

Implementasi Algoritma Support Vector Machine (SVM) pada Analisis Sentimen Opini Publik Tentang Larangan Penggunaan Obat Sirup bagi Kesehatan Ginjal

Implementation of The Support Vector Machine (SVM) Algorithm on Sentiment Analysis of Public Opinion on The Prohibition of the use of Syrupy Drugs for Kidney Health

¹Galih Purnomo, ²Rumini*, ³Tri Susanto

^{1,2,3}Informatika, Fakultas Ilmu Komputer, Universitas Amikom Yogyakarta

*e-mail: rumini@amikom.ac.id

(received: 5 August 2024, revised: 7 October 2024, accepted: 8 October 2024)

Abstrak

Pada tahun 2022, Kementerian Kesehatan Indonesia melaporkan beberapa kasus gagal ginjal akut pada anak (GGAPA), yang mengakibatkan angka kematian sebesar 59%, terutama di antara anak-anak berusia antara 1-5 tahun. Penyebab utama diidentifikasi oleh Menteri Kesehatan Budi Gunadi Sadikin sebagai tiga pelarut *etilen glikol* (EG), *dietilen glikol* (DEG), dan *etilen glikol butil eter* (EGBE). Sebagai tanggapan, pemerintah menerapkan pembatasan konsumsi zat kental tersebut, yang menyebabkan beragam reaksi publik yang teramati di bagian komentar YouTube. Tujuan dari penelitian ini adalah untuk mengevaluasi pendapat masyarakat tentang larangan sirup untuk kesehatan ginjal. Komentar-komentar tersebut akan diklasifikasi menggunakan metode *Support Vector Machine* (SVM), dan akan diidentifikasi dengan kernel yang paling efektif di antara *linear*, *sigmoid*, *polynomial*, dan *RBF*. Data dikumpulkan melalui *web scraping* dengan 5000 data awal, dan setelah dilakukan *preprocessing*, data yang diolah sebanyak 4794 data. Hasil analisis menunjukkan bahwa kernel *linear* memiliki akurasi tertinggi sebesar 75,63%, diikuti oleh kernel *sigmoid* 75,29%, *RBF* 74,79%, dan *polynomial* 71,09%. Sedangkan pengujian *K-Fold Cross Validation* dengan nilai $k = 10$, menghasilkan nilai sebesar 74,64% untuk kernel *linear*. Penelitian ini menyimpulkan bahwa algoritma *Support Vector Machine* (SVM) dengan *kernel linear* mencapai akurasi tertinggi dalam analisis sentimen.

Kata kunci: gagal ginjal akut, obat sirup, analisis sentimen, *support vector machine*, youtube

Abstract

In 2022, the Indonesian Ministry of Health reported several cases of pediatric acute renal failure (GGAPA), which resulted in a mortality rate of 59%, mainly among children aged between 1-5 years. The main causes were identified by Health Minister Budi Gunadi Sadikin as the three solvents ethylene glycol (EG), diethylene glycol (DEG), and ethylene glycol butyl ether (EGBE). In response, the government implemented restrictions on the consumption of these condensed substances, which led to mixed public reactions observed in the YouTube comments section. The purpose of this study is to evaluate public opinions on the syrup ban for kidney health. The comments will be classified using the *Support Vector Machine* (SVM) method, and the most effective kernel among *linear*, *sigmoid*, *polynomial*, and *RBF* will be identified. Data was collected through *web scraping* with 5000 initial data, and after *preprocessing*, 4794 data were processed. The analysis results show that the *linear* kernel has the highest accuracy of 75.63%, followed by the *sigmoid* kernel 75.29%, *RBF* 74.79%, and *polynomial* 71.09%. While the *K-Fold Cross Validation* test with a value of $k = 10$, produced a value of 74.64% for the *linear* kernel. This research concludes that the *Support Vector Machine* (SVM) algorithm with a *linear* kernel achieves the highest accuracy in sentiment analysis.

Keywords: acute renal failure, syrup medicine, sentiment analysis, *support vector machine*, youtube

1 Pendahuluan

Kementerian Kesehatan (Kemenkes) mengungkapkan terjadinya beberapa kasus gagal ginjal akut pada anak (GGGA) di Indonesia pada tahun 2022, yang menimbulkan berbagai reaksi dari masyarakat selain itu, kasus ini memiliki tingkat kematian yang meningkat secara signifikan, seperti yang ditunjukkan oleh Angka Kematian (*Case Fatality Rate/CFR*) sebesar 59%, Kasus GGAPA sebagian besar terjadi pada anak-anak berusia antara 1 hingga 5 tahun, dengan total 190 kasus yang dilaporkan dan 130 kematian [1]. Menteri Kesehatan Budi Gunadi Sadikin menyatakan bahwa tiga zat pelarut yaitu *ethylene glycol* (EG), *diethylene glycol* (DEG), dan *ethylene glycol butyl ether* (EGBE) sebagai penyebab utama meningkatnya kejadian gagal ginjal akut pada anak-anak [2]. Akibatnya, pemerintah segera menerapkan larangan sementara atas penggunaan obat sirup. Kejadian ini mengakibatkan tekanan psikologis, kekhawatiran, dan ketidaknyamanan yang mengancam kesejahteraan fisik dan mental masyarakat umum.

Larangan pemberian obat sirup kepada anak di bawah umur telah menimbulkan banyak reaksi dari masyarakat umum, seperti yang terlihat dari komentar-komentar yang diposting di YouTube. Youtube secara luas dianggap sebagai salah satu platform media sosial paling populer di Indonesia, yang memberikan akses kepada pengguna ke beragam informasi global. Platform video milik Google itu meraih persentase sebesar 65,41%, jika dibandingkan periode yang sama tahun lalu ada kenaikan 2,39% dengan penggunaanya yang mencapai 181,9 juta orang indonesia [3]. Menurut data yang dilakukan oleh *We Are Social* dan *Hootsuite*, mayoritas individu yang mendedikasikan waktu untuk menonton video di YouTube berada di antara kelompok usia 16 hingga 24 tahun [4]. YouTube berfungsi sebagai sumber informasi berharga yang dapat digunakan sebagai subjek studi untuk analisis sentimen pada subjek tertentu. Youtube adalah media sosial yang digunakan untuk memposting video, melihat berbagai video, dan juga dapat berbagi video yang dapat dilihat oleh semua orang [5]. YouTube berfungsi sebagai sumber informasi berharga yang dapat digunakan sebagai subjek studi untuk analisis sentimen pada subjek tertentu. Memanfaatkan data yang diekstrak dari komentar YouTube akan meningkatkan daya tarik pengetahuan saat ini, sehingga mendorong keinginan untuk melakukan eksplorasi lebih lanjut. Ada banyak cara untuk mempelajarinya tetapi analisis sentimen menjadi salah satu jenis analisis yang dapat digunakan untuk menganalisis opini publik. Analisis sentimen atau *opinion mining* adalah ilmu komputer yang memiliki tujuan untuk mengidentifikasi dan mengartikulasikan opini, perasaan, penilaian, sikap, emosi, subjektivitas, penilaian, atau perspektif yang diekspresikan dalam suatu teks [6]. Dalam analisis sentimen terdapat tiga kategori opini, yaitu positif, netral dan negatif sehingga analisis ini dapat digunakan untuk melihat pendapat atau kesamaan seseorang opini terhadap suatu hal atau objek tertentu [7].

Pada penelitian ini, peneliti menggunakan algoritma *Support Vector Machine* (SVM) karena metode ini sangat baik untuk klasifikasi teks dan tidak memerlukan kemampuan komputasi yang kuat, metode klasifikasi ini sangat mudah digunakan pada perangkat yang tidak terlalu kuat untuk pembelajaran mesin dan sering menjadi standar untuk metode pembelajaran mesin lainnya.

