Memori (RAM)	±12.6 GB
Penyimpanan Sementara	±100 GB (temporary runtime disk)
GPU (Opsional)	NVIDIA Tesla K80 / T4 (tergantung ketersediaan saat <i>runtime</i>)
Sistem Operasi Virtual	Ubuntu 22.04 LTS (dalam container)

Selain lingkungan eksekusi, berbagai perangkat lunak dan pustaka (*library*) pendukung juga digunakan untuk mendukung proses analisis data, pelatihan model dan visualisasi hasil.

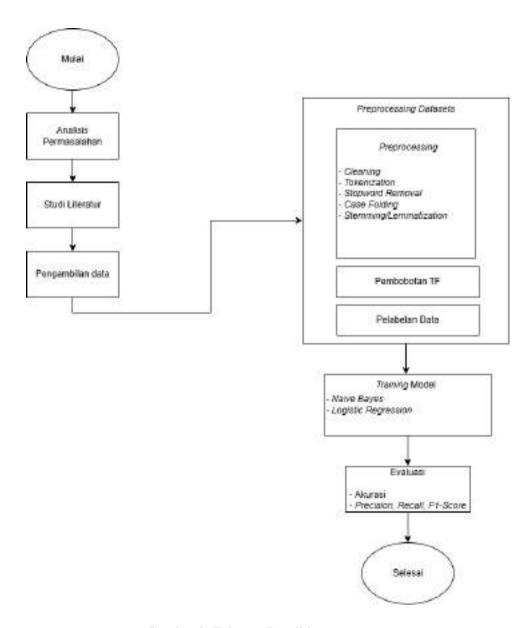
Rincian perangkat lunak yang digunakan disajikan dalam tabel 4 berikut:

Tabel 5. Perangkat Lunak dan Library yang Digunakan.

Komponen	Keterangan
Python Version	3.10 (default di Google Colab)
<i>Library</i> Utama	 nltk untuk preprocessing teks `Sastrawi` untuk stemming Bahasa Indonesia `scikit-learn` untuk klasifikasi (Naive Bayes & Logistic Regression) `matplotlib`, `seaborn` untuk visualisasi `pandas`, `numpy` untuk pengolahan data `transformers` dari HuggingFace (opsional, untuk model BERT) `wordcloud` untuk visualisasi kata

3.3 Tahapan Penelitian

Penelitian ini meliputi beberapa tahapan sistematis dimulai dari analisis permasalahan, studi literatur, dan pengumpulan data dari media sosial bertagar #kaburajadulu. *Datasets* yang dikumpulkan kemudian diproses melalui tahap *preprocessing*, seperti *cleaning*, tokenisasi, *stopword removal*, *case folding*, dan *stemming/lemmatization*. Setelah itu dilakukan pembobotan *Term Frequency* (TF) dan pelabelan sentimen secara otomatis. Data yang telah siap kemudian digunakan untuk melatih model menggunakan algoritma *Naïve Bayes* dan *Logistic Regression*, lalu diuji menggunakan metrik evaluasi seperti akurasi, *precision*, *recall*, *dan F1-score* guna menentukan performa terbaik dalam klasifikasi sentimen.



Gambar 2. Tahapan Penelitian.

3.3.1 Analisis Permasalahan

Pada tahap awal penelitian ini, dilakukan analisis terhadap unggahan media sosial yang menggunakan tagar #KABURAJADULU. Analisis ini bertujuan untuk memahami opini dan sentimen publik terkait isu yang dikaitkan dengan tagar tersebut. Data unggahan akan dikumpulkan dari platform media sosial seperti X dan TikTok guna mendapatkan informasi yang dapat digunakan dalam tahap selanjutnya. Topik penelitian ini berfokus pada membandingkan dua metode analisis sentimen, yaitu $Naive\ Bayes\ dan\ Logistic\ Regression$, dalam mengklasifikasikan sentimen unggahan dengan

tagar #KABURAJADULU. Perbandingan kedua metode ini bertujuan untuk menentukan metode yang lebih efektif dalam mengklasifikasikan sentimen positif, negatif, dan netral berdasarkan teks unggahan pengguna. Hasil penelitian ini diharapkan dapat memberikan wawasan mengenai akurasi dan kinerja kedua metode dalam analisis sentimen berbasis *Natural Language Processing* (NLP).

