Vol. 13 No. 3, pISSN: 2303-0577 eISSN: 2830-7062

http://dx.doi.org/10.23960/jitet.v13i3.7012

# ANALISIS SENTIMEN ULASAN PENGGUNA APLIKASI TIKET.COM DENGAN K-NEAREST NEIGHBOR

Paulenta Silvania Silitonga<sup>1</sup>, Angeline Riendra Tatipang<sup>2</sup>, Eva Salsabilla<sup>3</sup>, Anggraini Puspita Sari<sup>4\*</sup>

<sup>1,2,3,4</sup>Universitas Pembangunan Nasional Veteran Jawa Timur, Surabaya, Jawa Timur 60294, Telp. (031) 8706369

#### **Keywords:**

Analisis Sentimen; K-NN; SMOTE;

### Corespondent Email: anggraini.puspita.if@upnjati m.ac.id

untuk pemesanan tiket dan akomodasi. Ulasan dari pengguna di aplikasi Google Play Store berperan sebagai sumber data yang esensial dalam mengevaluasi kualitas layanan aplikasi tersebut. Penelitian ini dilakukan dengan tujuan untuk mengetahui performa algoritma K-Nearest Neighbor (K-NN) dalam menganalisis sentimen ulasan pengguna tiket.com, dengan dan tanpa penerapan SMOTE sebagai upaya penanganan ketidakseimbangan distribusi kelas positif dan negatif. Data diperoleh melalui teknik web scraping kemudian data melalui tahapan data pre-processing. Selanjutnya, data diklasifikasikan ke dalam dua jenis, yaitu positif dan negatif, berdasarkan nilai rating. Nilai k optimal diperoleh melalui pengujian dengan variasi nilai k yaitu k = 3, 5, 7, dan 9, baik pada data asli maupun data yang telah diolah menggunakan metode SMOTE. Hasil pengujian menunjukkan bahwa penerapan SMOTE secara konsisten meningkatkan akurasi model pada setiap nilai K yang diuji. Nilai K terbaik ditemukan pada k = 5, dengan akurasi sebesar 82,35% pada data tanpa SMOTE dan meningkat menjadi 84,56% setelah diterapkan SMOTE. Hal ini menunjukkan bahwa penggunaan SMOTE berpengaruh terhadap akurasi performa model.

Abstrak. Perkembangan teknologi digital telah mendorong peningkatan

penggunaan aplikasi daring, seperti tiket.com, yang umum dimanfaatkan



JITET is licensed under a Creative Commons Attribution-NonCommercial 4.0 International License. Abstract. The advancement of digital technology has contributed to the increased use of online applications, such as Tiket.com, commonly used for ticket booking and accommodation. User reviews on the Google Play Store are a crucial data source for evaluating the service quality of such applications. This study aims to analyze K-Nearest Neighbor (K-NN) algorithm's performance in sentiment analysis of Tiket.com user reviews, both with and without the application of SMOTE as a method for handling class imbalance between positive and negative sentiments. The data was collected through web scraping and subsequently underwent data pre-processing. The reviews were then classified into two categories—positive and negative based on their rating values. The optimal k value was determined by testing different k values, namely k = 3, 5, 7, and 9, on both the original dataset and the dataset processed using SMOTE. The test results showed that the application of SMOTE consistently improved model accuracy for each tested k value. The best performance was observed at k = 5, with an accuracy of 82.35% on the dataset without SMOTE, which increased to 84.56% after applying SMOTE. This indicates that the use of SMOTE has a positive impact on the model's performance accuracy.

#### 1. PENDAHULUAN

digital Kemaiuan teknologi telah memberikan kenyamanan dalam berbagai bidang kehidupan, termasuk sektor perjalanan dan pariwisata. Adanya platform seperti tiket.com sangat membantu masyrakat Indonesia dalam merencanakan serta memesan perjalanan mereka dengan cara yang lebih mudah dan efektif. Tingginya frekuensi pemakaian aplikasi tiket.com di kalangan masyarakat Indonesia menjadikannya sumber yang melimpah dengan data ulasan pengguna, yang menunjukkan pengalaman langsung pelanggan terhadap layanan disediakan[1]. Masukan serta komentar dari pengguna di Google Play Store menjadi sumber informasi yang penting bagi manajemen untuk seberapa puas pelanggan menilai meningkatkan kualitas pelayanan. Dengan pengalaman mengetahui dan harapan perusahaan dapat melakukan pengguna, perbaikan yang lebih fokus untuk memperbaiki layanan mereka. Analisis sentimen digunakan untuk memahami emosi yang terkandung dalam ulasan, baik positif, negatif, maupun netral, sehingga dapat memberikan wawasan berharga bagi perusahaan dalam menyusun strategi peningkatan layanan yang lebih efektif dan efisien[2].

Salah satu teknik yang sering dipakai dalam analisis sentimen adalah K-Nearest Neighbor (KNN), yang handal dalam mengklasifikasikan informasi berdasarkan kesamaan fitur dengan data yang sudah ada. Kelebihan utama dari **KNN** dalam kesederhanaannya penerapan dan kemampuannya untuk menangani pola data yang rumit tanpa memerlukan asumsi sebelumnya[3]. Metode ini juga menggunakan prinsip pembelajaran berbasis instance, yang memungkinkan hasil klasifikasi beradaptasi dengan perubahan data[4]. Meskipun banyak studi yang telah menunjukkan keefektifan algoritma machine learning dalam memprediksi sentimen pengguna, mayoritas masih lebih fokus pada pendekatan deep learning atau metode statistik, sementara penelitian mengenai pengguna KNN dalam analisis sentimen terhadap ulasan di Google Play Store masih tergolong terbatas[5].