Tujuan dari penelitian ini adalah untuk meningkatkan pemahaman dan kesadaran masyarakat Indonesia mengenai dampak buruk obat-obatan sirup terhadap kesehatan ginjal. Selain itu, penelitian ini juga bertujuan untuk mencegah atau mengurangi terjadinya gagal ginjal dan mengumpulkan tanggapan masyarakat terhadap larangan pemerintah terhadap penggunaan obat sirup. Penelitian ini juga bertujuan untuk mengevaluasi efektivitas teknik *Support Vector Machine* (SVM) dengan menggunakan empat kernel yang berbeda (*linear*, *sigmoid*, RBF, dan *polynomial*) dalam menganalisis sikap masyarakat terhadap larangan tersebut.

2 Tinjauan Literatur

Sebelum dilakukan penelitian, peneliti mengumpulkan informasi dari penelitian sebelumnya yang relevan dengan masalah penelitian, untuk digunakan sebagai titik acuan. Berikut berbagai hasil dari penelitian sebelumnya yang relevan dengan masalah penelitian ini.

Penelitian terakhir berjudul “Analisis Sentimen Terhadap Aplikasi Ruangguru Menggunakan Algoritma *Naive Bayes*, *Random Forest* Dan *Support Vector Machine*” jurnal ini mengeksplorasi pentingnya ulasan pengguna dalam meningkatkan kualitas sebuah aplikasi dan sebagai teknik untuk

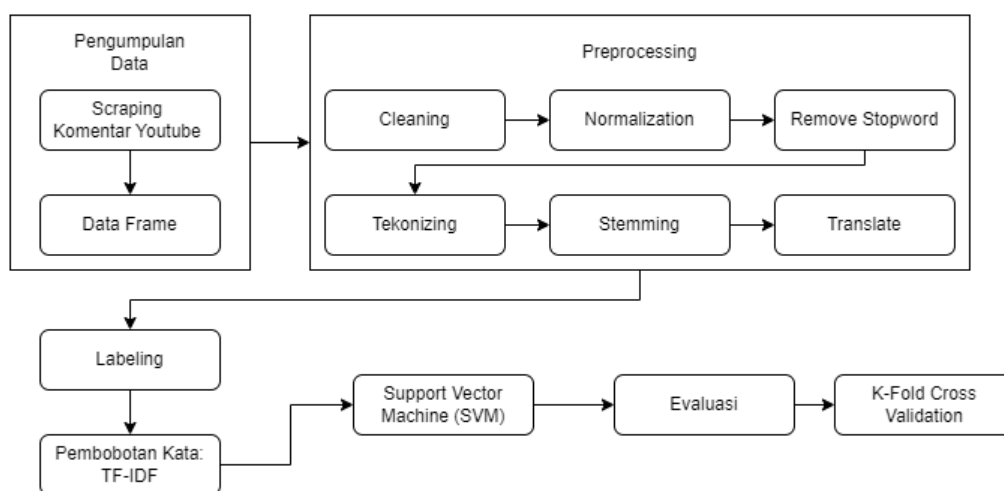
mengevaluasi kepuasan pengguna. Penelitian ini menyajikan temuan model klasifikasi *Random Forest*, yang mencapai tingkat akurasi 97,16% dengan menggunakan *Cross Validation* dan nilai AUC 0,996. Penelitian ini menunjukkan peningkatan akurasi sebesar 7,16% dibandingkan dengan penelitian sebelumnya, dengan model *Random Forest* yang menghasilkan kinerja tertinggi di antara semua model kategorisasi *Random Forest* [8].

Penelitian lainnya berjudul “Perbandingan Algoritma Support Vector Machine dan Random Forest untuk Analisis Sentimen Terhadap Kebijakan Pemerintah Indonesia Terkait Kenaikan Harga BBM Tahun 2022” menjelaskan Pentingnya bahan bakar minyak sebagai komoditas utama dan perannya yang sangat penting dalam operasi perusahaan telah dijelaskan. Kenaikan harga minyak dunia menjadi tantangan berat bagi negara-negara di seluruh dunia, termasuk Indonesia. Langkah ini menuai kritik luas dari masyarakat Indonesia dan mendapat sorotan tajam di berbagai platform media sosial. Berdasarkan temuan pengujian, kedua algoritma bekerja secara efektif seperti yang ditunjukkan oleh nilai akurasi masing-masing. Pendekatan *support vector machine* mencapai akurasi 77%, sedangkan teknik *Random Forest* mencapai akurasi 76%. [9].

Penelitian selanjutnya berjudul “Perbandingan Efektivitas Naïve Bayes dan SVM dalam Menganalisis Sentimen Kebencanaan di Youtube” menjelaskan bahwa Kemajuan di bidang *Natural Language Processing* (NLP) telah menciptakan prospek yang signifikan untuk analisis sentimen, khususnya dalam domain situasi krisis. Penelitian ini menyelidiki dan membandingkan dua metodologi analisis sentimen yang ada, khususnya *Naive Bayes* dan *Support Vector Machine*. Memanfaatkan data komentar YouTube yang berkaitan dengan bencana alam untuk mengevaluasi kemampuan algoritme dalam mendeteksi dan mengategorikan sentimen publik secara akurat sebagai netral, positif, atau negatif. Terkait dengan *Naive Bayes*, algoritma ini mencapai tingkat akurasi sebesar 79%. Selain itu, ia memiliki tingkat *presisi* 91% untuk prediksi negatif dan 73% untuk prediksi positif. Studi ini menawarkan wawasan berharga tentang pemanfaatan teknik pembelajaran mesin dalam analisis sentimen [10].

3 Metode Penelitian

Dalam melakukan penelitian, perlu adanya alur penelitian untuk memastikan bahwa penelitian yang dilakukan dapat berjalan dengan baik, sistematis, dan efektif. Pada penelitian ini terdiri dari tujuh proses. Proses-proses tersebut terdiri dari *scraping data*, *preprocessing data*, *labeling*, *TF-IDF*, *Modeling Support Vector Machine* (SVM), Evaluasi, *K-Fold Cross Validation*. Alur penelitian yang akan peneliti gunakan ditujukan sesuai Gambar 1.



Gambar 1. Alur penelitian

3.1 Pengumpulan Data

Pada awalnya, data komentar dikumpulkan dari platform media Youtube dengan memanfaatkan API youtube V3 dan modul *googleapiclient* dalam bahasa pemrograman *Python*. Dari proses *scraping* tersebut dataset yang diperoleh berjumlah 5000 data yang menggunakan beberapa *keyword*

<http://sistemasi.ftik.unisi.ac.id>

pencariannya youtube, seperti obat sirup bagi kesehatan ginjal, obat sirup buat gagal ginjal akut, penyebab gagal ginjal anak, larangan menggunakan obat sirup, komentar publik tentang larangan penggunaan sirup, dan sebagainya. Data yang telah diambil dari proses *scraping* disimpan dalam bentuk Data Frame/CSV yang nantinya akan diproses lebih lanjut.

3.2 Preprocessing

Sebelum proses data dimulai, komentar harus dibersihkan dari kata-kata yang tidak diperlukan yaitu kata hubung, tanda titik, koma, dan sebagainya. Ini akan membuat pembelajaran mesin yang dirancang lebih mudah dan lebih akurat..

3.2.1 Cleaning Data

Cleaning adalah proses melibatkan penghilangan elemen yang berlebihan, yaitu emoji, tanda baca, spasi ganda, angka, URL, jeda baris, dan fitur-fitur lain yang tidak penting. Selain itu, prosedur ini mencakup transformasi teks menjadi karakter huruf kecil, kadang-kadang disebut sebagai *case folding*. Pada tabel 1 merupakan contoh *cleaning data*.

Tabel 1. Contoh *cleaning data*

Sebelum	Sesudah
👤👤 KENAPA BARU SEKARANG DI KETAHUAN,,,, WOYYYYYYYYY??? BPOM KEMANA AJA, PERCUMA SKLH TINGGI2, KALO DAH BEGINI SAPA YG TANGGUNGJAWAB....moga anak2 gw selalu diberi kesehatan panjang umur ya TUHAN...	kenapa baru sekarang di ketahuan woyyyyyyyy bpom kemana aja, percuma sklh tinggi2, kalo dah begini sapa yg tanggungjawab moga anak2 gw selalu diberi kesehatan panjang umur ya tuhan

3.2.2 Normalization

Normalization adalah proses untuk merubah kata ulang yang awalnya tidak baku menjadi kalimat yang lebih baku. Pada tabel 1 merupakan contoh *normalization* pada kalimat yang diubah dapat berupa singkatan atau kalimat yang kurang jelas.