3.3.2 Studi Literatur

Studi ini bertujuan untuk mendapatkan fondasi teoritis dan praktis yang mendukung analisis sentimen terhadap unggahan media sosial bertagar #KABURAJADULU. Beberapa elemen utama dibahas dalam penelusuran literatur. Ini termasuk metode untuk menganalisis sentimen dalam konteks opini publik di media sosial, penggunaan teknik klasifikasi teks berbasis pembelajaran mesin, dan perbandingan efektivitas pengolahan data teks antara algoritma *Naive Bayes* dan *Logistic Regression*. Selain itu, penelitian sebelumnya juga melihat penggunaan tagar sebagai cara untuk menyampaikan aspirasi masyarakat terhadap masalah sosial-politik, seperti ketidakpuasan terhadap pemerintahan, pendidikan, dan kebijakan publik di Indonesia.

3.3.3 Pengambilan Data

Pengumpulan data dalam penelitian ini dilakukan dengan mengambil unggahan dari dua *platform* media sosial populer, yaitu X (sebelumnya *Twitter*) dan *TikTok*, yang menggunakan tagar #kaburajadulu. Tagar ini dipilih karena mencerminkan fenomena sosial yang sedang ramai dibicarakan, terutama berkaitan dengan keresahan generasi muda terhadap kondisi negara, seperti lemahnya sistem pendidikan, maraknya kasus korupsi, dan ketidakstabilan ekonomi. Data yang dikumpulkan berupa unggahan teks (*tweet* atau komentar) yang mengandung opini, keluhan, atau harapan masyarakat terhadap kehidupan sosial-politik di Indonesia. Proses pengambilan data dilakukan secara otomatis (*scraping*) dengan menggunakan bahasa pemrograman *Python*, yang didukung oleh pustaka *snscrape* untuk *platform* X. Sementara itu, pengambilan data komentar TikTok dilakukan melalui platform *cloud Apify*, menggunakan *bot* yang

tersedia di *marketplace* mereka, dengan metode input berupa *URL* konten TikTok yang relevan, agar sistem dapat mengakses dan mengunduh komentar-komentar dari pengguna.

Sebagai bagian dari pendekatan eksperimental, data yang dikumpulkan kemudian dibagi ke dalam tiga kelompok *dataset*, yaitu:

- a. Dataset X (Twitter): sebanyak 1.500 data unggahan
- b. *Dataset TikTok* A: sebanyak 3.000 data komentar dari berbagai konten dengan tagar terkait
- c. *Dataset TikTok* B: sebanyak 1.000 data komentar dari konten yang berbeda untuk keperluan validasi dan generalisasi model

Total keseluruhan data yang berhasil dikumpulkan mencapai 5.500 unggahan. Proses pengumpulan dilakukan dalam rentang waktu 3 bulan, mulai dari bulan Februari, Maret, dan April, yang bertujuan untuk menjaga relevansi data dan menangkap dinamika opini publik pada saat tagar #kaburajadulu sedang *viral*. *Dummy* data yang diperoleh dari ketiga *dataset* ini kemudian digunakan untuk eksperimen analisis sentimen, guna mengukur kinerja algoritma klasifikasi (*Naïve Bayes* dan *Logistic Regression*) secara terpisah pada tiap *dataset* serta mengevaluasi konsistensi hasil pelabelan sentimen lintas platform dan skala data.

3.3.4 Preprocessing Datasets

Setelah proses pengambilan data selesai, Tahap berikutnya adalah pengolahan kumpulan data, yang mencakup proses pembersihan dan normalisasi teks. Tujuan dari proses ini adalah untuk menghilangkan gangguan dan menyediakan data agar dapat diformat sesuai dengan model analisis sentimen. Menggabungkan teks menjadi huruf kecil (*case folding*), memisahkan kalimat menjadi kata-kata (*tokenisasi*), menghapus karakter yang tidak penting, dan menghapus kata-kata umum yang tidak penting (*stopword removal*).

3.3.5 Pelabelan Dataset

Tahap pelabelan *dataset* merupakan bagian penting dalam penelitian ini, karena model klasifikasi yang akan digunakan membutuhkan data teranotasi untuk proses pembelajaran. Sentimen dari masing-masing unggahan media sosial diklasifikasikan ke dalam tiga label: positif, negatif, dan netral, sesuai dengan konteks dan isi teks yang dikumpulkan melalui proses scraping dari *platform* X (*Twitter*) dan *TikTok*. Alur proses pelabelan dijelaskan sebagai berikut:

1. Persiapan Data

Sebelum proses pelabelan dilakukan, semua data telah melalui tahapan preprocessing seperti *cleaning*, tokenisasi, *stopword removal*, dan *stemming*. Data yang siap label hanya menyisakan kolom teks hasil *preprocessing* (*preprocessed_text*) yang akan dinilai makna sentimennya.

