

Terbit online pada laman web jurnal : <http://teknosi.fti.unand.ac.id/>

Jurnal Nasional Teknologi dan Sistem Informasi

| ISSN (Print) 2460-3465 | ISSN (Online) 2476-8812 |



Klik di sini dan tuliskan Kategori Artikel

Penerapan Algoritma *Decision Tree* untuk Ulasan Aplikasi Vidio di *Google Play*

Ivana Lucia Kharisma^a, Dhea Ayu Septiani^b, Anggun Fergina^c, Kamdan^d^{abcd} Program Studi Teknik Informatika, Universitas Nusa Putra, Kec. Cisaat, Kabupaten Sukabumi, Jawa Barat 43152, Indonesia

INFORMASI ARTIKEL

Sejarah Artikel:

Diterima Redaksi: 27 Juni 2023

Revisi Akhir: 11 September 2023

Diterbitkan Online: 20 September 2023

KATA KUNCI

Siaran Langsung,
Analisa Sentimen
Pohon Keputusan,
Matriks Kebingungan,
Streamlit.

KORESPONDENSI

E-mail: ivana.lucia@nusaputra.ac.id

A B S T R A C T

Aplikasi berbasis *video streaming* atau siaran langsung menjadi jenis aplikasi paling banyak digunakan di dunia. *Video On Demand* merupakan sistem interaktif yang memungkinkan kita memilih konten video yang akan ditonton. Vidio adalah portal *online* atau situs web *streaming* video yang didirikan pada tahun 2014. Situs web ini memungkinkan pengguna untuk menonton dan menikmati berbagai video dan layanan lain. Namun, berdasarkan ulasan di *Google Play*, Vidio mendapatkan rating rata-rata hanya sebesar 3.7 dari 623.000 lebih total ulasan. Hal tersebut yang mendorong dilakukannya penelitian ini. Data yang dikumpulkan adalah sebanyak 1000 data pada rentang waktu 2 Februari 2023 – 19 Februari 2023. Data tersebut diklasifikasikan ke dalam sentimen positif dan negatif menggunakan algoritma *Decision Tree* atau Pohon Keputusan. Berdasarkan 3 skenario pembagian data, didapatkan akurasi terbesar diperoleh dari pembagian data 80% data latihan dan 20% data uji yaitu sebesar 97.3%. sedangkan pada skenario pembagian data 70:30, akurasinya 96.8%, dan pembagian data 90:10 akurasinya sebesar 96.8%. Dari akurasi yang telah diperoleh, untuk evaluasi pengujian model, penelitian ini menggunakan *Confusion Matrix* atau Matriks Kebingungan. Agar prediksi dari model yang telah dilatih agar tersedia untuk orang lain, penelitian ini melakukan *model deployment* menggunakan Streamlit.

1. PENDAHULUAN

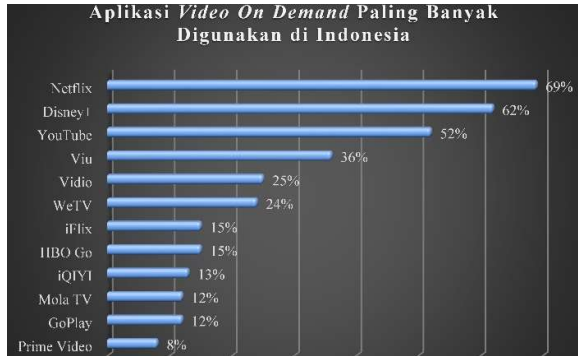
Aplikasi berbasis streaming video adalah aplikasi berbasis lalu lintas streaming yang paling banyak digunakan di dunia, terhitung 48,9% dari semua lalu lintas streaming Internet. [1]. Menariknya, Indonesia menjadi negara dengan total waktu yang dihabiskan untuk streaming tertinggi di dunia [2]. *Video On Demand* (selanjutnya disebut VOD) adalah sistem interaktif yang memungkinkan kita memilih konten video yang akan ditonton. Berbeda dengan televisi yang harus menunggu konten sesuai jadwal tayang, dengan VOD kita bisa menikmati, mengunduh dan memilih konten yang kita inginkan kapan saja dan di perangkat apa saja [3]. Berdasarkan temuan survei yang dirilis oleh Populix [4], alasan paling tinggi masyarakat Indonesia menggunakan layanan VOD adalah “Bisa nonton di mana saja” dengan persentase mencapai 84%, seperti yang terdapat pada Gambar 1.

Di dalam survei tersebut juga menunjukkan bahwa Netflix adalah aplikasi VOD yang paling banyak digunakan di Indonesia disusul Disney+, Youtube, dan Viu. Membuktikan bisa bersaing dengan platform VOD dari luar negeri, Vidio menjadi platform VOD lokal yang juga paling banyak digunakan oleh masyarakat Indonesia disusul oleh Mola TV dan GoPlay, seperti terdapat pada Gambar 2.



Gambar 1. Aplikasi *Video On Demand*. Paling Banyak Digunakan di Indonesia

Pada Gambar 1 dapat diketahui persentase alasan-alasan masyarakat Indonesia menggunakan aplikasi *video on demand*. Adapun aplikasi berbasis *video on demand* yang paling banyak digunakan di Indonesia terdapat pada Gambar 2.



Gambar 2. Alasan Masyarakat Menggunakan Aplikasi *Video On Demand*

Hingga saat ini (04/02/2023) Vidio sudah terunduh sebanyak 50 juta dengan *rating* 3.7 dan 623 ribu ulasan, serta menempati urutan pertama pada kategori *top grossing* segmentasi *entertainment* di situs *Google Play*. Selain memiliki konten serial