Tabel 2. Contoh *normalization*

Sebelum	Sesudah
kenapa baru sekarang di ketahuan woyyyyyyyy bpom kemana aja, percuma sklh tinggi2, kalo dah begini sapa yg tanggungjawab moga anak2 gw selalu diberi kesehatan panjang umur ya tuhan	kenapa baru sekarang di ketahuan woyyyyyyyy bpom kemana aja percuma sekolah tinggi kalo dah begini siapa yang tanggungjawab semoga anak saya selalu diberi kesehatan panjang umur ya tuhan

3.2.3 Stopword Remove

Suatu proses menghapus kata-kata yang sering muncul yang dianggap tidak penting dan tidak memiliki makna, contoh *stopword remove* pada tabel 3.

Tabel 3. Contoh *stopword remove*

Sebelum	Sesudah
kenapa baru sekarang di ketahuan woyyyyyyyyy bpom kemana aja percuma sekolah tinggi kalo dah begini siapa yang tanggungjawab semoga anak saya selalu diberi kesehatan panjang umur ya tuhan	kenapa baru sekarang ketahuan bpom kemana percuma sekolah tinggi begini tanggungjawab semoga anak selalu diberi kesehatan panjang umur tuhan

3.2.4 *Tokenizing*

Tokenizing adalah proses memisah kalimat menjadi kata per kata. Tujuan utama *tokenizing* adalah untuk menghasilkan struktur yang lebih terperinci dalam teks, dengan setiap kata berfungsi sebagai unit analisis yang terpisah, berikut contoh *tokenizing* pada tabel 4.

Tabel 4. Contoh *tokenizing*

Sebelum	Sesudah
kenapa baru sekarang ketahuan bpom kemana percuma sekolah tinggi begini tanggungjawab semoga anak selalu diberi kesehatan panjang umur tuhan	[kenapa, baru, sekarang, ketahuan, bpom, kemana, percuma, sekolah, tinggi, begini, tanggungjawab, semoga, anak, selalu, diberi, kesehatan, panjang, umur, tuhan]

3.2.5 *Stemming*

Stemming digunakan untuk menghilangkan imbuhan yang melekat pada awalan kata dan akhiran kata, sehingga kata tersebut menjadi kata dasar, pada tabel 5 contoh *stemming*.

Tabel 5. Contoh *stemming*

Sebelum	Sesudah
[kenapa, baru, sekarang, ketahuan, bpom, kemana, percuma, sekolah, tinggi, begini, tanggungjawab, semoga, anak, selalu, diberi, kesehatan, panjang, umur, tuhan]	[kenapa, baru, sekarang, ketahu, bpom, kemana, percuma, sekolah tinggi, begini, tanggungjawab, semoga, anak, selalu, beri, sehat, panjang, umur, tuhan]

3.2.6 *Translate*

Translate adalah proses mengubah kalimat bahasa Indonesia menjadi bahasa Inggris. Tujuannya untuk menjaga konsistensi bahasa dalam dataset, sehingga seluruh data memiliki format yang seragam, pada tabel 6 merupakan contoh *translate*.

Tabel 6. Contoh translate

Sebelum	Sesudah
[kenapa, baru, sekarang, ketahu, bpom, kemana, percuma, sekolah tinggi, begini, tanggungjawab, semoga, anak, selalu, beri, sehat, panjang, umur, tuhan]	why new now know bpom where in vain school high like this responsibility hopefully child always give healthy long age god

3.3 Labeling

Proses *Labeling* adalah proses untuk membedakan data menjadi dua kelas yaitu kelas positif dan kelas negatif dengan *textblob*[18], dimana label positif berisi kata-kata penyemangat, dukungan dan lainnya, sedangkan label negatif berisi penghinaan, cacian, dan cemoohan [11]. Tahap ini untuk menandai atau mengelompokkan setiap komentar berdasarkan tingkat sentimen yang ditunjukkan dalam skor komentar. Tujuan dari metode ini adalah untuk mempersiapkan data untuk diinterpretasikan dengan model analisis sentimen untuk memahami pendapat dan perspektif yang diungkapkan dalam komentar YouTube.

3.4 Pembobotan TF-IDF

TF-IDF adalah sebuah algoritma yang dapat digunakan untuk penghitungan atau pengekstrakan kata menjadi sebuah angka yang berbentuk *vector* [12]. Prinsip yang mendasari algoritma ini menyatakan bahwa jika sebuah istilah atau frasa lazim dalam satu kategori dokumen tetapi tidak dalam kategori lainnya, maka istilah atau frasa tersebut dapat dianggap sebagai pembeda yang efektif untuk klasifikasi hal ini dicapai dengan menghitung nilai TF-IDF, yang melibatkan perkalian *term frequency* (TF) dengan *inverse document frequency* (IDF) [13]. Berikut pada persamaan 1 :

$$w(x, y) = TF(x, y) \times \log \left(\frac{N}{df(y)} \right) \quad (1)$$

$TF(x,y)$ merepresentasikan *term frequency*, yaitu berapa kali term x muncul di dokumen y . N merujuk pada jumlah total dokumen, sedangkan $df(y)$ adalah jumlah dokumen yang mengandung term y .

3.5 Support Vector Machine

Support Vector Machine (SVM) adalah sebuah teknik yang digunakan untuk menentukan fungsi pemisahan atau *hyperplane* yang membagi data ke dalam beberapa kategori yang berbeda [14]. Proses ini banyak digunakan untuk mengklasifikasikan teks dengan menggunakan berbagai macam kernel yang mempercepat proses mendapatkan hasil dengan cepat, akurat, dan kuat. Terdapat empat fungsi rumus yang sering digunakan dalam SVM antara lain :

1. *Kernel Linear*, yaitu kernel paling sederhana yang digunakan untuk klasifikasi linier pada persamaan 2.

$$K(x_i, x) = x_i x \quad (2)$$

2. *Kernel Polynomial*, yaitu kernel yang menggunakan derajat, pada persamaan 3.

$$K(x_i, x_j) = (x_i x)^d \quad (3)$$

3. *Kernel Gaussian Radial Basis Function*, yaitu kernel untuk dataset yang tidak terpisah secara linier, pada persamaan 4.

$$K(x_i, x_j) = \exp\left(\frac{-\|x_i - x\|^2}{2\sigma^2}\right) \quad (4)$$

4. *Kernel Sigmoid*, yaitu kernel menggunakan proses pengembangan dari jaringan tiruan, pada persamaan 5.

$$K(x_i, x_j) = \tanh(\sigma(x_i \cdot x) + c) \quad (5)$$

Keterangan :

x_i dan x adalah data pada training, parameter δ merupakan sigma, c adalah *complexity*, d adalah *degree*, dan $-\|x_i - x\|^2$ merupakan kuadrat jarak vector x_i dan x .

3.6 Evaluasi

Pada tahap evaluasi ini menggunakan *Confusion Matrix* untuk menghitung nilai TP, FP, TN, dan FN. *Confusion Matrix* adalah pendekatan evaluasi berbasis matriks yang digunakan untuk menilai kinerja model klasifikasi, dimana *Confusion Matrix* memberikan perbandingan antara hasil klasifikasi yang dihasilkan oleh model dan hasil aktual [15]. Kriteria confusion matrix dapat dilihat pada Tabel 7.