2. Pembagian ke Volunteer

Total data yang berhasil dikumpulkan setelah melalui tahap *preprocessing* berjumlah 5.198 data teks. Untuk menjaga konsistensi sekaligus mengevaluasi perbedaan persepsi dalam proses pelabelan, data tersebut dilabelkan secara terpisah oleh dua pihak, yaitu seorang *volunteer* dan peneliti. Masing-masing melabelkan seluruh 5.198 data yang sama, sehingga memungkinkan dilakukan perbandingan terhadap hasil klasifikasi yang dihasilkan oleh masing-masing pelabel. Pendekatan ini bertujuan untuk mengukur tingkat kesepahaman (*agreement*) serta akurasi pelabelan sebelum data digunakan dalam tahap pelatihan model.

3. Petunjuk dan Panduan Penilaian

Masing-masing *volunteer* diberikan pedoman pelabelan, yang mencakup:

- a. Sentimen positif: berisi ekspresi keinginan pindah keluar negeri, antusiasme, atau semangat.
- b. Sentimen negatif: berisi kritik terhadap negara, pesimisme, atau ketidakpuasan yang kuat.
- c. Sentimen netral: berisi kalimat di luar topik, bersifat informatif, atau tanpa emosi yang jelas.

3.3.6 Klasifikasi Sentimen

Proses ini terdiri dari dua tahapan utama, penerapan *Naive Bayes* dan *Logistic Regression*. *Naive Bayes* dipilih karena sederhana dan efektif dalam menangani data teks dengan banyak fitur dan hubungan antarfitur yang bersifat independen. *Logistic Regression* digunakan sebagai pembanding karena dapat menghasilkan probabilitas kelas yang stabil dan memberikan hasil klasifikasi yang dapat diinterpretasikan. Untuk membandingkan kedua pendekatan ini pada *dataset* yang sama, masing-masing pendekatan digunakan.

3.3.7 Evaluasi

Hasil dari implementasi *Naive Bayes* dan *Logistic Regression* selanjutnya dilakukan proses evaluasi. Dilakukan evaluasi untuk mengetahui seberapa efektif masing-masing pendekatan dalam mengidentifikasi dan mengklasifikasikan dengan benar. Dalam penelitian ini, beberapa metrik evaluasi digunakan, termasuk *accuracy*, *precision*, *recall*, dan *f1-score*, yang dihitung berdasarkan *confusion matrix*. Metrik-metrik ini memberikan gambaran menyeluruh tentang kinerja model, baik dalam mengidentifikasi sentimen positif maupun negatif, serta seberapa baik model mencegah kesalahan klasifikasi. Hasil evaluasi akan digunakan untuk menentukan strategi mana yang paling efektif untuk menganalisis sentimen di unggahan media sosial dengan tagar #KABURAJADULU. Ini akan berfungsi sebagai dasar untuk penelitian terkait jenis ini di masa mendatang.

V. SIMPULAN DAN SARAN

5.1 Kesimpulan

Penelitian ini berhasil mengimplementasikan metode *Natural Language Processing* (NLP) dasar untuk menganalisis sentimen publik terhadap tagar #kaburajadulu pada *platform X (Twitter)* dan *TikTok*, dengan total data sebanyak 5.198 unggahan teks. Penelitian ini menjawab tantangan klasifikasi opini publik di media sosial serta membandingkan kinerja dua algoritma klasifikasi: *Naïve Bayes* dan *Logistic Regression*. Kesimpulan utama sebagai berikut:

- Tantangan utama dalam klasifikasi opini publik adalah penggunaan bahasa informal, konteks kalimat yang ambigu, dan ketidakseimbangan data, seperti dominasi sentimen positif sebesar 47,6% dan rendahnya sentimen negatif sebesar 14,8% pada dataset TikTok
 - a. Penggunaan bahasa informal ditangani melalui normalisasi teks dengan tokenisasi, *case folding*, dan *stemming*.
 - b. Konteks kalimat yang ambigu diantisipasi dengan pembobotan TF-IDF untuk memberikan bobot kata yang lebih relevan.
 - c. Ketidakseimbangan data diatasi melalui penggunaan Logistic Regression yang lebih stabil dibandingkan Naïve Bayes, terutama pada kelas netral dan negatif.
- 2) Proses *preprocessing* yang mencakup tokenisasi, *case folding*, *stopword removal*, *stemming*, dan pembobotan TF-IDF berhasil mengubah data mentah menjadi format terstruktur dan representatif, sehingga siap digunakan untuk proses klasifikasi secara optimal
- 3) Kinerja model menunjukkan bahwa *Logistic Regression* unggul dibandingkan *Naïve Bayes*, dengan akurasi tertinggi 85% pada *dataset TikTok* (2.940 data), *F1-score* lebih tinggi pada kelas netral dan negatif, serta kemampuan yang lebih baik dalam menangani distribusi data tidak seimbang. Sementara itu, *Naïve Bayes* cenderung melakukan *over*-prediksi pada sentimen positif akibat sensitivitas terhadap frekuensi kata, dengan akurasi tertinggi hanya 66%

4) Pemetaan opini publik mengindikasikan bahwa mayoritas pengguna menyuarakan keinginan untuk berpindah atau mencari kehidupan yang lebih baik di luar negeri, disertai pembahasan isu-isu sosial-politik di Indonesia.