Beberapa riset relevan telah dilaksanakan sebelumnya, seperti yang tertera dalam studi berjudul "Analisis Sentimen Publik di Media

Terhadap Sosial Twitter Tiket.com Menggunakan Algoritma Klasifikasi". Riset ini menerapkan beragam metode klasifikasi, antara lain Naive Bayes, K-Nearest Neighbor (KNN), Support Vector Machine (SVM), dan Random Forest (RF), untuk menganalisis perasaan masyarakat terhadap tiket.com. Temuan dari penelitian ini dapat dilihat bahwa algoritma Random Forest memiliki akurasi tertinggi bila dibandingkan dengan metode lainnva, sedangkan KNN dengan parameter k=11 berhasil mecapai akurasi sebesar 91% [6].

K-Nearest Neighbor Metode (KNN) memungkinkan penelitian ini dapat mengklasifikasikan sentimen ulasan dengan efektif berdasarkan kemiripan fitur dengan data historis. Dengan pendekatan ini, informasi yang dihasilkan dapat diandalkan tanpa memerlukan asumsi yang rumit. Keunggulan KNN dalam mengelola pola data yang beragam dan kemampuannya untuk menyesuaikan diri dengan perubahan tren sentimen menjadikannya pilihan yang tepat untuk mendapatkan pemahaman yang lebih akurat mengenai preferensi pelanggan. Selain itu, karena masih sedikit penelitian yang mengkaji penggunaan **KNN** dalam menganalisis sentimen pengguna di Google Play Store, studi memberikan sumbangan pengembangan metode klasifikasi yang lebih aplikatif dan efisien di bidang pariwisata digital. Oleh karena itu, penelitian ini mengarah pada eksplorasi reaksi dari sentimen pengguna tiket.com di Google Play Store dengan menggunakan algoritma K-Nearest Neighbor sekaligus menilai efektivitasnya dibandingkan teknik lainnya yang bisa dipakai. Selain itu, ini juga berupaya memberikan rekomendasi kepada pengembang aplikasi berdasarkan pola perasaan yang terindetifikasi. bertuiuan Penelitian ini untuk memberikan wawasan yang lebih mendalam bagi para pengembang aplikasi mengenai persepsi dan pandangan para pengguna. Dengan pemahaman yang lebih baik, mereka dapat mengembangkan rencana yang lebih efisien untuk memperbaiki pengalaman pelanggan dan meningkatkan mutu layanan yang diberikan.

# 2. TINJAUAN PUSTAKA

# 2.1. Analisis Sentimen

Analisis sentimen adalah metode yang digunakan untuk mengidentifikasi dan mengolah data berbentuk teks guna menemukan makna emosional dari teks tersebut[7]. Setelah melalui proses tersebut, data dibedakan menjadi dua jenis, yaitu data sentimen positif dan data sentimen negatif.

# 2.2. K-Nearest Neighbor

Algoritma K-Nearest Neighbor digunakan dalam penelitian ini sebagai salah satu teknik klasifikasi dalam pendekatan *supervised learning*. Algoritma ini bekerja dengan cara sederhana, dengan menghitung jarak antara data testing ke seluruh data training, kemudian memilih jarak terdekat dari semua nilai jarak yang ada [8].

# 2.3. API Google-Play-Scraper

API Google-Play-Scraper merupakan salah satu cara yang dapat diterapkan untuk melakukan pengambilan data menggunakan Python dari Google Play Store tanpa memerlukan pustaka eksternal tambahan [9].

#### 2.4. TF-IDF

TF-IDF (*Term Frequency-Inverse Document Frequency*) merupakan teknik yang bekerja dengan cara mengekstraksi fitur yang kemudian mengukur frekuensi kemunculan sebuah kata dalam dokumen tertentu (*term frequency*) dan diberi nilai bobot serta menghitung penyebaran suatu kata untuk mengetahui kelangkaannya pada koleksi dokumen (*inverse document frequency*) [10].

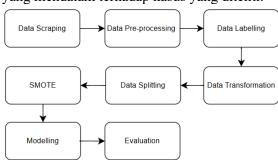
# 2.5. **SMOTE**

SMOTE (Synthetic Minority Oversampling Technique) adalah metode yang berguna untuk menyelesaikan masalah ketidakseimbangan kelas (class imbalance problem atau CIP). Metode ini dilakukan dengan menambahkan data sintetis pada kelas minoritas, agar distribusi merata antara jumlah data di kelas minoritas dan mayoritas [11].

### 3. METODE PENELITIAN

Penelitian ini dilakukan dengan pendekatan yang befokus pada data, melalui tahapan sistematis yang divisualisasikan dalam Gambar 1. Setiap langkah dalam proses penelitian dirancang untuk memastikan keakuratan dan validitas data yang diperoleh, sehingga hasil yang diperoleh dapat diandalkan dan relevan.

Selain itu, metode ini memungkinkan analisis yang mendalam terhadap kasus yang diteliti.



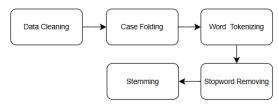
Gambar 1. Metode Penelitian

# 3.1. Data Scraping

Proses data scraping merupakan tahap awal dalam pelaksanaan penelitian, di mana data dikumpulkan dari sentimen pengguna aplikasi tiket.com di Google Play Store. Pengambilan data difokuskan pada ulasan yang diterbitkan sepanjang tahun 2024, yaitu mulai dari bulan Januari hingga Desember. Hal ini dilakukan untuk merepresentasikan pandangan dan pengalaman terbaru dari para pengguna terhadap aplikasi tersebut dalam satu tahun terakhir.

# 3.2. Data Pre-processing

Pada tahap ini, ulasan yang diperoleh dari data scraping dilakukan penghapusan elemenelemen seperti angka, tanda baca, emoji, dan duplikasi. Tahapan ini penting dalam analisis sentimen, sebagaimana dijelaskan oleh Attarik et al. [12], yang menerapkan preprocessing serupa sebelum mengklasifikasikan tweet menggunakan K-Nearest Neighbor.