Tabel 7. Contoh translate

Kelas Aktual	Kelas Prediksi	
	Positif	Negatif
Positif	TP	FN
Negatif	FP	TN

Keterangan:

TP = Kelas kata positif benar terprediksi positif

FN = Kelas kata positif terprediksi negatif

FP = Kelas kata negatif terprediksi positif

TN = Kelas kata negatif benar terprediksi negatif

Persamaan dan penjelasan singkat untuk akurasi, *precision*, *recall*, and *F1-score* ditunjukkan dibawah ini :

1. Akurasi

Akurasi digunakan untuk mengukur kemampuan suatu *classifier*, khususnya dalam algoritma SVM ini, dalam mengklasifikasi data dengan tepat. Rumus dari akurasi adalah sebagai berikut pada persamaan 6.

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \times 100 \quad (6)$$

2. *Precision*

Precision adalah proses pengukuran presentase dari data prediksi yang kelasnya sama dengan data aktual pada persamaan 7.

$$Precision = \frac{TP}{TP + FP} \times 100 \quad (7)$$

3. Recall

Recall atau bisa disebut dengan *sensitivity* digunakan untuk mengukur kemampuan model dengan presentase banyaknya data aktual yang di prediksi benar oleh sistem pada persamaan 8.

$$Recall = \frac{TP}{TP + FN} \times 100 \quad (8)$$

4. F1-Score

Suatu proses yang digunakan untuk mengukur keseimbangan antara *precision* dan *recall* dalam model klasifikasi, pada persamaan 9.

$$F1 - Score = 2 \times \frac{Precision \times Recall}{Precision + Recall} \quad (9)$$

3.7 K-Fold Cross Validation

K-Fold Cross Validation adalah tahapan terakhir pada penelitian yang berfungsi untuk melakukan validasi pengujian yang menilai keefektifan dan keakuratan proses algoritma [16]. Dalam proses validasi ini memiliki dua proses untuk melatih data, yaitu subproses *training* dan subproses *testing* [17]. Validasi silang *K-fold* memastikan bahwa tidak ada tumpang tindih dalam data pengujian. Pada penelitian ini nilai “k” yang digunakan adalah 10.

4 Hasil dan Pembahasan

Pada bagian ini diuraikan secara detail hasil penelitian yang mencakup pengumpulan data, *preprocessing data*, *labeling*, *tf-idf*, dan *modeling*.

4.1 Pengumpulan Data

Tahap pengumpulan data atau *scraping* data menjadi tahap pertama dalam penelitian ini yang nantinya data tersebut akan menjadi bahan untuk penelitian mengenai analisis sentimen masyarakat. Informasi tersebut bersumber dari komentar YouTube yang berkaitan dengan obat sirup yang digunakan untuk kesehatan ginjal, obat sirup untuk gagal ginjal akut, faktor-faktor yang berkontribusi terhadap gagal ginjal pada anak, pembatasan penggunaan obat sirup, umpan balik masyarakat mengenai pelarangan penggunaan obat sirup, pendapat masyarakat mengenai kebijakan mengenai obat sirup, dan topik-topik serupa. Hasil proses *scraping data* ditunjukkan pada Gambar 2 dibawah ini.

Unnamed: 0	author	published_at	like_count	text
0	@asiahaksiah5708	2022-11-02T01:57:38Z	0	Jika yang mengsumsi obat sirup obat hepelepsi ...
1	@suprisupri6670	2022-10-26T14:38:50Z	0	Tai kucengggg
2	@hudusiahspd8444	2022-10-24T21:14:26Z	0	Pak minuman instan juga sbiakny adstop beredar...
3	@hudusiahspd8444	2022-10-24T21:12:16Z	0	Klo memang betul btul brbhya, kn ank ank yg du...
4	@hudusiahspd8444	2022-10-24T21:11:23Z	0	Sy heran, kok obat berbhya bisa beredar dari d...
...
4995	@ulyred527	2022-10-27T07:45:42Z	0	Makasih dok PENJELASAN YANG SANGAT DIBUTUHKAN 🙏🙏
4996	@drix6301	2022-10-27T07:42:49Z	0	Itu Fungsi BPOM Apa Bosqu Untuk Mslh Ini Paham...
4997	@wafa4609	2022-10-27T07:32:16Z	0	Terimakasih dokter
4998	@rinihariani7625	2022-10-27T07:19:13Z	0	Dok..aman kah saya pake obat demam proris ibu ...
4999	@indarti2717	2022-10-27T07:02:02Z	0	Yaa kalau obat batuknya bikin gagal ginjal yaa...

5000 rows x 5 columns

Gambar 2. Hasil scraping data

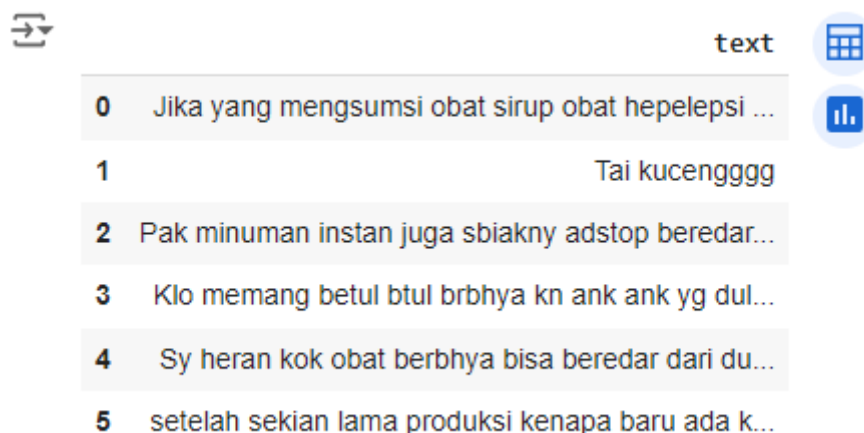
Berdasarkan Gambar 3, terdapat kolom berupa *author*, *published_at*, *like_count*, dan *text* dari komentar pengguna youtube. Data tersebut di ambil dari Video yang di unggah oleh kanal-kanal YouTube terpercaya seperti Kompas.com, CNBC Indonesia, CNN Indonesia, Metro TV, TV One, dan BBC News dengan total 5000 komentar. Pengambilan data tersebut menggunakan API YouTube V3 dan library *googleapiclient* dengan bahasa pemrograman python. Data yang telah diambil dari proses *scraping* disimpan dalam bentuk Data Frame/CSV yang akan diproses lebih lanjut.

4.2 Preprocessing Data

Tahap preprocessing data terdiri dari beberapa tahap termasuk *data cleaning*, *case folding*, *normalisasi*, *stopwords removal*, *tokenizing*, *stemming*, dan *translate*. Pada tahap ini peneliti hanya mengambil data teksnya saja untuk proses lebih lanjut. Tahap tersebut antara lain:

4.2.1 Cleaning Data

Cleaning data dilakukan untuk menghilangkan karakter yang tidak relevan atau mengganggu, sehingga menghasilkan kumpulan data yang lebih terorganisir dan dapat dikelola secara analisis. Dalam proses tersebut juga mencakup perubahan teks menjadi katakter huruf kecil. Hasil dari proses tersebut ditunjukkan pada Gambar 3 di bawah ini.

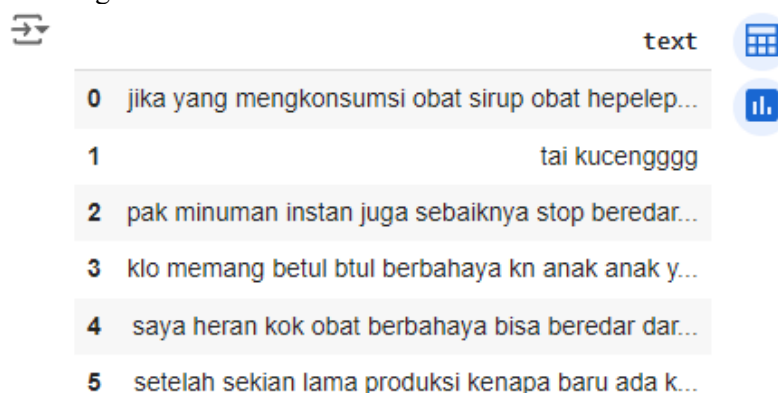


	text
0	Jika yang mengsumsi obat sirup obat hepelepsi ...
1	Tai kucengggg
2	Pak minuman instan juga sbiakny adstop beredar...
3	Klo memang betul btul brbhya kn ank ank yg dul...
4	Sy heran kok obat berbhya bisa beredar dari du...
5	setelah sekian lama produksi kenapa baru ada k...