Dengan demikian, penelitian ini tidak hanya berhasil menjawab rumusan masalah terkait tantangan dan perbandingan algoritma klasifikasi, tetapi juga membuktikan bahwa implementasi NLP dasar efektif untuk menganalisis dan memetakan opini publik secara otomatis dan skalabel.

5.2 Saran

Berdasarkan hasil dan temuan penelitian, beberapa saran yang dapat diberikan untuk pengembangan penelitian di masa depan adalah:

- 1) Penanganan ketidakseimbangan data diperlukan penerapan teknik penyeimbangan data seperti *oversampling* (misalnya SMOTE), *undersampling*, atau penyesuaian *class weight*, agar model tidak bias terhadap kelas mayoritas dan mampu mengenali sentimen minoritas secara lebih baik.
- 2) Eksplorasi model yang lebih kompleks penelitian selanjutnya dapat menggunakan algoritma klasifikasi lanjutan seperti *Support Vector Machine* (SVM), *Random Forest*, atau pendekatan *deep learning* seperti LSTM, GRU, atau model berbasis *transformer* seperti *BERT*, guna meningkatkan pemahaman konteks sentimen dan akurasi prediksi.
- 3) Peningkatan validitas pelabelan untuk meningkatkan kualitas label, disarankan melibatkan lebih dari satu *volunteer* dalam proses pelabelan, serta menggunakan metrik kesepakatan seperti *Cohen's Kappa* antar *annotator*, guna memastikan objektivitas dan konsistensi pelabelan data.
- 4) Analisis topik dan ekspansi konteks penelitian dapat dikembangkan dengan menerapkan teknik *topic modeling* (seperti LDA) untuk mengetahui topik-topik dominan dalam setiap kategori sentimen. Ini akan membantu dalam memahami konteks di balik ekspresi sentimen pengguna media sosial.
- 5) Penyajian visual interaktif disarankan membangun *dashboard* interaktif berbasis *web*, yang menyajikan visualisasi distribusi sentimen, kata kunci utama, serta tren sentimen dari waktu ke waktu. Hal ini akan meningkatkan pemanfaatan hasil analisis oleh pihak eksternal, seperti akademisi, peneliti kebijakan, atau media.

DAFTAR PUSTAKA

- Amal, I., & Jayanta. (2023). Perbandingan Pelabelan Otomatis Dan Manual Untuk Analisis Sentimen Terhadap Kenaikan Harga BBM Pertamina Pada Twitter Menggunakan Algoritma Support Vector Machine. *Seminar Nasional Mahasiswa Ilmu Komputer Dan Aplikasinya (SENAMIKA)*, 4(2), 473–487. https://conference.upnvj.ac.id/index.php/senamika/article/view/2562
- Amrizal, V. (2018). Penerapan Metode Term Frequency Inverse Document Frequency (Tf-Idf) Dan Cosine Similarity Pada Sistem Temu Kembali Informasi Untuk Mengetahui Syarah Hadits Berbasis Web (Studi Kasus: Hadits Shahih Bukhari-Muslim). In *Jurnal Teknik Informatika* (Vol. 11, Issue 2). repository.uinjkt.ac.id. https://doi.org/10.15408/jti.v11i2.8623
- Astiningrum, M., Saputra, P. Y., & Rohmah, M. S. (2018). Implementasi Nlp Dengan Konversi Kata Pada Sistem Chatbot Konsultasi Laktasi. *Jurnal Informatika Polinema*, *5*(1), 46–52. https://doi.org/10.33795/jip.v5i1.262
- Barupal, D. K., & Fiehn, O. (2019). Generating the blood exposome database using a comprehensive text mining and database fusion approach. In *Environmental Health Perspectives* (Vol. 127, Issue 9, pp. 2825–2830). jmlr.org. https://doi.org/10.1289/EHP4713
- Boateng, E. Y., & Abaye, D. A. (2019). A Review of the Logistic Regression Model with Emphasis on Medical Research. In *Journal of Data Analysis and Information Processing* (Vol. 07, Issue 04, pp. 190–207). scirp.org. https://doi.org/10.4236/jdaip.2019.74012
- Ceri, S., Bozzon, A., Brambilla, M., Della Valle, E., Fraternali, P., & Quarteroni, S. (2013). An Introduction to Information Retrieval. In *Web Information Retrieval*. edl.emi.gov.et. https://doi.org/10.1007/978-3-642-39314-3_1
- Fahrezi, M. A., Yazid, K. M. A., Laksono, I. L., Sa'adat, F., G, F. I. N., & Pribadi, M. R. (2022). Perancangan UI/UX Pada Aplikasi Daily Trade Dengan Menggunakan Metode Design Thinking. *MDP Student Conference*, *1*(1), 279–283. https://jurnal.mdp.ac.id/index.php/msc/article/view/1760
- Fransiska Vina Sari, & Arief Wibowo. (2019). Analisis Sentimen Pelanggan Toko Online Jd.Id Menggunakan Metode Naïve Bayes Classifier Berbasis Konversi Ikon Emosi. *Jurnal SIMETRIS*, *10*(2), 681–686. https://jurnal.umk.ac.id/index.php/simet/article/view/3487/1000
- Heydarian, M., Doyle, T. E., & Samavi, R. (2022). MLCM: Multi-label confusion matrix. *Ieee Access*. https://ieeexplore.ieee.org/abstract/document/9711932/
- Hindarto, D. (2023). a Comparative Study of Sentiment Classification: Traditional Nlp Vs. Neural Network Approaches. *Jurnal Teknologi Informasi Universitas Lambung Mangkurat (JTIULM)*, 8(2), 49–60. https://doi.org/10.20527/jtiulm.v8i2.178