Tabel 1. Penelitian Terkait

No	Penulis	Metode	Hasil Penelitian
1.	Aufar et al., 2020	<i>Decision Tree</i> dan <i>Random Forest</i>	Hasil yang didapatkan adalah <i>Decision Tree</i> menjadi algoritma dengan tingkat akurasi yang tinggi yaitu sebesar 89,4%, dibandingkan dengan <i>Random Forest</i> dengan akurasi 88,2%. Data yang digunakan adalah komentar pada <i>channel</i> Youtube 'Nokia Mobile' sebanyak 2000 komentar pada video-video berbeda [7].
2.	Putra et al., 2022	<i>Naïve Bayes</i> , <i>KNN</i> dan <i>Decision Tree</i>	Hasilnya adalah <i>Decision Tree</i> menghasilkan tingkat akurasi paling tinggi yaitu 61,92%, dibanding dengan <i>Naïve Bayes</i> dengan akurasi 55,49% dan akurasi <i>KNN</i> sebesar 61,47% [8].
3.	Batlayeri & Gatta, 2022	<i>Naïve Bayes</i> , <i>KNN</i> dan <i>Decision Tree</i>	Hasilnya adalah <i>Naïve Bayes</i> dengan akurasi 74,92%, <i>KNN</i> 76,22% dan yang paling tinggi adalah <i>Decision Tree</i> dengan akurasi 77,85%. Data yang dikumpulkan sebanyak 1000 <i>tweets</i> [9].
4.	Waluyan & Hartomo, 2022	<i>Support Vector Machine</i>	Penelitian ini menggunakan algoritma <i>Support Vector Machine</i> menghasilkan akurasi sebesar 86,04%. Data yang digunakan adalah komentar <i>Original Series</i> pada Vidio.com yang terdiri dari 1403 komentar [10].

Berdasarkan penelitian-penelitian terkait yang sudah diuraikan, algoritma *Decision Tree* menunjukkan performa yang cukup baik. Namun, ada penelitian lain yang menggunakan Vidio sebagai objek kajian. Perbedaan mendasar antara penelitian sebelumnya dan penelitian ini terletak pada sumber data yang digunakan serta algoritma yang diterapkan. Penelitian sebelumnya telah mengambil sampel data dari komentar-komentar yang terdapat pada *original series* di platform Vidio. Dalam kasus ini, algoritma yang digunakan adalah *Support Vector Machine* [10].

dan film, Vidio juga menayangkan program *live* termasuk program televisi lokal. Menurut laporan riset yang dilakukan oleh Media Partner Asia [5], Vidio berada di urutan pertama berdasarkan durasi nonton penggunaannya seperti pada Gambar 3. *Google Play* sendiri adalah sebuah toko *online* untuk pengguna android dimana pengguna bisa menemukan aplikasi, *game*, film, acara TV, buku dan konten lainnya. *Google Play* memiliki fitur yang menyertakan ulasan pengguna yang dapat menunjukkan pengalamannya dalam menggunakan aplikasi. Ulasan pengguna sering digunakan sebagai alat yang efektif untuk menemukan informasi tentang suatu produk atau layanan [6].

Banyaknya ulasan aplikasi Vidio di *Google Play* mendorong penelitian ini untuk dilaksanakan guna mengetahui sentimen positif dan negatif pengguna terhadap aplikasi Vidio.

2. METODE

2.1. Penelitian Terkait

Penelitian ini merujuk pada penelitian-penelitian yang sudah dilakukan sebelumnya. Tabel 1. berisi penelitian-penelitian yang memiliki kelebihan terkait algoritma *Decision Tree* serta penelitian yang memiliki topik yang sama.

Di sisi lain, penelitian ini mengambil pendekatan yang berbeda. Data yang digunakan berasal dari ulasan atau komentar yang berasal dari pengguna aplikasi Vidio di *Google Play Store*. Dalam hal ini, algoritma yang digunakan adalah *Decision Tree*. Dengan memilih data dari ulasan pengguna aplikasi Vidio di platform *Google Play*, penelitian ini dapat menganalisis bagaimana persepsi dan tanggapan pengguna terhadap aplikasi ini dari sudut pandang yang berbeda.

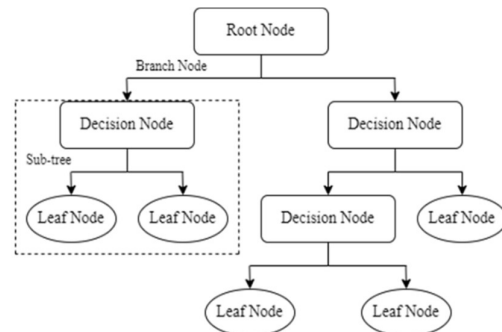
2.1.1. *Machine Learning*

Menurut Ibnu Daqiqil ID [11] pada bukunya yang berjudul “*Machine Learning: Teori, Studi Kasus dan Implementasi Menggunakan Python*”, Pembelajaran mesin (ML) adalah bidang penelitian yang berfokus pada desain dan analisis algoritma yang memungkinkan komputer untuk belajar. *Machine Learning* (ML) adalah bidang penelitian yang berfokus pada desain dan analisis algoritma yang memungkinkan komputer untuk belajar.

Klasifikasi mengukur kinerja berdasarkan perbandingan jumlah tebakan benar dan salah, sedangkan regresi dievaluasi berdasarkan seberapa dekat dengan nilai awal tebakan. Aplikasi atau model ML mengumpulkan pengalaman berdasarkan kumpulan data yang disediakan dalam proses pelatihan. Kumpulan data adalah kumpulan sampel yang harus dipelajari komputer untuk melakukan suatu tugas.

2.1.2. *Decision Tree*

Decision Tree dianggap sebagai metode yang paling menonjol dari metode-metode yang paling terkenal untuk representasi klasifikasi data. Peneliti yang berbeda dari berbagai bidang dan latar belakang telah mempertimbangkan masalah perluasan *Decision Tree* dari data yang tersedia, seperti studi mesin, pengenalan pola, dan statistik [12]. Ilustrasi dari algoritma *Decision Tree* dapat dilihat pada Gambar 3.



Gambar 3. Ilustrasi Algoritma *Decision Tree*

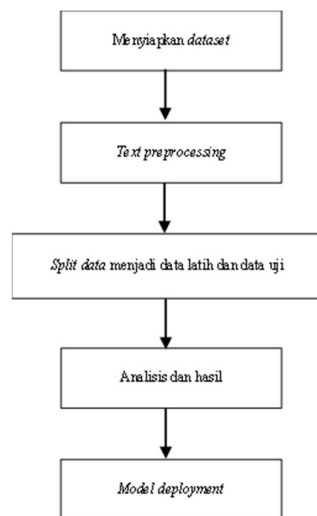
Pada Gambar 3 dapat dilihat setiap pohon terdiri dari simpul dan cabang. Setiap node mewakili fitur dalam kategori klasifikasi, dan setiap subset menentukan nilai yang dapat digunakan oleh node tersebut.