Gambar 2. Alur Data Pre-processing

Pra-pemrosesan teks yang digunakan dalam penelitian ini di visualisasikan dalam Gambar 2. Tahapan ini meliputi pembersihan data, pelipatan huruf (*case folding*), pemisahan kata (*tokenisasi*), penghilangan kata-kata umum (*stopword removal*), dan pengubahan kata ke bentuk dasar (*stemming*)[13]. Proses pembersihan data berguna untuk

menghilangkan komponen yang tidak relevan serta menghapus missing value. Case folding diterapkan untuk menyeragamkan teks dengan mengubah seluruh karakter menjadi lowercase. Tokenisasi berguna untuk memisahkan teks menjadi unit-unit kata individual. Stopword removing berguna untuk menghilangkan katakata umum yang kurang informatif, seperti "yang", "atau", dan "dan". Terakhir, stemming berguna untuk mereduksi kompleksitas data dengan memangkas imbuhan dari setiap kata, sehingga kata-kata dengan akar kata yang sama direpresentasikan secara seragam. Keseluruhan proses *pre-processing* ini bertujuan menyederhanakan representasi teks ulasan, menjadikannya lebih bersih dan terstruktur sehingga siap untuk dianalisis lebih lanjut dan meningkatkan akurasi model.

### 3.3. Data Labelling

Setelah data melalui pengolahan, pelabelan dilakukan dengan menggunakan pendekatan berbasis rating score, yaitu metode otomatis yang memanfaatkan nilai ulasan atau penilaian yang diberikan pengguna. Data yang terkumpul akan dikategorikan menjadi sentimen positif dan sentimen negatif. Ulasan dengan rating tinggi, yaitu 4 dan 5 akan dilabeli dengan label positif, sedangkan ulasan dengan rating rendah, yaitu 1 dan 2 akan dilabeli dengan label negatif. Ulasan dengan rating netral atau 3 diabaikan. Pendekatan ini umum digunakan dalam analisis sentimen berbasis data ulasan karena rating mencerminkan opini pengguna secara langsung dan bersifat eksplisit. Pendekatan serupa juga diterapkan oleh Pratama et al. (2023), yang menunjukkan bahwa pelabelan berbasis rating dapat mencapai akurasi hingga 87% dalam klasifikasi sentimen pengguna [14].

### 3.4. Data Transformation

Setelah proses pelabelan data, data teks ditransformasikan ke dalam format numerik melalui penerapan TF-IDF (*Term Frequency-Inverse Document Frequency*). Metode ini bekerja dengan cara mengekstraksi fitur yang kemudian menetapkan nilai bobot di setiap kata berdasarkan frekuensi kemunculannya dalam dokumen tertentu (*term frequency*) serta kelangkaannya di seluruh koleksi dokumen (*inverse document frequency*). Teknik ini berperan dalam menilai pentingnya sebuah kata dalam suatu dokumen [15].

### 3.5. Data Splitting

Pada tahapan ini, data dipecah ke dalam 2 bagian utama, yaitu data latih dan data uji[16]. Pembagian ini berguna untuk memastikan model dapat dilatih dan diuji secara adil. Pada penelitian ini, data dibagi dengan proporsi 80% untuk data latih dan 20% untuk data uji. Prosedur *splitting* ini bertujuan utama untuk mengukur kapabilitas model dalam menggeneralisasi dan memprediksi data baru yang belum pernah dikenali sebelumnya.

# 3.6. **SMOTE**

Pada tahap ini, setelah data dibagi menjadi set pelatihan dan pengujian, metode SMOTE (Synthetic Minority Oversampling Technique) diimplementasikan pada data latih guna menanggulangi masalah ketidakseimbangan distribusi antar kelas pada data latih. SMOTE menghasilkan sampel sintetis baru pada kelas minoritas dengan menginterpolasi nilai-nilai tetangga terdekat, bukan sekadar menduplikasi data yang ada. Tujuannya adalah untuk menyeimbangkan distribusi kelas, sehingga model dapat belajar secara lebih merata dan menghasilkan performa prediksi yang lebih akurat. Dengan demikian, performa prediksi model diharapkan dapat ditingkatkan, terutama dalam mengenali kelas minoritas yang sebelumnya kurang terwakili[17].

### 3.7. Modelling

Pada tahap ini, model dibangun dengan menggunakan algoritma K-Nearest Neighbor. Aldoritma pengklasifikasian ini bekerja dengan cara sederhana memilih tetangga terdekat berdasarkan nilai jarak antara data testing dengan seluruh data training[18]. Algoritma ini dimanfaatkan untuk menghasilkan analisis sentimen berupa kategori positif atau negatif dari data ulasan yang telah melalui proses pengolahan. Untuk menentukan nilai k yang paling optimal dalam proses analisis, peneliti mengevaluasi performa model berdasarkan akurasi dengan variasi nilai k, yaitu 3, 5, 7, dan 9."

#### 3.8. Evaluation

Proses evaluasi terhadap performa dari model KNN dilakukan menggunakan classification report dan confusion matrix. Classification report merangkum metrik performa penting seperti akurasi, presisi, recall,

dan F1-score secara detail untuk setiap kelas. Akurasi merepresentasikan proporsi prediksi yang tepat, sementara presisi mengukur ketepatan model dalam mengidentifikasi instance kelas positif. Recall mengevaluasi kapasitas model untuk mendeteksi semua kasus positif yang relevan, dan F1-score menyajikan ukuran keseimbangan antara presisi dan recall [19]. Di sisi lain, confusion matrix berfungsi dalam memvisualisasikan hasil evaluasi model dengan menampilkan jumlah data yang diprediksi benar maupun salah dalam format tabular, meliputi empat komponen utama, yaitu True Positive (TP), True Negative (TN), False Positive (FP), dan False Negative (FN)[20]. Evaluasi ini dilakukan untuk memastikan model mampu mengklasifikasikan data dengan secara efektif dan seimbang.