Gambar 3. Hasil cleaning data

4.2.2 Normalization

Dalam Gambar 4, hasil *Normalization* menunjukkan beberapa perubahan yang signifikan dalam komentar. Dalam beberapa kasus, normalisasi diperlukan untuk merubah kata singkatan dan kesalahan typo menjadi benar dan baku. Seperti kata singkat “sy” menjadi “saya” atau perbaikan typo seperti “bangett” menjadi “banget”.

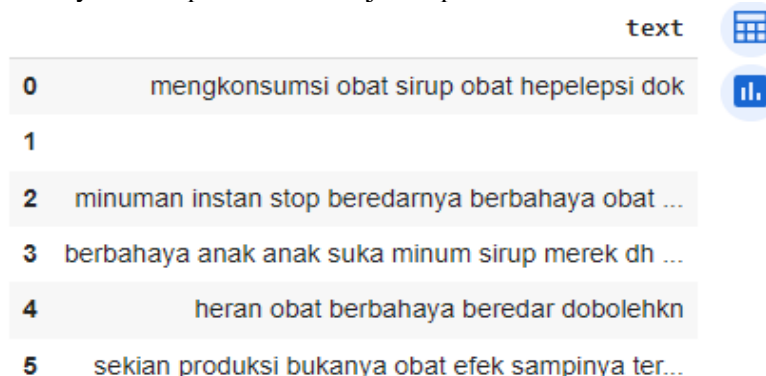


	text
0	jika yang mengkonsumsi obat sirup obat hepelep...
1	tai kucengggg
2	pak minuman instan juga sebaiknya stop beredar...
3	klo memang betul btul berbahaya kn anak anak y...
4	saya heran kok obat berbahaya bisa beredar dar...
5	setelah sekian lama produksi kenapa baru ada k...

Gambar 4. Hasil normalization

4.2.3 Stopword Remove

Proses ini digunakan untuk menghapus kata kata umum yang sering muncul dalam teks tetapi biasanya tidak memiliki nilai informasi yang signifikan, seperti “yang”, “atau”, “saat”, “dengan”, “dan”, “masih”, dan lainnya. Hasil proses ini ditunjukkan pada Gambar 5 dibawah ini.

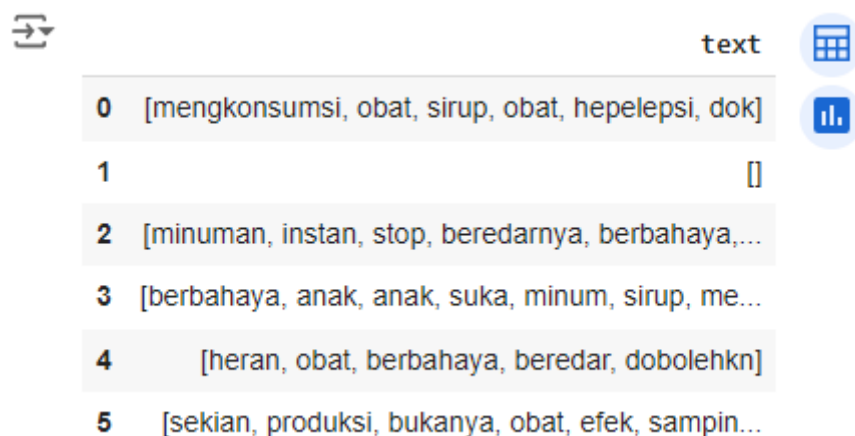


	text
0	mengkonsumsi obat sirup obat hepelepsi dok
1	
2	minuman instan stop beredarnya berbahaya obat ...
3	berbahaya anak anak suka minum sirup merek dh ...
4	heran obat berbahaya beredar dobolehkn
5	sekian produksi bukanya obat efek sampinya ter...

Gambar 5. Hasil stopword remove

4.2.4 Tokenizing

Dalam Gambar 6, teks komentar telah di pecah menjadi token token yang merupakan kata-kata atau frasa-frasa. Proses ini memungkinkan data menjadi lebih terstruktur dan siap untuk pengolahan lebih lanjut dalam analisis teks.

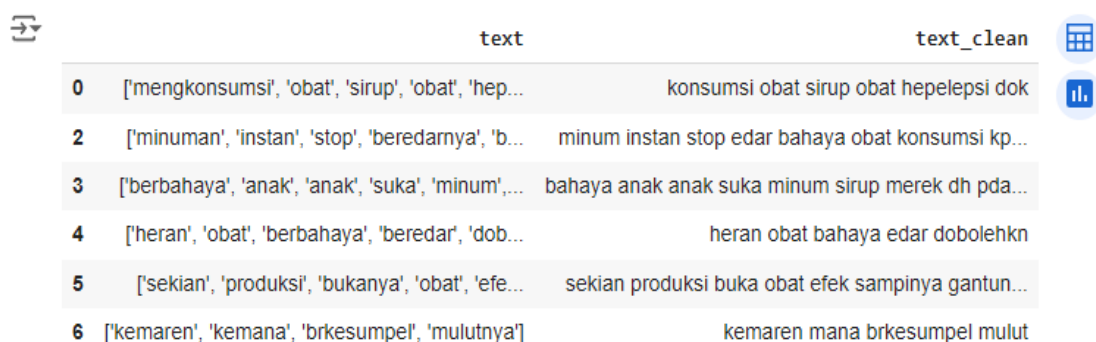


	text
0	[mengkonsumsi, obat, sirup, obat, hepelepsi, dok]
1	[]
2	[minuman, instan, stop, beredarnya, berbahaya,...]
3	[berbahaya, anak, anak, suka, minum, sirup, me...]
4	[heran, obat, berbahaya, beredar, dobolehn]
5	[sekian, produksi, bukanya, obat, efek, sampin...]

Gambar 6. Hasil tokenizing

4.2.5 Stemming

Dalam gambar 7, proses ini merubah kata-kata seperti “mengkonsumsi” diubah menjadi “konsumsi”, “minuman” menjadi “minum”, dan “berbahaya” menjadi “bahaya”. Setelah data di-*stemming*, data tersebut di proses kembali untuk menghilangkan data yang kosong dan data yang terduplikat. Jadi proses tersebut menyebabkan data berkurang menjadi 4794.

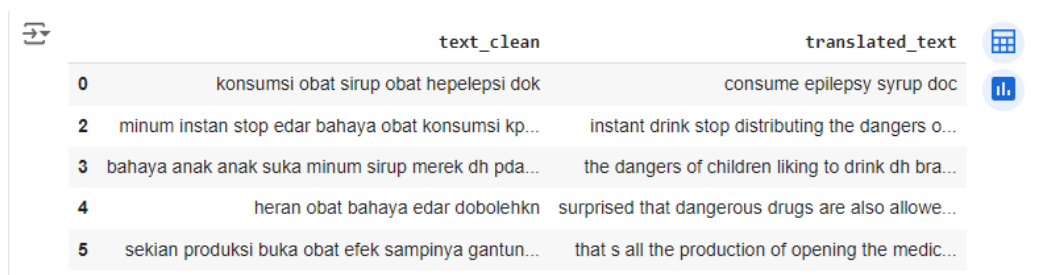


	text	text_clean
0	['mengkonsumsi', 'obat', 'sirup', 'obat', 'hep...]	konsumsi obat sirup obat hepelepsi dok
2	['minuman', 'instan', 'stop', 'beredarnya', 'b...]	minum instan stop edar bahaya obat konsumsi kp...
3	['berbahaya', 'anak', 'anak', 'suka', 'minum', '...]	bahaya anak anak suka minum sirup merek dh pda...
4	['heran', 'obat', 'berbahaya', 'beredar', 'dob...]	heran obat bahaya edar dobolehn
5	['sekian', 'produksi', 'bukanya', 'obat', 'efe...]	sekian produksi buka obat efek sampinya gantun...
6	['kemaren', 'kemana', 'brkesumpel', 'mulutnya']	kemaren mana brkesumpel mulut

Gambar 7. Hasil stemming

4.2.6 Translate

Proses *translate* ini dilakukan untuk merubah bahasa Indonesia menjadi bahasa inggris, dengan tujuan memaksimalkan proses selanjutnya. Hasil dapat dilihat dalam gambar 8 dibawah ini.



	text_clean	translated_text
0	konsumsi obat sirup obat hepelepsi dok	consume epilepsy syrup doc
2	minum instan stop edar bahaya obat konsumsi kp...	instant drink stop distributing the dangers o...
3	bahaya anak anak suka minum sirup merek dh pda...	the dangers of children liking to drink dh bra...
4	heran obat bahaya edar dobolehn	surprised that dangerous drugs are also allowe...
5	sekian produksi buka obat efek sampinya gantun...	that s all the production of opening the medic...