2.1.3. *Analisis Sentimen*

Analisis sentimen adalah teknik yang digunakan untuk mengekstraksi informasi tentang sikap seseorang terhadap suatu topik atau peristiwa dengan mengklasifikasikan polaritas teks. Pengelompokan tersebut menunjukkan apakah teks tersebut positif, negatif, atau netral [13]. Untuk itu diperlukan suatu tanda untuk mengklasifikasikan emosi linguistik berupa kosa kata yang terdapat dalam teks. menyukai misalnya pelajaran “baik” merupakan penanda emosi positif dan kata buruk 'buruk' adalah sebuah pertanda perasaan negatif.

2.2. *Tahapan Penelitian*

Flowchart yang menunjukkan tahapan-tahapan penelitian dari awal sampai akhir terdapat pada Gambar 4.



Gambar 4. Alur Penelitian

2.2.1. Menyiapkan dataset

Berikut dua tahap yang dilakukan dalam menyiapkan *dataset* pada penelitian ini:

Tabel 2. *Dataset Hasil Web Scraping*

<i>userName</i>	<i>Rating</i>	<i>At</i>	<i>Content</i>
taat triyanto	1	03/02/2023 07:14	Mengecewakan pelanggan!!!
Gigin Gunawan	2	08/02/2023 14:48	Setelah di perbarui.. Jadi mengecewakan..
fauzan	3	12/02/2023 14:25	Bagus, terimakasih
sutari ws	4	15/02/2023 11:30	sekarang vidio sudah mulai belajar memahami pelanggannya, sudah baguslah ttp terus ditingkatkan ya, jangan lipa buka link pengaduan bagaimana cara isi vocer langganan ya. Sy sdh beli tp lupa cara memasukkan
Alfin Nur latif	5	03/02/2023 09:12	Bagus ini bisa membantun Saya melihat Bola

1. *Web scraping*

Proses pengumpulan *dataset* dilakukan dengan cara *web scraping*. Proses *web scraping* menggunakan bahasa pemrograman *Python* dengan *Google Colab* sebagai *coding environment*-nya. Data yang diambil adalah data ulasan selama bulan Februari 2023 sebanyak 1000 ulasan mengenai aplikasi Vidio pada situs *Google Play* dengan kategori ‘terbaru’. Contoh ulasan yang dihasilkan dari proses *web scraping* dapat dilihat pada Tabel 2.

2. *Labeling*

Skema pelabelan yang digunakan adalah *average labeling*, dimana ulasan dengan rating 1, 2, dan 3 diberi label negatif (angka 0) dan ulasan dengan rating 4 & 5 diberi label positif (angka 1). Skema *average labeling* dapat dengan cepat diterapkan pada ukuran dataset yang besar. Karena dapat diterapkan pada dataset yang besar dengan waktu yang singkat, penelitian ini memilih menggunakan skema *average labeling* [14]. Setelah pelabelan dengan skema *average labeling*, tahap kedua adalah pelabelan manual yaitu *dataset* dikurasi kembali karena ada beberapa isi ulasan yang tidak sesuai *rating* [15]. Hal ini bertujuan untuk memastikan bahwa isi ulasan sesuai dengan *rating* yang diberikan.

2.2.2. *Text Preprocessing*

Preprocessing adalah tahap awal dalam pemrosesan data. *Preprocessing* dilakukan untuk mengubah teks menjadi data yang siap diproses [16], dengan bantuan Bahasa pemrograman *Python* dan *Google Colab*. Beberapa tahapan yang akan dilakukan diantaranya:

1. *Case Folding*: *Case folding* adalah teknik yang digunakan untuk mengubah semua teks dalam dataset menjadi huruf kecil (*lower case*) agar lebih mudah dibaca saat pengolahan data [17]. Tabel 3. menunjukkan contoh data ulasan yang belum dan sudah melalui proses *case folding*.

Tabel 3. Teks Saat Sebelum dan Sesudah Proses *Case Folding*

Sebelum <i>Case Folding</i>	Setelah <i>Case Folding</i>
Recomended sekali aplikasi ini. Nomer Bintang gak ada duanya. Yuk didownload untuk kalian yang ingin merasakan sensasi emosi yang luar biasa	recomended sekali aplikasi ini. nomer bintang gak ada duanya. yuk didownload untuk kalian yang ingin merasakan sensasi emosi yang luar biasa

2. *Cleaning*: *Cleaning* adalah proses menghapus karakter yang tidak perlu. Hal ini bertujuan untuk mengurangi *noise* yang dapat mengakibatkan prosedur komputasi klasifikasi menjadi kurang optimal [18]. Contohnya termasuk spasi ekstra, angka, simbol, tanda baca, emotikon, dan tautan. Tabel 4. memperlihatkan contoh ulasan sebelum dan sesudah proses *cleaning*.

Tabel 4. Teks Saat Sebelum dan Sesudah Proses *Cleaning*

Sebelum <i>Cleaning</i>	Sesudah <i>Cleaning</i>
Recomended sekali aplikasi ini. Nomer 1, Bintang 1, gak ada duanya, Yuk didownload , untuk kalian yang ingin merasakan sensasi emosi yang luar biasa..	Recomended sekali aplikasi ini. Nomer Bintang gak ada duanya. Yuk didownload untuk kalian yang ingin merasakan sensasi emosi yang luar biasa..

3. *Tokenizing*: *Tokenizing* adalah proses membagi kalimat menjadi kata per kata yang kemudian disebut token [19]. Contoh teks sebelum dan sesudah proses *tokenizing* terdapat pada Tabel 5.

Tabel 5. Teks Saat Sebelum dan Sesudah Proses *Tokenizing*

Sebelum <i>Tokenizing</i>	Setelah <i>Tokenizing</i>
recomended sekali aplikasi ini. Nomer bintang gak ada duanya. Yuk didownload untuk kalian yang ingin merasakan sensasi emosi yang luar biasa	['recomended', 'sekali', 'ini', 'nomer', 'bintang', 'gak', 'ada', 'duanya', 'yuk', 'didownload', 'untuk', 'kalian', 'yang', 'ingin', 'merasakan', 'sensasi', 'emosi', 'yang', 'luar', 'biasa']

4. *Normalization*: Proses *normalization* melibatkan pengubahan kata yang tidak baku menjadi baku, dan singkatan kembali ke kata aslinya [20]. Contoh sebelum dan sesudah proses *normalization* terdapat pada Tabel 6.