#### 4. HASIL DAN PEMBAHASAN

Pada bagian ini, akan diuraikan hasil yang diperoleh dari proses penelitian serta analisis mengenai penerapan metode yang telah digunakan. Di samping itu, diskusi ini juga akan mengevaluasi efektivitas dalam penerapan metode yang digunakan.

# 4.1. Data Scraping

Dataset ini berisi kumpul ulasan pengguna yang diperoleh melalui teknik pengambilan data menggunakan Google Play Scraper. Informasi yang terkumpul mencakup nama pengguna, isi komentar, nilai rating, serta tanggal dan waktu dengan rentang waktu 1 Januari 2024 hingga 31 Desember 2024. Proses ini menghasilkan total 2.743 ulasan seperti yang terlihat pada Gambar 3.

Total	data ulasan: 274	3		
	pengguna	ulasan	rating	tanggal
1	Pengguna Google	keren	5	2024-12-30 13:24:16
2	Pengguna Google	Reload terus, busuuuk	3	2024-12-30 08:38:59
3	Pengguna Google	mantap dan terpercaya.	5	2024-12-30 07:45:33
4	Pengguna Google	sangat membantu, mudah digunakan	5	2024-12-30 07:33:05
5	Pengguna Google	kerendapat tiket murah	5	2024-12-30 07:25:00
2739	Pengguna Google	Baru pertama kali melakukan booking ticket lew	1	2024-01-01 23:09:14
2740	Pengguna Google	antabs	5	2024-01-01 16:22:46
2741	Pengguna Google	Aplikasi tidak memuaskan, lanjutkan pembayaran	1	2024-01-01 09:42:00
2742	Pengguna Google	ok	5	2024-01-01 08:24:33
2743	Pengguna Google	Pertama kali pesan tiket langsung bermasalah,	1	2024-01-01 05:37:51

Gambar 3. Hasil data scraping

### 4.2. Data Pre-processing

Tahap pra-pemrosesan data berguna untuk menstandarisasi format data agar konsisten, sehingga dapat mengingkatkan kinerja dan akurasi dari model dalam proses analisis. Proses ini terdiri dari data cleaning, case folding, word tokenizing, stopword removing, dan stemming.

# 4.2.1. Data Cleaning

Pada tahapan ini, ada 2.158 data yang telah melalui proses *cleaning*. Elemen yang tidak relevan serta *missing value* sudah dihilangkan agar tidak mengganggu hasil analisis. Langkah ini penting untuk memastikan bahwa data yang digunakan dalam proses analisis adalah data yang memiliki validitas dan reliabilitas yang tinggi.

	ulasan	clean_review
1	keren	keren
2	Reload terus, busuuuk	Reload terus busuuuk
3	mantap dan terpercaya.	mantap dan terpercaya
4	sangat membantu, mudah digunakan	sangat membantu mudah digunakan
5	kerendapat tiket murah	keren dapat tiket murah
2154	mudah di pahami aplikasi nya	mudah di pahami aplikasi nya
2155	Baru pertama kali melakukan booking ticket lew	Baru pertama kali melakukan booking ticket lew
2156	antabs	antabs
2157	Aplikasi tidak memuaskan, lanjutkan pembayaran	Aplikasi tidak memuaskan lanjutkan pembayaran
2158	Pertama kali pesan tiket langsung bermasalah,	Pertama kali pesan tiket langsung bermasalah h

**Gambar 4.** Hasil *data cleaning* pada Data

### 4.2.2. Case Folding

Setelah data dibersihkan, kemudian seluruh karakter dalam teks dikonversi menjadi huruf kecil guna menyamakan format penulisan. Hal ini dilakukan untuk menghindari duplikasi makna akibat perbedaan penggunaan huruf besar dan kecil dalam analisis kata.

case folding	clean_review	ulasan	CDU.
keren	keren	keren	1
reload terus busuuuk	Reload terus busuuuk	Reload terus, busuuuk	2
mantap dan terpercaya	mantap dan terpercaya	mantap dan terpercaya.	3
sangat membantu mudah digunakan	sangat membantu mudah digunakan	sangat membantu, mudah digunakan	4
keren dapat tiket murah	keren dapat tiket murah	kerendapat tiket murah	5
mudah di pahami aplikasi nya	mudah di pahami aplikasi nya	mudah di pahami aplikasi nya	2154
baru pertama kali melakukan booking ticket lew	Baru pertama kali melakukan booking ticket lew	Baru pertama kali melakukan booking ticket lew	2155
antabs	antabs	antabs	2156
aplikasi tidak memuaskan lanjutkan pembayaran	Aplikasi tidak memuaskan lanjutkan pembayaran	Aplikasi tidak memuaskan, lanjutkan pembayaran	2157
pertama kali pesan tiket langsung bermasalah h	Pertama kali pesan tiket langsung bermasalah h	Pertama kali pesan tiket langsung bermasalah,	2158

Gambar 5. Hasil casefolding pada Data

### 4.2.3. Word Tokenizing

Setelah melalui proses *case folding*, data masuk ke tahapan *word* tokenizing yang dimana kalimat atau teks dipisahkan menjadi unit-unit kata (token). Dengan tokenisasi, kalimat utuh diurai menjadi kata-kata dasar untuk memudahkan proses analisis lanjutan seperti

penghitungan frekuensi, pembobotan, dan ekstraksi fitur.