Gambar 8. Hasil translate

4.3 Labeling

Labeling dilakukan untuk menentukan sentimen sebuah komentar, apakah komentar tersebut positif atau negatif dengan data yang sudah di *preprocessing*. Hasil labeling tersebut ditunjukkan pada Gambar 4.8 menggunakan *library TextBloom*. Pada Gambar 9, ditunjukkan bahwa data sentimen negatif memiliki 1496 sedangkan data positif memiliki 1477.

	text_clean	translated_text	subjektivitas	polaritas	sentimen
2	bahaya anak anak suka minum sirup merek dh pda...	the dangers of children liking to drink dh bra...	0.500000	0.5000	Positif
3	heran obat bahaya edar dobolehn	surprised that dangerous drugs are also allowe...	0.900000	-0.2500	Negatif
4	sekian produksi buka obat efek sampinya gantun...	that s all the production of opening the medic...	0.540064	0.1375	Positif
6	lantas guna bpom bubarkana sajah layak knpa pa...	then use bpom to disperse it why is it worth ...	0.216667	0.2750	Positif
8	minum sirup bahaya	drinking syrup is dangerous	0.900000	-0.6000	Negatif

Pada gambar 10 berikut merupakan detail jumlah data positif dan negatif.

Gambar 9. Hasil labeling

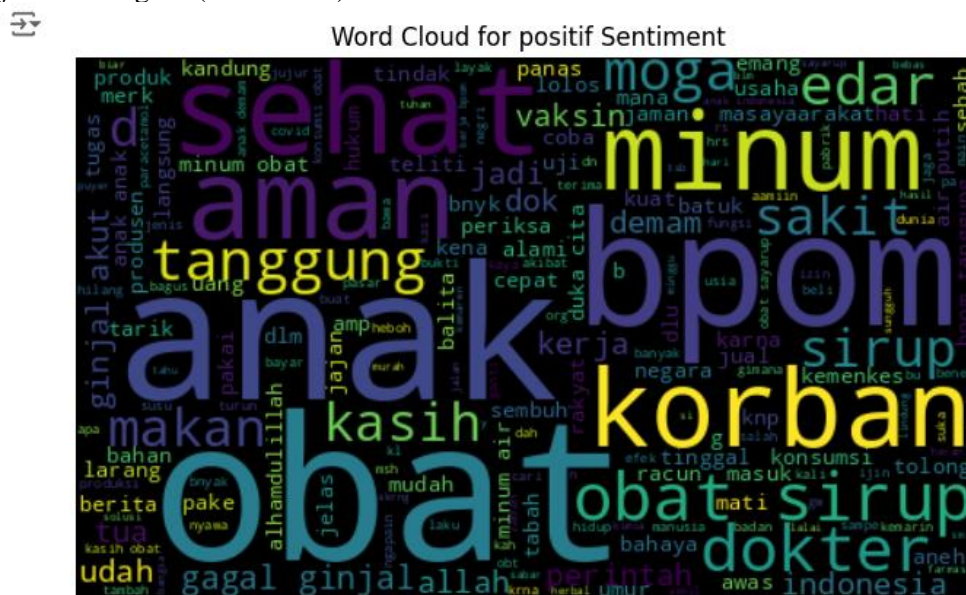
```

sentimen
Negatif      1496
Positif      1477
Name: count, dtype: int64
    
```

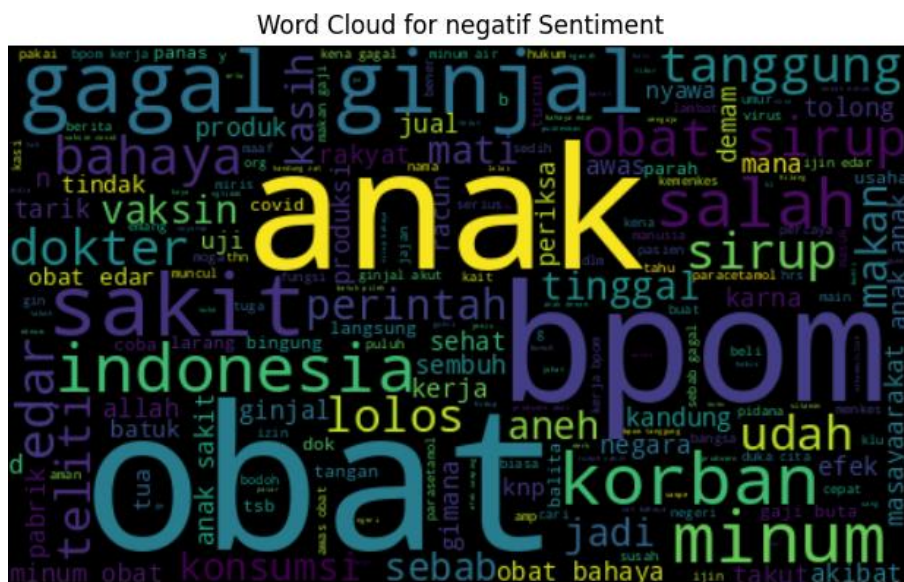
Gambar 10. Data positif dan data negatif

4.4 WordCloud

Wordcloud atau juga dikenal dengan sebagai klaster kata adalah secara visual menyajikan istilah yang paling sering muncul dalam set data pengujian yang berisi data positif dan negatif. Berdasarkan temuan ini, dapat disimpulkan bahwa istilah "sehat", "minum", dan "aman" sering muncul pada pengujian data positif (Gambar 11), sedangkan istilah "sakit", "gagal", dan "ginjal" ditampilkan pada pengujian data negatif (Gambar 12).



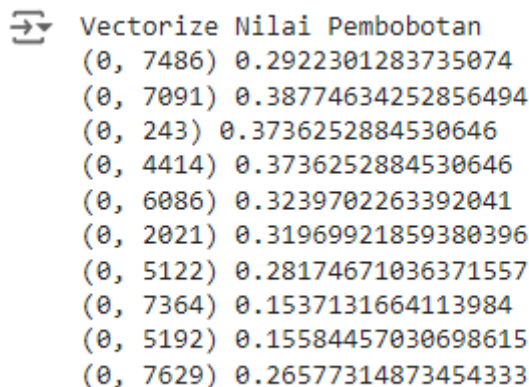
Gambar 11. WordCloud positif



Gambar 12. WordCloud negatif

4.5 Pembobotan Kata TF-IDF

Tahap TF-IDF digunakan untuk memberikan nilai pada arti penting pada kata-kata di dalam dokumen dan mengidentifikasi istilah-istilah yang sangat penting bagi dokumen. Hasil dari TF-IDF dapat dilihat pada Gambar 13 di bawah ini.