Tabel 6. Teks Saat Sebelum dan Sesudah Proses *Normalization*

Sebelum <i>Normalization</i>	Setelah <i>Normalization</i>
['recomended', 'sekali', 'aplikasi', 'ini', 'nomer', 'bintang', 'gak', 'ada', 'duanya', 'yuk', 'didownload', 'untuk', 'kalian', 'yang', 'ingin', 'merasakan', 'sensasi', 'emosi', 'yang', 'luar', 'biasa']	['direkomendasikan', 'sekali', 'aplikasi', 'ini', 'nomor', 'bintang', 'tidak', 'ada', 'duanya', 'yuk', 'diunduh', 'untuk', 'kalian', 'yang', 'ingin', 'merasakan', 'sensasi', 'emosi', 'yang', 'luar', 'biasa']

5. *Stopword*: *Stopword* adalah proses dimana kata-kata yang tidak penting atau memiliki informasi yang rendah pada sebuah teks dihilangkan [21]. Contoh *stopword* dalam Bahasa Indonesia adalah 'dan', 'yang', 'di', dan lain-lain. Teks saat sebelum dan sesudah proses *stopword* terdapat pada Tabel 7.

Tabel 7. Teks Saat Sebelum dan Sesudah Proses *Stopword*

Sebelum <i>Stopword</i>	Setelah <i>Stopword</i>
['direkomendasikan', 'sekali', 'aplikasi', 'ini', 'nomor', 'bintang', 'tidak', 'ada', 'duanya', 'yuk', 'diunduh', 'untuk', 'kalian', 'yang', 'ingin', 'merasakan', 'sensasi', 'emosi', 'yang', 'luar', 'biasa']	['direkomendasikan', 'aplikasi', 'nomor', 'bintang', 'duanya', 'diunduh']

6. *Stemming*: Kosakata berimbuhan yang ditemukan dalam dataset akan dipotong menjadi kata dasar selama tahap *stemming* [22]. Contoh terdapat pada Tabel 8. dimana terlampir teks sebelum dan sesudah melalui proses *stemming*.

Tabel 8 Teks Saat Sebelum dan Sesudah Proses *Stemming*

Sebelum <i>stemming</i>	Setelah <i>stemming</i>
['direkomendasikan', 'aplikasi', 'nomor', 'bintang', 'duanya', 'diunduh', 'merasakan', 'sensasi', 'emosi']	['rekomendasi', 'aplikasi', 'nomor', 'bintang', 'dua', 'unduh', 'rasa', 'sensasi', 'emosi']

2.2.3. Pembobotan Kata

Teknik pemberian nilai pada setiap kata yang berhasil melewati tahap *preprocessing* dikenal sebagai pembobotan kata. Pendekatan TF-IDF (*Term Frequency-Inverse Document Frequency*) digunakan dalam pembobotan kata pada penelitian ini, seperti yang ditunjukkan pada persamaan 1 dan persamaan 2. Kata-kata yang akan digunakan sebagai input dalam proses klasifikasi diberi bobot makna [23].

$$w_{ij} = t_{fij} \times idf_j \tag{1}$$

$$w_{ij} = t_{fij} \times \log(D/df_j) \tag{2}$$

2.2.4. Split Data

Pembagian data dilakukan setelah data selesai melewati tahapan *preprocessing*. Terdapat 2 jenis data pada tahapan pembagian ini yaitu data latih dan data uji. Skenario pembagian data pada

penelitian ini antara lain 70:30, 80:20, dan 90:10 [24], dimana setiap *dataset* akan dibagi menjadi 70% data latih dan 30% data uji, begitu juga untuk skenario yang lain.

2.2.5. Analisis dan Hasil

Pada tahapan ini dilakukan klasifikasi dengan pemodelan algoritma *Decision Tree*. Ulasan yang akan diklasifikasi berupa nilai dari pembobotan dari tiap kata melalui proses penghitungan TF-IDF. Contoh dari pembobotan kata ditunjukkan pada Tabel 9 berikut :

Tabel 9. Contoh pembobotan kata dengan TF-IDF

Stemming	aktif	eror	langgan
['eror', 'langgan', 'aktif']	0.6789	0.4169	0.6042

Hasil klasifikasi akan dilakukan pengujian menggunakan dataset yang sudah dibagi sebelumnya. Proses pertama yang akan dilakukan adalah pengujian dengan data latih terlebih dahulu untuk pelatihan terhadap sistem, kemudian pengujian dengan dataset uji agar memperoleh hasil dari akurasi.

2.3.5.1 Confusion Matrix

Untuk evaluasi pengujian model, penelitian ini menggunakan *Confusion Matrix*. Jumlah hasil tes yang benar dan jumlah hasil tes yang salah diklasifikasikan dalam sebuah tabel yang disebut *confusion matrix* [25]. Tabel *confusion matrix* bisa dilihat pada Tabel 10.

Tabel 10. *Confusion Matrix*

		<i>Confusion Matrix</i>	
		<i>Predicted Labels</i>	
<i>True Labels</i>	0 (negatif)	TP	FN
	1 (positif)	FP	TN
	0 (negatif)	0 (negatif)	1 (positif)
1 (positif)	0 (negatif)	1 (positif)	

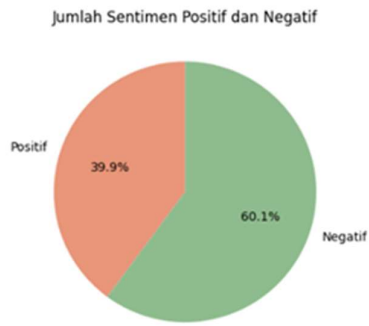
Keterangan:

- TP (*True Positive*): merujuk pada data positif yang diprediksi dengan benar atau diprediksi sebagai positif.
- TN (*True Negative*): merujuk pada data negatif yang diprediksi dengan benar atau diprediksi sebagai negatif.
- FP (*False Positive*): merujuk pada data negatif yang diprediksi salah, yaitu menjadi data positif.
- FN (*False Negative*): merujuk pada data positif yang diprediksi secara salah, yaitu menjadi data negatif.

Pengujian performa model menggunakan *confusion matrix* mempunyai beberapa metrik performa yang umum digunakan seperti di bawah ini [26]:

1. Accuracy

Seberapa akurat model dapat mengklasifikasikan objek dijelaskan oleh *accuracy*. *Accuracy* adalah proporsi positif dan negatif yang diprediksi benar terhadap *dataset*. *Accuracy* dapat dihitung dengan persamaan 3.