toker	case folding	clean_review	ulasan
[keren	keren	keren	keren
[reload, terus, busuuuk	reload terus busuuuk	Reload terus busuuuk	Reload terus, busuuuk
[mantap, dan terpercaya	mantap dan terpercaya	mantap dan terpercaya	mantap dan terpercaya.
[sangat, membantu mudah, digunakan	sangat membantu mudah digunakan	sangat membantu mudah digunakan	sangat membantu, mudah digunakan
[keren, dapat, tiket murah	keren dapat tiket murah	keren dapat tiket murah	kerendapat tiket murah
-			
[mudah, di, pahami aplikasi, nya	mudah di pahami aplikasi nya	mudah di pahami aplikasi nya	mudah di pahami aplikasi nya
[baru, pertama, kali melakukan, booking tick	baru pertama kali melakukan booking ticket lew	Baru pertama kali melakukan booking ticket lew	Baru pertama kali melakukan booking ticket lew
[antabs	antabs	antabs	antabs
[aplikasi, tidak memuaskan, lanjutkan pembay	aplikasi tidak memuaskan lanjutkan pembayaran	Aplikasi tidak memuaskan lanjutkan pembayaran	Aplikasi tidak memuaskan, lanjutkan pembayaran
[pertama, kali, pesan tiket, langsung bermas	pertama kali pesan tiket langsung bermasalah h	Pertama kali pesan tiket langsung bermasalah h	Pertama kali pesan tiket langsung bermasalah,

**Gambar 6**. Hasil word tokenizing pada Data

# 4.2.4. Stopwoord Removing

Pada tahap ini, data diproses lagi ke tahapan stopword removing untuk menghapus kata-kata umum (stopwords) yang tidak memiliki kontribusi makna signifikan, seperti "di", "langsung", atau "dapat". Penghapusan stopword ini dapat membantu mengurangi noise dalam data. Gambar 7 memperlihatkan hasil penerapan stopword removing pada baris ketiga, dengan perubahan kalimat dari "mantap dan terpercaya" menjadi "mantap terpercaya".

stopword	token	case folding	clean_review
keren	[keren]	keren	keren
reload busuuuk	[reload, terus, busuuuk]	reload terus busuuuk	Reload terus busuuuk
mantap terpercaya	[mantap, dan, terpercaya]	mantap dan terpercaya	mantap dan terpercaya
membantu mudah	[sangat, membantu, mudah, digunakan]	sangat membantu mudah digunakan	sangat membantu mudah digunakan
keren tiket murah	[keren, dapat, tiket, murah]	keren dapat tiket murah	keren dapat tiket murah
mudah pahami aplikasi	[mudah, di, pahami, aplikasi, nya]	mudah di pahami aplikasi nya	mudah di pahami aplikasi nya
booking ticket aplikasi yernyata mengecewakan	[baru, pertama, kali, melakukan, booking, tick	baru pertama kali melakukan booking ticket lew	Baru pertama kali melakukan booking ticket lew
antabs	[antabs]	antabs	antabs
aplikasi memuaskan lanjutkan pembayaran disitu	[aplikasi, tidak, memuaskan, lanjutkan, pembay	aplikasi tidak memuaskan lanjutkan pembayaran	Aplikasi tidak memuaskan lanjutkan pembayaran
pesan tiket bermasalah kampung	[pertama, kali, pesan, tiket, langsung, bermas	pertama kali pesan tiket langsung bermasalah h	Pertama kali pesan tiket langsung bermasalah h

**Gambar 7**. Hasil *stopword removing* pada Data

#### 4.2.5. Stemming

Setelah tahap stopword removing, data diproses ke tahapan stemming yang berguna untuk mengurangi kompleksitas data dengan menghapus bentuk imbuhan, sehingga makna yang sama tidak dikategorikan sebagai entitas yang berbeda. Proses ini juga membantu menyederhanakan variasi kata agar lebih

mudah dianalisis oleh algoritma pemrosesan teks.

stemming	stopword	token	case folding
keren	keren	[keren]	keren
reload busuuuk	reload busuuuk	[reload, terus, busuuuk]	reload terus busuuuk
mantap percaya	mantap terpercaya	[mantap, dan, terpercaya]	mantap dan terpercaya
bantu mudah	membantu mudah	[sangat, membantu, mudah, digunakan]	sangat membantu mudah digunakan
keren tiket murah	keren tiket murah	[keren, dapat, tiket, murah]	keren dapat tiket murah
mudah paham aplikasi	mudah pahami aplikasi	[mudah, di, pahami, aplikasi, nya]	mudah di pahami aplikasi nya
booking ticket aplikasi yernyata kecewa resche	booking ticket aplikasi yernyata mengecewakan	[baru, pertama, kali, melakukan, booking, tick	baru pertama kali melakukan booking ticket lew
antabs	antabs	[antabs]	antabs
aplikasi muas lanjut bayar situ situ geser men	aplikasi memuaskan lanjutkan pembayaran disitu	[aplikasi, tidak, memuaskan, lanjutkan, pembay	aplikasi tidak memuaskan lanjutkan pembayaran
pesan tiket masalah kampung	pesan tiket bermasalah kampung	[pertama, kali, pesan, tiket, langsung, bermas	pertama kali pesan tiket langsung bermasalah h

Gambar 8. Hasil Stemming pada Data

# 4.3. Data Labelling

Pelabelan data dilakukan secara otomatis berdasarkan nilai *rating score* yang terdapat pada masing-masing ulasan. Sebelum tahap pelabelan, data dibersihkan dari *missing value*, sehingga tersisa 2.117 data yang dapat diolah. Selanjutnya, ulasan dengan label netral atau rating 3 dihilangkan karena analisis difokuskan pada sentimen positif dan negatif saja, sehingga jumlah data yang digunakan menjadi 2.036.