- Jovanica, C., Rahmintaningrum, D. D., Nuradni, H. A., & Salsabila, A. (2022). Analisis Pengaruh Aktor Pada Tagar #Roketchina Di Media Sosial Twitter Menggunakan Social Network Analysis (Sna). In *Jurnal Ilmiah Komunikasi Makna* (Vol. 10, Issue 1, p. 43). jurnal.unissula.ac.id. https://doi.org/10.30659/jikm.v10i1.15644
- Kim, Y., Hsu, S. H., & de Zúñiga, H. G. (2013). Influence of social media use on discussion network heterogeneity and civic engagement: The moderating role of personality traits. *Journal of Communication*, 63(3), 498–516. https://doi.org/10.1111/jcom.12034
- Mäntylä, M. V., Graziotin, D., & Kuutila, M. (2018). The evolution of sentiment analysis—A review of research topics, venues, and top cited papers. *Computer Science Review*, 27, 16–32. https://doi.org/10.1016/j.cosrev.2017.10.002
- Prasetya, M. A., Wulandari, M., & Nikmah, S. A. (2024). Implementasi NLP(Natural Language Processing) Dasar pada Analisis Sentiment Review Spotify. *Seminar Nasional Teknologi & Sains*, 3(1), 145–153. https://proceeding.unpkediri.ac.id/index.php/stains/article/view/4166
- Rahmad, F., Suryanto, Y., & Ramli, K. (2020). Performance Comparison of Anti-Spam Technology Using Confusion Matrix Classification. *IOP Conference Series:*Materials Science and Engineering, 879(1). https://doi.org/10.1088/1757-899X/879/1/012076
- Ratnawati, F. (2018). Implementasi Algoritma Naive Bayes Terhadap Analisis Sentimen Opini Film Pada Twitter. In *INOVTEK Polbeng Seri Informatika* (Vol. 3, Issue 1, p. 50). download.garuda.kemdikbud.go.id. https://doi.org/10.35314/isi.v3i1.335
- Syahputra, H., & Wibowo, A. (2023). Comparison of Support Vector Machine (SVM) and Random Forest Algorithm for Detection of Negative Content on Websites. In *Jurnal Ilmiah Teknik Elektro Komputer dan Informatika (JITEKI)* (Vol. 9, Issue 1, pp. 165–173). researchgate.net. http://journal.uad.ac.id/index.php/JITEKI
- Xu, J., Zhang, Y., & Miao, D. (2020). Three-way confusion matrix for classification: A measure driven view. *Information Sciences*, *507*, 772–794. https://doi.org/10.1016/j.ins.2019.06.064
- Yulianto, L., Rochim, A. I., & Hakim, L. (2018). Pelanggaran Kode Etik Pada Pemberitaan Media Sosial Intagram (Konflik Etnis Rohingnya). In *Representamen* (Vol. 4, Issue 02). Universitas 17 Agustus 1945 Surabaya. https://doi.org/10.30996/.v4i02.1742