Gambar 7. Diagram Jumlah Sentimen Positif dan Negatif

Pada Gambar 7 dapat diketahui bahwa jumlah data positif sebanyak 39.9% atau 399, dan jumlah data negatif sebanyak 60.1% atau 601. Data yang sudah dikumpulkan, akan melalui tahapan *preprocessing* guna untuk membersihkan data. Hasil *dataset* yang sudah melewati tahapan *preprocessing* dapat dilihat pada Tabel 13.

Tabel 11 *Dataset* Sebelum Dan Sesudah Tahapan *Preprocessing*

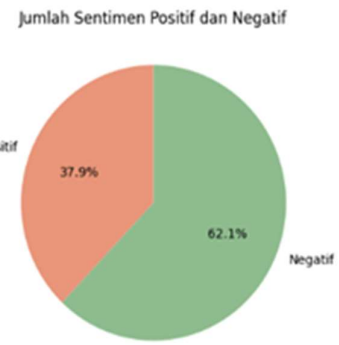
Data ulasan sebelum <i>preprocessing</i>	Data ulasan sesudah <i>preprocessing</i>
Ini kenapa Bein sport streamingnya ngefrezz tapi suaranya lancar	['siar langsung', 'beku', 'suara', 'lancar']
Bayar nya mahal tapi nonton lag banget padahal sinyal bagus, coba perbaiki lagi biar penonton merasa puas	['mahal', 'nonton', 'lambat', 'sinyal', 'bagus', 'coba', 'baik', 'biar', 'tonton', 'puas']
Lumayan lah lebih bagus lagi ya program nya banyakin	['lumayan', 'bagus', 'program', 'banyak']

Setelah melewati tahapan *preprocessing*, data ulasan mengalami pengurangan dikarenakan terdapat data yang kosong (*null*). Jumlah masing-masing data berdasarkan *rating* yang sudah mengalami pengurangan ditunjukkan pada Tabel 14.

Tabel 12. Jumlah Data Berdasarkan *Rating* Setelah Pengurangan

<i>Rating</i>	Jumlah data
1	472
2	58
3	54
4	54
5	302
Jumlah	940

Adapun jumlah sentimen positif dan negatif setelah mengalami pengurangan dapat dilihat pada Gambar 8.



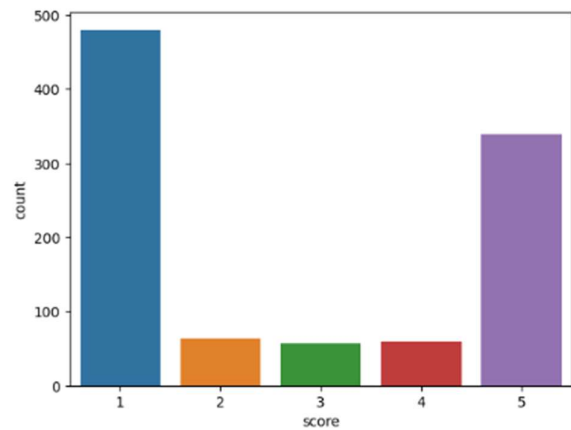
Gambar 8. Diagram Jumlah Sentimen Positif dan Negatif Setelah Pengurangan

Pada Gambar 8 diketahui jumlah data dengan sentimen positif setelah mengalami pengurangan adalah 37.9% atau sebanyak 356 data, sedangkan jumlah sentimen negatif menjadi 62.1% atau sebanyak 584 data.

Tahapan selanjutnya adalah proses pembagian data. Dalam penelitian ini terdapat 3 skenario pembagian data yaitu 70:30, 80:20 dan 90:10, atau 70% data latih dan 30% data uji. Setelah proses pembagian data, selanjutnya adalah tahap klasifikasi menggunakan algoritma *Decision Tree*. Hasil sentimen pada *dataset* pada penelitian ini hanya diklasifikasikan ke dalam kelas positif dan negatif saja. Setelah proses pelatihan model *Decision Tree* berhasil dilakukan, selanjutnya akan dilakukan pengujian. Pada prosesi pengujian dengan algoritma *Decision Tree* pada penelitian ini akan menghasilkan hasil prediksi terhadap data uji.

4. PEMBAHASAN

Tahap awal pada penelitian ini adalah mengumpulkan data dengan cara *web scraping* situs web *Google Play* untuk mendapatkan ulasan aplikasi Vidio. Gambar 9 menunjukkan jumlah data berdasarkan *rating*.



Gambar 9. Diagram Jumlah Data Berdasarkan *Rating*

Kemudian, pelabelan dilakukan secara manual berdasarkan rating dan kata-kata yang terkandung di dalamnya. Pada tabel 15, terdapat jumlah label positif dan negatif.

Tabel 13. Jumlah Sentimen Positif dan Negatif Sebelum Pengurangan

	Positif	Negatif
	399	601

Tabel 15. merupakan jumlah sentimen positif dan negatif sebelum data mengalami pengurangan dikarenakan penghapusan data kosong. Untuk jumlah sentimen positif dan negatif setelah penghapusan data kosong, terlihat pada Tabel 16.

Tabel 14. Jumlah Sentimen Positif dan Negatif Sebelum Pengurangan

	Positif	Negatif
	356	584

Tahapan preprocessing pada penelitian ini menggunakan metode casefolding, cleaning, tokenizing, stopword, normalization, dan stemming. Selanjutnya, data dikonversikan ke dalam bentuk angka dengan metode TF-IDF. Setelah itu, proses pembagian data. Berikut beberapa skenario pembagian data yang dilakukan pada penelitian ini:

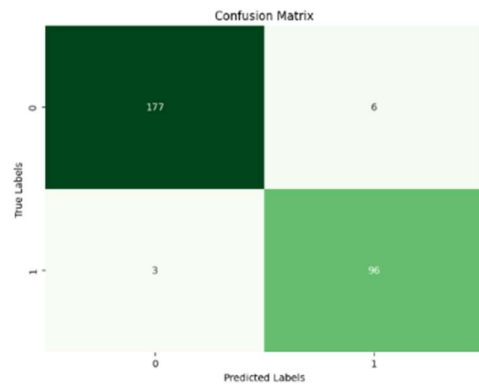
1. Pembagian data 70:30

Pada skenario ini, 940 jumlah data akan dibagi menjadi 70% data latih dan 30% data uji. Pada Gambar 10 data latih pada skenario ini adalah sebanyak 658 data, sedangkan data ujinya sebanyak 282 data.



Gambar 10. Pembagian Data 70:30

Akurasi yang didapatkan pada pembagian data 70:30 adalah sebesar 0.9680 atau 96.8%. Adapun bentuk confusion matrix pada skenario ini ditunjukkan pada Gambar 11.