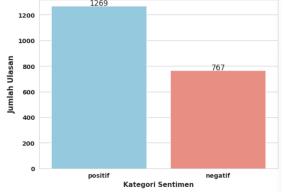
Data setelah pelabelan sentimen:

	stemming	label
1	keren	positif
2	mantap percaya	positif
3	bantu mudah	positif
4	keren tiket murah	positif
5	mantap	positif
2032	mudah paham aplikasi	positif
2033	booking ticket aplikasi yernyata kecewa resche	negatif
2034	antabs	positif
2035	aplikasi muas lanjut bayar situ situ geser men	negatif
2036	pesan tiket masalah kampung	negatif

Gambar 9. Hasil dari data labelling

Dari hasil pelabelan tersebut, terdapat 1.269 ulasan yang tergolong sentimen positif dan 767 ulasan sebagai sentimen negatif. Gambar 10 menyajikan distribusi sentimen dalam dataset dan menunjukkan bahwa sentimen positif paling mendominasi, mencerminkan

kecenderungan persepsi pengguna terhadap hal yang dinilai. Dominasi ini dapat menjadi indikator bahwa sebagian besar pengguna memiliki pengalaman yang memuaskan terhadap layanan atau produk yang digunakan.



**Gambar 10**. Distribusi Ulasan Positif dan Negatif Setelah Pelabelan

# 4.4. Data Transformation

Tahap transformasi data dilakukan dengan metode TF-IDF. Hasil dari transformasi data ini didapatkan nilai representasi numerik dari katakata dalam data berdasarkan tingkat kepentingannya, sehingga didapatkan informasi yang paling relevan dari data. Hasil TF-IDF pada 2 baris pertama divisualisasikan pada Gambar 11.

Baris ke-1:

Kata    :	Hasil Perhitungan TF-IDF
keren	1.0000
Baris ke-2:	
Kata    :	Hasil Perhitungan TF-IDF
mantap     percaya	0.5770   0.8167

**Gambar 11**. Hasil 2 Baris Pertama dari *Data Transformation* 

### 4.5. Data Splitting

Sebelum masuk pada tahap pemodelan, dilakukan pembagian dataset ulasan pengguna tiket.com menggunakan rasio 80:20, yang menghasilkan 1.628 data latih dan 408 data uji. Sebelum itu, data dibersihkan terlebih dahulu dengan menghapus nilai yang hilang (missing values) untuk memastikan kualitas data yang diolah. Gambar 12 menunjukkan hasil dari proses eksekusi kode pada tahap data splitting.

Jumlah data latih: 1628 Jumlah data uji: 408

Gambar 12. Hasil dari Data Splitting

#### 4.6. **SMOTE**

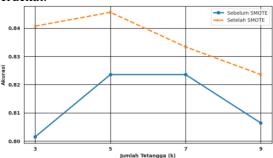
Setelah dilakukan tahapan data preparation, data ulasan menunjukkan distribusi kelas yang tidak seimbang pada data latih, dengan 1015 positif dan 613 data negatif. Ketidakseimbangan ini dapat menyebabkan bias terhadap kelas mayoritas. Untuk mengatasi hal tersebut, digunakan metode SMOTE yang untuk berguna menangani masalah ketidakseimbangan distribusi antar kelas pada datalatih. Setelah diterapkan SMOTE, data latih menjadi seimbang dengan masing-masing 1015 data untuk kelas positif dan negatif, sehingga total menjadi 2030 data.

**Tabel 1**. Distribusi Kelas Sebelum dan Sesudah Penerapan SMOTE pada Data Latih

	Sebelum	Sesudah
	SMOTE	SMOTE
Positif	1015	1015
Negatif	613	1015
Jumlah	1628	2030

# 4.7. Modelling

Tahap pemodelan dalam penelitian ini dilakukan dengan menerapkan algoritma K-Nearest Neighbor guna mengelompokkan sentimen pengguna. Pengujian dilakukan pada nilai k=3, 5, 7, dan 9. untuk mengevaluasi performa model pada berbagai jumlah tetangga terdekat.



**Gambar 13**. Grafik Perbandingan Akurasi KNN terhadap Nilai k

Berdasarkan hasil yang ditampilkan pada Gambar 13, terlihat bahwa model KNN yang dilatih dengan data yang telah diproses menggunakan SMOTE menunjukkan akurasi yang lebih tinggi pada seluruh nilai *k* dibandingkan dengan model yang tidak menggunakan SMOTE. Akurasi terbaik setelah SMOTE dicapai pada nilai k = 5 dengan nilai akurasi sebesar 0,8456 atau 84,56%. Sementara itu, akurasi tertinggi pada data sebelum SMOTE tercapai pada k = 5 dan k = 7, dengan akurasi sebesar 0,8235 atau 82,35%.

### 4.8. Evaluation

Evaluasi dilakukan untuk mengukur performa model K-Nearest Neighbor (KNN) dengan nilai k = 5 dalam mengklasifikasikan sentimen. Penilaian performa model menggunakan empat metrik utama, yaitu accuracy, precision, recall, dan F1-score. Pengujian dilakukan terhadap dua skenario, yaitu menggunakan data tidak seimbang (tanpa penerapan SMOTE) dan data yang telah diseimbangkan menggunakan metode SMOTE. Hasil evaluasi awal ditunjukkan melalui classification report dari kedua skenario. Visualisasi metrik klasifikasi untuk data tanpa penerapan SMOTE ditampilkan pada Gambar 13 dan hasil setelah penerapan SMOTE ditampilkan pada Gambar 14.