Gambar 13. Hasil TF-IDF

4.6 Support Vector Machine (SVM)

Sebelum mengklasifikasi data, data yang diambil dari langkah sebelumnya perlu dipartisi menjadi dua segmen, dengan distribusi 80%:20%. Porsi 20% akan digunakan sebagai dataset pengujian, sedangkan 80% sisanya akan digunakan sebagai dataset pelatihan. Tujuan dari pembagian ini adalah untuk memastikan bahwa model dapat dievaluasi dengan menggunakan data yang tak terlihat yang akan menghasilkan hasil pengujian yang lebih akurat dan tidak bias. Setelah pembagian dataset, langkah selanjutnya adalah melatih mode menggunakan empat kernel yaitu *linier*, *polynomial*, RBF, dan *sigmoid*. Penggunaan berbagai kernel ini bertujuan untuk menemukan akurasi paling tinggi dari empat kernel tersebut yang nantinya akan digunakan untuk proses evaluasi.

```

kernels = ['linear', 'poly', 'rbf', 'sigmoid']

# Dictionary untuk menyimpan akurasi setiap kernel
accuracy_dict = {}

# Melakukan pelatihan dan evaluasi untuk setiap kernel
for kernel in kernels:
    svm = SVC(kernel=kernel)
    svm.fit(X_train, y_train)
    svm_prediction = svm.predict(X_test)
    accuracy = accuracy_score(y_test, svm_prediction)
    accuracy_dict[kernel] = accuracy
    print(f"Akurasi untuk kernel {kernel}: {accuracy:.4f}")

# Menemukan kernel dengan akurasi tertinggi
best_kernel = max(accuracy_dict, key=accuracy_dict.get)
best_accuracy = accuracy_dict[best_kernel]

print(f"\nKernel dengan akurasi tertinggi: {best_kernel} dengan akurasi {best_accuracy:.4f}")

```

Akurasi untuk kernel linear: 0.7563
 Akurasi untuk kernel poly: 0.7109
 Akurasi untuk kernel rbf: 0.7479
 Akurasi untuk kernel sigmoid: 0.7529

Kernel dengan akurasi tertinggi: linear dengan akurasi 0.7563

Gambar 14. Hasil pengujian empat kernel SVM

Berdasarkan Gambar 14, hasil pengujian akurasi tersebut terlihat bahwa kernel *linear* mengungguli tiga jenis kernel lainnya dengan akurasi 75,63%, menjadikannya pilihan yang paling efektif untuk dataset ini. Kernel *sigmoid* menunjukkan akurasi yang sebanding dengan kernel linear, mencapai nilai 75,29%. Sebaliknya, kernel *RBF* menunjukkan akurasi 74,79%, dan kernel *polynomial* menunjukkan akurasi terendah pada 71,09%. Hasil ini menunjukkan bahwa pemisahan data secara *linear* menghasilkan hasil prediksi yang paling tepat untuk model klasifikasi yang digunakan.

4.7 Evaluasi

Setelah proses modeling SVM selesai dan kernel *linear* dipilih, selanjutnya data akan memasuki tahap *evaluasi* menggunakan metode *confusion matrix* yang akan menghasilkan nilai *precision*, *recall*, *F1-score*, dan akurasi. Hasil *evaluasi confusion matrix* ditunjukkan pada Gambar 15, Untuk kategori *negatif* menunjukkan nilai *precision* sebesar 77%, *recall* sebesar 74%, dan *F1-score* sebesar 75 %. Sedangkan untuk kategori *positif* menunjukkan nilai *precision* sebesar 75%, *recall* sebesar 77%, dan *F1-score* sebesar 76%.

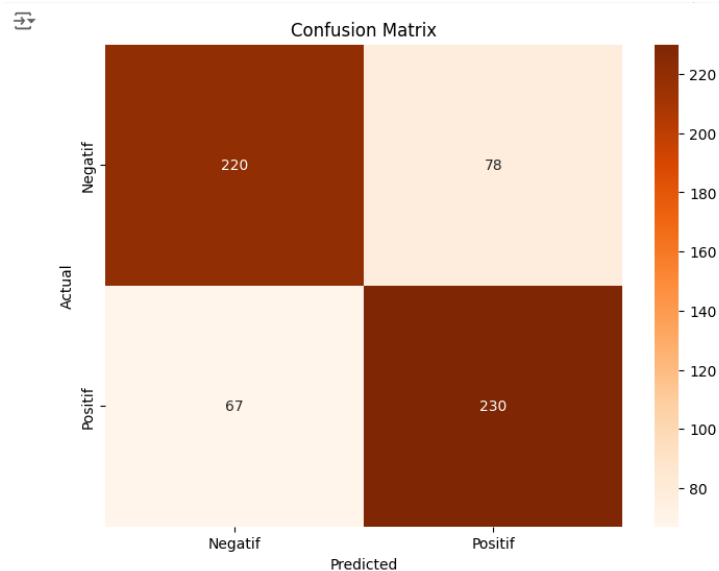
Evaluasi SVM dengan kernel terbaik:
Classification Report:

	precision	recall	f1-score	support
Negatif	0.77	0.74	0.75	298
Positif	0.75	0.77	0.76	297
accuracy			0.76	595
macro avg	0.76	0.76	0.76	595
weighted avg	0.76	0.76	0.76	595

Accuracy: 0.7563025210084033

Gambar 15. Hasil evaluasi svm kernel terbaik (linear)

Pada Gambar 16 menunjukkan hasil sebagai berikut *True Negative* = 220, *True Positive* = 230, *False Positive* = 68, dan *False Negative* = 78.



Gambar 16. Hasil confusion matrix

4.8 K-Fold Cross Validation

Selanjutnya akan dilakukan pengujian ulang menggunakan *K-Fold Cross Validation* dengan tujuan agar hasil evaluasi kerja algoritma memperoleh hasil maksimal. Pada penelitian ini menggunakan nilai k sebesar 10 dengan hasil pengujian rata-rata yang dapat dilihat pada Gambar 17 sebesar 74,64%.

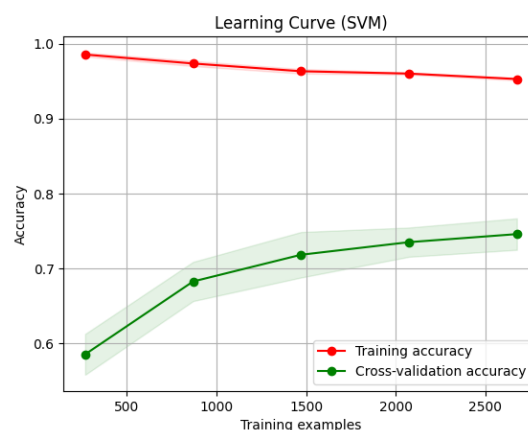
```
[ ] # k-fold Cross Validation untuk SVM
k = 10
svm_cv_scores = cross_val_score(svm, X_final, y, cv=k, scoring='accuracy')

# Print hasil cross validation
print("Mean Accuracy (SVM):", np.mean(svm_cv_scores))
```

Mean Accuracy (SVM): 0.7464002440512507

Gambar 17. Hasil rata-rata akurasi SVM kernel linear

Pada penelitian ini, memanfaatkan *learning curve* untuk menganalisis kesesuaian data yang menghubungkan performa model dengan jumlah data yang digunakan, serta untuk memahami proses belajar algoritma. Gambar 18, menunjukkan bahwa *learning curve* pada model SVM tersebut belajar dengan baik dan tidak mengalami tumpang tindih (*overlap*)/menyilang. Hal ini membuktikan bahwa SVM memiliki performa yang optimal pada data baru.



Gambar 18. Hasil learning curve SVM kernel linear

5 Kesimpulan

Berdasarkan data dan hasil penelitian yang sudah dilakukan peneliti, bahwa data yang digunakan pada penelitian ini di dapat dari hasil data *crawling/scraping* sebanyak 5000 data dan dilakukan preprocessing sehingga tersisa 4794 data yang dapat diolah. Data diambil dengan menggunakan keyword dari pencarian youtube seperti "obat sirup untuk kesehatan ginjal", "obat sirup untuk gagal ginjal akut", "penyebab gagal ginjal pada anak", "larangan penggunaan obat sirup", "komentar publik tentang larangan penggunaan obat sirup", "perspektif publik tentang kebijakan obat sirup". Data tersebut diambil pada tahun 2022.