Gambar 11. Confusion Matrix Data Uji 30%

Pada Gambar 11 didapatkan hasil sebagai berikut:

- TN (*True Negative*) sebanyak 177 data. Dari 282 data yang diuji, 177 diantaranya diprediksi benar yaitu sebagai sentimen negatif.
- FP (*False Positive*) sebanyak 6. Artinya, hanya 6 data yang diprediksi salah. Seharusnya negatif tetapi hasil prediksinya menjadi positif.
- FN (*False Negative*) sebanyak 3. 3 data yang seharusnya positif, diprediksi salah yaitu menjadi negatif.
- TP (*True Positive*) sebanyak 96. Artinya dari 282 data uji, 96 data diuji dengan benar, yaitu sebagai sentimen positif.
- Untuk menghasilkan *accuracy*, *precision*, *recall* dan *f-1 score* pada pembagian data 70:30 adalah sebagai berikut:

- Accuracy

$$\frac{TP+TN}{Total} = \frac{96+177}{282} = 0.968$$

- Precision

$$\frac{TP}{FP+TP} = \frac{96}{6+96} = 0.941$$

- Recall

$$\frac{TP}{FN+TP} = \frac{96}{3+96} = 0.969$$

- F-1 Score

$$2 \times \frac{precision \times recall}{precision+r} = 2 \times \frac{0.941 \times 0.969}{0.941+0.969} = 0.953$$

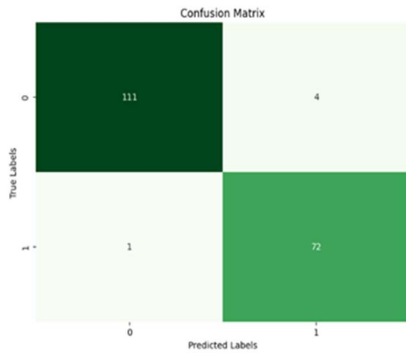
2. Pembagian data 80:20

Skenario ini menghasilkan data latih sebanyak 752, dan data uji sebanyak 188. Penggambaran pembagian data ini dapat dilihat pada Gambar 12.



Gambar 12. Pembagian Data 80:20

Akurasi yang didapatkan pada pengujian dengan pembagian data 80:20 adalah 0.973 atau 97.3%. Adapun bentuk *confusion matrix*-nya ditunjukkan pada Gambar 13.



Gambar 13. *Confusion Matrix* Data Uji 20%

Pada Gambar 13 didapatkan hasil sebagai berikut:

- TN (*True Negative*) sebanyak 111 data. Artinya dari 188 data yang diuji, 111 diantaranya diprediksi benar yaitu sebagai sentimen negatif.
- FP (*False Negative*) sebanyak 4. Artinya, hanya 4 data yang diprediksi salah. Seharusnya negatif tetapi hasil prediksinya menjadi positif.
- FN (*False Negative*) sebanyak 1. Hanya ada 1 data yang seharusnya diprediksi positif, tetapi diprediksi secara salah yaitu menjadi negatif.
- TP (*True Positive*) sebanyak 72. Dari 188 data uji, 72 data diuji dengan benar, yaitu sebagai sentimen positif.
- Untuk menghasilkan *accuracy*, *precision*, *recall* dan *f-1 score* pada pembagian data 80:20 adalah sebagai berikut:

- *Accuracy*

$$\frac{TP+TN}{Total} = \frac{72+111}{188} = 0.973$$

- *Precision*

$$\frac{TP}{FP+TP} = \frac{72}{4+72} = 0.947$$

- *Recall*

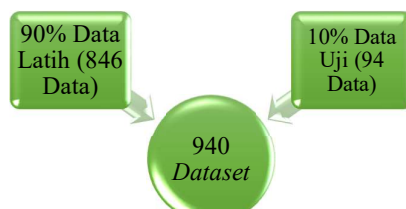
$$\frac{TP}{FN+TP} = \frac{72}{1+72} = 0.986$$

- *F-1 Score*

$$2 \times \frac{precision \times recall}{precision+recall} = 2 \times \frac{0.947 \times 0.986}{0.947+0.986} = 0.966$$

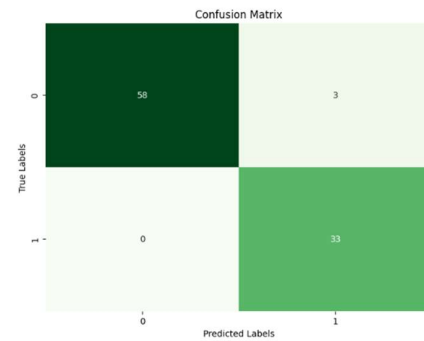
3. Pembagian data 90:10

Pada skenario pembagian data 90:10, didapatkan 846 data latih dan 94 data uji. Penggambaran pembagian data 90:10 dapat dilihat pada Gambar 14.



Gambar 14. Pembagian Data 90:10

Akurasi yang didapatkan skenario ini adalah sebesar 0.968 atau 96.8%. *Confusion matrix* pada skenario ini dapat dilihat pada Gambar 15.



Gambar 15. *Confusion Matrix* Data Uji 10%

Pada Gambar 15 didapatkan hasil sebagai berikut:

- TN (*True Negative*) sebanyak 58 data. Dari 94 data yang diuji, 58 diantaranya diprediksi benar yaitu sebagai sentimen negatif.
- FP (*False Negative*) sebanyak 3. Disimpulkan bahwa hanya ada 3 data yang diprediksi salah. Seharusnya negatif tetapi hasil prediksinya menjadi positif.
- FN (*False Negative*) sebanyak 0. Artinya, tidak ada data yang seharusnya diprediksi salah.
- TP (*True Positive*) sebanyak 33. Itu berarti 33 data diuji dengan benar, yaitu sebagai sentimen positif.

Untuk menghasilkan *accuracy*, *precision*, *recall* dan *f-1 score* pada pembagian data 90:10 adalah sebagai berikut:

- *Accuracy*

$$\frac{TP+T}{Total} = \frac{33+58}{94} = 0.968$$

- *Precision*

$$\frac{TP}{FP+TP} = \frac{33}{3+33} = 0.916$$

- *Recall*

$$\frac{TP}{FN+TP} = \frac{33}{0+33} = 1$$

- *F-1 Score*

$$2 \times \frac{precision \times recall}{precision+recall} = 2 \times \frac{0.916 \times 1}{0.916+1} = 0.956$$

Pada penelitian ini, proses membuat *model deployment*-nya menggunakan *framework* Streamlit. Pada Gambar 16 terdapat tampilan Streamlit yang digunakan.