Classification	n Report			
	precision	recall	f1-score	support
negatif	0.78	0.75	0.76	154
positif	0.85	0.87	0.86	254
accuracy			0.82	408
macro avg	0.81	0.81	0.81	408
weighted avg	0.82	0.82	0.82	408

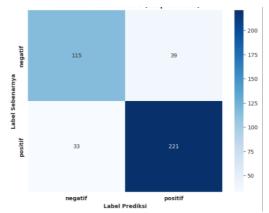
**Gambar 13.** Classification Report Tanpa Penerapan SMOTE

Classificatio	n Report precision	recall	f1-score	support
negatif positif	0.73 0.95	0.94 0.79	0.82 0.86	154 254
accuracy macro avg weighted avg	0.84 0.87	0.86 0.85	0.85 0.84 0.85	408 408 408

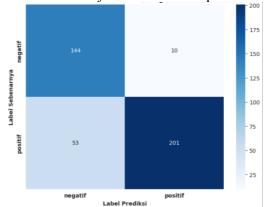
**Gambar 13.** Classification Report Dengan Penerapan SMOTE

Model KNN yang dilatih menggunakan data dengan penerapan SMOTE menunjukkan peningkatan pada hampir seluruh metrik evaluasi dibandingkan model yang dilatih pada data tidak seimbang. Peningkatan ini terlihat jelas pada *recall* untuk kelas negatif, yang naik dari 0,75 menjadi 0,94, serta *F1-score* yang meningkat dari 0,76 menjadi 0,82. Meskipun *precision* untuk kelas negatif menurun dari 0,78 menjadi 0,73.

Untuk memberikan gambaran yang lebih rinci mengenai hasil klasifikasi pada masingmasing kelas, digunakan *confusion matrix*. Visualisasi *confusion matrix* untuk model tanpa SMOTE ditampilkan pada Gambar 15 dan hasil dengan penerapan SMOTE ditampilkan pada Gambar 16.



Gambar 14. Confusion Matrix Tanpa SMOTE



**Gambar 15.** *Confusion Matrix* Dengan SMOTE

Berdasarkan Gambar 14 model mengindikasikan 115 true negatives dan 39 false positives, serta 221 true positives dan 33 false negatives. Konfigurasi ini mencerminkan tantangan dalam mendeteksi kelas minoritas, yang kemungkinan besar diakibatkan oleh ketidakseimbangan data. Setelah implementasi SMOTE yang ditunjukkan pada Gambar 15, performa model mengalami perbaikan, dengan akurasi yang meningkat dari 82,35% menjadi 84,56%. Peningkatan ini didukung oleh perubahan pada confusion matrix, yaitu pada true negatives yang meningkat menjadi 144, sementara false positives berkurang menjadi 10. Di sisi lain, jumlah true positives sedikit berkurang menjadi 201, dengan peningkatan false negatives menjadi 53. Secara keseluruhan, hasil ini menunjukkan pengurangan kesalahan klasifikasi dan peningkatan kemampuan model dalam mengidentifikasi kelas negatif secara lebih akurat, meskipun peningkatan ini diiringi oleh sedikit penurunan pada metrik true positives.

Untuk mempermudah perbandingan secara kuantitatif, hasil evaluasi pada masing-masing metrik dirangkum dalam Tabel 2.

**Tabel 2**. Perbandingan *Classification Report*Dengan dan Tanpa Penerapan SMOTE

	Accuracy (%)	Precision (%)	Recall (%)	F1- score (%)
Sebelum	82,35	81,35	80,84	81,08
Setelah	84,56	84,18	86,32	84,25

Berdasarkan Tabel 2, model dengan SMOTE mencapai akurasi sebesar 84,56%, lebih tinggi dibandingkan model tanpa SMOTE yang hanya mencapai 82,35%. Peningkatan juga terlihat pada precision (84,18%), recall (86,32%),dan F1-score (84,25%),menunjukkan bahwa model dengan SMOTE lebih mampu menangani ketidakseimbangan data, terutama dalam mengklasifikasikan kelas minoritas. Peningkatan recall yang signifikan menandakan bahwa model dengan SMOTE lebih efektif dalam mendeteksi instance dari kelas yang sebelumnya kurang terwakili. Dengan demikian, teknik SMOTE terbukti meningkatkan akurasi dan keandalan model KNN dalam analisis sentimen, menjadikannya lebih optimal untuk aplikasi praktis.

### 5. KESIMPULAN

Penelitian ini memiliki tujuan untuk mengkaji emosi pengguna terhadap aplikasi tiket.com pada Google Play Store dan menilai kinerja algoritma K-Nearest Neighbor dalam mengklasifikasikan sentimen. Dengan menjalani beberapa langkah pemrosesan data awal dan penerapan algoritma klasifikasi, studi berusaha memberikan ini pemahaman mengenai pola emosi pengguna keberhasilan metode yang digunakan dalam analisis data teks dari ulasan.

Sebanyak 2.743 ulasan pengguna tiket.com dari Google Play Store tahun dikumpulkan melalui metode scraping. Setelah melalui proses pembersihan data, total ulasan yang tersisa mencapai 2.158. Kemudian, processing dilakukan proses case folding, tokenization, penghapusan stopword, dan stemming, sehingga tersisa 2.117 ulasan. proses pelabelan sentimen Setelah berdasarkan rating, berhasil diidentifikasi

- 2.036 ulasan yang terlabeli, terdiri dari 1.269 ulasan positif dan 767 ulasan negatif.
- b. Algoritma K-Nearest Neighbor mampu menghasilkan model dengan performa klasifikasi yang baik dalam analisis ulasan pengguna aplikasi tiket.com, dengan akurasi tertinggi mencapai 84,56% pada k=5 dan dengan menggunakan SMOTE.
- Penerapan SMOTE yang membantu menyeimbangkan jumlah data antar kelas, berpengaruh terhadap peningkatan akurasi model klasifikasi.
- d. Penelitian ini memiliki keterbatasan berupa data mungkin tidak yang sepenuhnya representatif karena pengambilan hanya dari Google Play Scraper dan penghilangan ulasan netral yang berpotensi menghilangkan informasi penting. Meskipun K-Nearest Neighbor dan SMOTE meningkatkan performa, metode ini memiliki risiko menghasilkan data sintetis yang kurang akurat serta keterbatasan dalam menangkap nuansa emosi pengguna.