Hasil akurasi menggunakan algoritma *Support Vector Machine* (SVM) menggunakan empat kernel menunjukkan bahwa kernel linear memiliki akurasi tertinggi sebesar 75,63%, diikuti oleh kernel sigmoid 75,29%, RBF 74,79%, dan Polynomial 71,09%.

Berdasarkan hasil evaluasi terhadap metode *Support Vector Machine* (SVM) dengan menggunakan kernel *linear*, dapat disimpulkan bahwa model ini menunjukkan performa yang memuaskan dalam mengklasifikasi data secara akurat ke dalam kelas negatif dan positif. *Evaluasi* terhadap sentimen opini publik tentang larangan penggunaan obat sirup bagi kesehatan ginjal menunjukkan bahwa Untuk kategori *negatif* menunjukkan nilai *precision* sebesar 77%, *recall* sebesar 74%, dan *F1-score* sebesar 75 %. Sedangkan untuk kategori *positif* menunjukkan nilai *precision* sebesar 75%, *recall* sebesar 77%, dan *F1-score* sebesar 76%. Hasil ini menunjukkan bahwa metode SVM dengan kernel linear memiliki kemampuan yang cukup baik dalam mengklasifikasikan sentiment opini publik ke dalam kategori negatif dan positif.

Hasil validasi yang diperoleh dengan *K-Fold Cross Validation* dengan K=10 menunjukkan akurasi rata-rata sebesar 74,64%. Selain itu, *learning curve* menunjukkan bahwa model tidak mengalami tumpang tindih (*overlap*) dan kurvanya cenderung baik serta tidak menyilang.

Referensi

- [1] A. Putri Riani, N. Sulistyowati, T. Ridwan, and A. Voutama, "METHOMIKA: Jurnal Manajemen Informatika & Komputerisasi Akuntansi," vol. 7, no. 2, 2023, doi: 10.46880/jmika.Vol7No2.pp325-339.
- [2] Diva Lufiana Putri and Sari Hardiyanto, "3 Zat Berbahaya Temuan Kemenkes pada Pasien Gagal Ginjal Akut, Apa Saja?," Kompas. Accessed: May 28, 2024. [Online]. Available: <https://www.kompas.com/tren/read/2022/10/20/132900065/3-zat-berbahaya-temuan-kemenkes-pada-pasien-gagal-ginjal-akut-apa-saja-?page=all>
- [3] Fitri Wulandari, Elin Haerani, Muhammad Fikry, and Elvia Budianita, "Analisis sentimen larangan penggunaan obat sirup menggunakan algoritma naive bayes classifier," *Jurnal CoSciTech (Computer Science and Information Technology)*, vol. 4, no. 1, pp. 88–96, May 2023, doi: 10.37859/coscitech.v4i1.4781.
- [4] Irham Zharfan and Ferlangga, "Youtube Jadi Raja Media Sosial Di Indonesia, Diakses 94 Persen Warga ," Binus. Accessed: May 28, 2024. [Online]. Available: <https://student-activity.binus.ac.id/himti/2022/08/25/youtube-jadi-raja-media-sosial-di-indonesia-diakses-94-persen-warga/>
- [5] T. M. Tinambunan and C. Siahaan, "Tresia Monica Tinambunan, dan Chontina Siahaan Pemanfaatan Youtube Sebagai Media Komunikasi Massa Di Kalangan Pelajar," 2022. [Online]. Available: www.youtube.com
- [6] Z. Nanda Aulia, G. Kuncoro Jati, and I. Santoso, "Analisis Sentimen Tanggapanpublic Mengenai E-Tilang Melalui Media Sosial Youtube Menggunakan Algoritma Naive Bayes," Feb. 2023. [Online]. Available: <https://journals.upi-yai.ac.id/index.php/ikraith-informatika/issue/archive>

- [7] O. I. Gifari, M. Adha, I. Rifky Hendrawan, F. Freddy, and S. Durrand, "Analisis Sentimen Review Film Menggunakan TF-IDF dan Support Vector Machine," *JIFOTECH (Journal Of Information Technology)*, vol. 2, no. 1, 2022.
- [8] E. Fitri, Y. Yuliani, S. Rosyida, and W. Gata, "Analisis Sentimen Terhadap Aplikasi Ruangguru Menggunakan Algoritma Naive Bayes, Random Forest Dan Support Vector Machine," *TRANSFORMTIKA*, vol. 18, no. 1, pp. 71–80, 2020, [Online]. Available: www.nusamandiri.ac.id,
- [9] M. Samantri, "Perbandingan Algoritma Support Vector Machine dan Random Forest untuk Analisis Sentimen Terhadap Kebijakan Pemerintah Indonesia Terkait Kenaikan Harga BBM Tahun 2022," *Jurnal Teknologi Informasi dan Komunikasi*, vol. 8, no. 1, p. 2024, 2024, doi: 10.35870/jti.
- [10] T. Aura Azzahra *et al.*, "Jurnal Media Informatika Budidarma Perbandingan Efektivitas Naive Bayes dan SVM dalam Menganalisis Sentimen Kebencanaan di Youtube," 2024, doi: 10.30865/mib.v8i1.7186.
- [11] D. R. Manalu, M. C. L. Tobing, and M. Yohanna, "Analisis Sentimen Twitter Terhadap Wacana Penundaan Pemilu Dengan Metode Support Vector Machine," *METHOMIKA Jurnal Manajemen Informatika dan Komputerisasi Akuntansi*, vol. 6, no. 6, pp. 149–156, Oct. 2022, doi: 10.46880/jmika.Vol6No2.pp149-156.
- [12] M. Hafizh Mahendra, D. Triantoro Murdiansyah, and K. Muslim Lhaksana, "Dike : Jurnal Ilmu Multidisiplin Analisis Sentimen Tweet COVID-19 Menggunakan Metode K-Nearest Neighbors dengan Ekstraksi Fitur TF-IDF dan CountVectorizer," 2023.
- [13] P. Widyantara *et al.*, "Analisis Sentimen pada Teks Berbahasa Bali Menggunakan Metode Multinomial Naive Bayes dengan TF-IDF dan BoW," *JNATIA*, vol. 2, no. 1, 2023.
- [14] A. A. Nurkhaliza and A. W. Wijayanto, "Perbandingan Algoritma Klasifikasi Support Vector Machine dan Random Forest pada Prediksi Status Indeks Mitigasi dan Kesiapsiagaan Bencana (IMKB) Satuan Kerja BPS di Indonesia Tahun 2020," *Jurnal Informatika Universitas Pamulang*, vol. 7, no. 1, pp. 2622–4615, 2022, doi: 10.32493/informatika.v7i1.16117.
- [15] M. Raffi, A. Suharso, and I. Maulana, "Analisis Sentimen Ulasan Aplikasi Binar Pada Google Play Store Menggunakan Algoritma Naive Bayes Sentiment Analysis Of Binar Application Reviews On Google Play Store Using Naive Bayes Algorithm," *Journal of Information Technology and Computer Science (INTECOMS)*, vol. 6, no. 1, 2023.
- [16] S. Alpin Rizaldi, S. Alam, and I. Kurniawan, "Analisis Sentimen Pengguna Aplikasi Jmo (Jamsostek Mobile) Pada Google Play Store Menggunakan Metode Naive Bayes 1)," vol. 2, no. 3, pp. 109–117, 2023, doi: 10.55123.
- [17] Y. Widyaningsih, G. P. Arum, and K. Prawira, "Aplikasi K-Fold Cross Validation Dalam Penentuan Model Regresi Binomial Negatif Terbaik," *BAREKENG: Jurnal Ilmu Matematika dan Terapan*, vol. 15, no. 2, pp. 315–322, Jun. 2021, doi: 10.30598/barekengvol15iss2pp315-322.
- [18] Nugroho, Adytyo Wahyu. Dan Norhikmah., 2024. "Analisis Sentimen menggunakan Algoritma Support Vector Machine pada Covid_19, Jurnal SISTEMASI : Jurnal Sistem Informasi Volume 13, No. 4, 2024, ISSN:2302-8149.