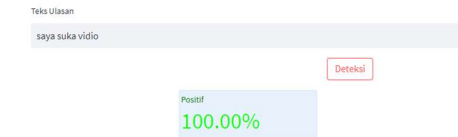
Prediksi Ulasan Aplikasi Vidio



Gambar 16. Tampilan Streamlit

Setelah ulasan dimasukkan dan dideteksi, tampilan Streamlit akan seperti pada Gambar 17. Pada Gambar 17, terlihat jika dimasukkan ulasan 'saya suka vidio', maka teks dideteksi sebagai sentimen positif dengan akurasi 100%. Contoh masukkan lain yang mengandung kata-kata negatif dapat dilihat pada Gambar 18.

Prediksi Ulasan Aplikasi Vidio



Gambar 17. Tampilan Streamlit Terhadap Sentimen Positif

Pada Gambar 18 terlihat jika masukkan 'aplikasi eror' terdeteksi sebagai sentimen negatif dengan akurasi 100%.

Prediksi Ulasan Aplikasi Vidio



Gambar 18. Tampilan Streamlit Terhadap Sentimen Negatif

5. KESIMPULAN

Dari pengujian yang sudah dilakukan, terdapat beberapa kesimpulan yang dapat diambil dalam analisis sentimen terhadap ulasan aplikasi Vidio di situs *Google Play* dengan menggunakan algoritma *Decision Tree*. Setelah mengalami pengurangan dikarenakan terdapat data kosong, jumlah *dataset* yang sebelumnya 1000, menjadi 940. Adapun dari *dataset* 940, jumlah ulasan kelas positif sebanyak 356 (37.9%), sedangkan untuk kelas negatif sebanyak 584 (62.1%).

Berdasarkan perbandingan pembagian data latih dan data uji 70:30, 80:20 dan 90:10, didapatkan akurasi terbaik adalah dengan skenario 80:20, yaitu 97.3% dengan *precision* 0.947 (94.7%), *recall* 0.986 (98.6%), dan *f-1 score* 0.966 (96.6%). Sedangkan pembagian 70:30 mendapatkan akurasi sebesar 96.8% dengan *precision* 0.941 (94.1%), *recall* 0.969 (96.9%), dan *f-1 score* 0.953 (95.3%). Pembagian data 90:10 mendapatkan akurasi sebesar 96.8% dengan *precision* 0.916 (91.6%), *recall* 1 (100%), dan *f-1 score* 0.956 (95.6%). Hasil akurasi tersebut memberikan hasil yang lebih baik dari penelitian terdahulu yang sejenis yaitu analisa ulasan sentimen komentar pada channel Youtube yaitu sebesar 89,4%.

Dengan besarnya akurasi yang dihasilkan oleh ketiga skenario pembagian pengujian data, serta kecilnya angka ketidakakuratan prediksi, maka dapat disimpulkan bahwa algoritma *Decision Tree* berjalan dengan baik pada penelitian ini. Algoritma ini mampu

memberikan klasifikasi yang tepat terhadap ulasan sesuai dengan pelabelan yang diberikan serta mampu mempermudah pada proses analisa setiap ulasan dengan fasilitas *streamlit* yang digunakan. Dari hasil klasifikasi yang telah dilakukan, dapat disimpulkan bahwa ulasan aplikasi Vidio mendapat sentimen negatif lebih besar daripada sentimen positif.

DAFTAR PUSTAKA

- [1] Sandvine, "The Mobile Internet Phenomena Report," 2021. Accessed: Jan. 27, 2023. [Online]. Available: https://www.sandvine.com/hubfs/Sandvine_Redesign_2019/Downloads/2021/Phenomena/MIPR%20Q1%202021%2020210510.pdf
- [2] data.ai, "STATE OF MOBILE," 2022. Accessed: Feb. 27, 2023. [Online]. Available: <https://www.data.ai/en/go/state-of-mobile-2022/>
- [3] N. A. Yusuf and Indrawati, "Analisis Faktor yang Memengaruhi Pembentukan Minat Berlangganan di Industri Video-On-Demand di Indonesia.," vol. 3, pp. 161–173, 2019.
- [4] Populix, "Indonesian Video Entertainment on Demand Consumption," Jul. 2022. Accessed: Feb. 03, 2023. [Online]. Available: <https://info.populix.co/report/indonesian-video-entertainment-on-demand-consumption/>
- [5] Media Partner Asia, "Southeast Asia Online Consumer Insights and Analytics." pp. 1–2, Nov. 22, 2022. [Online]. Available: https://media-partners-asia.com/AMPD/Q1_2022/SEA/PR.pdf
- [6] S. A. Saputra, D. Rosiyadi, W. Gata, and S. M. Husain, "Analisis Sentimen E-Wallet Pada Google Play Menggunakan Algoritma Naive Bayes Berbasis Particle Swarm Optimization," *RESTI (Rekayasa Sistem dan Teknologi Informasi)*, vol. 3, pp. 377–382, 2019.
- [7] M. Afar, R. Andreswari, and D. Pramesti, "Sentiment Analysis on Youtube Social Media Using Decision Tree and Random Forest Algorithm: A Case Study," *2020 IEEE International Conference on Data Science and Its Applications (ICoDSA)*, pp. 1–7, Oct. 2020, doi: [10.1109/ICoDSA50139.2020.9213078](https://doi.org/10.1109/ICoDSA50139.2020.9213078).
- [8] T. W. Putra, A. Triayudi, and Andrianingsih, "Analisis Sentimen Pembelajaran Daring menggunakan Metode Naïve Bayes, KNN, dan Decision Tree," *JTIK (Jurnal Teknologi Informasi dan Komunikasi)*, pp. 20–26, Jan. 2022.
- [9] P. D. Batlayeri and W. Gatta, "Analisis Sentimen Pejualan Jafra Dalam Pandemi Covid-19 Dengan Algoritma Klasifikasi," *JIRE (Jurnal Informatika & Rekayasa Elektronika)*, vol. 5, pp. 11–18, Apr. 2022.
- [10] M. T. Waluyan and K. D. Hartomo, "Analisis Sentimen Kebutuhan Fast Track Pada Originals Vidio Menggunakan Support Vector Machine," *JATISI (Jurnal Teknik Informatika dan Sistem Informasi)*, vol. 9, no. 3, pp. 2153–2162, Sep. 2022, [Online]. Available:

- <https://jurnal.mdp.ac.id/index.php/jatisi/article/view/2348>
- [11] Ibnu Daqiqil ID, *Machine Learning: Teori, Studi Kasus dan Implementasi Menggunakan Python*, 1st ed. Riau: UR Press, 2021.
- [12] B. T. Jijo and A. M. Abdulazeez, "Classification Based on Decision Tree Algorithm for Machine Learning," *Journal Of Applied Science And Technology Trends*, vol. 02, no. 01, pp. 20–28, 2021.
- [13] M. I. Aditama, R. I. Pratama, K. H. U. Wiwaha, and N. A. Rakhmawati, "Analisis Klasifikasi Sentimen Pengguna Media Sosial Twitter Terhadap Pengadaan Vaksin COVID-19," *JIEET (Journal Information Engineering and Educational Technology)*, vol. 04, no. 02, pp. 90–92, 2020, [Online]. Available: <https://journal.unesa.ac.id/index.php/jieet/article/view/11018>
- [14] H. Nguyen, A. Veluchamy, M. Diop, and R. Iqbal, "Comparative Study of Sentiment Analysis with Product Reviews Using Machine Learning and Lexicon-Based Approaches," *SMU Data Science Review*, vol. 1, no. 4, 2018.
- [15] B. A. Adli, "Analisis Sentimen Customer Review Produk Tokopedia Menggunakan Algoritma Decision Tree," 2022.
- [16] N. K. Widyasanti, I. K. G. D. Putra, and N. K. D. Rusjayanthi, "Seleksi Fitur Bobot Kata dengan Metode TFIDF untuk Ringkasan Bahasa Indonesia," *MERPATI*, vol. 6, no. 7, pp. 119–126, Aug. 2018.
- [17] E. Undamayanti, T. I. Hermanto, and I. Kaniawulan, "Analisis Sentimen Menggunakan Metode Naive Bayes Berbasis Particle Swarm Optimization Terhadap Pelaksanaan Program Merdeka Belajar Kampus Merdeka," *Jurnal Sains Komputer & Informatika (J-SAKTI)*, vol. 6, no. 2, pp. 916–930, Sep. 2022.
- [18] D. A. N. Krisna and U. Salamah, "Perbandingan Algoritma Naïve Bayes Dan K-Nearest Neighbor Untuk Klasifikasi Berita Hoax Kesehatan Di Media Sosial Twitter," *Jurnal Teknik Informatika Kaputama (JTIK)*, vol. 6, no. 2, pp. 836–845, Jul. 2022.
- [19] D. P. Santoso and W. Wibowo, "Analisis Sentimen Ulasan Aplikasi Buzzbreak Menggunakan Metode Naïve Bayes Classifier pada Situs Google Play Store," *JURNAL SAINS DAN SENI ITS*, vol. 11, no. 2, pp. 190–196, 2022.
- [20] S. Khairunnisa, Adiwijaya, and S. Al Faraby, "Pengaruh Text Preprocessing terhadap Analisis Sentimen Komentar Masyarakat pada Media Sosial Twitter (Studi Kasus Pandemi COVID-19)," *JURNAL MEDIA INFORMATIKA BUDIDARMA*, vol. 5, no. 2, pp. 406–414, Apr. 2021.
- [21] C. Cahyaningtyas, Y. Nataliani, and I. R. Widiyari, "Analisis sentimen pada rating aplikasi Shopee menggunakan metode Decision Tree berbasis SMOTE," *AITI: Jurnal Teknologi Informasi*, vol. 18, no. 2, pp. 174–185, Aug. 2021.
- [22] I. Mahayani, A. D. R., and M. E. Supriyadi, "Analisis Sentimen Twitter Terhadap Pembayaran ShopeePayLater Pada Aplikasi Belanja Online (Shopee) Menggunakan Metode Lexicon Based dan Naïve Bayes Classifier," *Jurnal Ilmiah KOMPUTASI*, vol. 19, no. 4, pp. 545–558, Dec. 2020.
- [23] R. Kosasih and A. Alberto, "Analisis Sentimen Produk Permainan Menggunakan Metode TF-IDF Dan Algoritma K-Nearest Neighbor," *InfoTekJar: Jurnal Nasional Informatika dan Teknologi Jaringan*, vol. 6, no. 1, pp. 134–139, 2021.
- [24] S. W. Iriananda, R. P. Putra, and K. S. Nugroho, "Analisis Sentimen Dan Analisis Data Eksploratif Ulasan Aplikasi Marketplace Google Playstore," *The 4th Conference on Innovation and Application of Science and Technology (CIASTECH)*, pp. 473–482, Dec. 2021.
- [25] D. Normawati and S. A. Prayogi, "Implementasi Naïve Bayes Classifier Dan Confusion Matrix Pada Analisis Sentimen Berbasis Teks Pada Twitter," *Jurnal Sains Komputer dan Informatika (J-SAKTI)*, vol. 5, no. 2, pp. 697–711, Sep. 2021.
- [26] I. W. Saputro and B. W. Sari, "Uji Performa Algoritma Naïve Bayes untuk Prediksi Masa Studi Mahasiswa," *Citec Journal*, vol. 6, no. 1, pp. 1–11, Jan. 2019.
- [27] Terence Shin, "What Does it Mean to Deploy A Machine Learning Model?," Towards Data Science. Accessed: May 21, 2023. [Online]. Available: <https://towardsdatascience.com/what-does-it-mean-to-deploy-a-machine-learning-model-dddb983ac416#:~:text=Deploying%20a%20machine%20learning%20model,input%20and%20return%20an%20output.>
- [28] A. B. Prasetyo and T. G. Laksana, "Optimasi Algoritma K-Nearest Neighbors dengan Teknik Cross Validation Dengan Streamlit (Studi Data: Penyakit Diabetes)," *Journal of Applied Informatics and Computing (JAIC)*, vol. 6, no. 2, pp. 194–204, Dec. 2022.

NOMENKLATUR

w_{ij}	bobot term t_j
d_i	bobot dokumen
tf_{ij}	jumlah kemunculan term t_j dalam dokumen d_i .
D	jumlah total dokumen dalam database
df_j	jumlah total dokumen yang mengandung term t_j (paling tidak ada satu kata, yaitu term t_j)