#### **DAFTAR PUSTAKA**

- [1] H. Utami, "Analisis Sentimen dari Aplikasi Shopee Indonesia Menggunakan Metode Recurrent Neural Network," *Indonesian Journal of Applied Statistics*, vol. 5, no. 1, 2022.
- [2] Budiman, et al., "Analisis Sentimen Publik pada Media Sosial Twitter Terhadap Tiket.com Menggunakan Algoritma Klasifikasi," *Jurnal Informatika*, vol. 11, no. 1, 2024.
- [3] F. Bei and S. Saepudin, "Analisis Sentimen Aplikasi Tiket Online di Play Store Menggunakan Metode Support Vector Machine (SVM)," in *Prosiding Seminar* Nasional Sistem Informasi dan Manajemen Informatika Universitas Nusa Putra, vol. 1, no. 3, 2021.
- [4] R. Kurniawan, et al., "Sentiment Analysis of Google Play Store User Reviews," *Sistem Informasi dan Komputer*, vol. 13, no. 2, 2024.
- [5] M. Riski, et al., "Klasifikasi Sentimen Ulasan Aplikasi WhatsApp di Play Store Menggunakan Metode K-Nearest Neighbor," *KLIK: Kajian Ilmiah Informatika* dan Komputer, vol. 4, no. 1, 2023.
- [6] Budiman B., et al., "Analisis Sentimen Publik pada Mendia Sosial Twitter Terhadap Tiket.com Menggunakan Algoritma Klasifikasi," *Jurnal Informatika*, vol.11, no.1, 2024.

- [7] O. I. Gifari, et al., "Analisis Sentimen Review Film Menggunakan TF-IDF dan Support Vector Machine," *Journal of Information Technology*, vol.2, no.1, 2022.
- [8] R. Sari, "Analisis Sentimen pada Review Objek Wisata Dunia Fantasi Menggunakan Algoritma K-Nearest Neighbor (K-NN)," Jurnal Sains dan Manajemen, vol.8, no.1, 2020.
- [9] G. Riansyah, "Analisis Sentimen Ulasan Aplikasi DANA di Google Play Store Menggunakan Algoritma Naïve Bayes," *Politeknik TEDC Bandung*, vol.8, no.5, 2024.
- [10] I. Apriani, et al., "Perbandingan Pembobotan FItur TF-IDF dan TF-ABS Dalam Klasifikasi Berita Online Menggunakan Support Vector Machine (SVM)," e-Proceeding of Engineering, vol.10, no.3, 2023.
- [11] C. Christian, et al., "Analisis Sentimen pada Rating Aplikasi Shopee Menggunakan Metode Decision Tree Berbasis SMOTE," *Jurnal Teknologi Informasi*, vol. 18, no. 2, 2021.
- [12] K. D. ATTARIK, N. Safriadi, and Y. Yulianti, "ANALISIS SENTIMEN KEBIJAKAN PEMBERLAKUAN PEMBATASAN KEGIATAN MASYARAKAT (PPKM) TERHADAP PERTUMBUHAN EKONOMI SEKTOR E-COMMERS DI INDONESIA MENGGUNAKAN METODE K-NEAREST NEIGHBOR (KNN)", *JITET*, vol. 12, no. 3S1, Oct. 2024.
- [13] R. A. Prasetyo and R. Hidayat, "Analisis Sentimen Ulasan Aplikasi Mobile Menggunakan Metode Naïve Bayes dan K-Nearest Neighbor," *Jurnal RESTI (Rekayasa Sistem dan Teknologi Informasi)*, vol. 5, no. 2, 2021.
- [14] R. K. Pratama, M. N. Anwar, and A. R. Salam, "Analisis Sentimen terhadap Aplikasi Indodana: Paylater & Pinjaman menggunakan Naive Bayes," *Jurnal Ilmu Manajemen Terapan (JIMAT)*, vol. 4, no. 3, 2023.
- [15] R. A. P. Ricky, N. Rudiman, and A. V. Naufal, "Metode Pembobotan TF-IDF untuk Klasifikasi Teks Quick Count Pemilihan Wakil Presiden Indonesia 2024 pada X Twitter dengan Metode SVM," Jurnal Teknologi Informasi, vol. 18, no. 2, 2024.
- [16] Wijiyanto, A. I. P. Sopingi, "Perbandingan Data Untuk Memprediksi Ketepatan Studi Berdasarkan Atribut Keluarga Menggunakana Machine Learning," *Jurnal of Informatics*, vol. 8, no. 2, 2024.
- [17] A. Setiawan and B. Nugroho, "Optimalisasi Pra-Pemrosesan Teks untuk Analisis Sentimen Berbasis Bahasa Indonesia Menggunakan Pendekatan Text Mining," Jurnal Teknologi

- Informasi dan Ilmu Komputer, vol. 7, no. 4, 2020.
- [18] A. P. Sari, A. N. Sihananto, and D. A. Prasetya, "Implementasi Metode K-NN dalam Klasterisasi Kasus Kesehatan Jantung," *ALINIER: Jurnal Teknik Elektro dan Informatika*, vol. 3, no. 2, 2022.
- [19] M. Fadli and A. S. Rizal, "Klasifikasi dan Evaluasi Performa Model Random Forest Untuk Prediksi Stroke," *Jurnal Teknik*, vol. 12, no. 2, 2023.
- [20] A. R. Adinata, et al., "Implementasi Algoritma Convolutional Neural Network dan YOLOV8 Untuk Klasifikasi Ras Kucing," *Building of Informatics, Technology and Science*, vol. 6, no. 3, 2024.